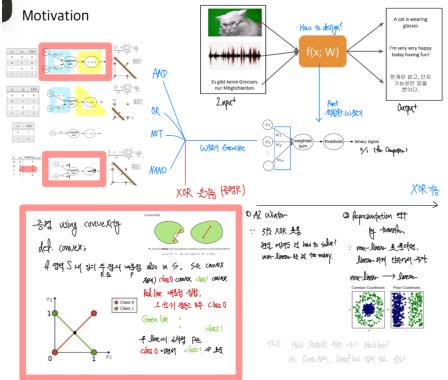


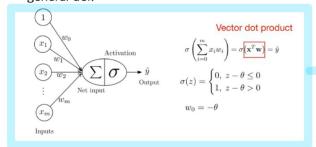
# Lec2. Perceptrons

## Motivation

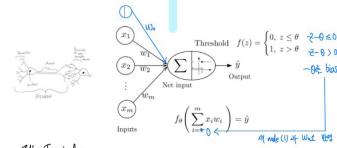


## Multi-layer perceptrons

Motivation: Single-layer perceptrons  
general def.



naive def.



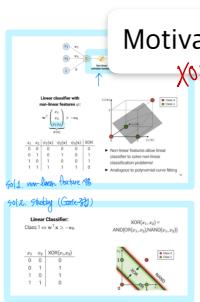
$$\text{New input } \vec{x} = w_0 + w_1x_1 + w_2x_2 + \dots + w_mx_m + b$$

$$= \vec{w}^T \vec{x} + b = \vec{x}^T \vec{w} + b$$

Activations  $0 < \delta_j(x_j) < b$  s.t.  $b$ : activation function (= threshold function in MLP)  
Input weight  $\vec{w}_j = \vec{w}$

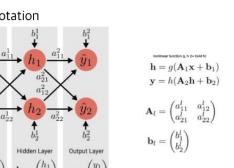
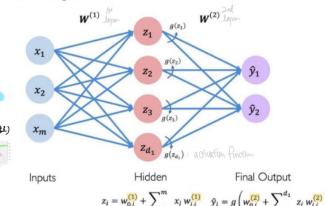
## Motivation

$XOR^{2D \rightarrow 1}$



## Multi-layer Perceptrons

### MLP: diagram

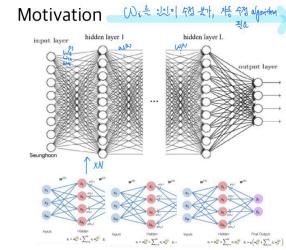


► Example MLP with  $L = 2$  layers of width 2 and parameters  $\{A_1, A_2, b_1, b_2\}$

## Learning MLP

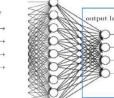
### Motivation

$(W_k \text{는 } \text{次의 } \text{特征, } \text{推 } \text{의 } \text{特征})$



### Learning MLP

### Perceptron Learning Algorithm



Let  $\text{error} = \text{target} - \text{output}$

- i)  $f = 0$ :  $0 - 0 = 0$
- ii)  $f = 1$ :  $1 - 0 = 1$
- iii)  $f = -1$ :  $0 - (-1) = 1$
- iv)  $f = 1$ :  $1 - 1 = 0$

Let  $B = \{x^{(1)}, x^{(2)}, \dots, x^{(N)}\} \in \mathbb{R}^{n \times N}$

- Initialization:  $W \leftarrow 0^{n \times 1}$  (ensure zero bias weight incl. bias)
- For every training example:
  - For every  $x^{(i)} \in B$ :  $\text{target} = y^{(i)}$
  - Calculate  $f = g(W^T x^{(i)})$
  - Compute error:  $\text{error} = (\text{target} - f)^2$
  - Update weight:  $W \leftarrow W - \text{error} \cdot x^{(i)}$

