# Data Analysis(3)

**Dept. of Mechanical System Design Engineering,**
**Seoul National University of Science and Technology**
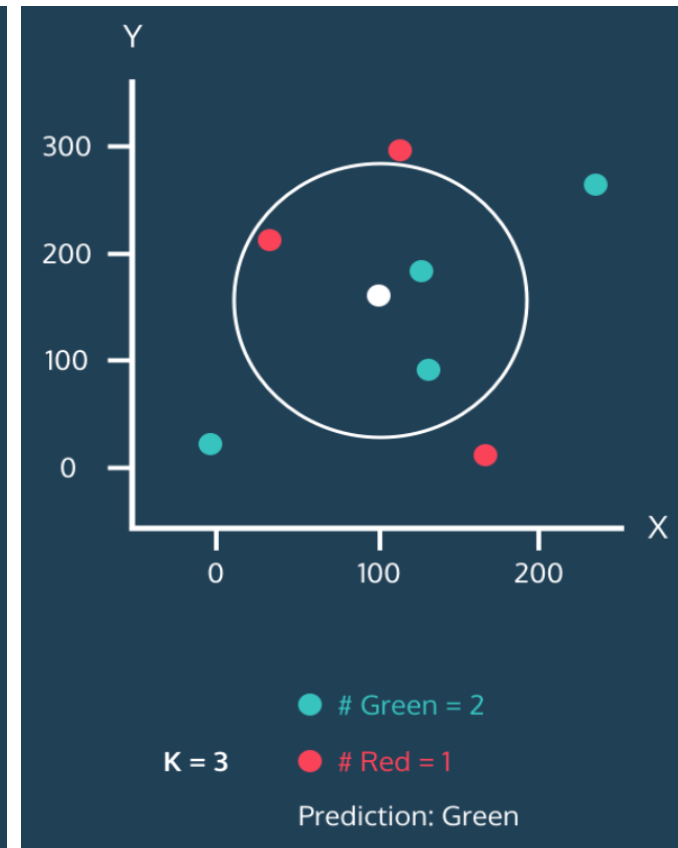