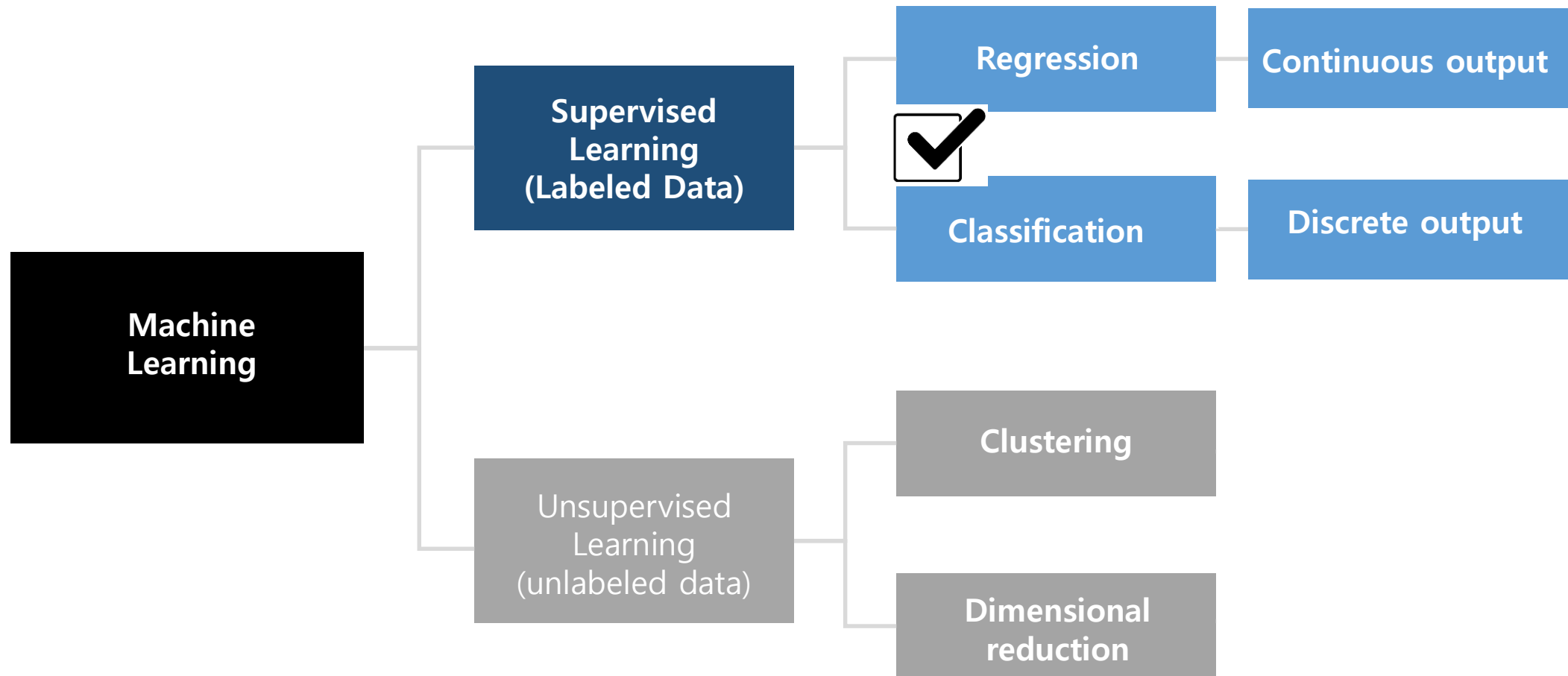# Linear Classification

**Prof. Je-Won Kang**
**Electronic & Electrical Engineering**
**Ewha Womans University**

# Machine learning problems

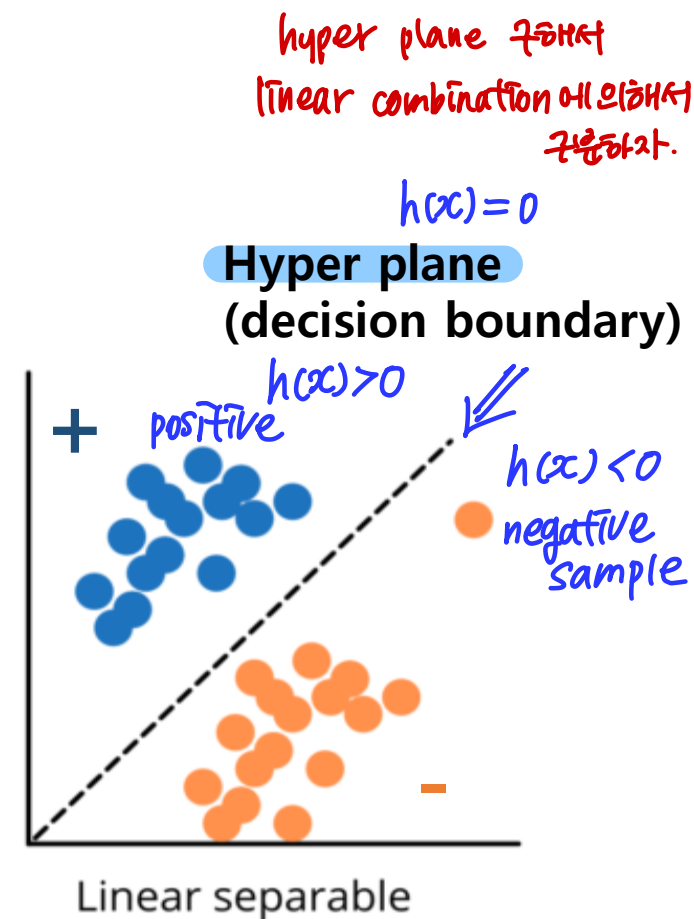# Linear classification

- Predict a discrete output $y$ (classification ID) from $x$ when $D = (x, y)$ is given
  - ID = 0 or 1 (binary classification)
  - ID = 0, 1, ..., $N$-1 (multi classification)

- Hypothesis set $\mathcal{H}$ : a set of lines

model parameter
feature

$$h_w(x) = w_0 + w_1 x_1 + \cdots + w_d x_d = \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}$$

$w$: model parameter (learnable parameter)

$$h_w(x) = w_0 + w_1 \phi(x_1) + \cdots + w_d \phi(x_d) = \boldsymbol{w}^{\mathrm{T}} \phi(\boldsymbol{x})$$

Linear model with a set of features

hyper plane 구해서
linear combination 에 의해서
구분하자.

