

BACKGROUND OF DATA MINING

Week02

Essential Math for Data Mining

통계 배우기

widely used Algorithm

to find best parameter 예측 변수를 가장 좋게 만들고자 하는

방법

to (minimize) objective function value

목표함수

어떤 파라미터가 좋지 않은지를 가늠하는

ex) MSE, J(SBM)

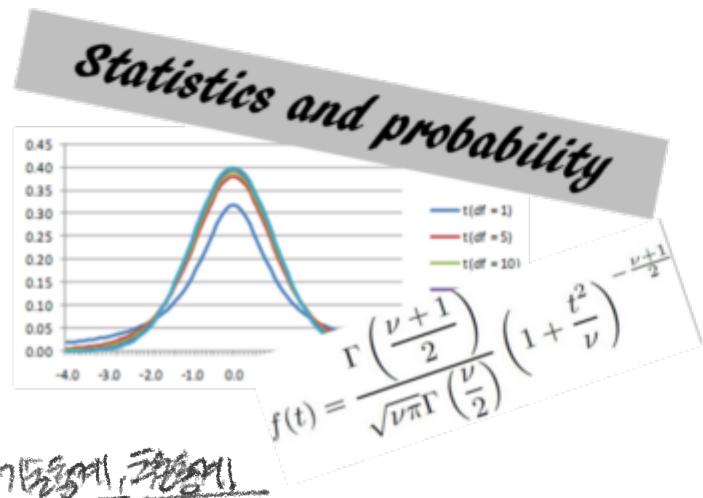
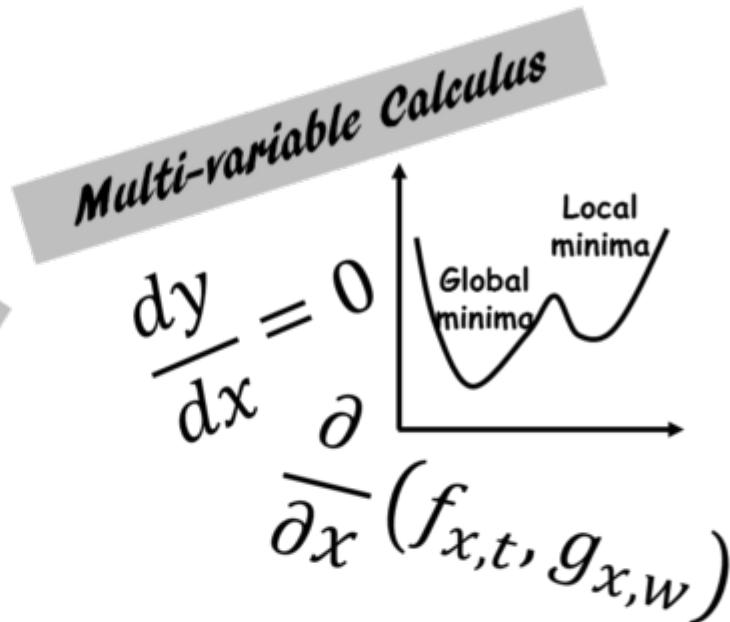
require optimization to find it

최적화

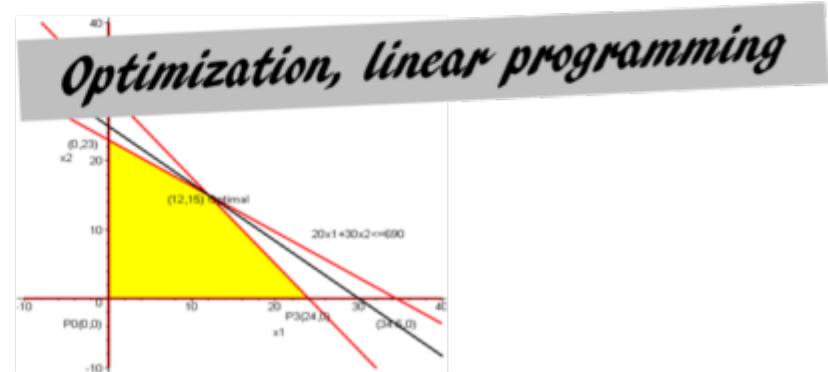
Statistics

Linear Algebra

$$\vec{a} \rightarrow \begin{bmatrix} x_1 & x_2 \\ x_3 & x_4 \end{bmatrix} \vec{b}$$



Descriptive Inferential



Essential Math for Data Mining: Statistics

□ Two main branches of statistics

▣ Descriptive statistics 기술통계

- Describe the basic features of data 데이터의 기본적인 특성(크기,형태)을 설명
- Data summaries and descriptive statistics, central tendency, variance, covariance, correlation 통계량
평균·중앙값·분산
분산
상관

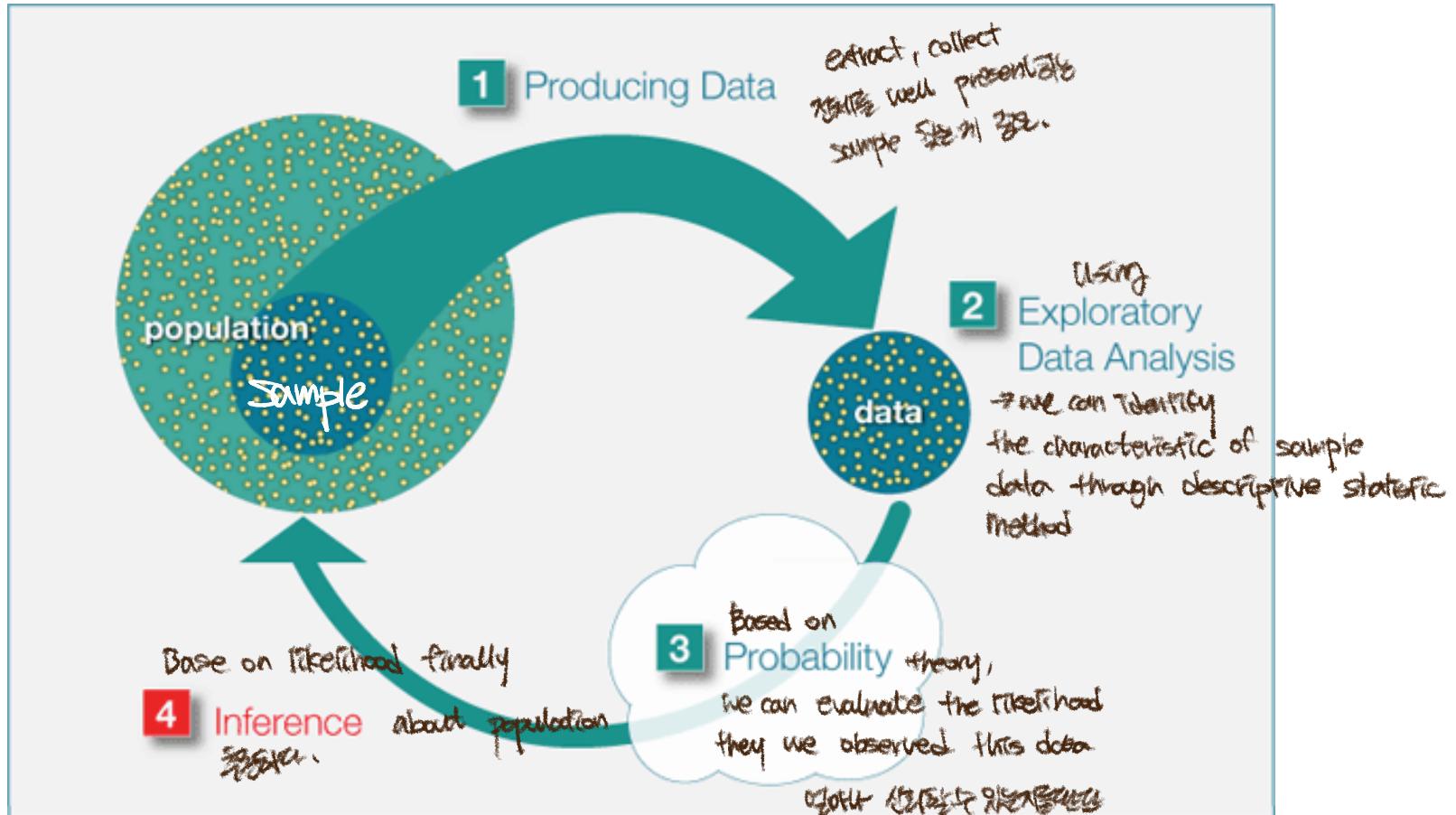
▣ Inferential statistics 판정통계

- Deduce properties of an underlying distribution of probability Probability model