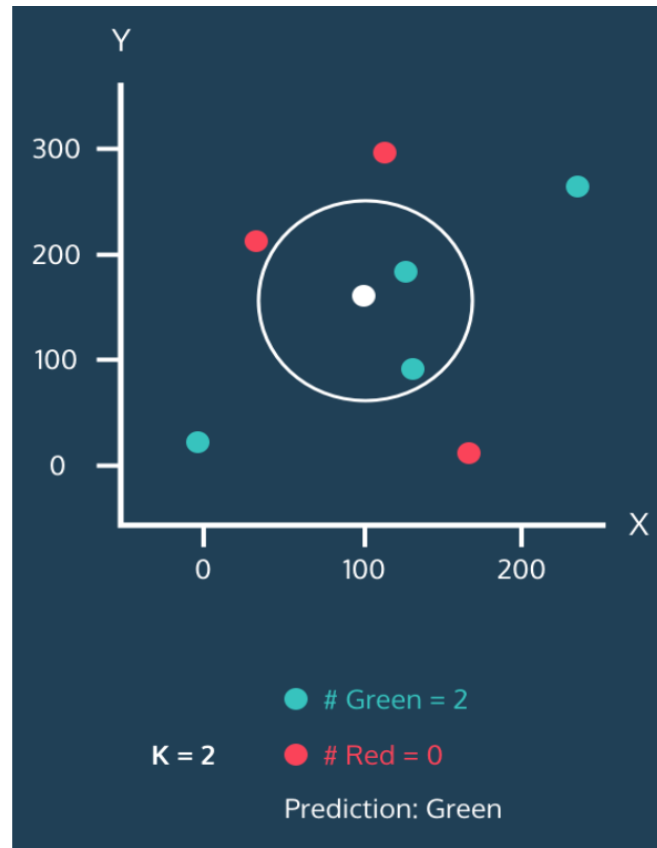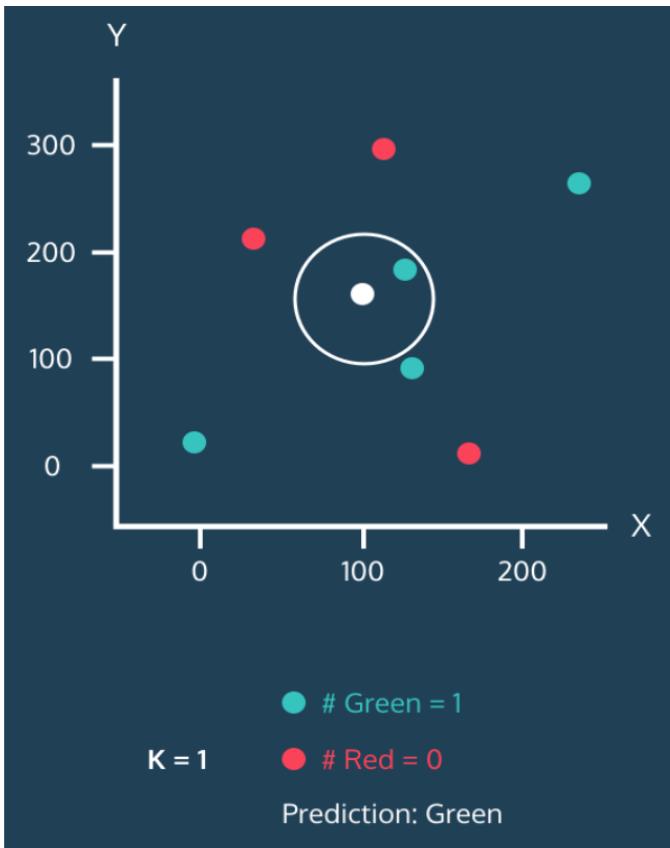
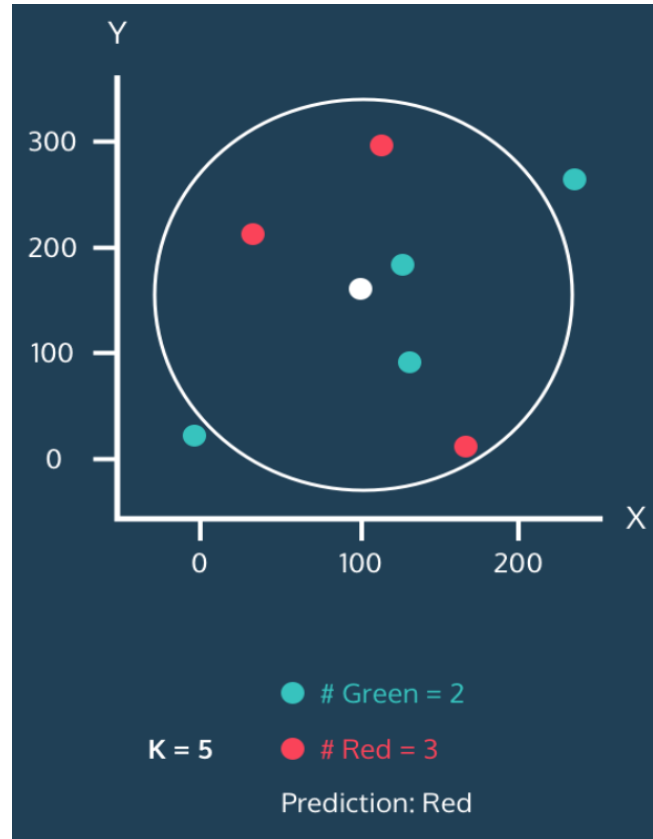**Prof. Ju Yeon Lee**
**(jylee@seoultech.ac.kr)**

