# Neuron Network Basis

## Chun-Hsiang Chan

Undergraduate Program in Intelligent Computing and Big Data, Chung Yuan Christian University
Master Program in Intelligent Computing and Big Data, Chung Yuan Christian University
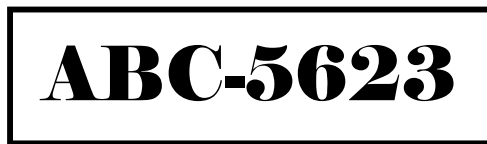
# Outline

- Model Formulation

- Single Neuron

- Hidden Layer

- Activation Function

- Linearity and Non-linearity
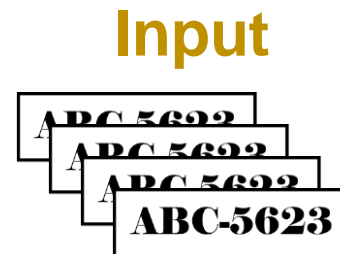
- Model Parameters

- Loss Function

# Model Formulation

- Imagine that you have a problem, such as a digit recognition, object detection, or numerical prediction.

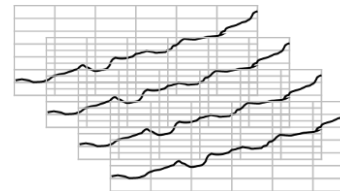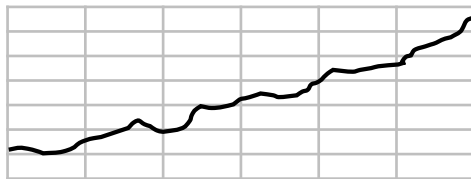**Digit recognition**          **Input**          **Function**          **Output**



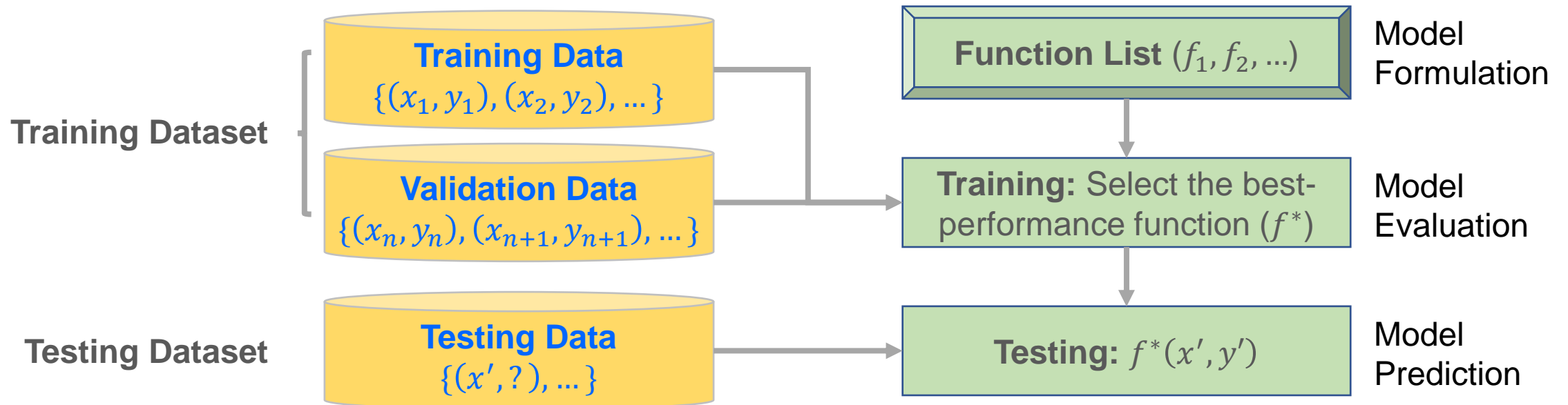$f(x)$          > "ABC-5623"

**Numerical prediction**



$f(x)$          > 234, 235, …
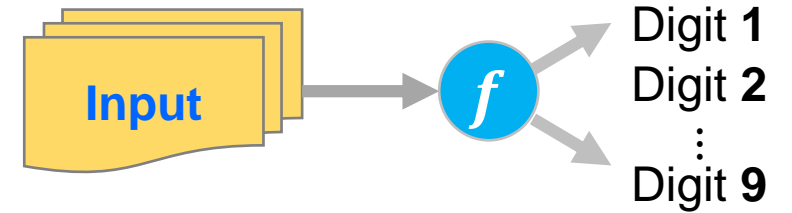
# Model Formulation

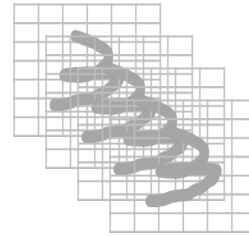- In machine learning and deep learning, we usually divide our input data into training and validation datasets for the training model.



**Training Dataset**

**Training Data**
$\{(x_1, y_1), (x_2, y_2), \dots\}$

**Validation Data**
$\{(x_n, y_n), (x_{n+1}, y_{n+1}), \dots\}$

**Testing Dataset**

**Testing Data**
$\{(x', ?), \dots\}$

**Function List** $(f_1, f_2, \dots)$

Model Formulation

**Training:** Select the best-performance function $(f^*)$

Model Evaluation

**Testing:** $f^*(x', y')$

Model Prediction

# Model Formulation

- **Prediction Tasks**
  - **Handwritten Recognition**



Input $\rightarrow$ $f$ $\rightarrow$ Digit **1** Digit **2** ⋮ Digit **9**

  - **Abnormal Chip Testing**

Input $\rightarrow$ $f$ $\rightarrow$ Normal / Abnormal

  - **Stock Prediction**

Input $\rightarrow$ $f$ $\rightarrow$ $y_{t+1}, y_{t+2}, \ldots$
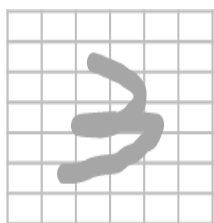
# Model Formulation