□ Probability 이유로 데이터의 특성을 예측하는 확률.

- Sampling, measurement, error, random number generation 표본
측정
오류
난수생성
- Basic probability: basic idea, expectation, probability calculus, Bayes theorem, conditional probability 기초 확률
기대치
확률계산
베이즈 정리
조건부 확률
- Probability distribution functions—uniform, normal, binomial, chi-square, student's t-distribution, Central limit theorem 제공 확률
정규분포
이항 분포
хи-квадрат
t-분포
중심극한정리

Essential Math for Data Mining: Statistics



Essential Math for Data Mining: Linear Algebra

□ Linear algebra

- The study of linear sets of equations and their transformation properties
- Concern linear equations, linear functions and their representations in vector spaces and through matrices
- Used in most areas of science and engineering, because it allows modeling many natural phenomena, and efficiently computing with such models

final model은 왜 정답??

선형방정식의 집합 그것의 반환을 다루는 수학 분야
선형 방정식, 선형 함수/ 그리고 벡터 공간이나 행렬을 통해 이
들을 표현하는데에 관심있음

선형방정식 = 형식의 표현

$$\begin{aligned}3x + 5y &= 7 \\x - 2y &= 6\end{aligned}$$

2개의 방정식을 가진 linear equations



matrix operation

$$\begin{bmatrix} 3 & 5 \\ 1 & -2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 7 \\ 6 \end{bmatrix}$$

matrix를 표현

In data mining, data can be represented
in vector, matrix form, and desired value
derived by matrix operation



$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 3 & 5 \\ 1 & -2 \end{bmatrix}^{-1} \begin{bmatrix} 7 \\ 6 \end{bmatrix}$$

finding inverse matrix

제거 행렬의 역행렬을 구한 뒤, 결과행렬과 곱해
마지막 차수를 구함.

Essential Math for Data Mining: Calculus

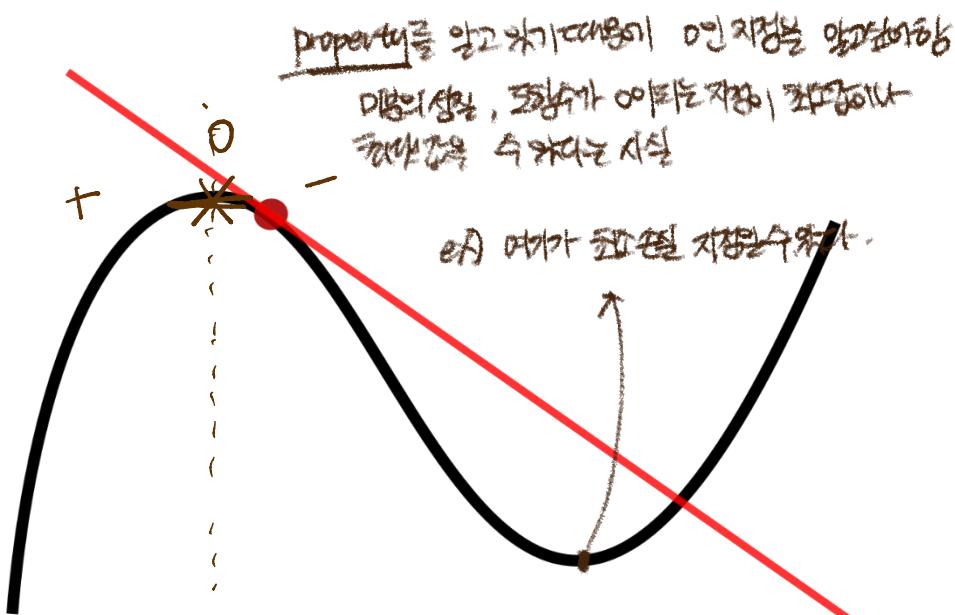
미적분학

□ Calculus

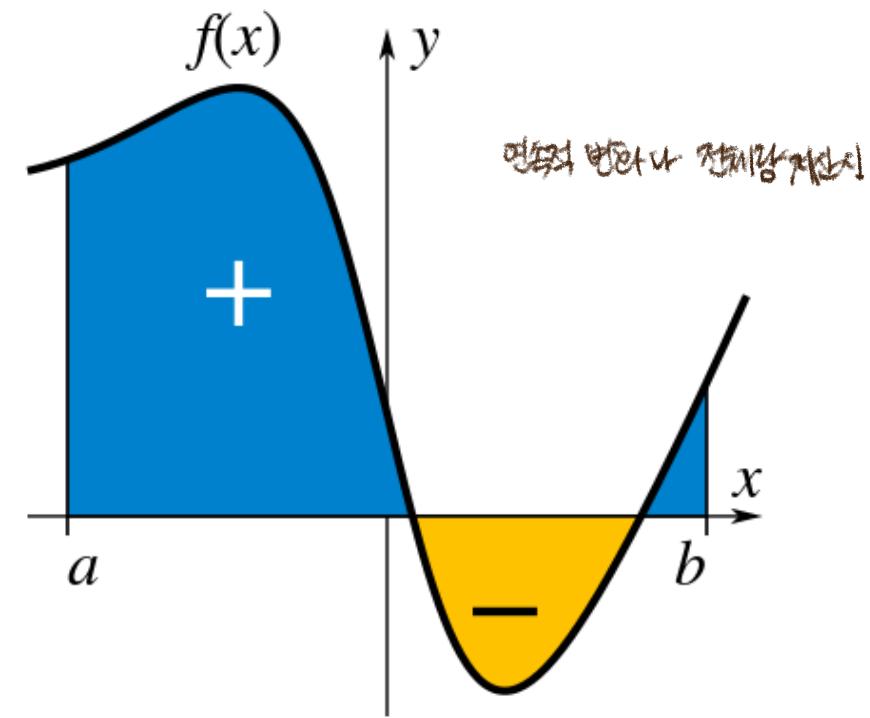
- Branch of mathematics concerned with the calculation of instantaneous rates of change (differential calculus) and the summation of infinitely many small factors to determine some whole (integral calculus)

순차적인 변화율을 구하는 수학

이루어진 조각들을 더해서 전체를 계산
적분
기법 사용



employed lowest error value
마이너스에서 최소값을 찾거나 하는 방법을 찾거나 하는 방법
differential calculus
optimization의 핵심



integral calculus

Essential Math for Data Mining: Optimization

- Optimization 최적화: 어떤 조건에서 가장
 - ▣ Collection of mathematical principles and methods used for solving optimization problems
 - ▣ Optimization problem is the problem of finding the best solution from all feasible solutions
 - In the simplest case, an optimization problem consists of maximizing or minimizing a real function

Different data, different optimization
→ to find best parameter,

Many problems in data mining,
finding best function to explain the given data well
Various optimization techniques are used to find
the optional function.