국립**서울과학기술대학교**
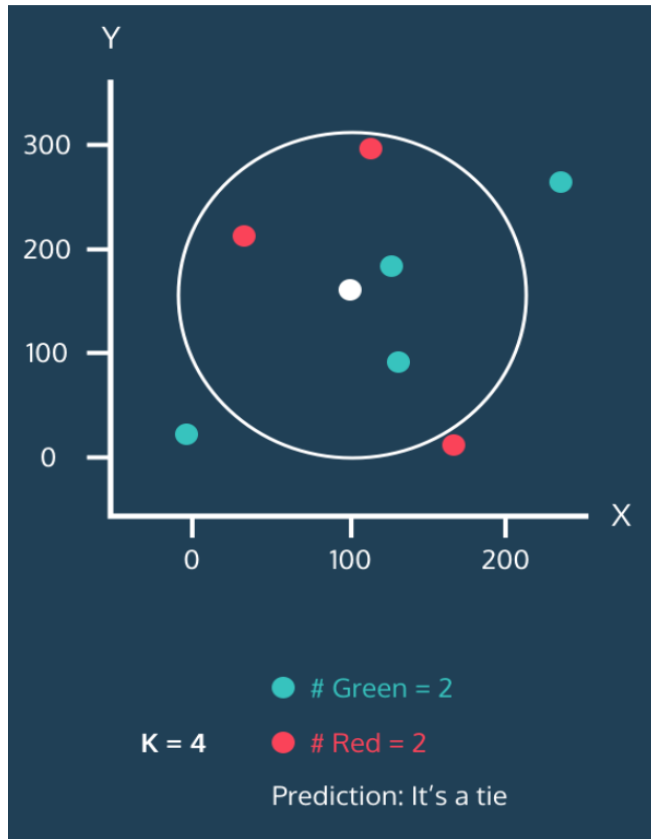
# Review

# KNN : K-Nearest Neighbor

# KNN : K-Nearest Neighbor

# Distance Formula

## Euclidean Distance



$$d = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$$

$$\sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + \ldots + (a_n - b_n)^2}$$

## Manhattan Distance



$$d = \mid a_1 - b_1 \mid + \mid a_2 - b_2 \mid$$

$$\mid a_1 - b_1 \mid + \mid a_2 - b_2 \mid + \ldots + \mid a_n - b_n \mid$$

서울과학기술대학교

# Distance Formula

## Minkowski Distance

$$D(X, Y) = \left( \sum_i^n (|x_i - y_i|)^p \right)^{\frac{1}{p}}$$

- **P = 1, Manhattan Distance**

- **P = 2, Euclidean Distance**

서울과학기술대학교

# Decision Tree



**A**                         **B**



| | |
|---|---|
| #1 | #2 |
| P : 0.9 | P : 0.1 |

| | |
|---|---|
| #1 | #2 |
| P : 0.5 | P : 0.5 |

- *A : 0.9 * 0.1 = 0.09*
- *B : 0.5 * 0.5 = 0.25*

- **Gini Impurity :** $Gini = 1 - \sum_{i=1}^{C} (p_i)^2$

- **Entropy Index :** $H(X) = -\sum_{i=1}^{n} p_i \log_2 p_i$

서울과학기술대학교

$$Ent(D) = -\sum_{k=1}^{|Y|} p_k \, log_2 p_k$$

$$Gain(D, a) = Ent(D) - \sum_{v=1}^{V} \frac{|D^v|}{|D|} Ent(D^v)$$

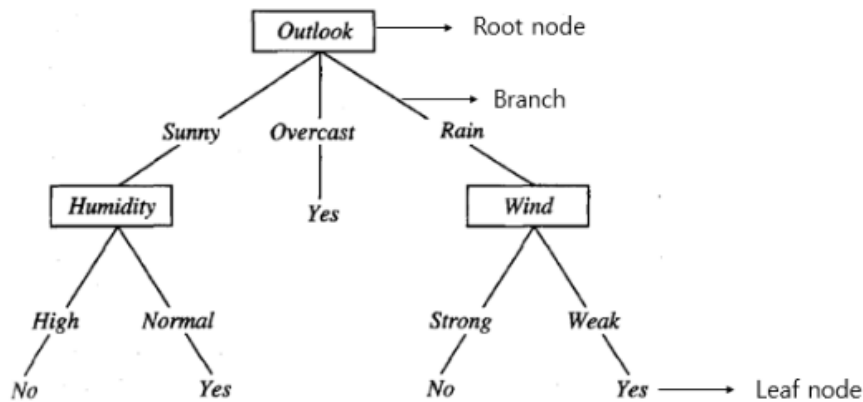| Day | Outlook | Temperature | Humidity | Wind | PlayTennis |
|-----|---------|-------------|----------|------|------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