$h(x)=0$

**Hyper plane**
**(decision boundary)**

$h(x)>0$

+ positive

$h(x)<0$
negative sample

−

Linear separable

# Example: image recognition



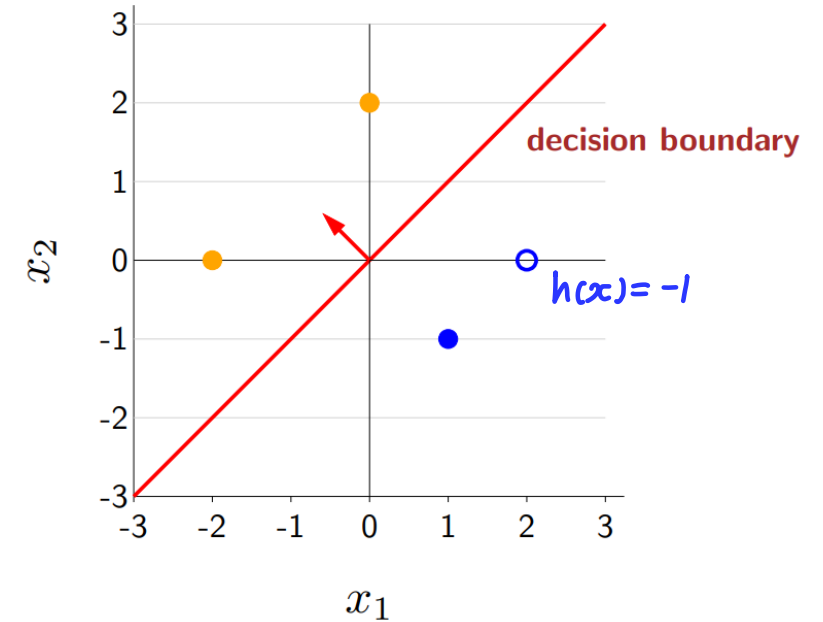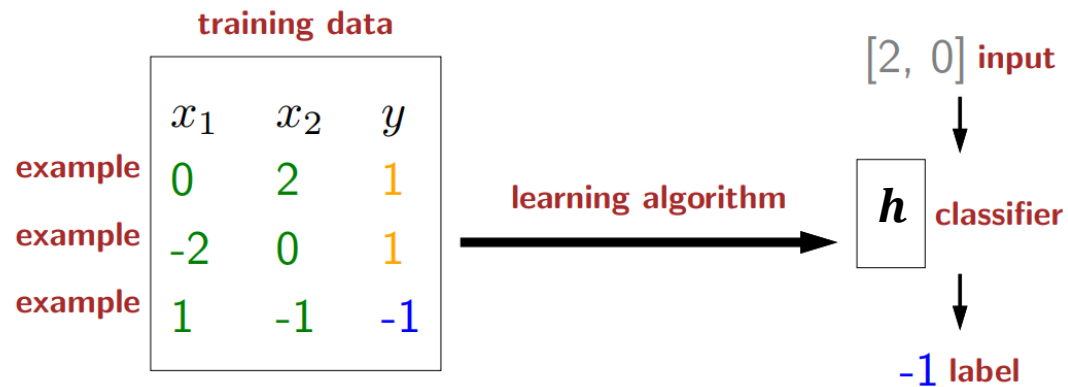hyper plane이 다수 존재

airplane classifier
car classifier
deer classifier

0

# Problem formulation

- $X = \mathbf{R}^d$ is an input space
  - $\mathbf{R}^d$ : a *d*-dimensional Euclidean space
  - input vector $x \in X$: $x = (x_1, x_2, \ldots, x_d)$, e.g. *d = 2*
- $Y = \{+1, -1\}$ is an output space
  - Binary (yes/no) decision
- Now, we want to approximate a target function *f*
  - $f: X \to Y$ (unknown ideal function)
  - Data $(x^1, y^1), \ldots, (x^N, y^N)$ ; dataset where $y^N = f(x^N)$
  - Correct label is ready for a training set
  - **Hypothesis** $g: X \to Y$ (ML model to approximate *f*) : $g \in H$

