

# Statistical Machine Learning

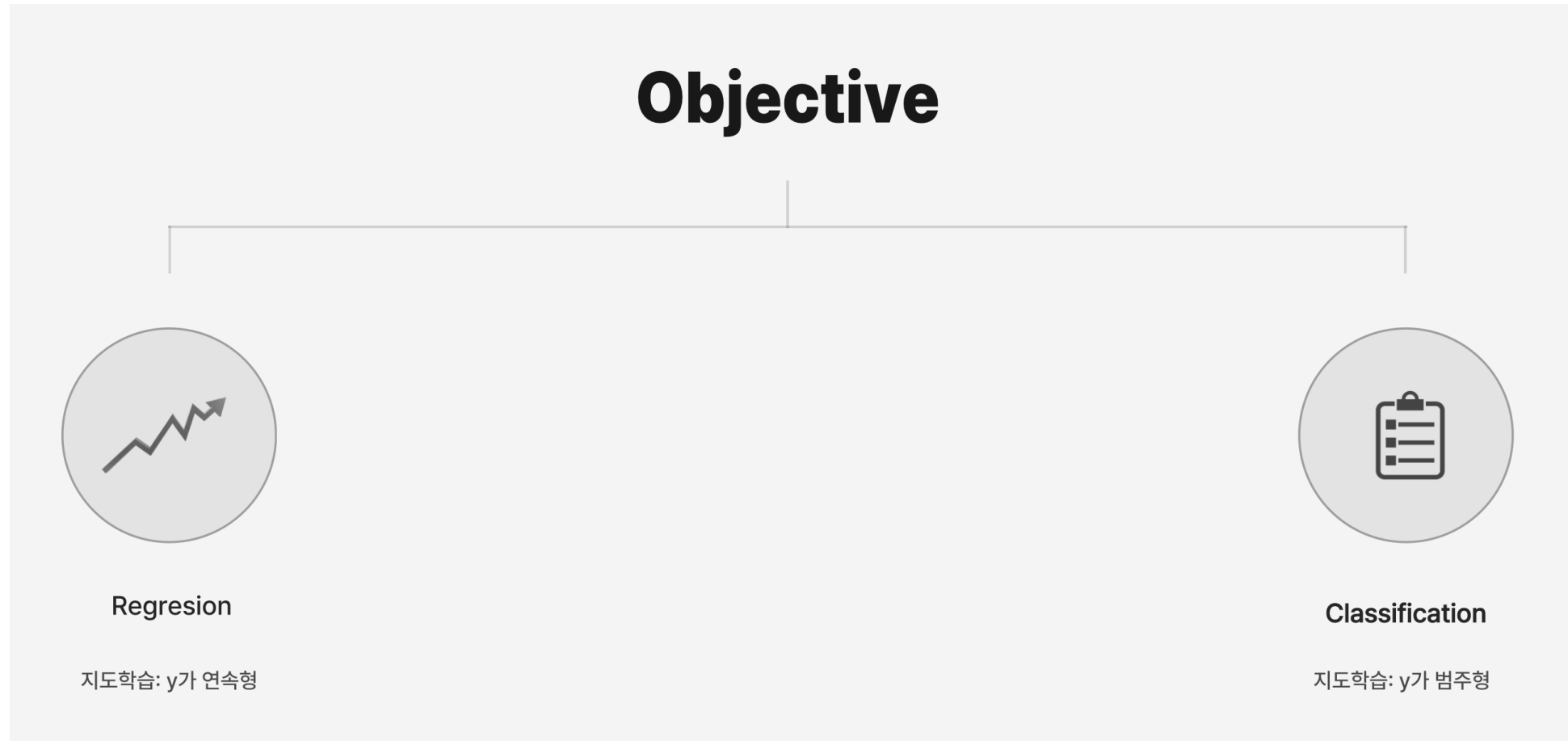
1주차

담당: 15기 염윤석

# Classification

# Classification

---



**1. Bayesian Decision Theory**

**2. Parametric Method**

**3. Non-parametric Method**

**4. Model Evaluation**

# 1. Bayesian Decision Theory

# Bayes' Rule

$$\textit{posterior} \rightarrow P(C | \mathbf{x}) = \frac{\overset{\textit{prior}}{P(C)} \overset{\textit{likelihood}}{p(\mathbf{x} | C)}}{\underset{\textit{evidence}}{p(\mathbf{x})}}$$

$$P(C=0) + P(C=1) = 1$$

$$p(X) = p(X | C=1)P(C=1) + p(X | C=0)P(C=0)$$

$$p(C=0 | X) + p(C=1 | X) = 1$$

$$X = \{x_1, x_2\}$$

$$\text{choose} \begin{cases} C=1 & \text{if } P(C=1 | x_1, x_2) > 0.5 \\ C=0 & \text{otherwise} \end{cases}$$

or

$$\text{choose} \begin{cases} C=1 & \text{if } P(C=1 | x_1, x_2) > P(C=0 | x_1, x_2) \\ C=0 & \text{otherwise} \end{cases}$$

# Bayes' Rule ( $K > 2$ classes)

---

$$P(C_i | \mathbf{x}) = \frac{p(\mathbf{x} | C_i)P(C_i)}{p(\mathbf{x})}$$
$$= \frac{p(\mathbf{x} | C_i)P(C_i)}{\sum_{k=1}^K p(\mathbf{x} | C_k)P(C_k)}$$

$$P(C_i) \geq 0 \text{ and } \sum_{i=1}^K P(C_i) = 1$$

Choose  $C_i$  if  $P(C_i | X) = \max_k P(C_k | X)$

## 2. Parametric Method



## 2-1. Naïve Bayes Classifier

# Parametric Estimation

---

$$P(C_{i.}|X) = \frac{P(X|C_{i.})P(C_{i.})}{P(X)} = \frac{P(X|C_{i.})P(C_{i.})}{\sum_{k=1}^K P(X|C_{k.})P(C_{k.})}$$