- **Handwritten Recognition**: alphabet and number

$$x: image$$



**7 x 6 array**

$$f$$

$$y: label$$

$$f: R^N \rightarrow R^M$$

1 indicates the element contains ink
0 indicates the element contains nothing

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \Rightarrow 7 \times 6 = 42 \; elements$$

10+26 dimension for label recognition

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{matrix} a \\ b \\ c \\ \vdots \\ 8 \\ 9 \end{matrix} \quad \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$
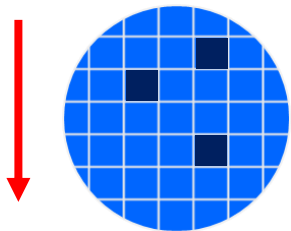
Detect          Ground Truth

# Model Formulation

- **Abnormal Chip Testing**

$$x: image$$



1 indicates a defect chip
0 indicates a normal chip

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \Rightarrow dimensions = size(chips)$$

$$f$$

$$y: label$$
2 dimensions
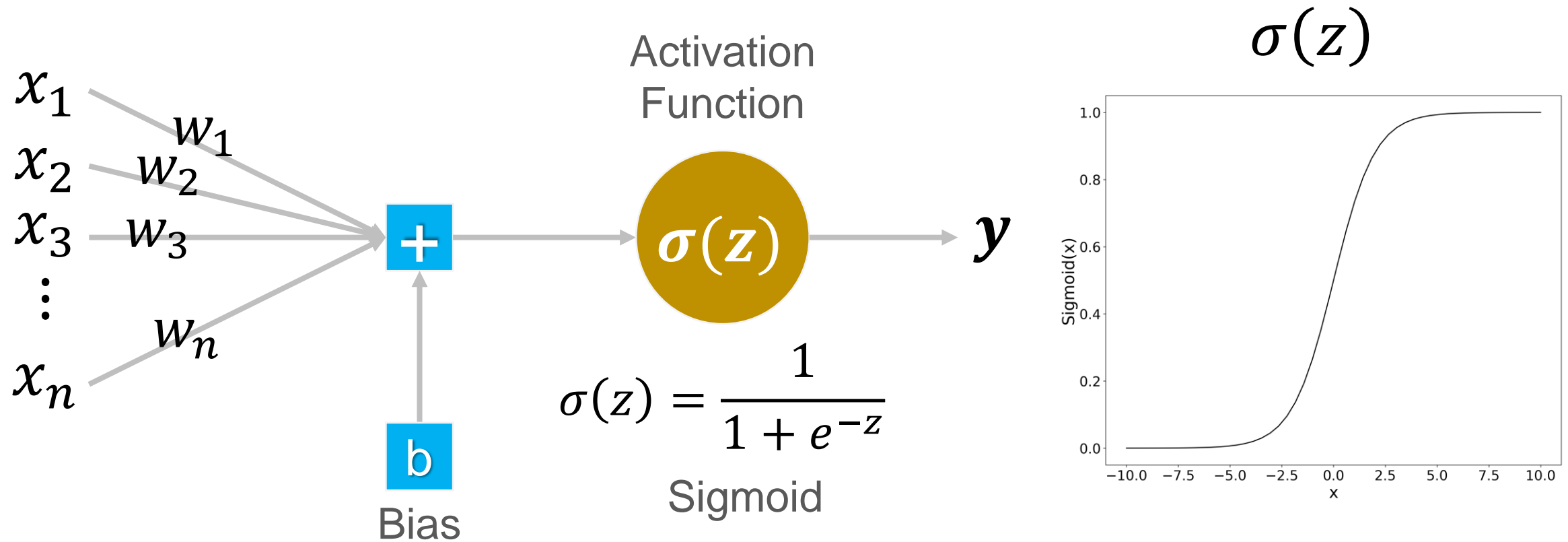Normal chip or Defect chip

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} V \\ V \\ \times \\ \vdots \\ V \\ V \end{matrix} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$  Confirm

Detect   Ground Truth

# Single Neuron

$x_1$

$w_1$

$x_2$

$w_2$

$x_3$

$w_3$

$\vdots$

$w_n$

$x_n$

**+**

**b**

Bias

Activation
Function

$\sigma(z)$

$y$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Sigmoid

$\sigma(z)$

# Single Neuron



$x_1$
$w_1$
$x_2$
$w_2$
$x_3$
$w_3$
$\vdots$
$w_n$
$x_n$
$b$

$+$

Activation
Function

$\sigma(z)$

$y$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Sigmoid

$1$

Bias

$\sigma(z)$

# Single Neuron

$x_1$

$w_1$

$x_2$　$w_2$

$x_3$　$w_3$

$\vdots$

$w_n$

$x_n$　$b$

$+$

Activation
Function

$\sigma(z)$

$y$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Sigmoid

1

Bias

$$y = h_{w,b}(x)$$
$$= \sigma(w^T x + b)$$

# A Layer of Neurons – Perceptron



$x_1$
$x_2$
$x_3$
$\vdots$
$x_n$

$+$　　$\sigma(z)$　　$y_1 \; (is\ a?)$

$+$　　$\sigma(z)$　　$y_2 \; (is\ b?)$

$+$　　$\sigma(z)$　　$y_3 \; (is\ c?)$

$\Big\}$ Choose the max one

**1**

Bias

Three neurons for three classes/ labels

# A Layer of Neurons – Perceptron



$$x_1, x_2, x_3, \ldots, x_n$$
$$w_1, w_2, w_3, w_n, b$$
$$\boxed{1}$$
$$+ \quad \sigma(z) \quad y$$

$$\begin{cases} true & y \geq \theta \\ false & y < \theta \end{cases}$$

$$\theta \; is \; threshold$$

$y = \sigma(w_1 x_1 + w_2 x_2 + b)$ is a linear formula...