# Linear classification framework

**Hypothesis function to build a decision boundary**

training data

|  | $x_1$ | $x_2$ | $y$ |
|---|---|---|---|
| example | 0 | 2 | 1 |
| example | -2 | 0 | 1 |
| example | 1 | -1 | -1 |

learning algorithm ⟶

[2, 0] input

↓

$h$ classifier

↓

-1 label

decision boundary

$h(x) = -1$

| Which predictor? Hypothesis class | How good is a predictor? Loss function | How to compute the best predictor? Optimization algorithm |
|---|---|---|

$$h(x) = \text{sign}\,(w^\mathrm{T}x)$$

**Zero-one loss**
**Hinge loss**
**Cross-entropy loss**
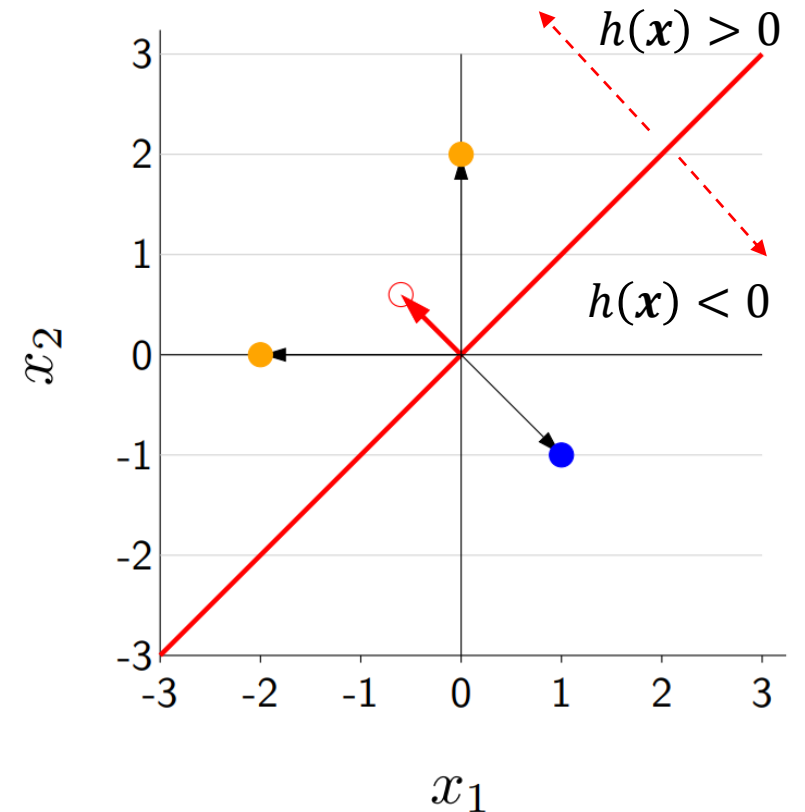
**Gradient descent algorithm**

# Linear classification model

- The linear formula $g \in \boldsymbol{H}$ can be written as

$$h(x) = \text{sign}\left(\left(\sum_{i=1}^{d} w_i x_i\right) + w_0\right)$$

$$= \text{sign}\left(\left(\sum_{i=0}^{d} w_i x_i\right)\right), \ x_0 = 1$$

$$= \text{sign}(w^T x)$$

$w_o : a \ bias \ term$

$$sign(x) = \begin{cases} 1, & if \ x > 0 \\ -1, & if \ x < 0 \end{cases}$$

$x_o : 1$

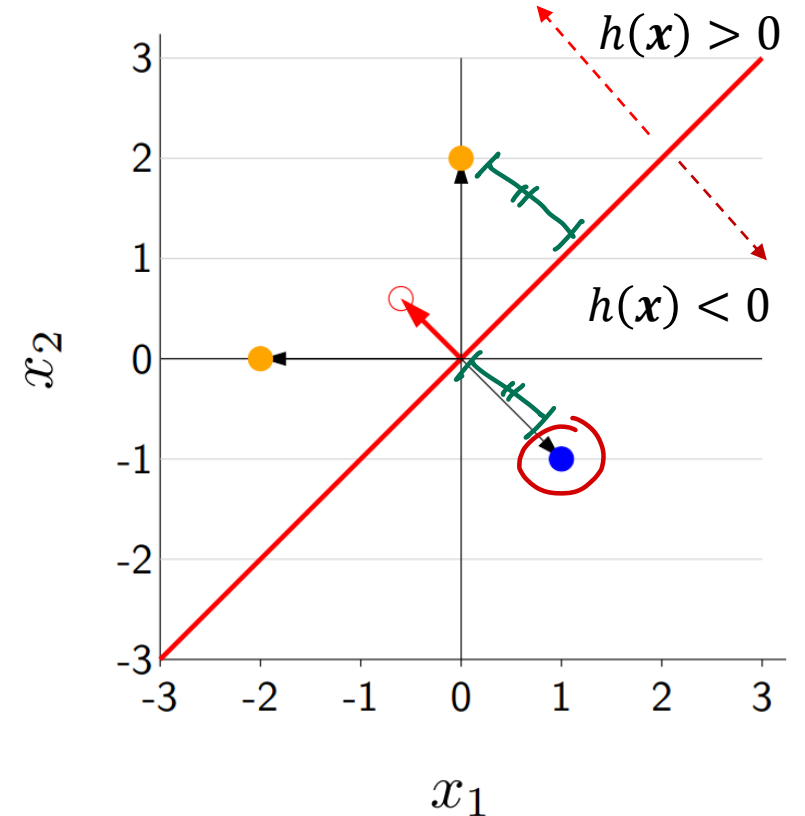$$h(\boldsymbol{x}) = w_0 + w_1 x_1 + w_2 x_2 = 0$$