是 算法 問題。

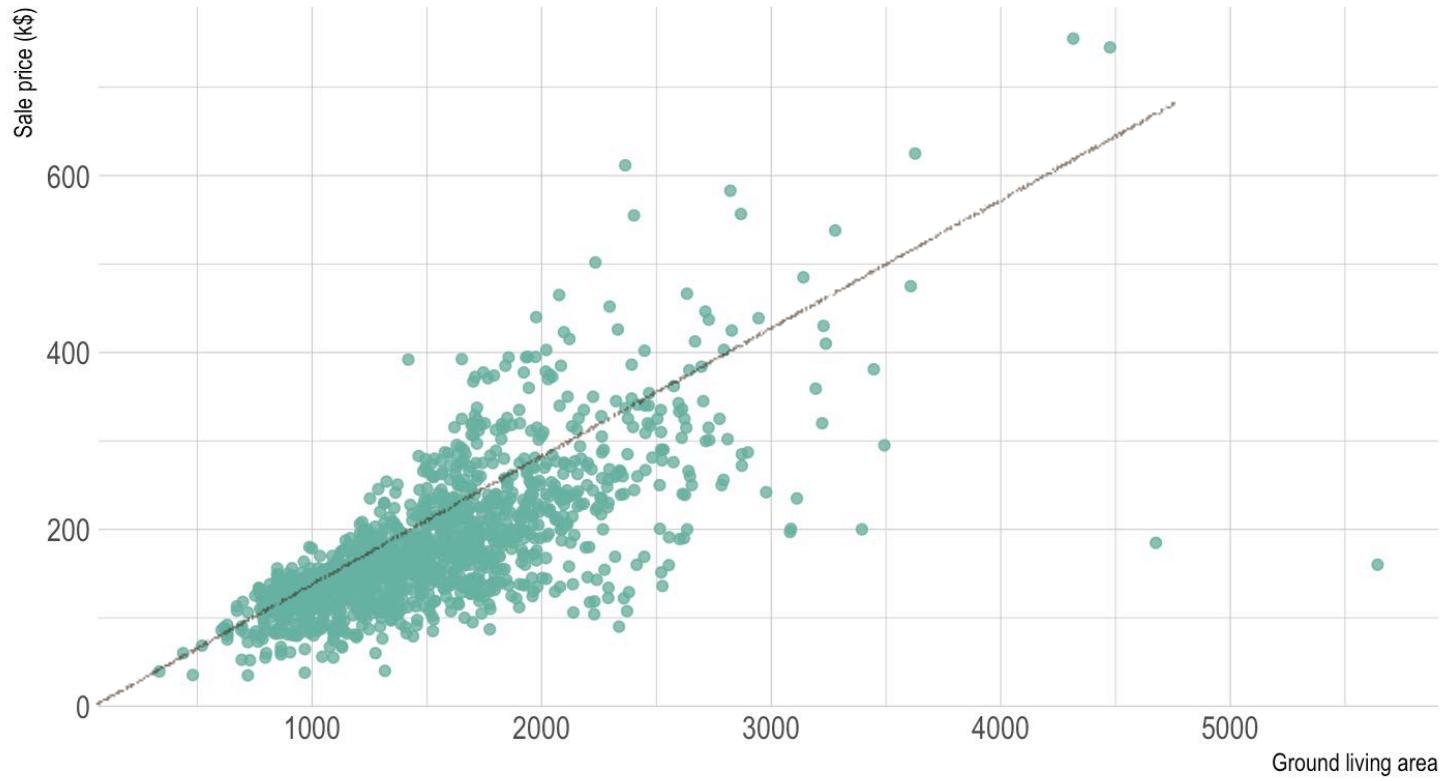


Statistics



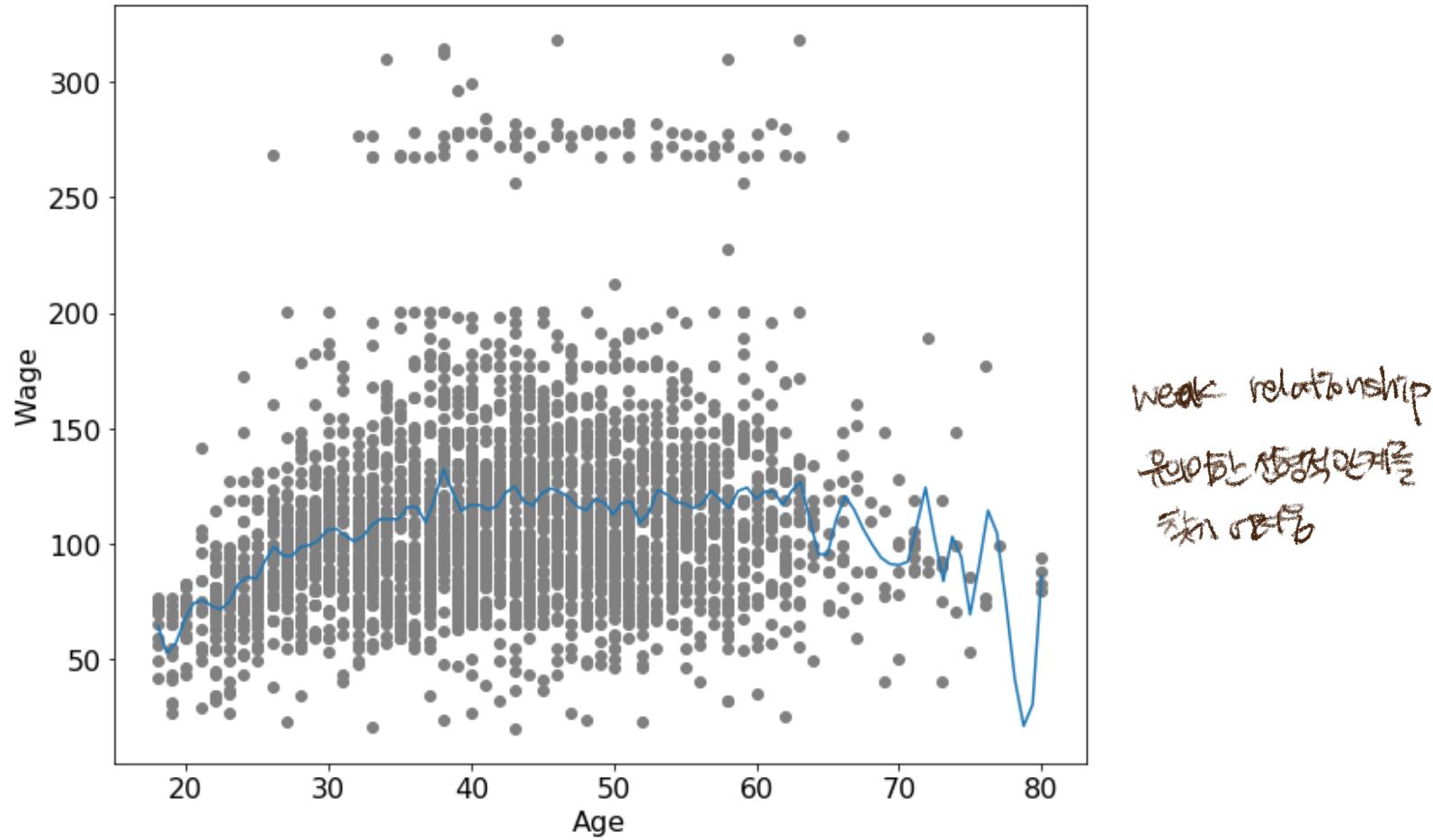
- A vast set of tools for understanding data