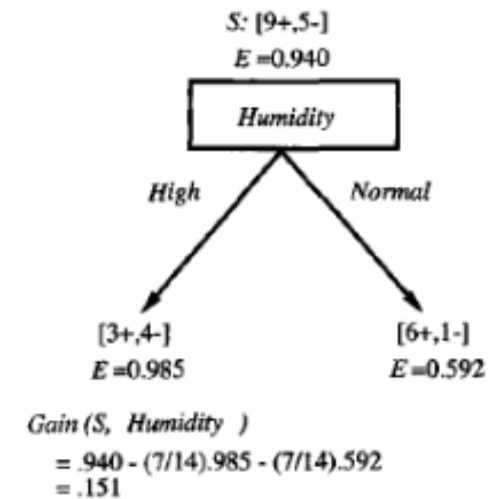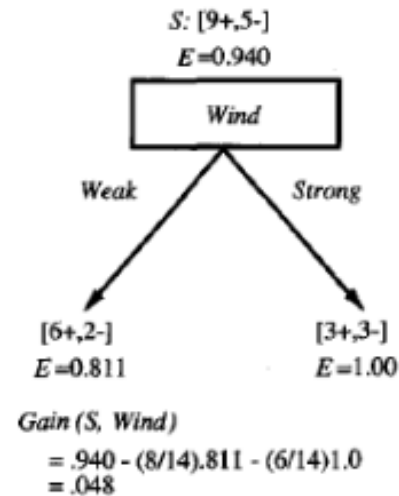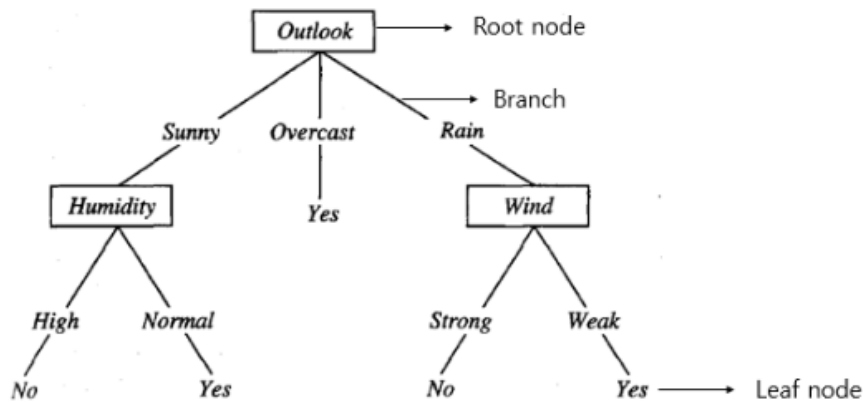$$Values(Wind) = Weak, Strong$$
$$S = [9+, 5-]$$
$$S_{Weak} \leftarrow [6+, 2-]$$
$$S_{Strong} \leftarrow [3+, 3-]$$

$$Gain(S, Wind) = Entropy(S) - \sum_{v \in \{Weak, Strong\}} \frac{|S_v|}{|S|} Entropy(S_v)$$
$$= Entropy(S) - (8/14)Entropy(S_{Weak})$$
$$- (6/14)Entropy(S_{Strong})$$
$$= 0.940 - (8/14)0.811 - (6/14)1.00$$
$$= 0.048$$