# Example of linear classifier

$$h(\boldsymbol{x}) = w_0 + w_1 x_1 + w_2 x_2 = 0$$

$$\overbrace{\phantom{xx}}^{\boldsymbol{w}} \overbrace{\phantom{xx}}^{\boldsymbol{\phi}(\boldsymbol{x})}$$

$$h(\boldsymbol{x}) = \text{sign}\left([-1, 1]\, [x_1, x_2]^{\mathrm{T}}\right)$$

$$\text{sign}(z) = \begin{cases} +1 & \text{if } z > 0 \\ -1 & \text{if } z < 0 \\ 0 & \text{if } z = 0 \end{cases}$$

$$h([0,2]) = \text{sign}\left([-1, 1]\, [0, \overset{2}{2}]^{\mathrm{T}}\right) = 1$$

$$h([1,-1]) = \text{sign}\left(\underset{w^{\mathsf{T}}x}{[-1, 1]\, [1, -1]^{\mathrm{T}}}\right) = -1$$

# Hypothesis class : which classifier?

$$h(x) = w_0 + w_1 x_1 + w_2 x_2 = 0$$

$$h(x) = \text{sign}\,([-1, 1]\, [x_1, x_2]^{\text{T}})$$

$$h(x) = \text{sign}\,([0.5, 1]\, [x_1, x_2]^{\text{T}})$$

## For optimization
### Define a metric and compute an error

$$\text{Loss}_{0\text{-}1}(x, y, \mathbf{w}) = \mathbf{1}[f_{\mathbf{w}}(x) \neq y] \quad \text{zero-one loss}$$
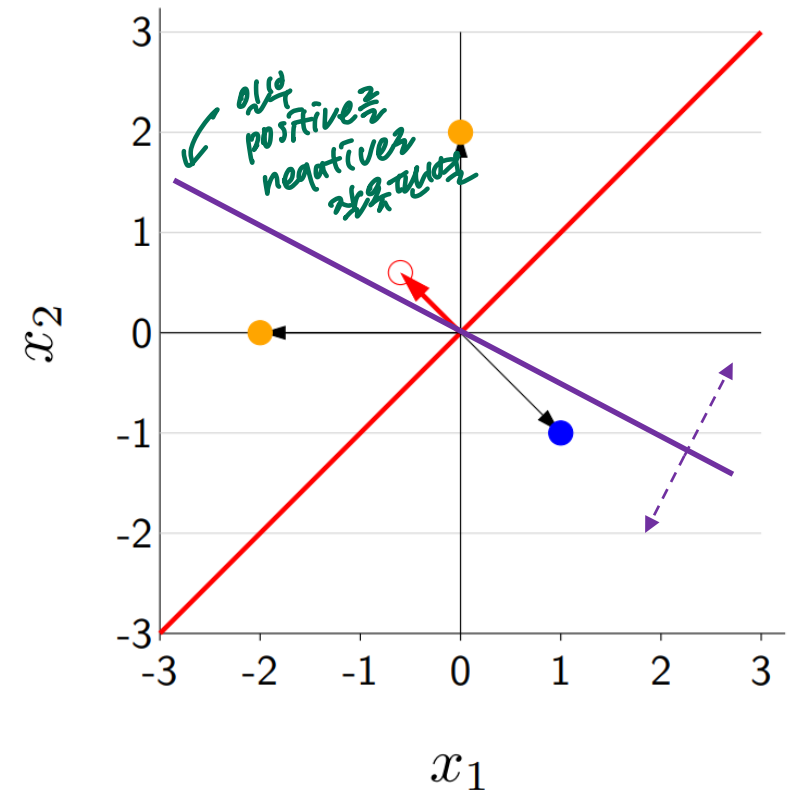
*Zero-one-function*

$$\text{Loss}([0, 2], 1, [0.5, 1]) = \mathbf{1}[\text{sign}([0.5, 1] \cdot [0, 2]) \neq 1] = 0$$

$$\text{Loss}([-2, 0], 1, [0.5, 1]) = \mathbf{1}[\text{sign}([0.5, 1] \cdot [-2, 0]) \neq 1] = 1$$

$$\text{Loss}([1, -1], -1, [0.5, 1]) = \mathbf{1}[\text{sign}([0.5, 1] \cdot [1, -1]) \neq -1] = 0$$



*old positives, negatives 각각 판단점*

# Score and margin

- Input data : $x$
- Predicted label : $h(\boldsymbol{x}) = \text{sign}\,(w^T \phi(\boldsymbol{x}))$
- Target label: $y$

✓ Score : the score on an example $(x, y)$ is $w \cdot \phi(x)$, how **confident** we are in predicting +1.