Ground living area partially explains sale price of apartments



Statistics

- A vast set of tools for understanding data

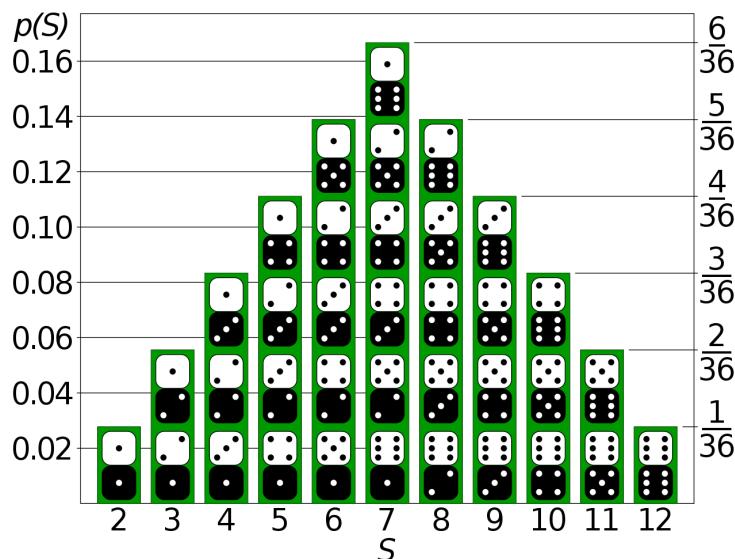


Statistics

- Descriptive statistics 기술통계
 - A summary statistic that quantitatively describes or summarizes features of a collection of information ↔ qualitative
 - Univariate 하나의 변수에 대한 Single variable
 - $\frac{\sum x_i}{n}$ Mean, Median, Mode center applicable!
 - Variance, standard deviation, Percentile
 - Skewness, kurtosis asymmetry
 - Bivariate or multivariate 2개 이상 multiple variables (relationship)
 - Cross-tabulations and contingency tables
 - Graphical representation via scatterplots
 - Quantitative measures of dependence (covariance, correlation)

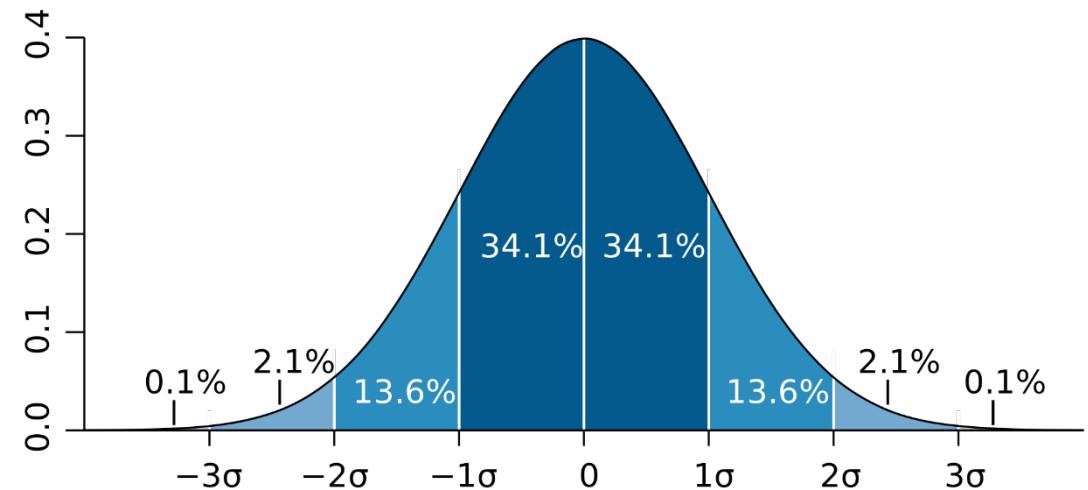
Statistics

- Probability distribution
 - A mathematical function that provides the probabilities of occurrence of different possible outcomes
 - The probabilities of occurrence of the specific observations



Discrete random variable →
probability mass function

비교적이 확률, 이동평균



Continuous random variable →
probability density function

비교적 평균

Statistics: Discrete Probability Distributions

이제의 베르누이 분포

비탄생률이 특별한 경우

Bernoulli distribution

n=1의 시점

2 outcome
0 or 1

discrete

- The discrete probability distribution of a random variable which takes the value 1 with probability p and the value 0 with probability $q = 1 - p$

$$\Pr(X = 1) = p = 1 - \Pr(X = 0) = 1 - q$$

- A special case of the binomial distribution where a single trial is conducted (so n would be 1 for such a binomial distribution)

- Probability mass function

$$f(X = x; p) = p^k(1 - p)^{1-k}$$

$$\begin{array}{ll} k=1 & p \\ k=0 & 1-p \end{array}$$

Binomial distribution

n개의 투표 결과

- The discrete probability distribution of the number of successes in a sequence of n independent experiments, each asking a yes-no question, and each with its own Boolean-valued outcome: success (with probability p) or failure (with probability $q = 1 - p$).

- Probability mass function

$$f(X = k; n, p) = \Pr(X = k) = \frac{n!}{k!(n - k)!} p^k(1 - p)^{n-k}$$

성공

실패

성공수

성공 (k) n개 중에 k개 성공하는 경우의 수

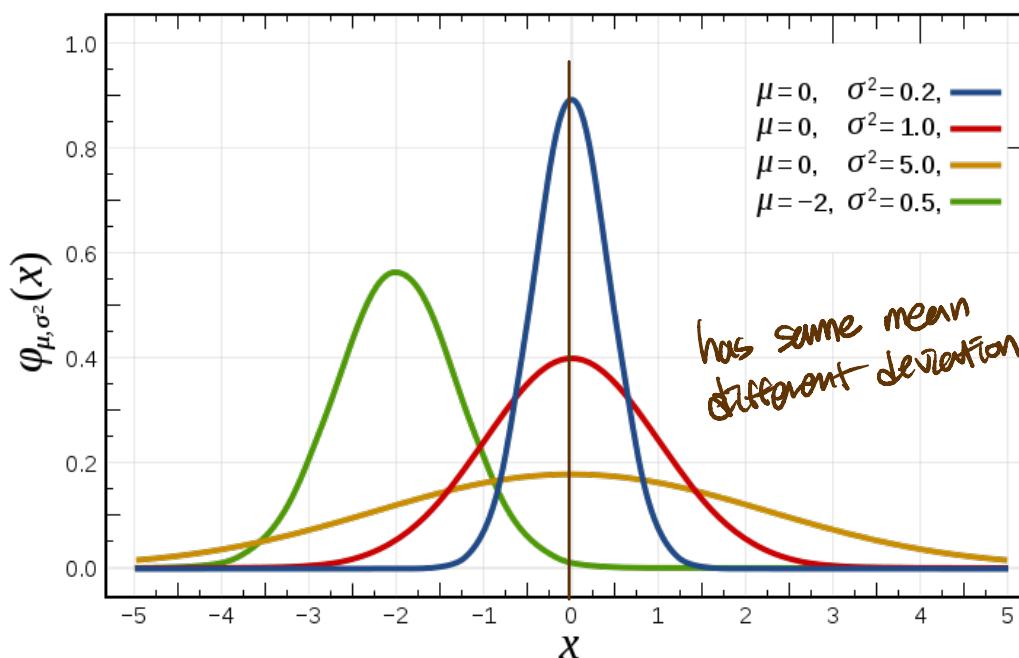
예) 동전을 n번 던져서 앞면이 k번 나온 경우

Statistics: Continuous Distributions

- Normal (Gaussian) distribution
 - ▣ Very common continuous probability distribution
 - ▣ Bell-shaped