**Discriminant** :  $g_i(x) = P(X|C_{i.})P(C_{i.}) \rightarrow g_i(x) = \log_2 P(X|C_{i.}) + \log_2 P(C_{i.})$

Do know about the exact distribution?  $\rightarrow$  **need Estimation!**

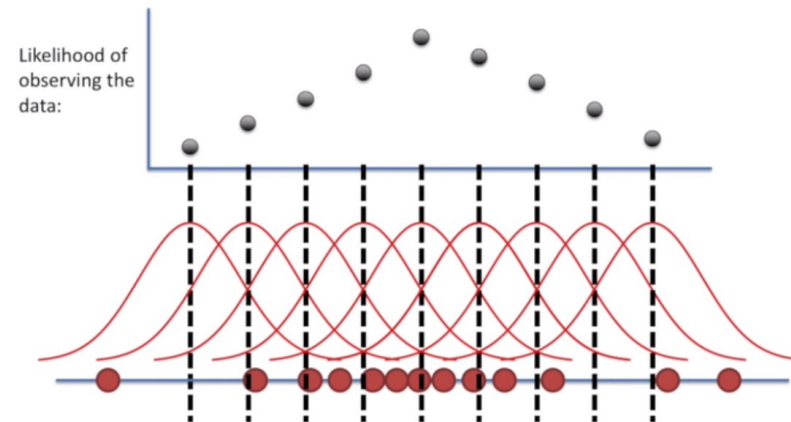
# Back to MLE

$$P(C_{i.}|X) = \frac{P(X|C_{i.})P(C_{i.})}{P(X)} = \frac{P(X|C_{i.})P(C_{i.})}{\sum_{k=1}^K P(X|C_{k.})P(C_{k.})}$$

**Discriminant** :  $g_i(x) = P(X|C_{i.})P(C_{i.}) \rightarrow g_i(x) = \log_2 P(X|C_{i.}) + \log_2 P(C_{i.})$

Do know about the exact distribution?  $\rightarrow$  **need Estimation!**

$$\begin{aligned}\theta_{MLE} &= \arg \max_{\theta} \log P(X|\theta) \\ &= \arg \max_{\theta} \log \prod_i P(x_i|\theta) \\ &= \arg \max_{\theta} \sum_i \log P(x_i|\theta)\end{aligned}$$



# Log Likelihood Function

---

- **Bernoulli distribution**

$$\log L(p) = \sum_{i=1}^n (y_i \log p + (1 - y_i) \log (1 - p))$$

- **Binomial distribution**

$$\log L(p) = \log \binom{n}{c} + \sum_{i=1}^n (y_i \log p + (1 - y_i) \log (1 - p))$$

- **Multinomial distribution**

$$\log L(p) = \sum_{i=1}^n \sum_{j=1}^c y_{ij} \log p_j$$

- **Normal distribution**

$$\log L(\mu) \approx - \frac{\sum_{i=1}^n (y_i - \mu)}{\sigma^2}$$

# Parametric Estimation

$$P(C_i|X) = \frac{P(X|C_i)P(C_i)}{P(X)} = \frac{P(X|C_i)P(C_i)}{\sum_{k=1}^K P(X|C_k)P(C_k)}$$

**Discriminant** :  $g_i(x) = P(X|C_i)P(C_i) \rightarrow g_i(x) = \log_2 P(X|C_i) + \log_2 P(C_i)$

Example >  $P(X|C_i) \sim$  Gaussian Distribution

$$P(X|C_i) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

→ **MLE** for  $\mu$  &  $\sigma$

- $m = \frac{\sum_t x^t}{N}$
- $s^2 = \frac{\sum_t (x^t - m)^2}{N}$

$$g_i(x) = -\frac{1}{2} \log 2\pi - \log \sigma_i - \frac{(x - \mu_i)^2}{2\sigma_i^2} + \log P(C_i)$$



$$g_i(x) = -\frac{1}{2} \log 2\pi - \log s_i - \frac{(x - m_i)^2}{2s_i^2} + \log \hat{P}(C_i)$$

Choose  $C_i$  if  $P(C_i | X) = \max_k P(C_k | X) = \max_k g_k(x)$

# Parametric Estimation

$$P(C_i|X) = \frac{P(X|C_i)P(C_i)}{P(X)} = \frac{P(X|C_i)P(C_i)}{\sum_{k=1}^K P(X|C_k)P(C_k)}$$

**Discriminant** :  $g_i(x) = P(X|C_i)P(C_i) \rightarrow g_i(x) = \log_2 P(X|C_i) + \log_2 P(C_i)$

Example >  $P(X|C_i) \sim$  Gaussian Distribution

$$P(X|C_i) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

→ **MLE** for  $\mu$  &  $\sigma$

- $m = \frac{\sum_t x^t}{N}$
- $s^2 = \frac{\sum_t (x^t - m)^2}{N}$

$$g_i(x) = -\frac{1}{2} \log 2\pi - \log \sigma_i - \frac{(x - \mu_i)^2}{2\sigma_i^2} + \log P(C_i)$$



$$g_i(x) = -\frac{1}{2} \log 2\pi - \log s_i - \frac{(x - m_i)^2}{2s_i^2} + \log \hat{P}(C_i)$$

Choose  $C_i$  if  $P(C_i | X) = \max_k P(C_k | X) = \max_k g_k(x)$

# Parametric Estimation

---

$$P(C_i|X) = \frac{P(X|C_i)P(C_i)}{P(X)} = \frac{P(X|C_i)P(C_i)}{\sum_{k=1}^K P(X|C_k)P(C_k)}$$

**Discriminant** :  $g_i(x) = P(X|C_i)P(C_i) \rightarrow g_i(x) = \log_2 P(X|C_i) + \log_2 P(C_i)$

Example >  $P(X|C_i) \sim \text{Bernoulli}, X = \{0,1\}$

$$P(X|C_i) = p^X(1-p)^{(1-X)}$$

→ **MLE** for  $p$

- $p = \frac{\sum_t x^t}{N}$

$$g_i(x) = \log \prod_t p^{x^t}(1-p)^{(1-x^t)} + \log_2 P(C_i)$$



Choose  $C_i$  if  $P(C_i | X) = \max_k P(C_k|X) = \max_k g_k(x)$

# Parametric Estimation

$$P(C_i|X) = \frac{P(X|C_i)P(C_i)}{P(X)} = \frac{P(X|C_i)P(C_i)}{\sum_{k=1}^K P(X|C_k)P(C_k)}$$

**Discriminant**:  $g_i(x) = P(X|C_i)P(C_i) \rightarrow g_i(x) = \log_2 P(X|C_i) + \log_2 P(C_i)$

Example >  $P(X|C_i) \sim$  Multinomial,  $X_j = \{0,1\}$   
( $X = \{X_1, X_2, X_3, \dots, X_K\} \mid K > 2$ )

$$P(X_1, X_2, X_3, \dots, X_K | C_i) = \prod_j p_j^{X_j}$$

$\rightarrow$  **MLE** for  $p_i$

- $p_j = \frac{\sum_t X_j^t}{N}$

$$g_i(x) = \log \prod_t \prod_j p_j^{X_j^t} + \log_2 P(C_i)$$



Choose  $C_i$  if  $P(C_i | X) = \max_k P(C_k | X) = \max_k g_k(x)$



# Naïve Bayes Classifier

---

Assume **Independent** among attributes  $X_j$  when class  $C_i$  is given

**Discriminant** :  $g_i(x) = P(X|C_i)P(C_i) = P(C_i) \prod_j P(X_j|C_i)$   
 $\rightarrow \log_2 P(C_i) + \sum_j P(X_j|C_i)$