✓ Margin : the margin on an example $(x, y)$ is $(w \cdot \phi(x))y$, how **correct** we are.

score x y

score · y

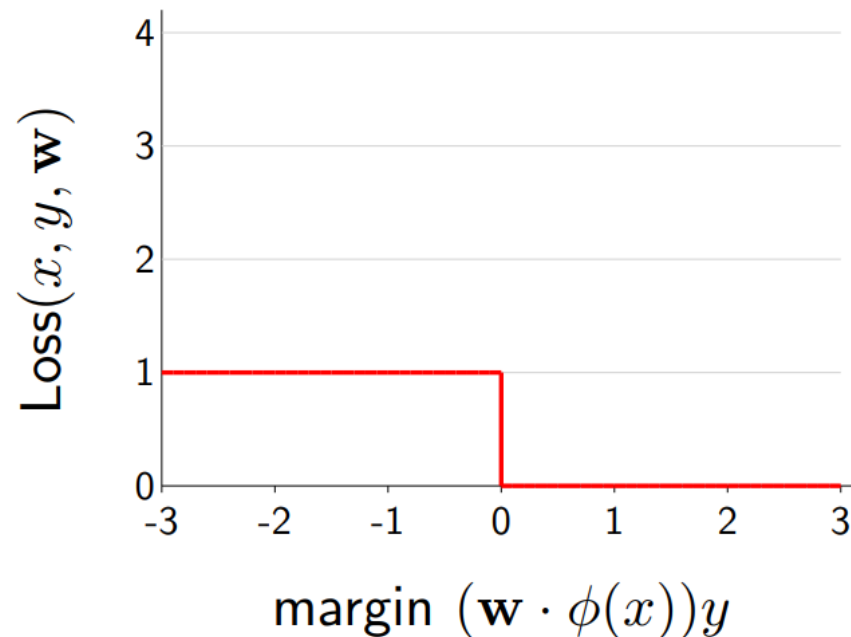| ⊕ | y=1 | → margin 大 → 정답맞춤. |
| ⊖ | y=-1 | → margin 大 → " |
| ⊕ | y=-1 | → margin 小 → 틀림. |

# Zero-one loss

$$\text{Loss}_{0\text{-}1}(x, y, \mathbf{w}) = \mathbf{1}[\underbrace{(\mathbf{w} \cdot \phi(x))y}_{\text{margin}} \le 0]$$



The goal is to minimize the loss

To run gradient descent, compute the gradient:

$$\nabla_{\mathbf{w}} \text{TrainLoss}(\mathbf{w}) = \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{(x,y) \in \mathcal{D}_{\text{train}}} \nabla \text{Loss}_{0\text{-}1}(x, y, \mathbf{w})$$

Gradient is zero almost everywhere!

학습불가

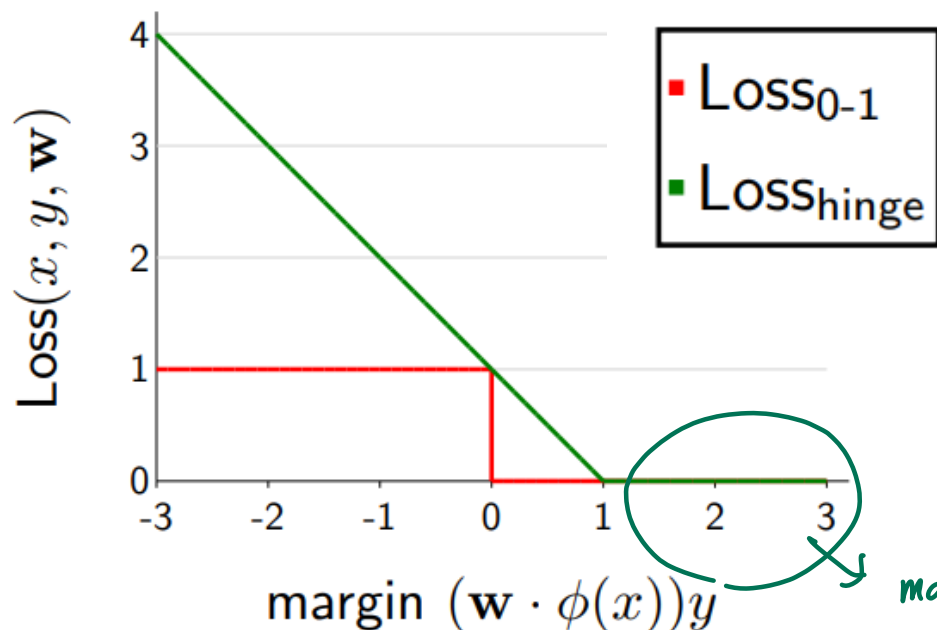$$\nabla_{\mathbf{w}} \text{Loss}_{0\text{-}1}(x, y, \mathbf{w}) = \nabla \mathbf{1}[(\mathbf{w} \cdot \phi(x))y \le 0]$$