Root node

Branch

Leaf node

$S: [9+,5-]$
$E=0.940$

Wind

Weak        Strong

$[6+,2-]$        $[3+,3-]$
$E=0.811$        $E=1.00$

Gain $(S, Wind)$
$= .940 - (8/14).811 - (6/14)1.0$
$= .048$

$S: [9+,5-]$
$E=0.940$

Humidity

High        Normal

$[3+,4-]$        $[6+,1-]$
$E=0.985$        $E=0.592$

Gain $(S, Humidity)$
$= .940 - (7/14).985 - (7/14).592$
$= .151$

| Day | Outlook | Temperature | Humidity | Wind | PlayTennis |
|-----|---------|-------------|----------|------|------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

$$Gain(S, Outlook) = 0.246$$
$$Gain(S, Humidity) = 0.151$$
$$Gain(S, Wind) = 0.048$$
$$Gain(S, Temperature) = 0.029$$
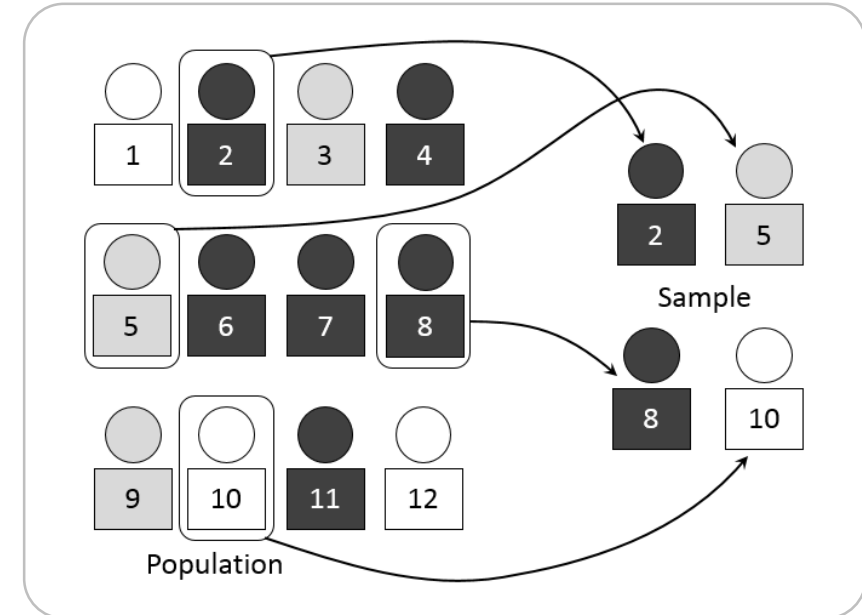
# Sampling

# Sampling?

*"Sampling" is the selection of a subset or a statistical sample (termed sample for short) of individuals from within a statistical population to estimate characteristics of the whole population*

Source : Wikipedia

- **Sampling Type :**

  - ✓ Balanced sampling : Balance between data classes
    → **Simple random sampling (단순 임의 샘플링)**