Perceptron output

# A Layer of Neurons – Perceptron

$x_1$
$x_2$
$x_3$
$\vdots$
$x_n$

$w_1$
$w_2$
$w_3$
$w_n$
$b$

$+$

$\sigma(z)$

$y$

$1$

$$\begin{cases} true & y \geq \theta \\ false & y < \theta \end{cases}$$

$\theta \ is \ threshold$



(a) $x_1$ **and** $x_2$

(b) $x_1$ **or** $x_2$

(c) $x_1$ **xor** $x_2$

Minsky & Papert (1969) pricked the neural network balloon

# A Layer of Neurons – Perceptron



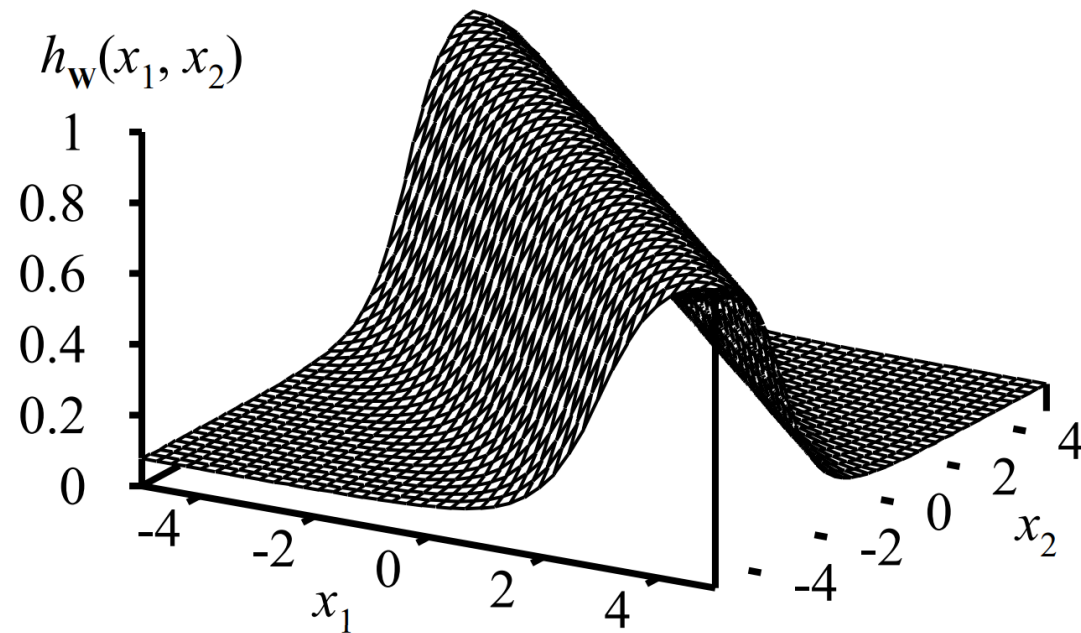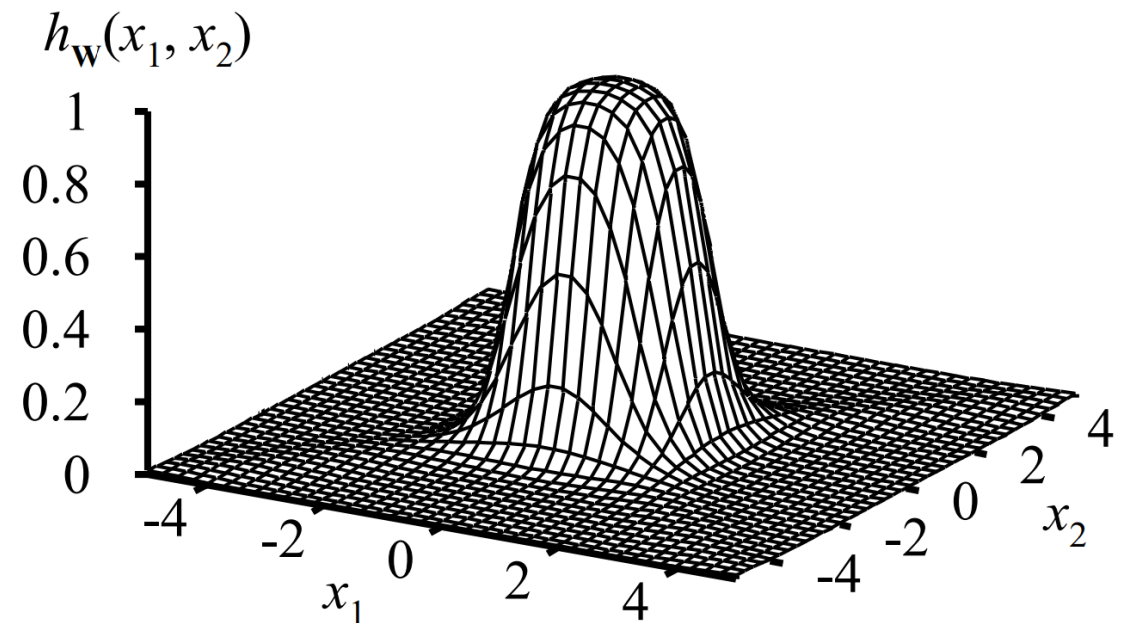**Hidden unit**