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

① μ: mean
② σ: standard deviation
spread



Statistics: Continuous Distributions

Normal distribution

Student's t -distribution (t -distribution)

- Continuous probability distributions that arises when estimating the mean of a normally distributed population in situations where the sample size is small and population standard deviation is unknown

- Let X_1, \dots, X_n be independent and identically distributed (iid) as $N(\mu, \sigma^2)$

- Sample mean

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$$

- Sample variance

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

- The random variable $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ has a standard normal distribution
- The random variable $\frac{\bar{X} - \mu}{S/\sqrt{n}}$ has a Student's t -distribution with $n - 1$ degrees of freedom

표본 분산을 사용하여 표준화

standard deviation을
보통 표준 편차라고 하며
표본의 표준 편차는 표본 표준 편차라고 한다

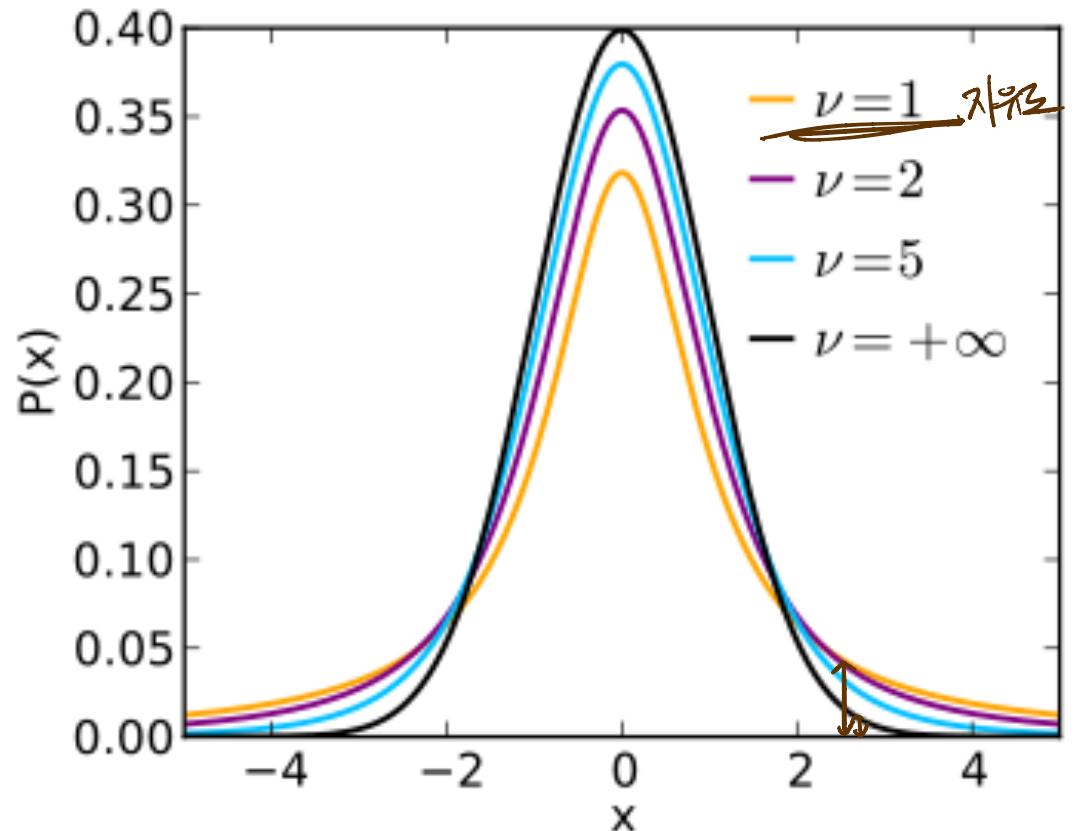
sample → error

정체성
statistical task

Statistics: Continuous Distributions

- The probability density function of t -distribution with varying degree of freedom *자유도가 바뀔 때 확률 밀도 분포가 어떻게 바뀌는가?*

자유도 ↑ thicker tail



자유도
observe extreme.
accuracy ↓

Statistics: Student's t -distribution

~~t-分布概要~~

□ Probability density function

$$f(x; \nu) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi}\Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}}$$

- ν : degree of freedom
- Γ : gamma function

$$\Gamma(n) = (n-1)! \quad \text{if } n \text{ is positive integer}$$

$$\Gamma(z) = \int_0^{-\infty} x^{z-1} e^{-x} dx$$

Statistics: Continuous Distributions

- Chi-squared distribution (χ^2)
 - The distribution of a sum of the squares of k independent standard normal random variables

표준정규분포를 가진 k 개의 독립적인 정규분포를 합한 결과의 분포 .

- Let X_1, \dots, X_k be independent, standard normal random variables

$$Y = \sum_{i=1}^k X_i^2$$

is distributed according to the chi-squared distribution with k degrees of freedom

자유도가 k 인 치아제곱분포를 갖는다 .

$$Y \sim \chi^2(k) \text{ or } Y \sim \chi_k^2$$

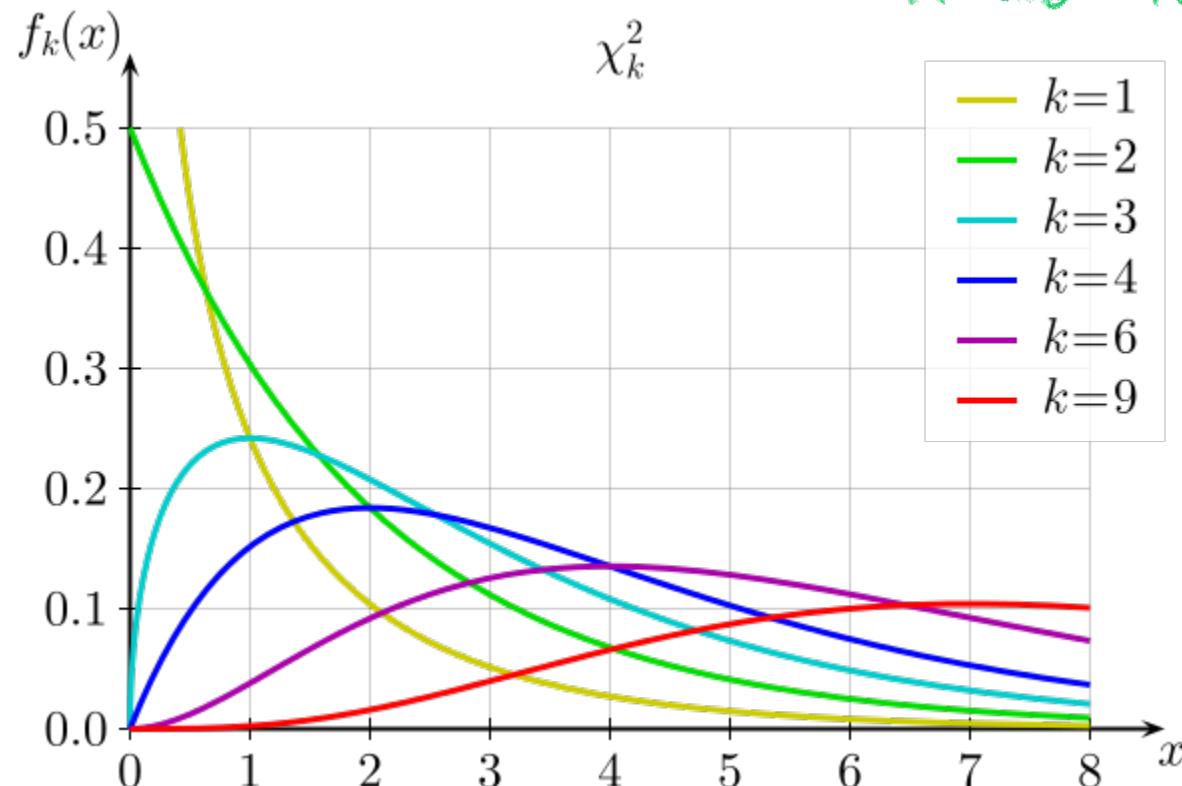
- Probability density function

$$f(x; k) = \frac{1}{2^{k/2} \Gamma(k/2)} x^{\frac{k}{2}-1} e^{-\frac{x}{2}}$$

Statistics: Continuous Distributions

- The probability density function of chi-squared distribution with varying degree of freedom