# Hinge loss

1-margin ⇒ 정답 값 맞추고 있는 상태면
1-margin 음수없일것.

margin

$$\text{Loss}_{\text{hinge}}(x, y, \mathbf{w}) = \max\{1 - (\mathbf{w} \cdot \phi(x))y, 0\}$$

- Zero loss if it is classified confidently and correctly
- Misclassification incurs a linear penalty w.r.t. confidence



$$\nabla\text{Loss}_{\text{hinge}}(x, y, \mathbf{w}) = \begin{cases} -\phi(x)y & \text{if } 1 > \{(\mathbf{w} \cdot \phi(x))y\} \\ 0 & \text{otherwise} \end{cases}$$

margin

repeat until convergence {
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$
}

margin 이 ⊕ 일때

Slides from regression and classification, CS221 Stanford

# Cross-entropy loss

- Considers two probability mass functions (pmf) $\{p, 1-p\}$ and $\{q, 1-q\}$ with a binary outcomes

- Cross entropy for these two pmfs : defined by

$$D(p \,\|\, q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)}$$

$$p \log \frac{1}{q} + (1-p) \log \frac{1}{1-q}$$

**Cf.** $\sum_{x \in X} p(x) \log \frac{1}{q(x)} = \left(H(p)\right) + D(p \,\|\, q)$ ← P 와 Q의 유사도에 따라

↳ 거의 안바뀜.

↳ K-L divergence

값 바뀜.

P 와 Q 유사하면

loss 줄기들다

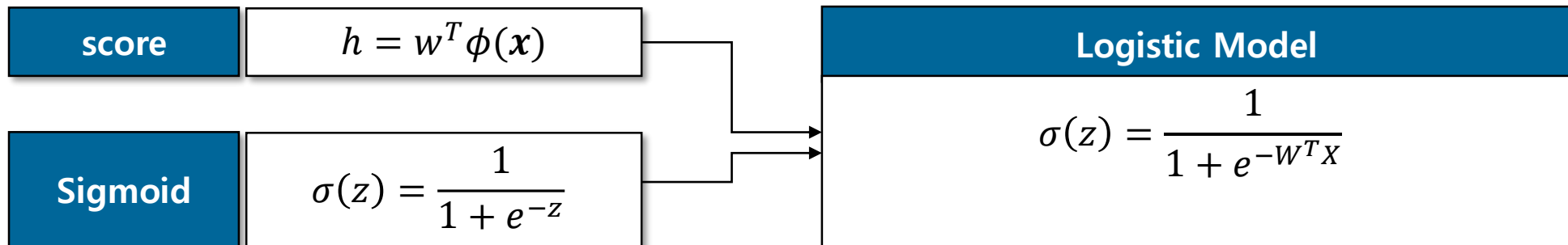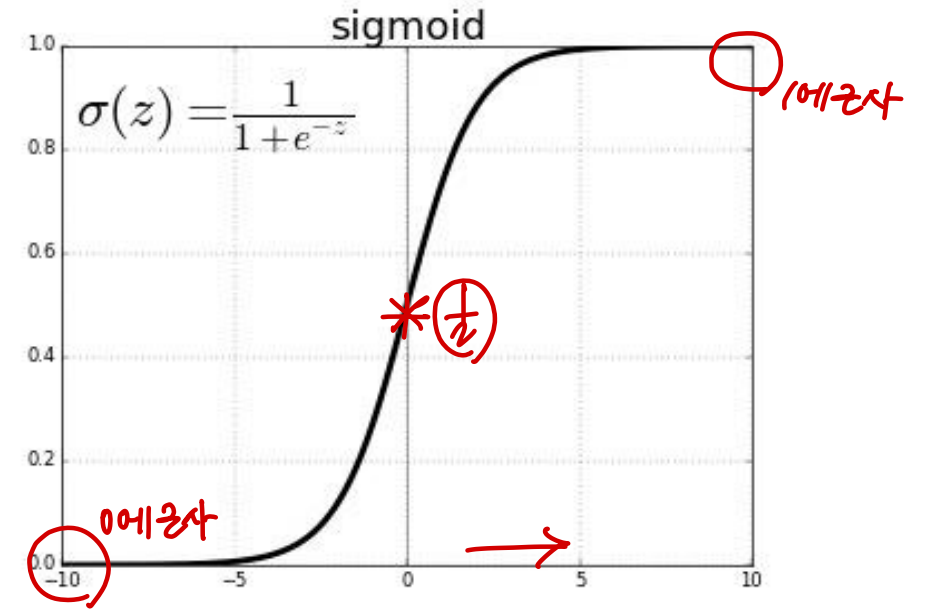Kullback-Leibler (K-L) divergence is a measure of dissimilarity of two distributions