  - ✓ Imbalanced/unbalanced sampling : Imbalance between data classes
    → **Stratified sampling (층화 추출)**
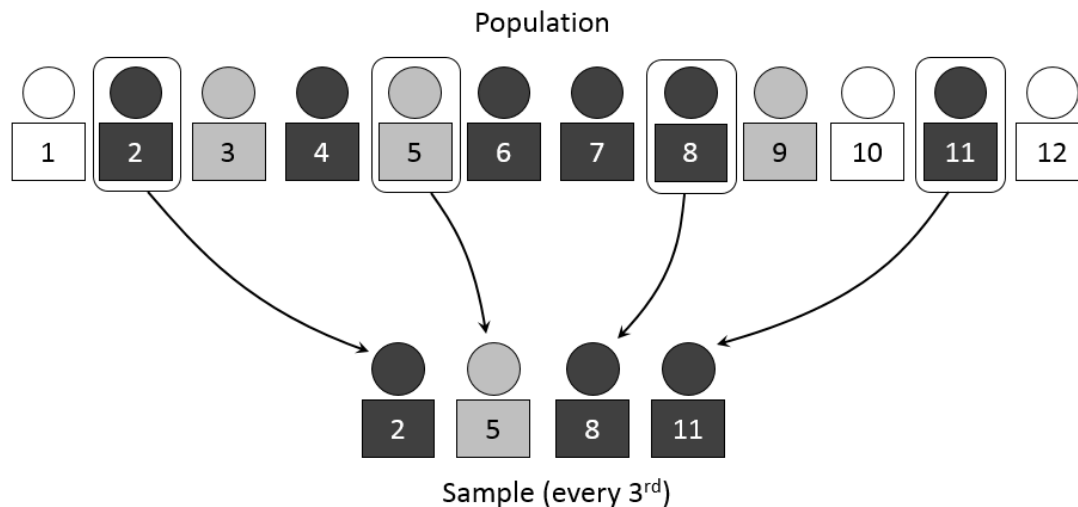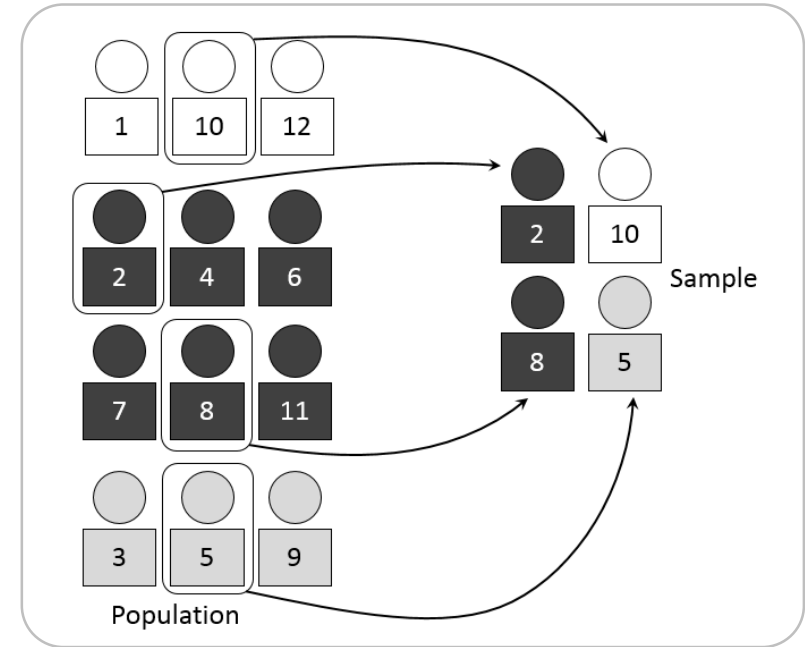    → **Systematic sampling (계통 추출)**

- **Simple random sampling (단순 임의 샘플링)**
  All subsets of a sampling frame have an equal probability of being selected. Each element of the frame thus has an equal probability of selection