Discrete  $X_j$

$\rightarrow$  Bernoulli or Multinomial

Continuous  $X_j$

$\rightarrow$  Gaussian (Normal) distribution

- Robust to isolated noise points
- Handle missing values by ignoring the instance during estimation
- Robust to irrelevant attributes
- Independence assumption may not hold for some attributes  $\rightarrow$  BBN(Bayesian Belief Networks)

## 2-2. Linear Discriminant

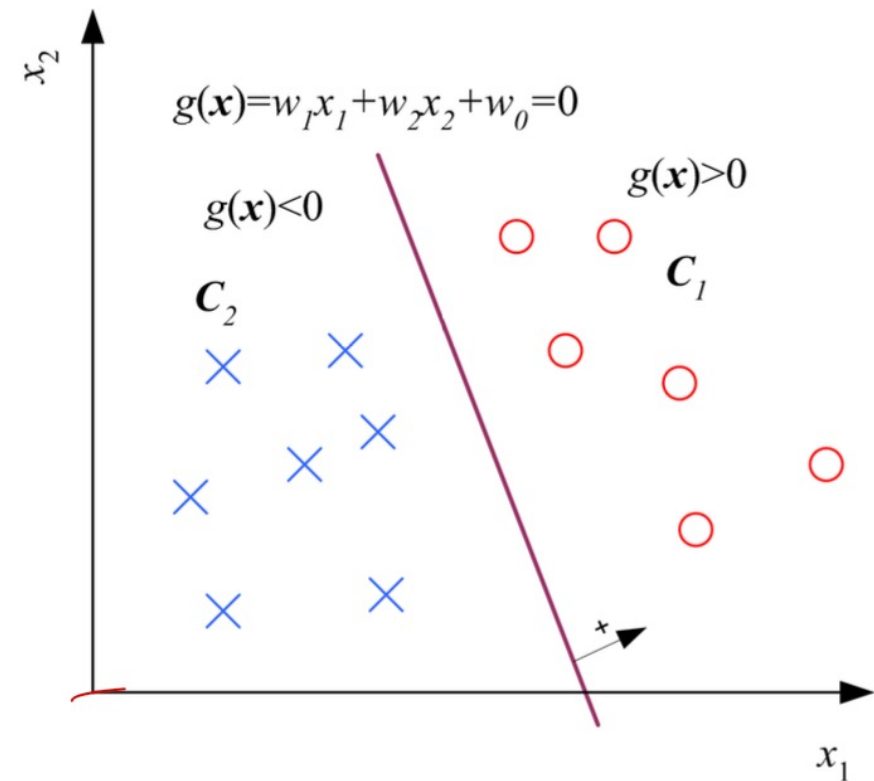
# Likelihood - vs Discriminant -based Classification

## Likelihood-based

- Use Bayes' Rule to calculate  $P(C_i|X)$
- Need Parametric estimation for  $P(X|C_i)$
- Purpose :  $g_i(x) = \log_2 P(X|C_i) + \log_2 P(C_i)$

## Discriminant Method

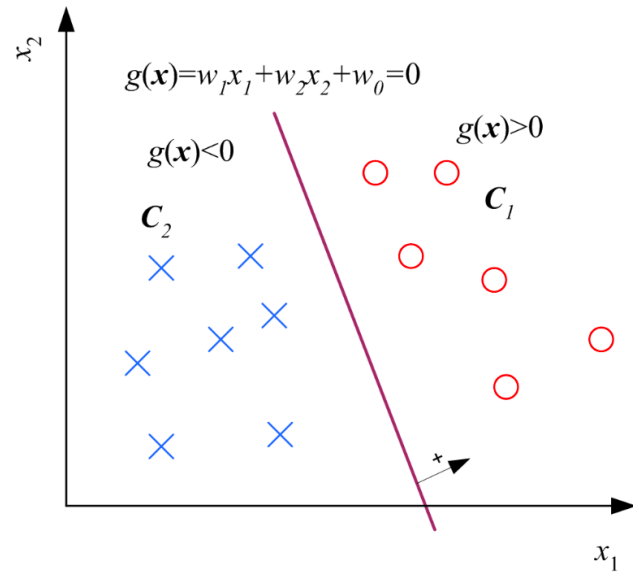
- Assume model  $g_i(x)$  directly, **no density estimation**
- Estimate boundary  $g_i(x)$  from data  $x$



# Linear Discriminant

**Discriminant** :  $g_i(x) = \sum_j^d w_{ij}x_j + w_{i0} = \mathbf{w}_i^T \mathbf{x} + w_{i0}$

If Two classes



$$\begin{aligned} g(\mathbf{x}) &= g_1(\mathbf{x}) - g_2(\mathbf{x}) \\ &= (\mathbf{w}_1^T \mathbf{x} + w_{10}) - (\mathbf{w}_2^T \mathbf{x} + w_{20}) \\ &= (\mathbf{w}_1 - \mathbf{w}_2)^T \mathbf{x} + (w_{10} - w_{20}) \\ &= \mathbf{w}^T \mathbf{x} + w_0 \end{aligned}$$

$$\text{choose } \begin{cases} C_1 & \text{if } g(\mathbf{x}) > 0 \\ C_2 & \text{otherwise} \end{cases}$$

Multi-classes ( $k > 2$ )

Choose  $C_i$  if  $P(C_i | X) = \max_k P(C_k | X) = \max_k g_k(x)$

# From Discriminant to Posterior

---

This is optimal solution... why?

Let assume  $P(X|C_i) \sim$  Gaussian Distribution

$$g_i(x) = w_i^T x + w_{i0} \qquad g_i(x) = -\frac{1}{2} \log 2\pi - \log \sigma_i - \frac{(x - \mu_i)^2}{2\sigma_i^2} + \log P(C_i)$$

$$\mathbf{w}_i = \Sigma^{-1} \mu_i \quad w_{i0} = -\frac{1}{2} \mu_i^T \Sigma^{-1} \mu_i + \log P(C_i)$$

$$y \equiv P(C_1 | \mathbf{x}) \text{ and } P(C_2 | \mathbf{x}) = 1 - y$$

$$\text{choose } C_1 \text{ if } \begin{cases} y > 0.5 \\ y/(1-y) > 1 \\ \log [y/(1-y)] > 0 \end{cases} \text{ and } C_2 \text{ otherwise}$$