- Cross entropy measures the error when approximating an observed pmf $\{p, 1-p\}$ between a fitted pmf $\{q, 1-q\}$

# Cross-entropy loss

Real value $h = w^T \phi(x)$    0 or 1

| Height (cms) | Weight (kg) | Fitness |
|---|---|---|
| 150 | 50 | Fit |
| 187 | 75 | Fit |
| 156 | 80 | Not Fit |
| 163 | 60 | Fit |
| 170 | 49 | Not Fit |
| 179 | 70 | Fit |

**sigmoid**

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



| **score** | $h = w^T \phi(x)$ |
|---|---|
| **Sigmoid** | $\sigma(z) = \dfrac{1}{1+e^{-z}}$ |

**Logistic Model**

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

# Sigmoid function

- Squash the output of the linear function

$$\sigma(-\mathrm{w}^T x) = \frac{1}{1 + e^{-\mathrm{w}^T x}}$$



sigmoid

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

- A better approach : interpret as a probability

$$P_{\mathrm{w}}(y = 1|x) = \sigma(-\mathrm{w}^T x) = \frac{1}{1 + e^{-\mathrm{w}^T x}}$$

$$P_{\mathrm{w}}(y = 0|x) = 1 - \sigma(-\mathrm{w}^T x) = \frac{e^{-\mathrm{w}^T x}}{1 + e^{-\mathrm{w}^T x}}$$

# Cross-entropy loss

| $\bar{Y}$ (estimated) |
|:---:|
| 0.7 |
| 0.2 |
| 0.1 |

*더 정확한 의미에서는* (handwritten, red)

| 0.9 |
|:---:|
| 0.1 |
| 0.0 |

$$\text{CE}(S, L) = -\sum_{\forall i} L_i \log(S_i)$$

$\sum y_i = 1$ (handwritten, red)

| $L = Y$ (real) |
|:---:|
| 1.0 |
| 0.0 |
| 0.0 |

## Understanding this Cost Function

- Suppose that $L = [1,0,0]$,
  - If $\bar{Y} = [1,0,0]$, then $D = -1 \cdot \log 1 - 0 \cdot \log 0 - 0 \cdot \log 0 = -1 \cdot 0 - 0 \cdot (-\infty) - 0 \cdot (-\infty) = 0$ (no cost).
  - If $\bar{Y} = [0,1,0]$, then $D = -1 \cdot \log 0 - 0 \cdot \log 1 - 0 \cdot \log 0 = -1 \cdot (-\infty) - 0 \cdot 0 - 0 \cdot (-\infty) = \infty$ (huge cost).
  - If $\bar{Y} = [0,0,1]$, then $D = -1 \cdot \log 0 - 0 \cdot \log 0 - 0 \cdot \log 1 = -1 \cdot (-\infty) - 0 \cdot (-\infty) - 0 \cdot 0 = \infty$ (huge cost).

## Gradient Descent Method

$$W \leftarrow W - \alpha \frac{\partial}{\partial W} \text{CE}$$

# Training a linear classifier

- Iterative optimization using gradient descent

1. **Initialize weights at time step $t = 0$**
2. **Compute the gradients**

$$\nabla E_{Train}(w_t) = -\frac{1}{N}\sum_{n=1}^{N}\frac{y_n x_n}{1 + e^{-y_n w_t^T x_n}}$$