Bias

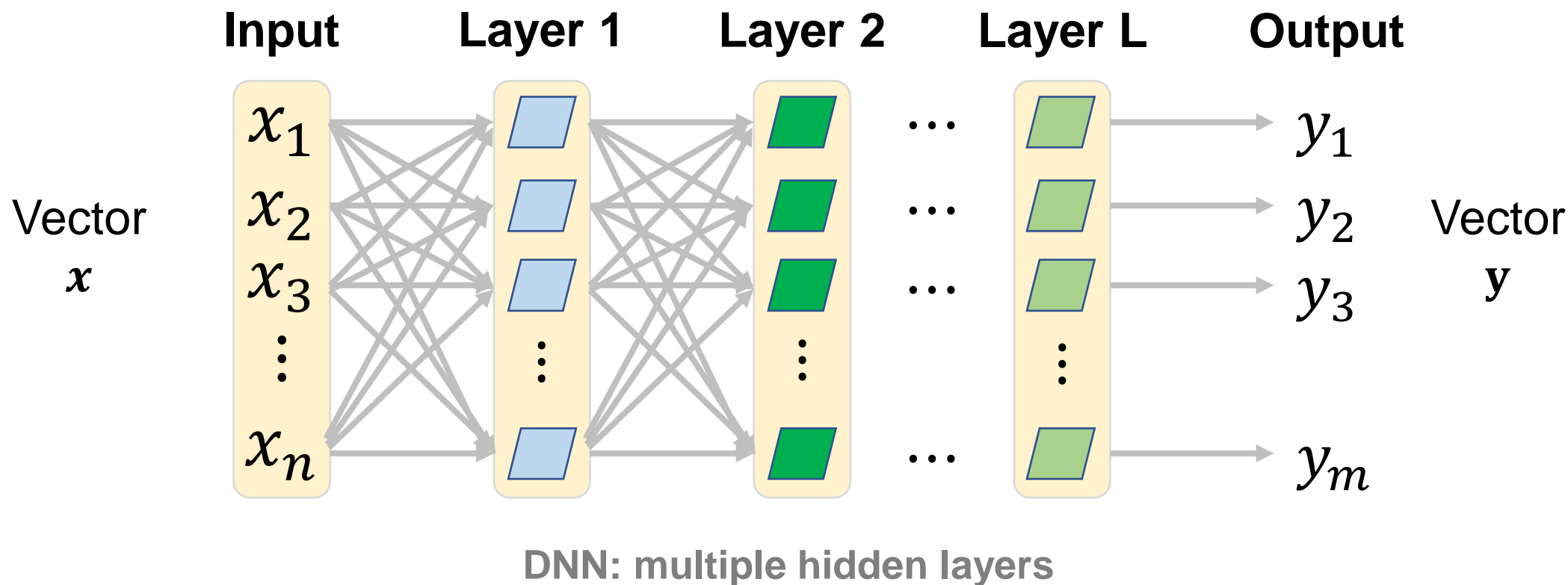# A Layer of Neurons – Perceptron



Combine two opposite-facing threshold functions to make a ridge

Combine two perpendicular ridges to make a bump

# Fully Connected Feedforward Network

**Deep Neural Network (DNN)**



DNN: multiple hidden layers

# Hidden Layer – Layer

**Layer** $l-1$   **Layer** $l$



| 1 | $a_1^{l-1}$ |
| 2 | $a_2^{l-1}$ |
| 3 | $a_3^{l-1}$ |
| ⋮ | |
| j | $a_j^{l-1}$ |

$N_{l-1}$ **nodes**

| 1 | $a_1^{l}$ |
| 2 | $a_2^{l}$ |
| 3 | $a_3^{l}$ |
| ⋮ | |
| i | $a_i^{l}$ |

$N_l$ **nodes**

**A Neuron Output**

$$a_1^l \quad \longrightarrow \text{Layer } l$$
$$\quad \longrightarrow \text{neuron } i$$

$$a^l = \begin{bmatrix} \vdots \\ a_i^l \\ \vdots \end{bmatrix}$$

One layer ➔ One Vector

# Hidden Layer – Weights

**A Weight between two neurons/nodes**

**Layer $l-1$**          **Layer $l$**



$w_{ij}^l$ → Layer $l-1$ to Layer $l$

$w_{ij}^l$ → From neuron $j$ (layer $l-1$) to neuron $i$ (layer $l$)
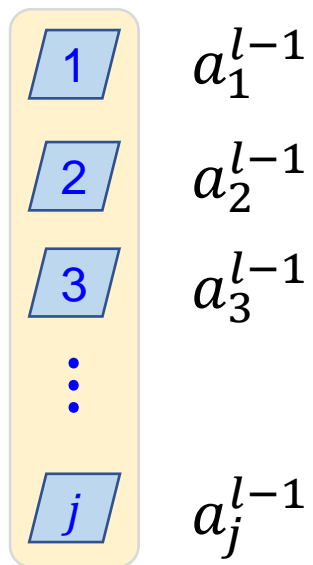
$$w^l = \begin{bmatrix} w_{11}^l & w_{12}^l & \cdots \\ w_{21}^l & w_{22}^l & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