# From Discriminant to Posterior

---

$$\begin{aligned}\text{logit}(P(C_1 | \mathbf{x})) &= \log \frac{P(C_1 | \mathbf{x})}{1 - P(C_1 | \mathbf{x})} = \log \frac{P(C_1 | \mathbf{x})}{P(C_2 | \mathbf{x})} \\ &= \log \frac{p(\mathbf{x} | C_1)}{p(\mathbf{x} | C_2)} + \log \frac{P(C_1)}{P(C_2)} \\ &= \log \frac{(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp\left[-(1/2)(\mathbf{x} - \mu_1)^T \Sigma^{-1} (\mathbf{x} - \mu_1)\right]}{(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp\left[-(1/2)(\mathbf{x} - \mu_2)^T \Sigma^{-1} (\mathbf{x} - \mu_2)\right]} + \log \frac{P(C_1)}{P(C_2)} \\ &= \mathbf{w}^T \mathbf{x} + w_0\end{aligned}$$

$$\text{where } \mathbf{w} = \Sigma^{-1}(\mu_1 - \mu_2) \quad w_0 = -\frac{1}{2}(\mu_1 + \mu_2)^T \Sigma^{-1}(\mu_1 - \mu_2)$$

The inverse of logit

$$\log \frac{P(C_1 | \mathbf{x})}{1 - P(C_1 | \mathbf{x})} = \mathbf{w}^T \mathbf{x} + w_0$$

$$P(C_1 | \mathbf{x}) = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + w_0) = \frac{1}{1 + \exp\left[-(\mathbf{w}^T \mathbf{x} + w_0)\right]}$$

# Logistic Regression (K = 2)

---

**Discriminant** :  $g_i(x) = w_i^T \mathbf{x} + w_{i0} = \text{score} = z$

Odds =  $\frac{P(C_1 | X)}{P(C_2 | X)} = \frac{y}{1-y} \rightarrow \text{한계가 있다(?) } \rightarrow \log(\text{odds}) = \text{logit} = z$  (실수 전체 범위)

$\log(\text{odds}) = \log + \text{probit} = \text{logit}$

$$\log \frac{P(C_1 | X)}{P(C_2 | X)} = \log \frac{y}{1-y} = z$$

$y \equiv P(C_1 | \mathbf{x})$  and  $P(C_2 | \mathbf{x}) = 1 - y$

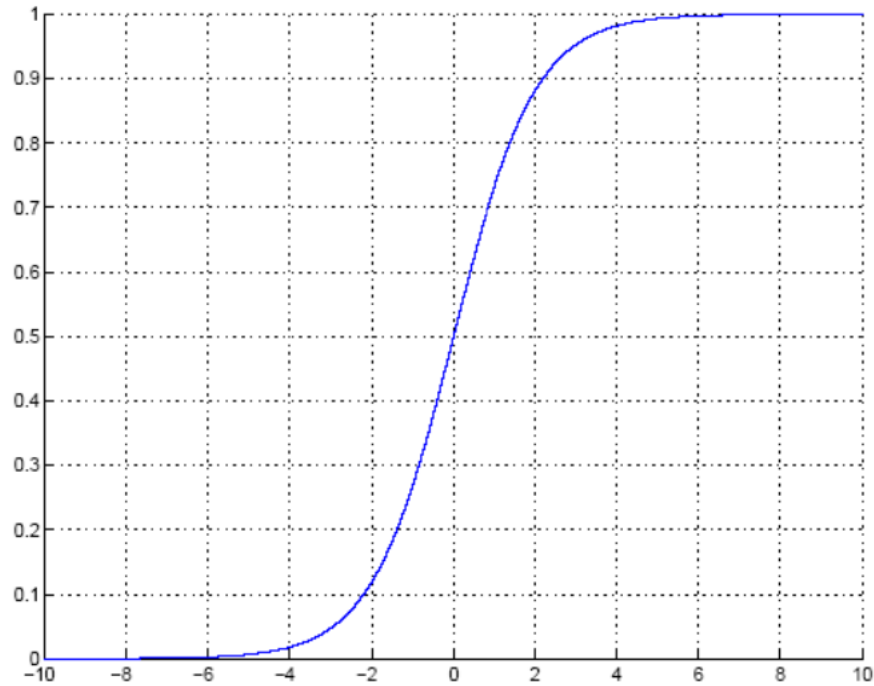
choose  $C_1$  if  $\begin{cases} y > 0.5 \\ y/(1-y) > 1 \\ \log[y/(1-y)] > 0 \end{cases}$  and  $C_2$  otherwise

$$y = \frac{1}{1+e^{-z}} = \sigma(z) = \text{sigmoid function}$$

# Logistic Regression (K = 2)

---

$$y = \frac{1}{1+e^{-z}} = \sigma(z) = \textit{sigmoid function}$$



Choose  $C_1$  when  $z > 0, y > 0.5$

Q. But why sigmoid function?



# Logistic Regression ( $K > 2$ )

---

**Discriminant** :  $g_i(x) = w_i^T x + w_{i0} = \text{score} = z_i$

$$\text{Odds} = \frac{P(C_i | X)}{P(C_k | X)} = e^{z_i}$$

$$\sum_1^{K-1} \frac{P(C_i | X)}{P(C_k | X)} = \sum_1^{K-1} e^{z_i} = \frac{1 - P(C_k | X)}{P(C_k | X)} \quad P(C_k | X) = \frac{1}{1 + \sum_1^{K-1} e^{z_i}}$$

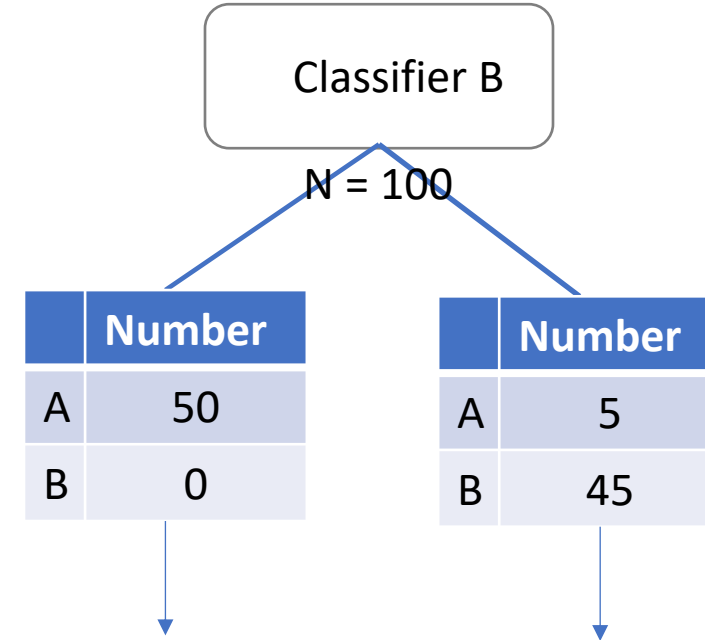
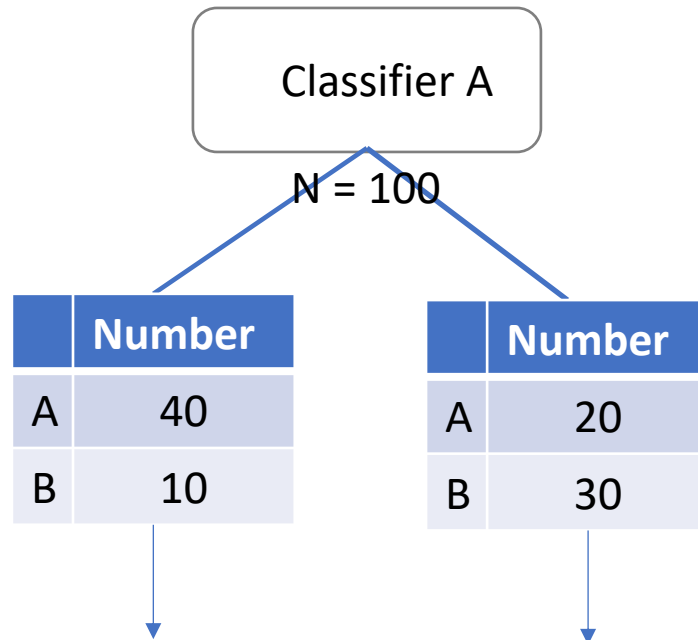
$$P(C_i | X) = P(C_k | X) \times e^{z_i} = \frac{1}{1 + \sum_1^{K-1} e^{z_i}} \times e^{z_i} = \frac{e^{z_i}}{\sum_1^K e^{z_i}}$$