## TMI: vectorization (행렬화)

### (a) Scalar Operation

$$\begin{array}{c} A_0 + B_0 = C_0 \\ A_1 + B_1 = C_1 \\ A_2 + B_2 = C_2 \\ A_3 + B_3 = C_3 \end{array}$$

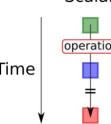
행렬화 44

### (b) SIMD Operation

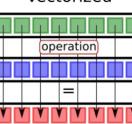
$$\begin{array}{c} A_0 + B_0 = C_0 \\ A_1 + B_1 = C_1 \\ A_2 + B_2 = C_2 \\ A_3 + B_3 = C_3 \end{array}$$

행렬화 14

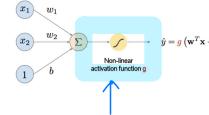
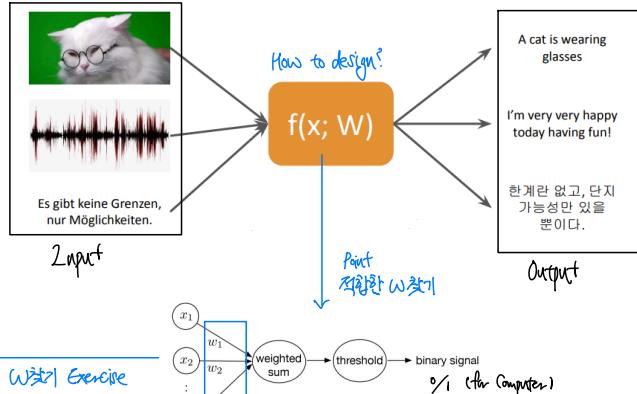
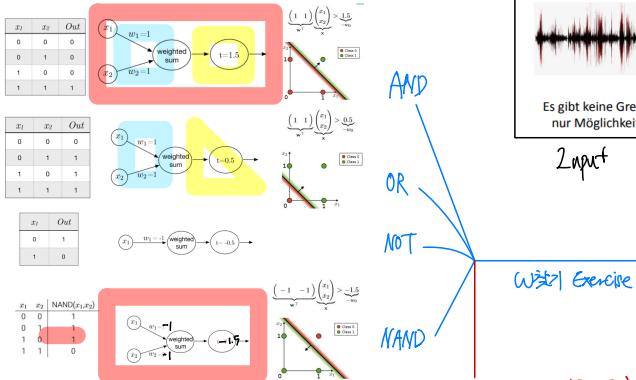
### Scalar



### Vectorized



# Motivation



Linear classifier with non-linear features  $\psi$ :

$$\mathbf{w}^T \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} > -w_0$$

$$\psi(\mathbf{x})$$

$x_1$	$x_2$	$\psi_1(\mathbf{x})$	$\psi_2(\mathbf{x})$	$\psi_3(\mathbf{x})$	XOR
0	0	0	0	0	0
0	1	0	1	0	1
1	0	1	0	0	1
1	1	1	1	1	0

- Non-linear features allow linear classifier to solve non-linear classification problems!
- Analogous to polynomial curve fitting

25

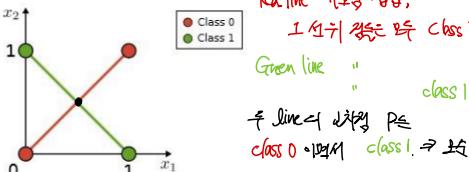
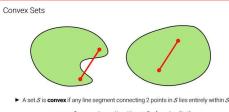
Sol 1. non-linear feature  $\psi$

Sol 2. Stacking (Gate-23)

ZigZag using convexity

def. convex:

if 영역 S 내부 두 점의连线 also in S.  $\Rightarrow$  S convex  
P.P.  $\Rightarrow$  class 0 convex class 1 convex.