자유도 = 차원의 수
k↑ 대형 k↓ 확장 확장



Statistics: Continuous Distributions

- *F-distribution require 2 degree of freedom*
 - A random variate of the *F*-distribution with parameters d_1 and d_2 arises as the ratio of two appropriately scaled chi-squared variates
- Let X_1 and X_2 be two independent random variables and $X_1 \sim \chi^2(d_1)$ and $X_2 \sim \chi^2(d_2)$

$$Y = \frac{X_1/d_1}{X_2/d_2}$$

is distributed according to the *F*-distribution with d_1 and d_2 degrees of freedom

$$Y \sim F(d_1, d_2)$$

- Probability density function

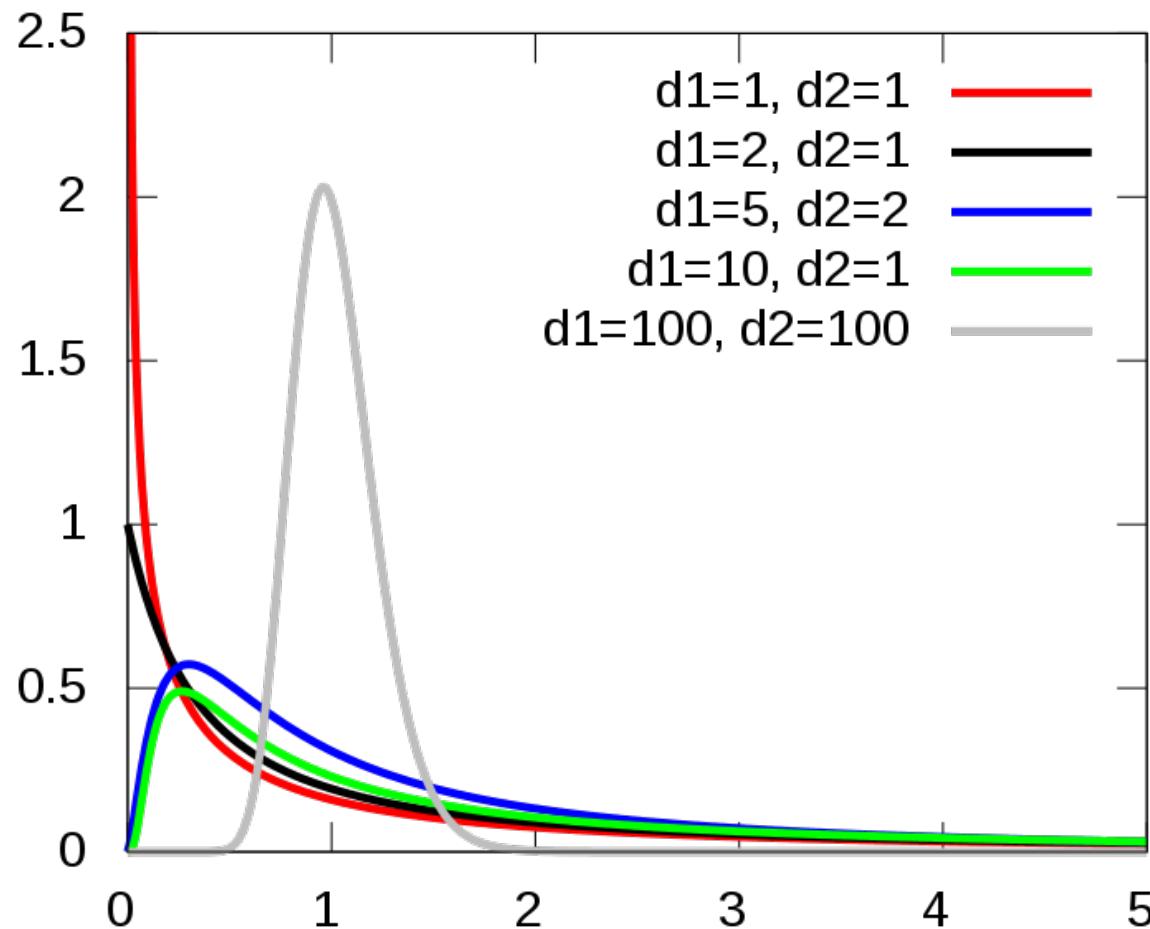
$$f(x; d_1, d_2) = \frac{\sqrt{\frac{(d_1 x)^{d_1} d_2^{d_2}}{(d_1 x + d_2)^{d_1 + d_2}}}}{x B\left(\frac{d_1}{2}, \frac{d_2}{2}\right)}$$

- B: beta function

$$B(x, y) = \int_0^1 t^{x-1} (1-t)^{y-1} dt$$

Statistics: Continuous Distributions

- The probability density function of F -distribution with varying degree of freedom



Linear Algebra

find ~~matrix~~

- Linear Algebra
 - ▣ Basic properties of matrix and vectors—scalar multiplication, linear transformation, transpose, conjugate, rank, determinant
 - ▣ Inner and outer products, matrix multiplication rule and various algorithms, matrix inverse
 - ▣ Special matrices—square matrix, identity matrix, triangular matrix, idea about sparse and dense matrix, unit vectors, symmetric matrix, Hermitian, skew-Hermitian and unitary matrices
 - ▣ Gaussian/Gauss-Jordan elimination, solving $Ax=b$ linear system of equation
 - ▣ Matrix factorization and decomposition
 - ▣ Vector space, basis, span, orthogonality, orthonormality, linear least square
 - ▣ Eigenvalues, eigenvectors, and diagonalization, singular value decomposition (SVD)

Linear Algebra

- Linear algebra is the study of vectors and linear functions
 - ▣ Scalar
 - A scalar is a number
 - ▣ Vector
 - A vector is a list of numbers
- ▣ Matrix
 - A matrix is also a collection of numbers *row / column*
 - The difference is that a matrix is a table of numbers rather than a list
- ▣ Linear equation
$$a_1x_1 + \cdots + a_nx_n = b$$
- ▣ Linear function
$$(x_1, \dots, x_n) \mapsto a_1x_1 + \cdots + a_nx_n$$

Linear Algebra

- Vectors

- ▣ Addition

$$\mathbf{v} + \mathbf{w}$$

- Example

$$\mathbf{v} + \mathbf{w} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} + \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 4 \\ 6 \end{bmatrix}$$

- ▣ Linear combination

$$a\mathbf{v} + b\mathbf{w}$$

- Example

$$3\mathbf{v} + 4\mathbf{w} = 3 \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 4 \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 15 \\ 22 \end{bmatrix}$$

$$\begin{bmatrix} 3 \\ 6 \end{bmatrix} \quad \begin{bmatrix} 12 \\ 16 \end{bmatrix}$$