$$P(C_i | X) = \frac{e^{z_i}}{\sum_1^K e^{z_i}} = \text{softmax}(z_i)$$

## 2-3. Learning Classifier

# Entropy

---



# Entropy

---

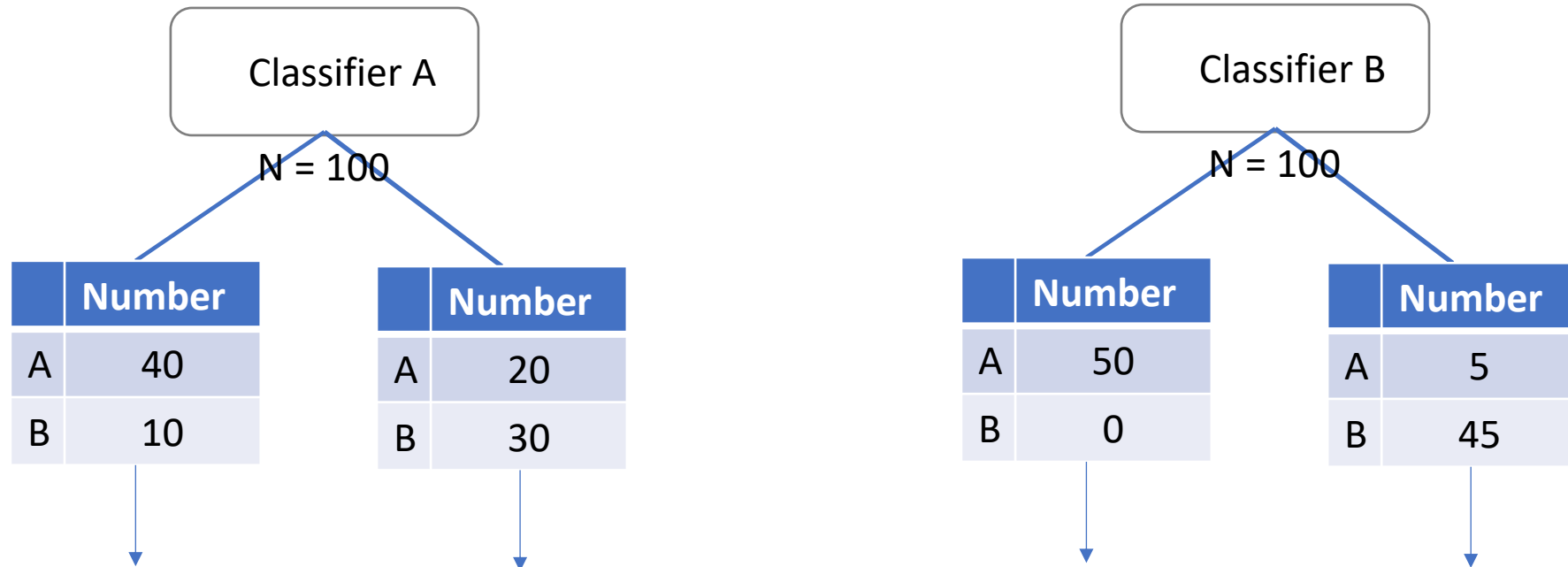
## Entropy (불균형도)

- 특정 node  $t$  에서 불순도
- 데이터 분포의 purity를 측정하는 척도, 여기서는 클래스의 분포의 purity를 측정
- Entropy가 낮을 수록 purity가 높은 것
- Max :  $\log_2 n_c$  ( $n_c$ : 클래스 총 개수)
- Min : 0 (클래스가 1개 밖에 없을 경우)

$$Entropy(t) = - \sum_{j=\text{class}} p(j|t) \cdot \log_2 p(j|t)$$

# Entropy

---



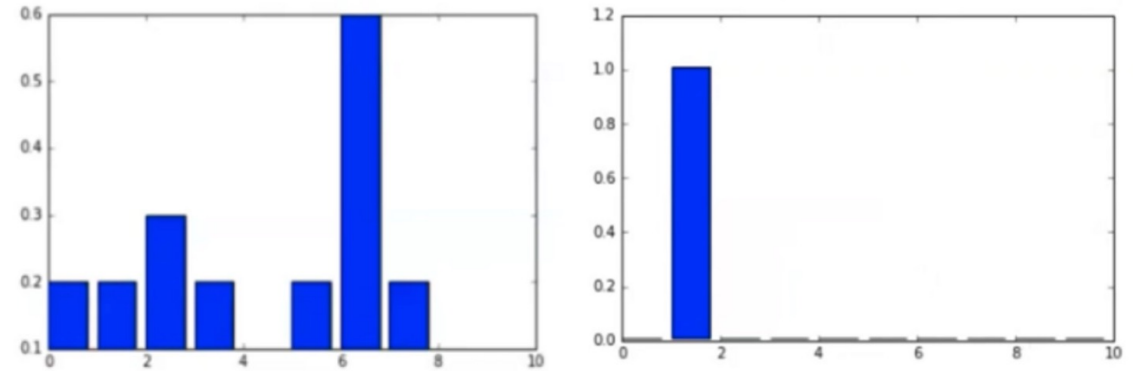
**Which one is better?**

# Cross-Entropy

두 분포의 차이의 척도

$$\text{Cross-entropy} = - \sum_{i=1}^N p_i \log q_i$$

p: 실제 정답의 분포  
q: 모델을 통해 구한 답의 분포



Minimize Cross-Entropy!