A2 Winter

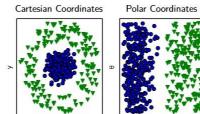
$\therefore$  5% XOR

$\Rightarrow$  5% XOR how to solve?  
non-linear is too many,

② representation off by transform

non-linear is too many,  
linear is not many.

non-linear  $\rightarrow$  linear



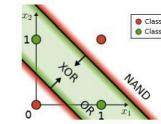
TMI : A2은 대안은 어떤 것인가?

16. Data 3D, Dataflow 3D로 해보니!

Linear Classifier:  
Class 1  $\Leftrightarrow \mathbf{w}^T \mathbf{x} > -w_0$

$XOR(x_1, x_2) =$   
 $\text{AND}(\text{OR}(x_1, x_2), \text{NAND}(x_1, x_2))$

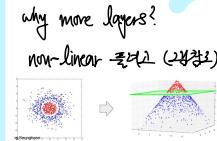
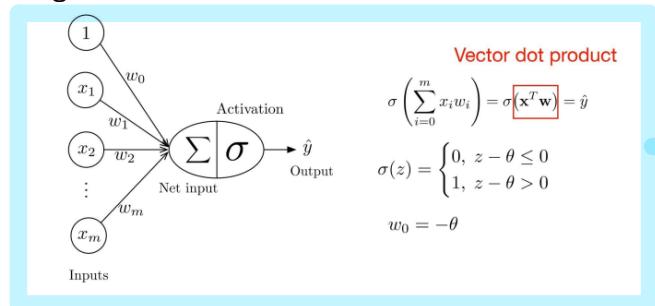
$x_1$	$x_2$	$\text{XOR}(x_1, x_2)$
0	0	0
0	1	1
1	0	1
1	1	0



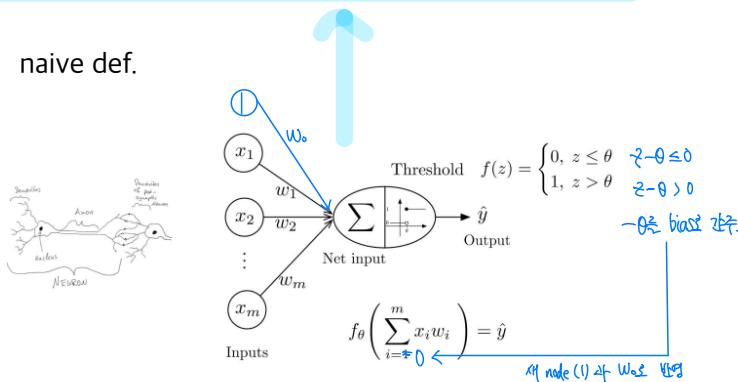
# Multi-layer perceptrons

Motivation: Single-layer perceptrons

general def.



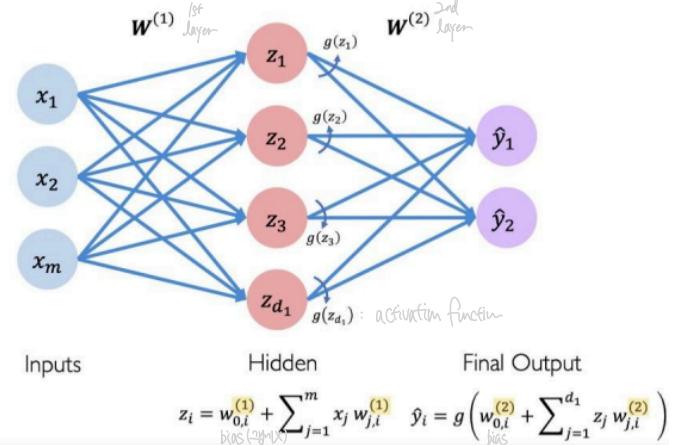
naive def.