Linear Algebra

□ Vectors

▣ Transpose

- **column** vector \leftrightarrow **row** vector
- Example

$$\mathbf{v} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \rightarrow \mathbf{v}^T = [1 \quad 2]$$

▣ Dot product, inner product

$$\mathbf{v} \cdot \mathbf{w}$$

- Example

$$\mathbf{v} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \mathbf{w} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$$
$$\mathbf{v} \cdot \mathbf{w} = (1)(3) + (2)(4) = 9$$

▣ Length

$$\|\mathbf{v}\| = \sqrt{\mathbf{v} \cdot \mathbf{v}}$$

- Example

$$\mathbf{v} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \|\mathbf{v}\| = \sqrt{(1)(1) + (2)(2)} = \sqrt{5}$$

Linear Algebra

- Matrix

- Addition

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} + \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} = \begin{bmatrix} a_{11} + b_{11} & a_{12} + b_{12} \\ a_{21} + b_{21} & a_{22} + b_{22} \end{bmatrix}$$

- Multiplication

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \cdot \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} 1 \cdot 5 + 2 \cdot 7 & 1 \cdot 6 + 2 \cdot 8 \\ 3 \cdot 5 + 4 \cdot 7 & 3 \cdot 6 + 4 \cdot 8 \end{bmatrix}$$

- Linear equation

$$Ax = \mathbf{b}$$

- Example

$$\begin{array}{rcl} x_1 & = & b_1 \\ -x_1 + x_2 & = & b_2 \\ -x_2 + x_3 & = & b_3 \end{array}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

Linear Algebra

□ Matrix

□ Inverse matrix

- An n -by- n square matrix, A is called invertible (or nonsingular) if there exists an n -by- n square matrix, B such that

$$AB = BA = I_n$$

where I_n denotes the n -by- n identity matrix which is a square matrix with ones on the main diagonal and zeros elsewhere

- B is the inverse of A (A^{-1})
- If A has no inverse, A is singular or non-invertible
- Example

$$A = \begin{bmatrix} -1 & \frac{3}{2} \\ 1 & -1 \end{bmatrix}, A^{-1} = \begin{bmatrix} 2 & 3 \\ 2 & 2 \end{bmatrix}$$

□ Solution of a linear equation

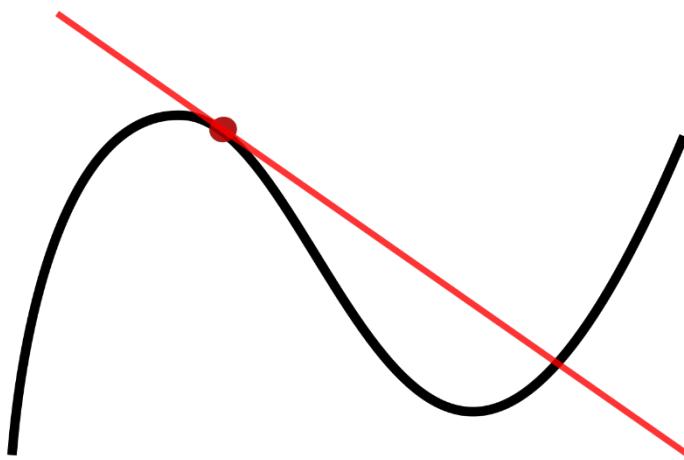
$$\mathbf{x} = A^{-1}\mathbf{b}$$

Calculus

□ Derivative

- A function of a real variable measures the sensitivity to change of the function value (output value) with respect to a change in its argument (input value)