Minimize Loss Function!

# How to find parameters

---

## Classification

- Binary Cross Entropy

$$BCE = -\frac{1}{N} \sum_{i=0}^N y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i)$$

- Categorical Cross Entropy

$$CCE = -\frac{1}{N} \sum_{i=0}^N \sum_{j=0}^J y_j \cdot \log(\hat{y}_j) + (1 - y_j) \cdot \log(1 - \hat{y}_j)$$

# MLE? → Loss function

---

If K=2 (Binary Classification)

- **Bernoulli distribution**

$$\log L(p) = \sum_{i=1}^n (y_i \log p + (1 - y_i) \log (1 - p))$$

Maximize Log Likelihood

We know about p (output of model)

$$p = \frac{1}{1+e^{-z}} = \sigma(z) = \text{sigmoid function}$$

$$P(C_i | X) = \frac{e^{z_i}}{\sum_1^K e^{z_i}} = \text{softmax}(z_i) \text{ if } K > 2$$

- **Binary Cross Entropy**

$$BCE = -\frac{1}{N} \sum_{i=0}^N y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i)$$

Minimize Loss Function



# Gradient Descent

---

## Minimize Loss Function

We know about  $p$  (output of model)

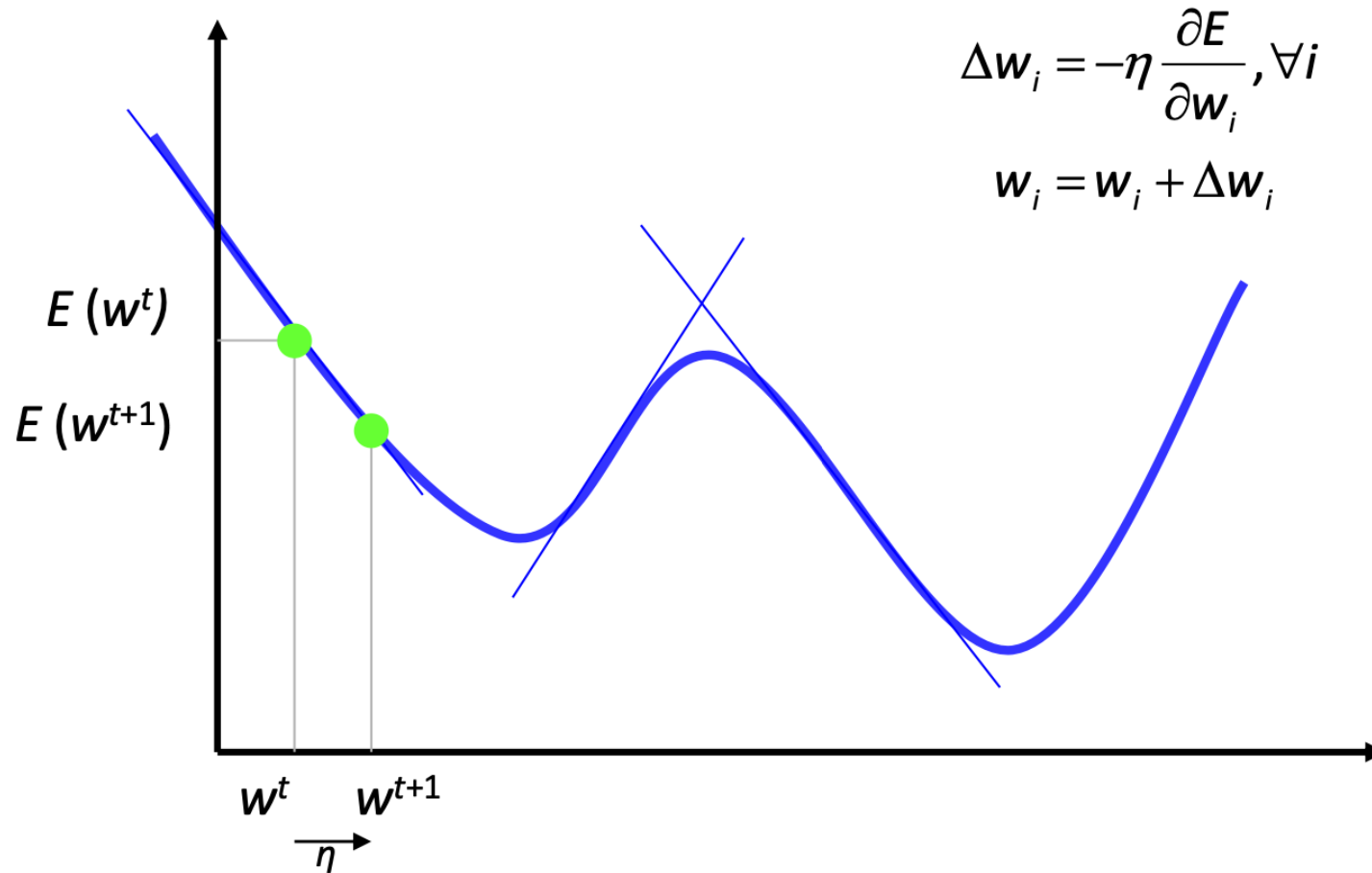
$$p = \frac{1}{1+e^{-z}} = \sigma(z) = \textit{sigmoid function}$$

$$P(C_i | X) = \frac{e^{z_i}}{\sum_1^K e^{z_i}} = \text{softmax}(z_i) \text{ if } K > 2$$

1. **model** :  $g_i(x) = w_i^T x + w_{i0} = \text{score} = z_i$
2. **Loss function** :  $E(w | X) = \text{Cross-Entropy}$
3. **Optimization** :  $w^* = \text{argmin}_w E(w|X)$