서울과학기술대학교

# Sampling?
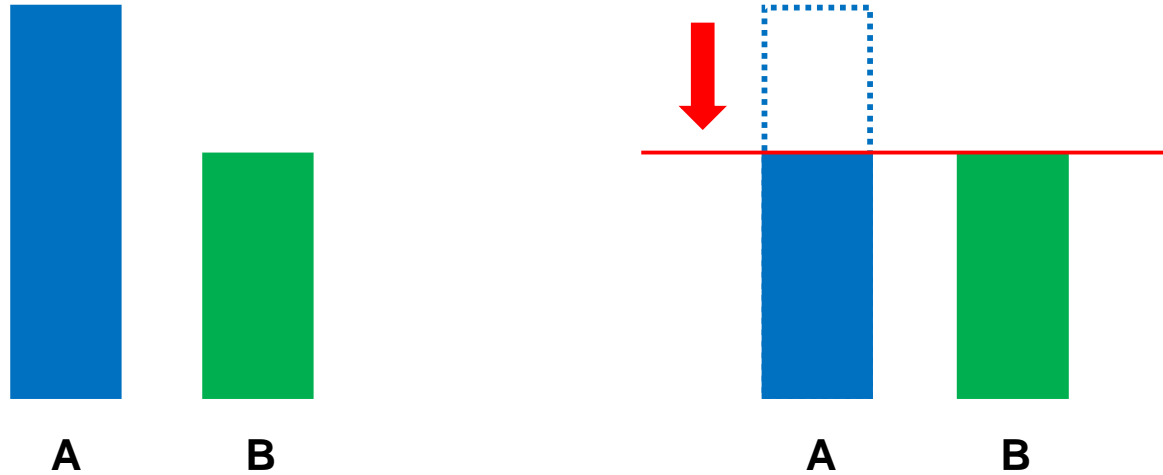
- **Stratified sampling (층화 추출)**
  When the population embraces a number of distinct categories, the frame can be organized by these categories into separate "strata." Each stratum is then sampled as an independent sub-population, out of which individual elements can be randomly selected
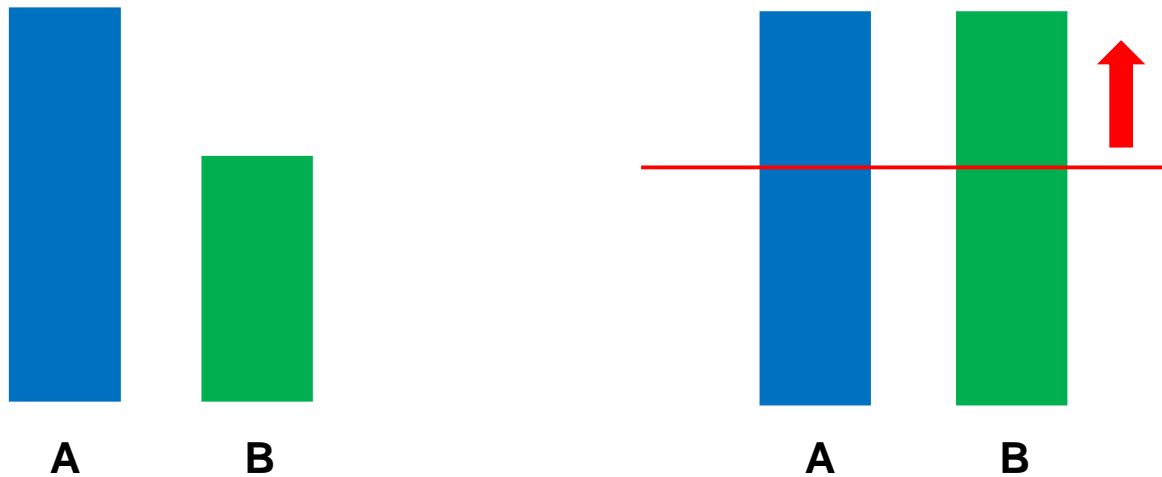


- **Systematic sampling (계통 추출)**
  Systematic sampling (also known as interval sampling) relies on arranging the study population according to some ordering scheme and then selecting elements at regular intervals through that ordered list