$N_{l-1}$

$N_l$

Weights between layers ➔ One Matrix

# Hidden Layer – Bias

**Layer** $l-1$                    **Layer** $l$



$a_1^{l-1}$

$a_2^{l-1}$

$a_3^{l-1}$

$a_j^{l-1}$

$N_{l-1}$ **nodes**

$b_1^l$

$b_2^l$

$b_3^l$

$b_i^l$

$a_1^l$

$a_2^l$

$a_3^l$

$a_i^l$

Bias

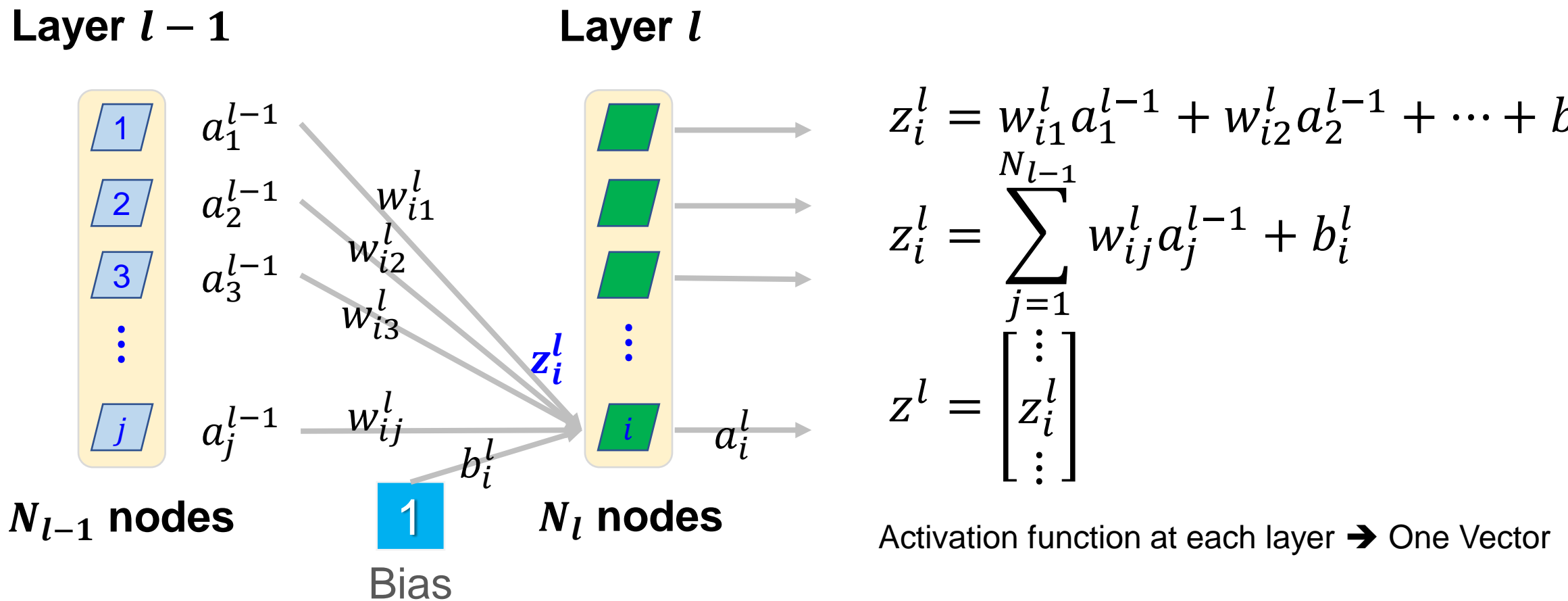$N_l$ **nodes**

$b_i^l$  Bias for neuron $i$ at layer $l$

$$b^l = \begin{bmatrix} \vdots \\ b_i^l \\ \vdots \end{bmatrix}$$

Bias of each layer ➜ One Vector

# Hidden Layer – Activation Function

**Layer $l-1$**                    **Layer $l$**



$$z_i^l = w_{i1}^l a_1^{l-1} + w_{i2}^l a_2^{l-1} + \cdots + b_i^l$$

$$z_i^l = \sum_{j=1}^{N_{l-1}} w_{ij}^l a_j^{l-1} + b_i^l$$

$$z^l = \begin{bmatrix} \vdots \\ z_i^l \\ \vdots \end{bmatrix}$$

**$N_{l-1}$ nodes**          **$N_l$ nodes**

Bias

Activation function at each layer ➔ One Vector

# Hidden Layer – Overall

**Layer** $l - 1$          **Layer** $l$



$$z_1^l = w_{11}^l a_1^{l-1} + w_{12}^l a_2^{l-1} + \cdots + b_1^l$$
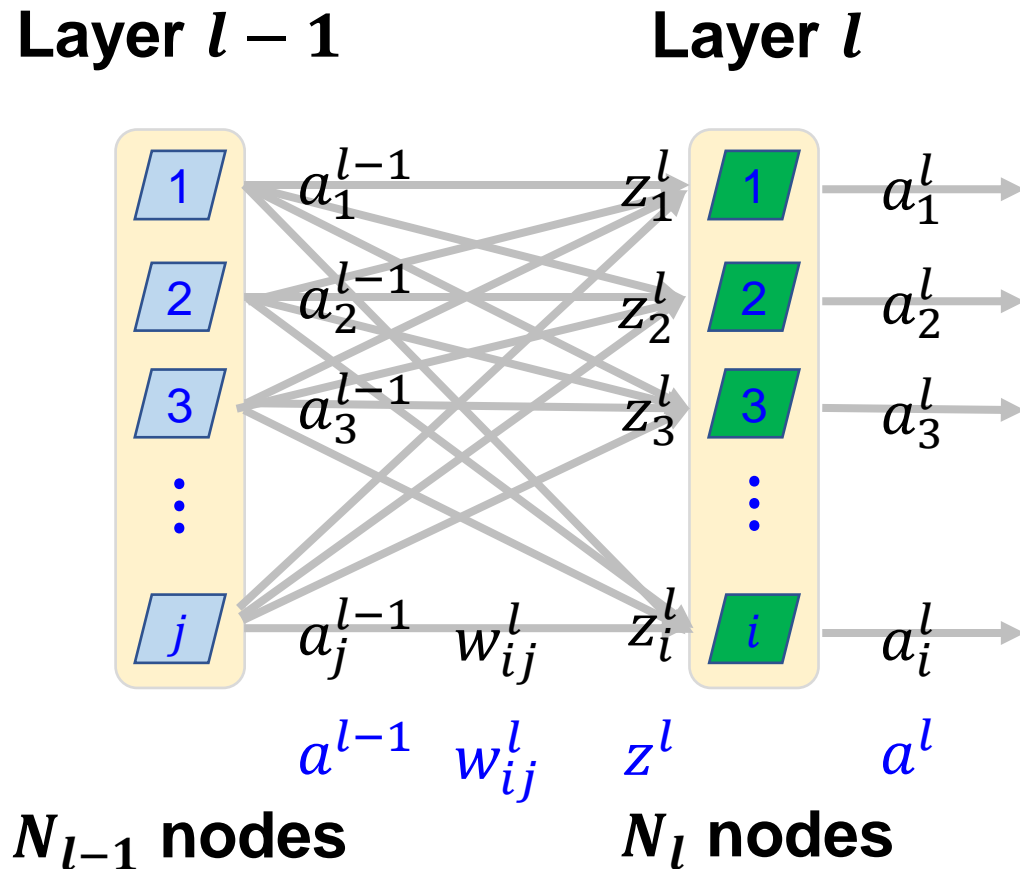
$$\vdots$$

$$z_i^l = w_{i1}^l a_1^{l-1} + w_{i2}^l a_2^{l-1} + \cdots + b_i^l$$

$$\begin{bmatrix} \vdots \\ z_i^l \\ \vdots \end{bmatrix} = \begin{bmatrix} w_{11}^l & w_{12}^l & \cdots \\ w_{21}^l & w_{22}^l & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} \vdots \\ a_i^{l-1} \\ \vdots \end{bmatrix} + \begin{bmatrix} \vdots \\ b_i^l \\ \vdots \end{bmatrix}$$