Gradient :  $\nabla_w E$

# Gradient Descent

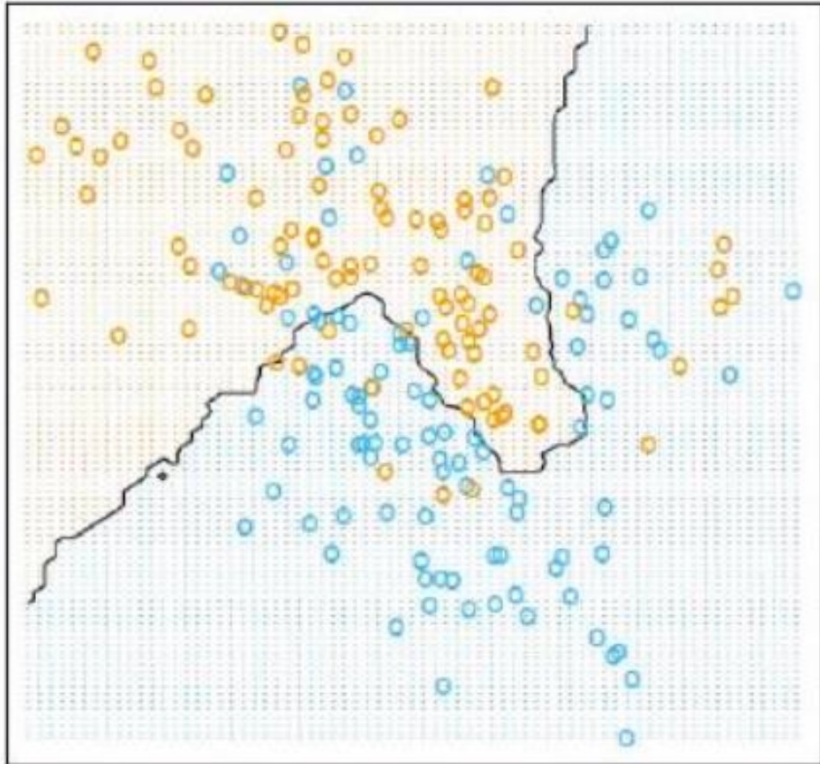


## 4. Non-parametric Method

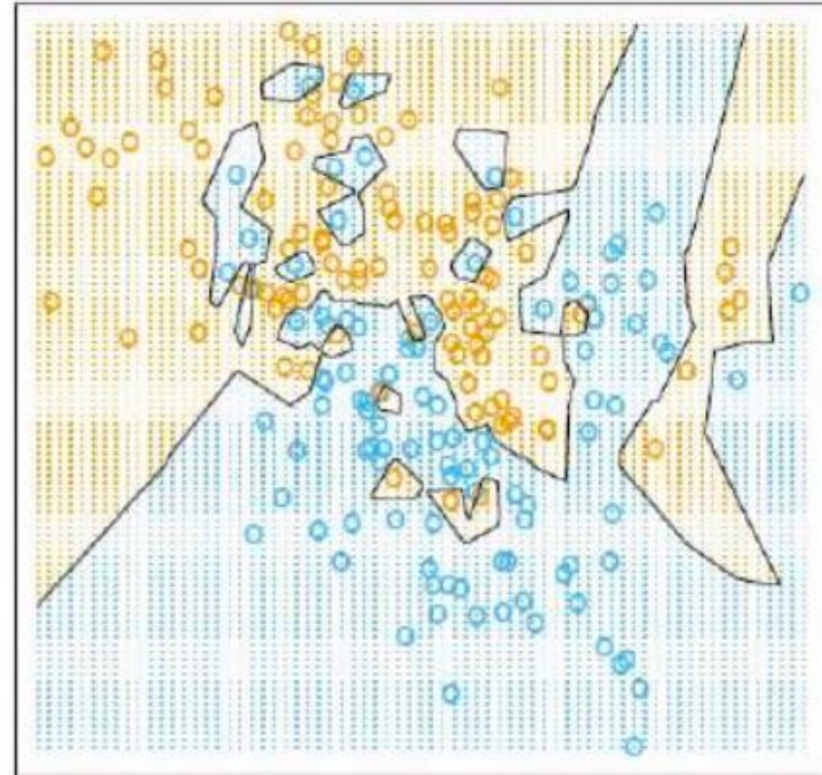
# KNN (K- Nearest Neighborhood)

전형적인 non-parametric method

15-Nearest Neighbor Classifier



1-Nearest Neighbor Classifier



# KNN (K- Nearest Neighborhood)

---

- Distance measure

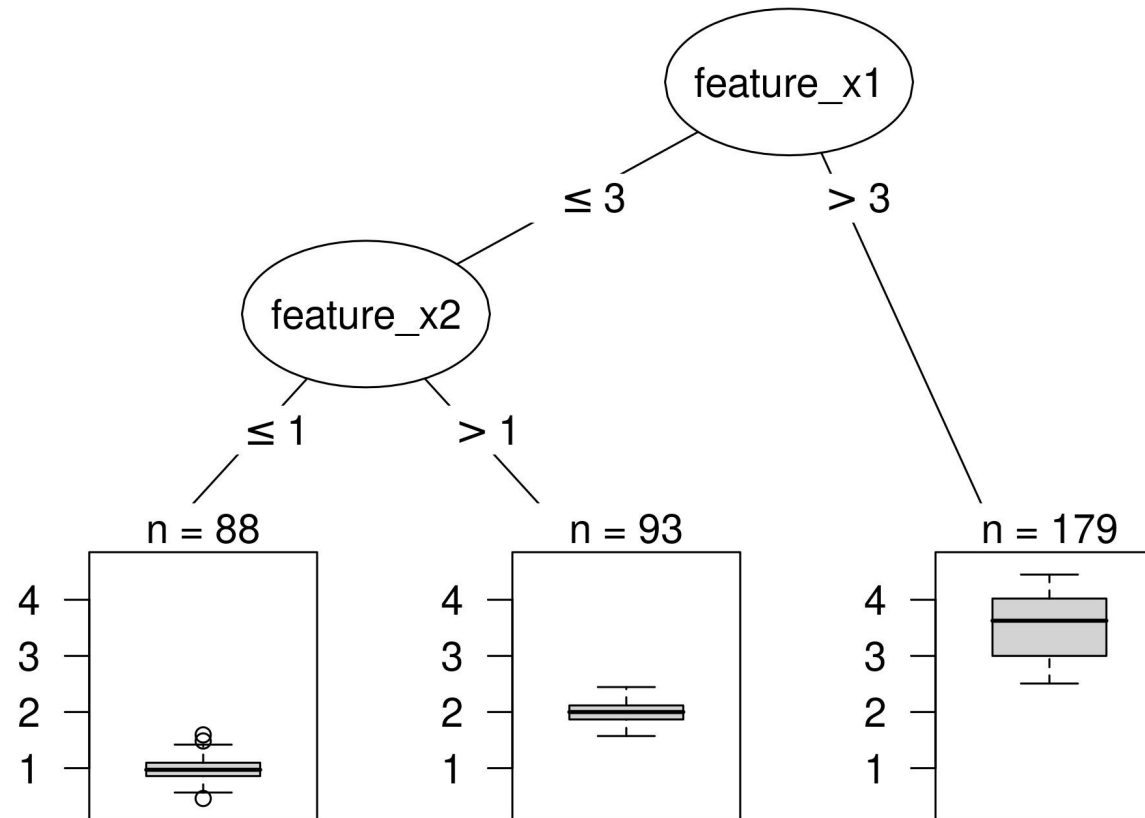
$$d(\mathbf{u}, \mathbf{v}) = (\sum |u_i - v_i|^2)^{\frac{1}{2}} = ||\mathbf{u} - \mathbf{v}||_2 \quad \text{Euclidean (L2 norm)}$$

$$d(\mathbf{u}, \mathbf{v}) = \sum |u_i - v_i| = ||\mathbf{u} - \mathbf{v}||_1 \quad \text{Manhattan (L1 norm)}$$

$$d(\mathbf{u}, \mathbf{v}) = (\sum |u_i - v_i|^p)^{\frac{1}{p}} = ||\mathbf{u} - \mathbf{v}||_p \quad \text{Minkowski (Lp norm)}$$

$$d(\mathbf{u}, \mathbf{v}) = \sqrt{(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})} \quad \text{Mahalanobis Distance}$$