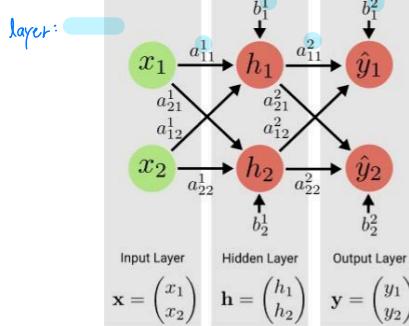
Activations  $a = g(z)$  s.t.  $g$ : activation function (= threshold function in MLP)

Total output  $\hat{y} = f(a)$  s.t.  $f$ : threshold function  $f = g$

MLP: diagram



MLP: matrix notation



$$\mathbf{A}_l = \begin{pmatrix} a_{11}^l & a_{12}^l \\ a_{21}^l & a_{22}^l \end{pmatrix}$$

$$\mathbf{b}_l = \begin{pmatrix} b_1^l \\ b_2^l \end{pmatrix}$$

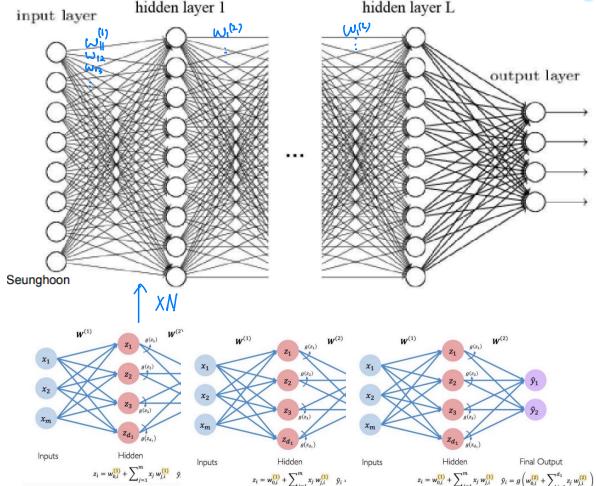
$\times g, h$ : activation function

► Example MLP with  $L = 2$  layers of width 2 and parameters  $\{\mathbf{A}_1, \mathbf{A}_2, \mathbf{b}_1, \mathbf{b}_2\}$

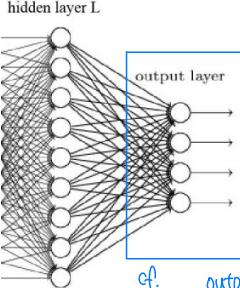
# Learning MLP

## Motivation

W는 뭇 절차로 학습하거나, 자동 학습 algorithm  
작동



## Perceptron Learning Algorithm



cf.	output	target	correct	To do
i)	$\hat{y} = 0$	0	Yes	X
ii)	$\hat{y} = 0$	1	No	$w \leftarrow w + \alpha$
iii)	$\hat{y} = 1$	0	No	$w \leftarrow w - \alpha$
iv)	$\hat{y} = 1$	1	Yes	X

일반화

```

let error = target - output
    i) / - 0 = /
    ii) 1 - 0 = 1
    iii) 0 - 1 = -1
    iv) 1 - 1 = 0

```

$w \leftarrow w + error \times \alpha$

Let

$$\mathcal{D} = (\langle \mathbf{x}^{[1]}, y^{[1]} \rangle, \langle \mathbf{x}^{[2]}, y^{[2]} \rangle, \dots, \langle \mathbf{x}^{[n]}, y^{[n]} \rangle) \in (\mathbb{R}^m \times \{0, 1\})^n$$

1. Initialize  $w := 0^m$  (assume notation where weight incl. bias)

2. For every training epoch:

- A. For every  $\langle \mathbf{x}^{[i]}, y^{[i]} \rangle \in \mathcal{D}$ :  
 (a)  $\hat{y}^{[i]} := \sigma(\mathbf{x}^{[i]^\top} w)$   
 (b) err :=  $(y^{[i]} - \hat{y}^{[i]})$   
 (c)  $w := w + err \times \mathbf{x}^{[i]}$

target :  $y^{[i]}$

output :  $\hat{y}^{[i]}$