$$z^l = W^l a^{l-1} + b^l$$

$a^{l-1}$   $w_{ij}^l$   $z^l$        $a^l$

$N_{l-1}$ **nodes**        $N_l$ **nodes**

# Hidden Layer – Overall
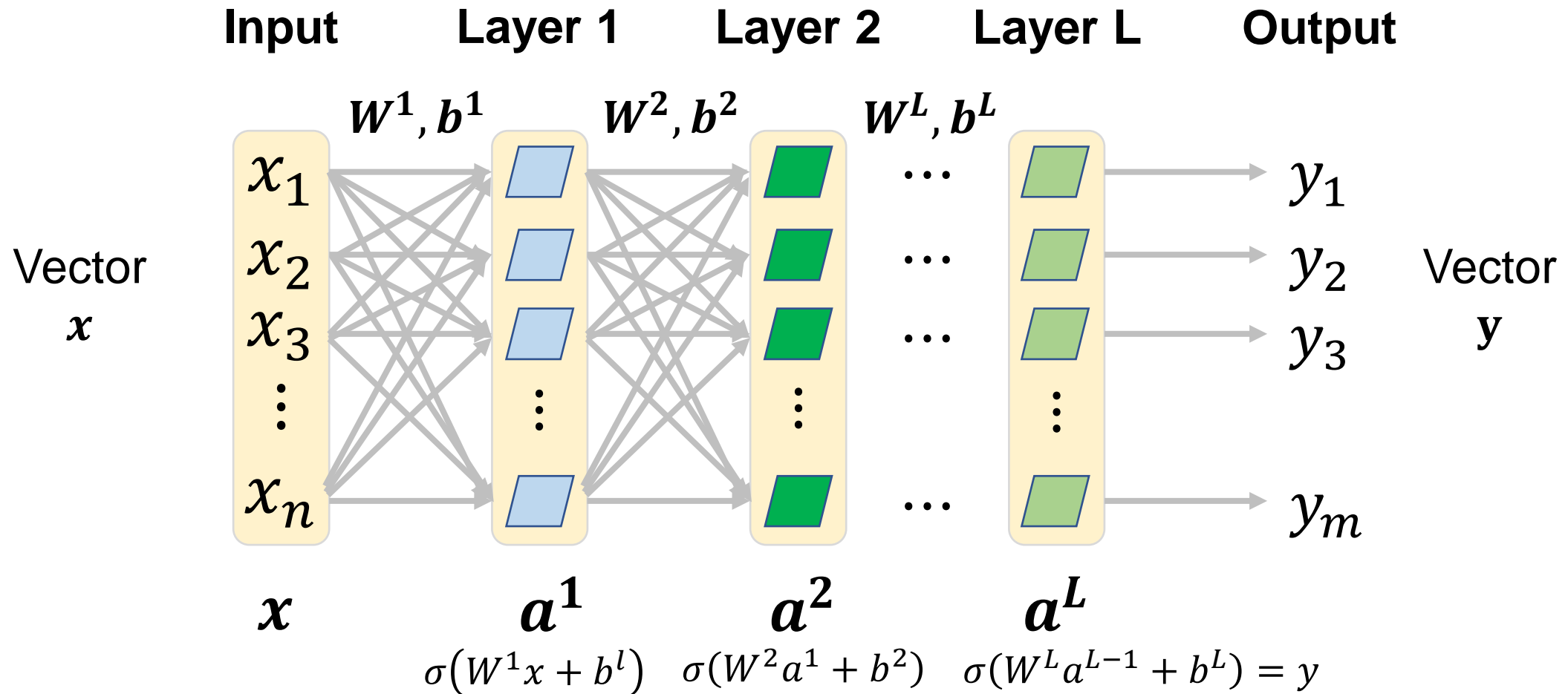
**Layer $l-1$**          **Layer $l$**
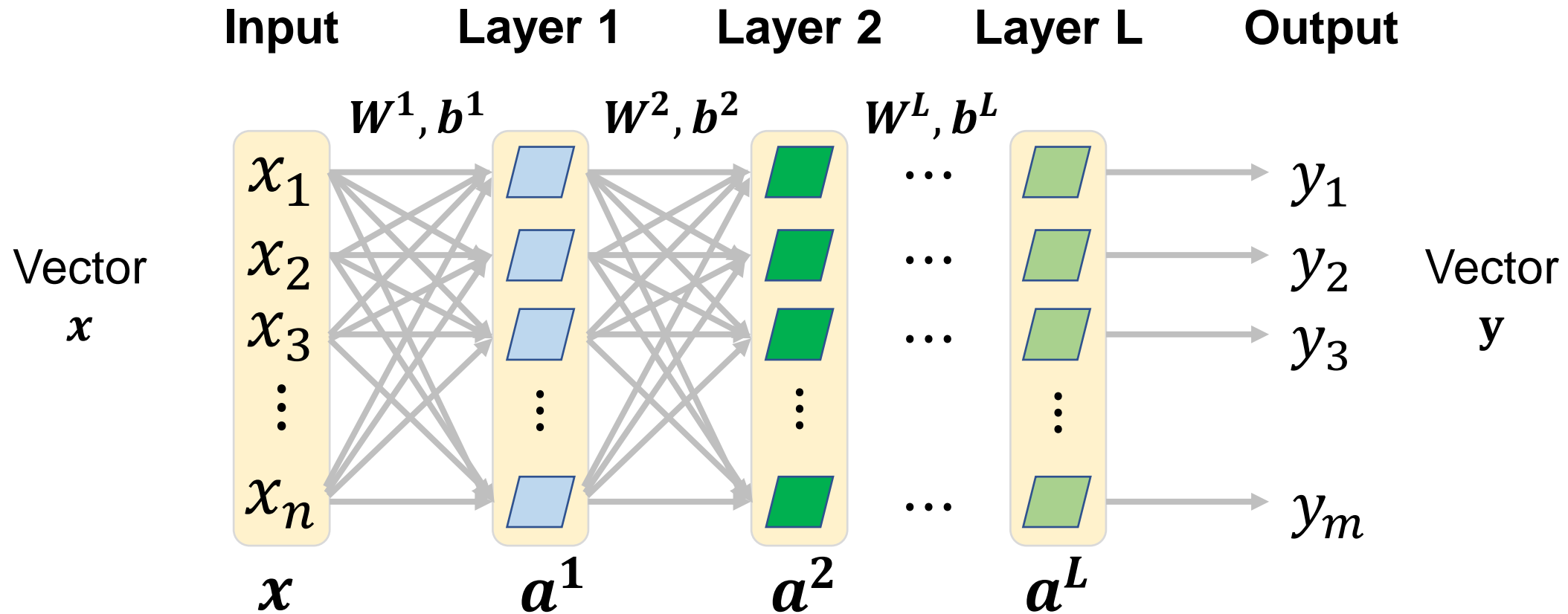