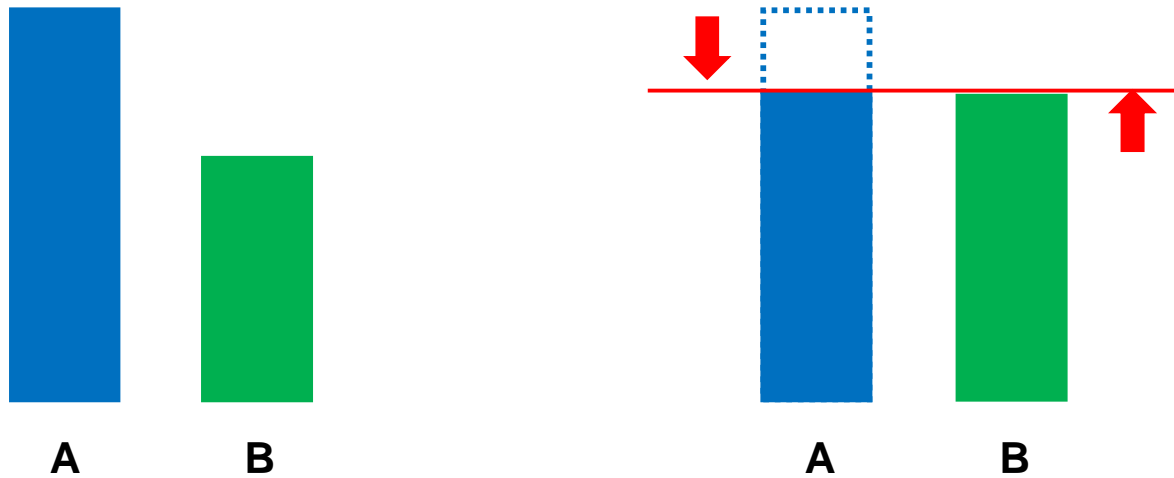서울과학기술대학교

# Imbalanced sampling

## 1. Under/Down Sampling



A    B          A    B

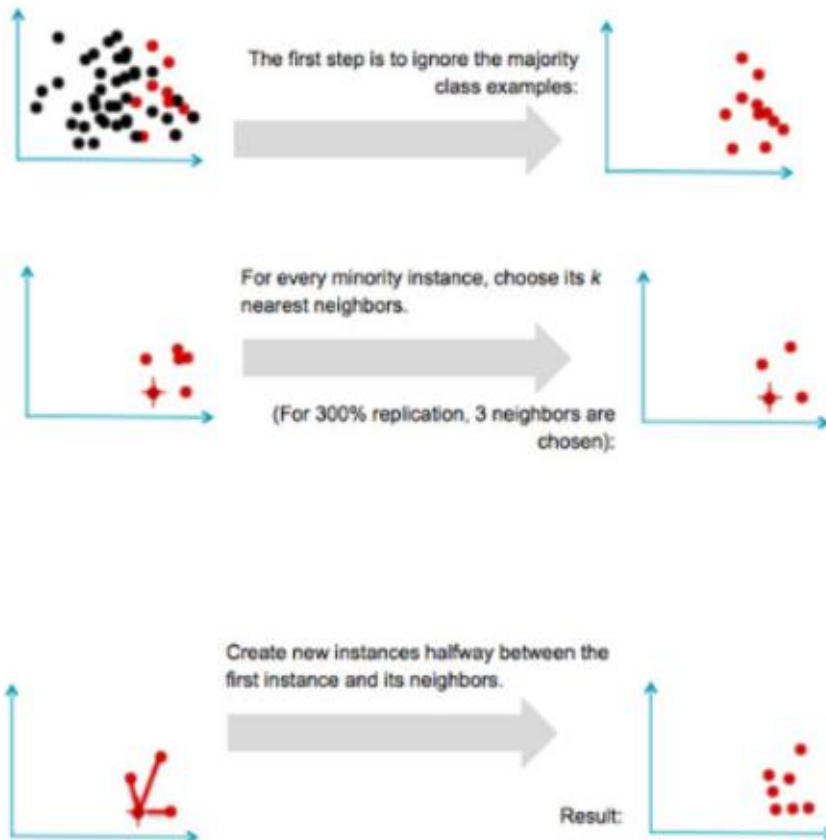## 2. Over/Up Sampling



A    B          A    B

국립서울과학기술대학교

## 3. Combination Sampling : Under + Over

서울과학기술대학교

# SMOTE

**SMOTE(Synthetic Minority Over-Sampling Technique) :**
generate new data between neighboring minority classes from random minority class data
For numerical features



**SMOTENC(Synthetic Minority Over-Sampling Technique for Nominal and Continuous) :**

For dataset containing numerical and categorical features

However, it is not designed to work with only categorical features

https://imbalanced-learn.org/stable/references/index.html

서울과학기술대학교

# Thank you

## Q & A