# Decision Tree



## 5. Model Evaluation

# Confusion Matrix

---

Actual class	Predicted Class		
		Positive	Negative
	Positive	<b>True Positive(TP)</b>	<b>False Negative(FN)</b>
	Negative	<b>False Positive(FP)</b>	<b>True Negative(TN)</b>

$$\text{Accuracy} : \frac{TP+TN}{TP+TN+FP+FN}$$

Q. What is Limitation of Accuracy?



# Confusion Matrix

Actual class	Predicted Class		
		Positive	Negative
	Positive	True Positive(TP)	False Negative(FN)
	Negative	False Positive(FP)	True Negative(TN)

Precision(정밀도) :  $\frac{TP}{TP+FP}$  → 양성 예측 중, 실제로 맞은 비율

Recall(sensitivity, 재현율, 민감도) :  $\frac{TP}{TP+FN}$  → 실제 양성 중, 맞은 비율

Specificity(특이도) :  $\frac{TN}{TN+FP}$  → 실제 음성 중, 맞은 비율

# F1-score

---

What was limitation of Accuracy?

## Precision & recall Trade-off

$$\text{Precision(정밀도)} : \frac{TP}{TP+FP}$$

$$\text{Recall(sensitivity, 재현율, 민감도)} : \frac{TP}{TP+FN}$$

→ 둘 다 높이는 것이 가능한가...?

→ 좋은 모델은 positive한 것을 모두 제대로 분류하고, positive한 것만 제대로 분류하면 된다.

About precision

About recall

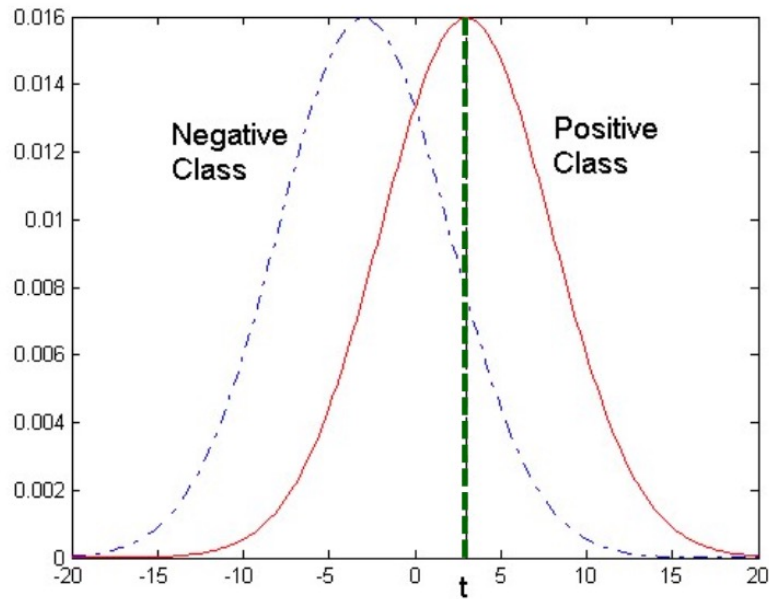
[Harmonized mean(조화 평균)]

$$\frac{1}{F1\ score} = 0.5 \left( \frac{1}{precision} + \frac{1}{recall} \right)$$

$$F1\ score = \frac{2 * precision * recall}{(precision + recall)}$$

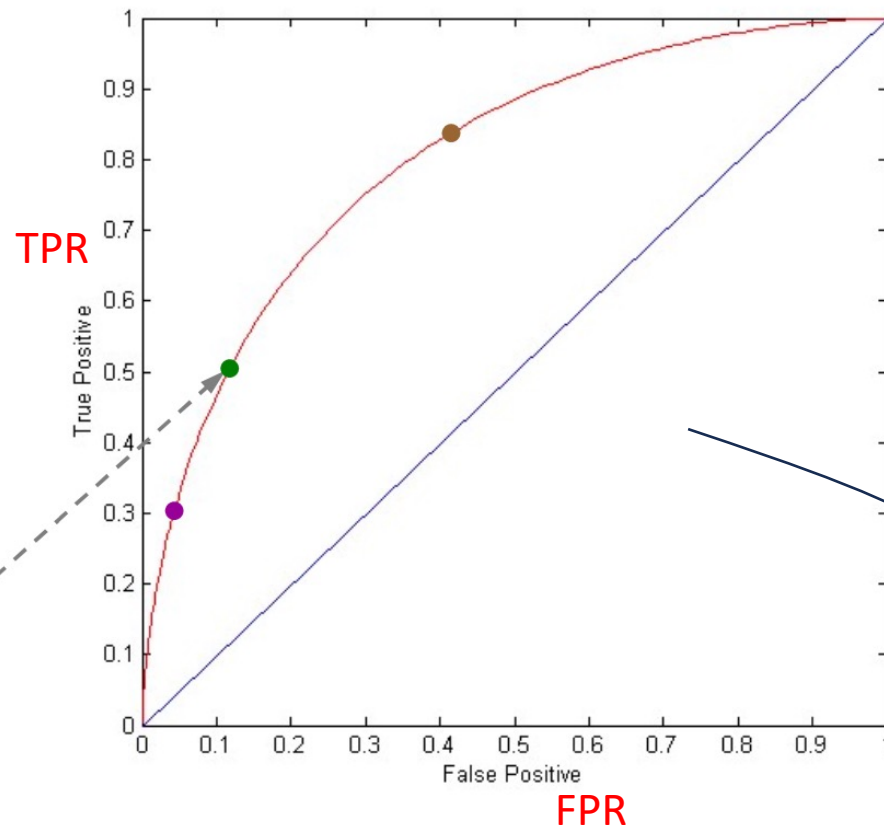
# ROC Curve

ROC(Reviewer Operating Characteristics) -Curve



At threshold  $t$ :

**TP=0.5, FN=0.5, FP=0.12, TN=0.88**



TPR(True positive Rate) = recall

FPR(False positive Rate) = specificity

AUC(Area Under Curve)

Bad= 0.5

Ideal = 1

# 수고하셨습니다!

해당 세션자료는 KUBIG Github에서 보실 수 있습니다!  
다음은 이번 주차 과제 설명이 있습니다!