$$a_i^l = \sigma(z_i^l)$$

$$\begin{bmatrix} a_1^l \\ a_2^l \\ \vdots \\ a_i^l \end{bmatrix} = \begin{bmatrix} \sigma(z_1^l) \\ \sigma(z_2^l) \\ \vdots \\ \sigma(z_i^l) \end{bmatrix}$$
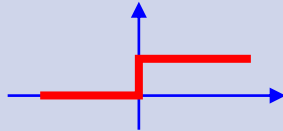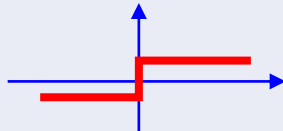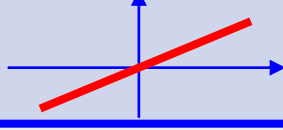
$$a^l = \sigma(z^l)$$

$a^{l-1}$  $w_{ij}^l$  $z^l$     $a^l$

$N_{l-1}$ **nodes**       $N_l$ **nodes**

# Hidden Layer – Overall



**Input      Layer 1      Layer 2      Layer L      Output**

$W^1, b^1$      $W^2, b^2$      $W^L, b^L$

Vector
$x$

$x_1$
$x_2$
$x_3$
$\vdots$
$x_n$

$\cdots$
$\cdots$
$\cdots$
$\vdots$
$\cdots$

$y_1$
$y_2$
$y_3$

$y_m$

Vector
$y$

$x$              $a^1$              $a^2$              $a^L$

$\sigma(W^1 x + b^l)$      $\sigma(W^2 a^1 + b^2)$      $\sigma(W^L a^{L-1} + b^L) = y$

# Hidden Layer – Overall



$$y = f(x) = \sigma(W^L \dots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \dots + b^L)$$

# Activation Function $\sigma(\cdot)$

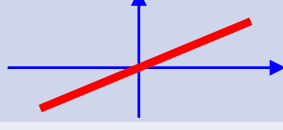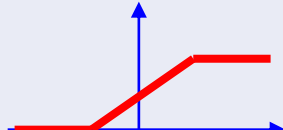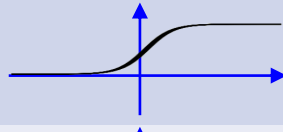| | Activation Function | Equation | Example | 1D Graph |
|---|---|---|---|---|
| **Bounded function** | Unit step (Heatviside) | $\phi(z) = \begin{cases} 0 & z < 0 \\ 0.5 & z = 0 \\ 1 & z > 0 \end{cases}$ | Perceptron variant | |
| | Sign (Signum) | $\phi(z) = \begin{cases} -1 & z < 0 \\ 0 & z = 0 \\ 1 & z > 0 \end{cases}$ | Perceptron variant | |
| | Linear | $\phi(z) = z$ | Adaline, linear regression | |
| **Bounded function** | Piece-wise Linear | $\phi(z) = \begin{cases} 1 & z \geq \frac{1}{2} \\ z + \frac{1}{2} & -\frac{1}{2} < z < \frac{1}{2} \\ 0 & z \leq -\frac{1}{2} \end{cases}$ | Support vector machine | |
| | Logistic (sigmoid) | $\phi(z) = \dfrac{1}{1 + e^{-z}}$ | Logistic regression, multi-layer NN | |
| | Hyperbolic Tangent | $\phi(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$ | Multi-layer NN | |

# Activation Function $\sigma(\cdot)$

| Activation Function | Equation | Example | 1D Graph |
|---|---|---|---|
| **Boolean** Unit step (Heatviside) | $\phi(z) = \begin{cases} 0 & z < 0 \\ 0.5 & z = 0 \\ 1 & z > 0 \end{cases}$ | Perceptron variant | |
| Sign (Signum) | $\phi(z) = \begin{cases} -1 & z < 0 \\ 0 & z = 0 \\ 1 & z > 0 \end{cases}$ | Perceptron variant | |
| **Linear** Linear | $\phi(z) = z$ | Adaline, linear regression | |
| Piece-wise Linear | $\phi(z) = \begin{cases} 1 & z \geq \frac{1}{2} \\ z + \frac{1}{2} & -\frac{1}{2} < z < \frac{1}{2} \\ 0 & z \leq -\frac{1}{2} \end{cases}$ | Support vector machine | |
| **Non-linear** Logistic (sigmoid) | $\phi(z) = \dfrac{1}{1 + e^{-z}}$ | Logistic regression, multi-layer NN | |
| Hyperbolic Tangent | $\phi(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$ | Multi-layer NN | |

# Activation Function (non-linearity)

- **Sigmoid**

$$sigmoid(x) = \frac{1}{1 + e^{-x}}$$

- **Tanh**

$$\tanh(x) = \frac{\sinh(x)}{\cosh(x)} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

- **Rectified Linear Unit (ReLU)**

$$ReLU(x) = \max(x, 0)$$

# Linearity and Non-linearity



https://vitalflux.com/how-know-data-linear-non-linear/

# Linearity and Non-linearity

- With linearity, the deep neural network is the same as linear transform.

$$W_1(W_2 \cdot x) = (W_1 W_2)x = Wx$$

- With non-linearity, the deep neural network with multiple layers could have a more complex function.