3. Set the direction to move :

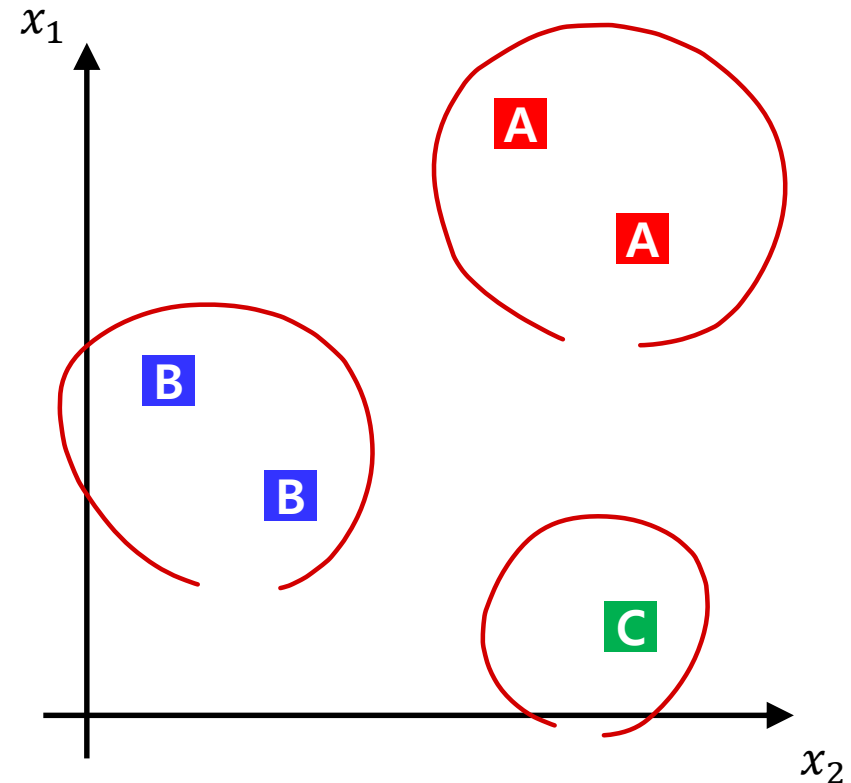$$v_t = -\nabla E_{Train}(w_t)$$

4. Update weights

$$w_{t+1} = w_t + \alpha v_t$$

5. Iterate to next step until converging

# Multiclass classification

- Not all classification predictive models support multi-class classification.
- split the multi-class classification dataset into multiple binary classification datasets and fit a binary classification model on each.
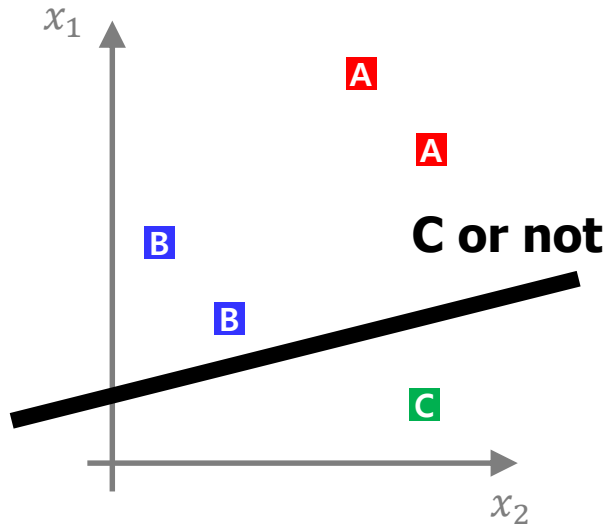
# Multiclass classification
## One-VS-All



multiclass를 binary로

score

$x_1$

A
A
B
**C or not**
B
C

$x_2$

$X \rightarrow \boxed{W} \rightarrow \boxed{f(t) = \frac{1}{1+e^{-t}}} \rightarrow \bar{Y}$

$x_1$

A
A
B
**B or not**
B
C

$x_2$

$X \rightarrow \boxed{W} \rightarrow \boxed{f(t) = \frac{1}{1+e^{-t}}} \rightarrow \bar{Y}$

$x_1$

A
A
B
**A or not**
B
C

$x_2$

$X \rightarrow \boxed{W} \rightarrow \boxed{f(t) = \frac{1}{1+e^{-t}}} \rightarrow \bar{Y}$

$$\begin{array}{ccc} A & B & C \end{array}$$

$$\begin{array}{ccc} A & & \\ \end{array}$$

$$\begin{pmatrix} x_1 & x_2 \end{pmatrix} \cdot \begin{pmatrix} W_{A1} & W_{B1} & W_{C1} \\ W_{A2} & W_{B2} & W_{C2} \end{pmatrix} = \left( \boxed{x_1 \cdot W_{A1} + x_2 \cdot W_{A2}} \quad \boxed{x_1 \cdot W_{B1} + x_2 \cdot W_{B2}} \quad \boxed{x_1 \cdot W_{C1} + x_2 \cdot W_{C2}} \right)$$

# Multiclass classification

$$WX = \begin{pmatrix} \bar{Y}_A \\ \bar{Y}_B \\ \bar{Y}_C \end{pmatrix}$$

$$f(t) = \frac{1}{1+e^{-t}}$$

**A** $sigmoid(A)$
**B** $sigmoid(B)$
**C** $sigmoid(C)$

proportions

**PRPBABILITY**

**A** $p_A = 0.7$
**B** $p_B = 0.2$
**C** $p_C = 0.1$

- $0 \leq p_A, p_B, p_C \leq 1$
- $p_A + p_B + p_C = 1$

One-Hot Encoding

**A** 1.0
**B** 0.0
**C** 0.0

**A**

# Advantage of linear classification

- Simple!
- Interpretability (example in Murphy 2012) 해석가능성 증가
  - $x_1$: the number of cigarettes per day , $x_2$: minutes of exercise per day
  - The goal is to predict $P(Y = \text{lung cancer})$
  - Assume we have estimated the best parameter $w = (1.3, -1.1)$ to have $h(x) = 1.3x_1 - 1.1x_2$

  $\longrightarrow$ For every cigarettes per day, the risk increased by a factor of $e^{1.3}$

  $$\frac{p(y=+1\,|\,x)}{p(y=-1\,|\,x)} = e^{w^T x} = e^{w_1 x_1 + w_2 x_2}$$

# Reference

- Book: Pattern Recognition and Machine Learning (by Christopher M. Bishop)
- Book: Machine Learning: a Probabilistic Perspective (by Kevin P. Murphy)
- https://www.andrewng.org/courses/