# Model Parameters

$$y = f(x) = \sigma(W^L \ldots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \ldots + b^L)$$

Function set                    Different parameters $W$ and $b$ → different functions

## Formal definition

$$f(x; \theta) \Rightarrow model\ parameter\ set$$
$$\theta = \{W^1, b^1, W^2, b^2, \ldots, W^L, b^L\}$$

# Loss Function & Reward Function

- Define a function to measure the quality of parameters set $\theta$
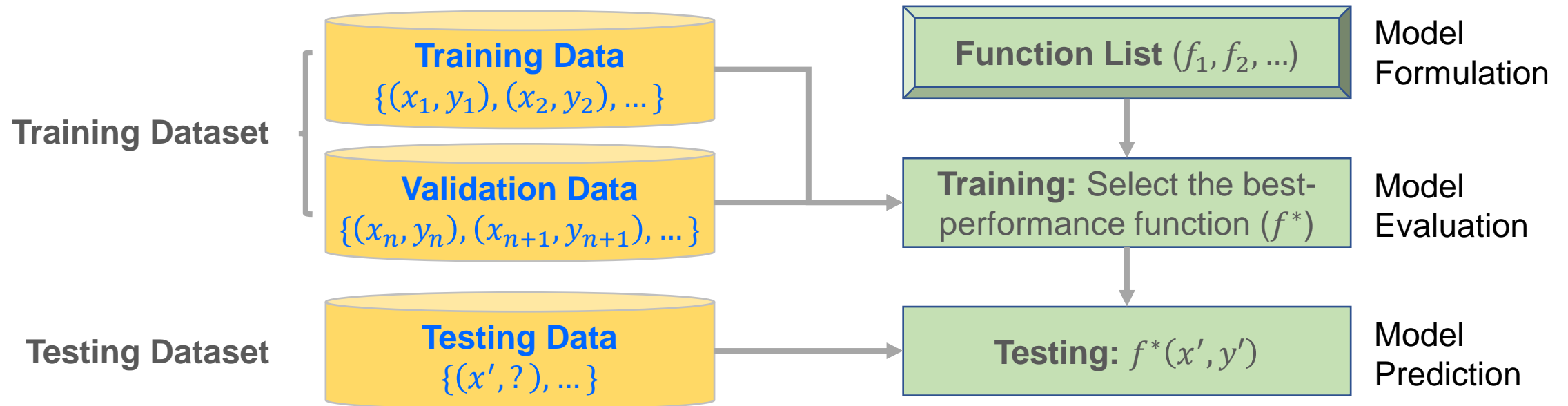  - Evaluating by a loss/cost/error function $C(\theta)$ ➜ how bad

$$\theta^* = \arg \min_{\theta} C(\theta)$$

  - Evaluating by an objective/reward function $O(\theta)$ ➜ how good

$$\theta^* = \arg \max_{\theta} O(\theta)$$

# Loss Function



**Best function:** $f(x; \theta) \sim \hat{y}$ ➜ $\|\hat{y} - f(x; \theta)\| \approx 0$

**Define a loss function:** $C(\theta) = \sum_k \|\hat{y}_k - f(x_k; \theta)\|$

# Loss Function

| Symbol | Name | Equation |
|---|---|---|
| $\mathcal{L}_1$ | L$_1$ loss | $\|y - o\|_1$ |
| $\mathcal{L}_2$ | L$_2$ loss | $\|y - o\|_2^2$ |
| $\mathcal{L}_1 \circ \sigma$ | Expectation loss | $\|y - \sigma(o)\|_1$ |
| $\mathcal{L}_2 \circ \sigma$ | Regularized expectation loss | $\|y - \sigma(o)\|_2^2$ |
| $\mathcal{L}_\infty \circ \sigma$ | Chebyshev loss | $max_j \left| \sigma(o)^j - y^j \right|$ |
| $Hinge$ | Hinge (margin) los | $\sum_j \max\left(0, \frac{1}{2} - \hat{y}^j o^j\right)$ |
| $Hinge^2$ | Squared hinge (margin) loss | $\sum_j \max\left(0, \frac{1}{2} - \hat{y}^j o^j\right)^2$ |
| $Hinge^3$ | Cubed hinge (margin) loss | $\sum_j \max\left(0, \frac{1}{2} - \hat{y}^j o^j\right)^3$ |

**Table 1** List of losses analyzed in this paper. $y$ is true label as one-hot encoding, $\hat{y}$ is true label as $+1/-1$ encoding, $o$ is the output of the last layer of the network, $\cdot$ $(j)$ denotes $j - th$ dimension of a given vector, and $\sigma(\cdot)$ denotes probability estimate.

# Loss Function

| | | |
|---|---|---|
| **Table 1** List of losses analyzed in this paper. $y$ is true label as one-hot encoding, $\hat{y}$ is true label as $+1/-1$ encoding, $o$ is the output of the last layer of the network, $\cdot\ (j)$ denotes $j - th$ dimension of a given vector, and $\sigma(\cdot)$ denotes probability estimate. | | |
| **Symbol** | **Name** | **Equation** |
| log | Log (cross entropy) loss | $-\sum_{j} y^j log\sigma(o)^j$ |
| $\log^2$ | Squared log loss | $-\sum_{j} \left[y^j log\sigma(o)^j\right]^2$ |
| tan | Tanimoto loss | $\dfrac{-\sum_j \sigma(o)^j y^j}{\|\sigma(o)\|_2^2 + \|y\|_2^2 - \sum_j \sigma(o)^j y^j}$ |

# References

- 臺大電機系李宏毅教授講義
- 臺大資工系陳縕儂教授講義
- NVIDIA Deep Learning Tutorial
- https://aima.cs.berkeley.edu/slides-pdf/chapter20b.pdf
- Janocha and Czarnecki (2017) On Loss Functions for Deep Neural Networks in Classification. *Computer Science ArXiv*. abs/1702.05659.

# Thank you for your attention!