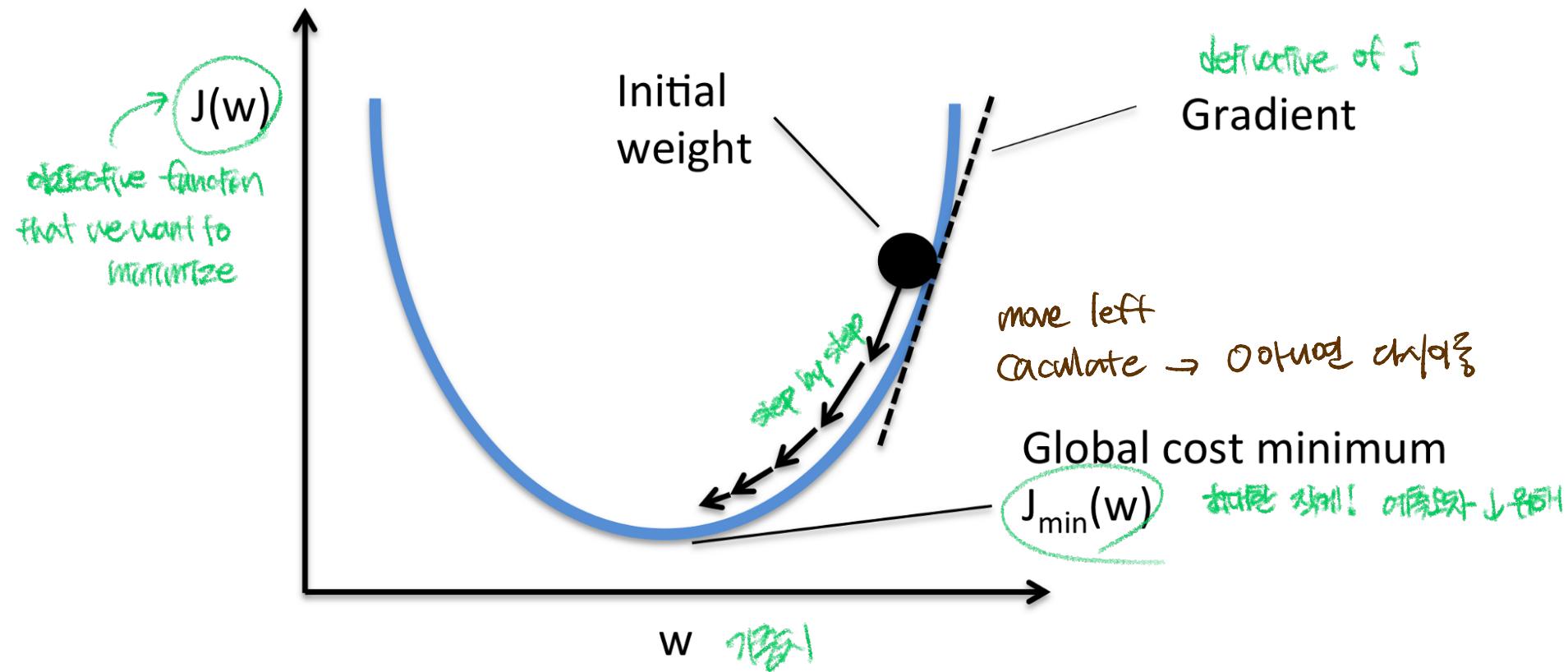
$$\frac{dy}{dx}$$



Calculus

Gradient Descent 경사하강법

Cost function을 미만화하는 법



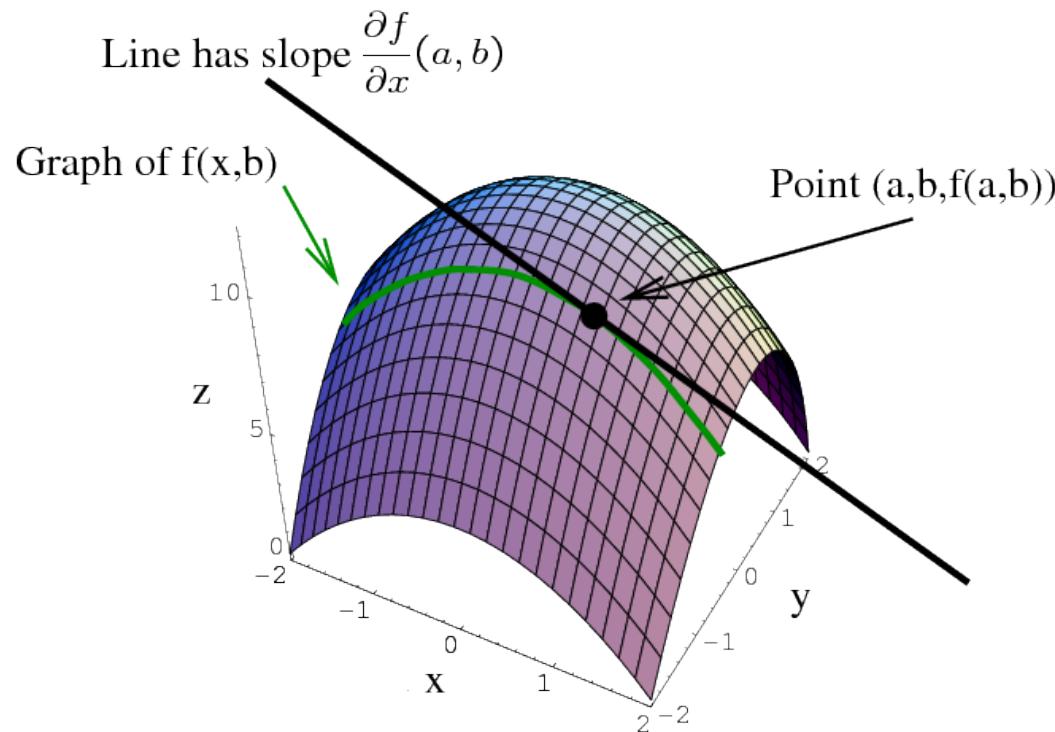
https://en.wikipedia.org/wiki/File:Gradient_Descent_in_2D.webm

Calculus

- Partial derivative ~~defn~~

- A function of several variables is its derivative with respect to one of those variables, with the others held constant

$$\frac{\partial f}{\partial x}$$



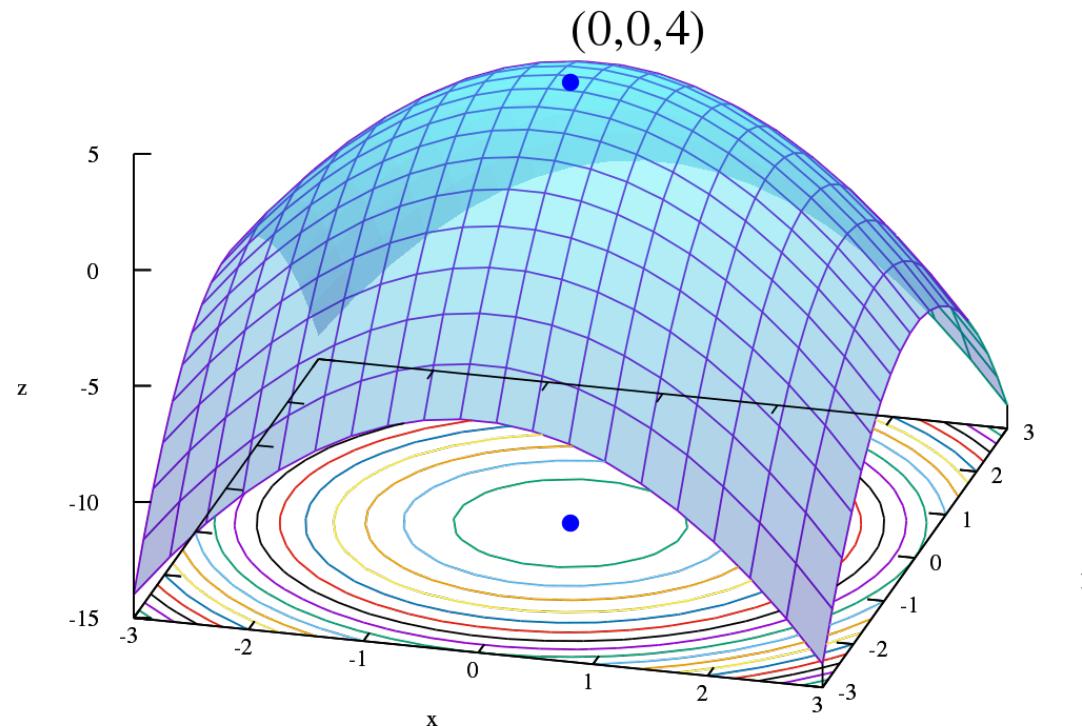
Optimization

- Optimization
 - ▣ Basics of optimization —how to formulate the problem
 - ▣ Linear programming, simplex algorithm
 - ▣ Integer programming *solution should be Integer*
 - ▣ Constraint programming, knapsack problem
 - ▣ Randomized optimization techniques—hill climbing, simulated annealing, Genetic algorithms

연속적이고 미지수의 수가 적은 경우

Optimization

- Optimization problem *to find max/min variable.*
 - Maximizing or minimizing a real function by systematically choosing input values from within an allowed set and computing the value of the function



Optimization

- Example **linear optimization problem**
 - For materials, the manufacturer has 750 m^2 of cotton textile and $1,000 \text{ m}^2$ of polyester. Every pair of pants (1 unit) needs 1 m^2 of cotton and 2 m^2 of polyester. Every jacket needs 1.5 m^2 of cotton and 1 m^2 of polyester.
 - The price of the pants is fixed at \$50 and the jacket, \$40.
 - **What is the number of pants and jackets that the manufacturer must give to the stores so that these items obtain a maximum sale?**

- Variables to be determined

$$x = \text{number of pants}$$

$$y = \text{number of jackets}$$

- Objective function

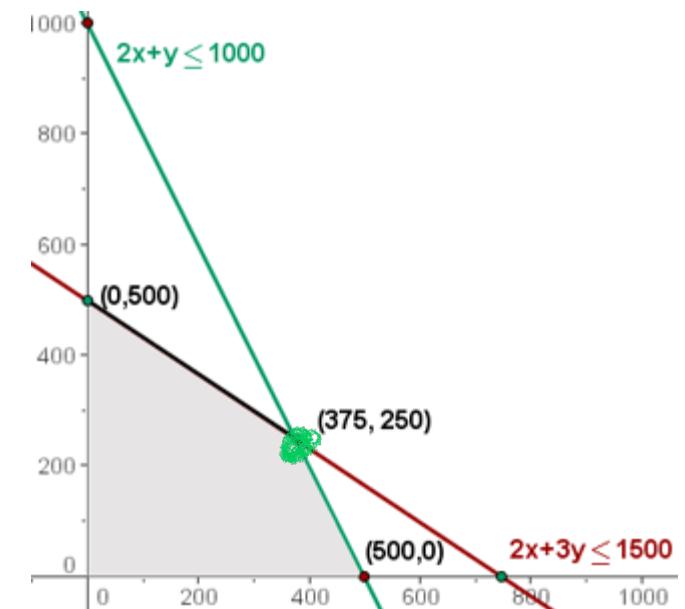
graph $f(x, y) = 50x + 40y$

- Constraints

$$x + 1.5y \leq 750$$

$$2x + y \leq 1000$$

zu 1500



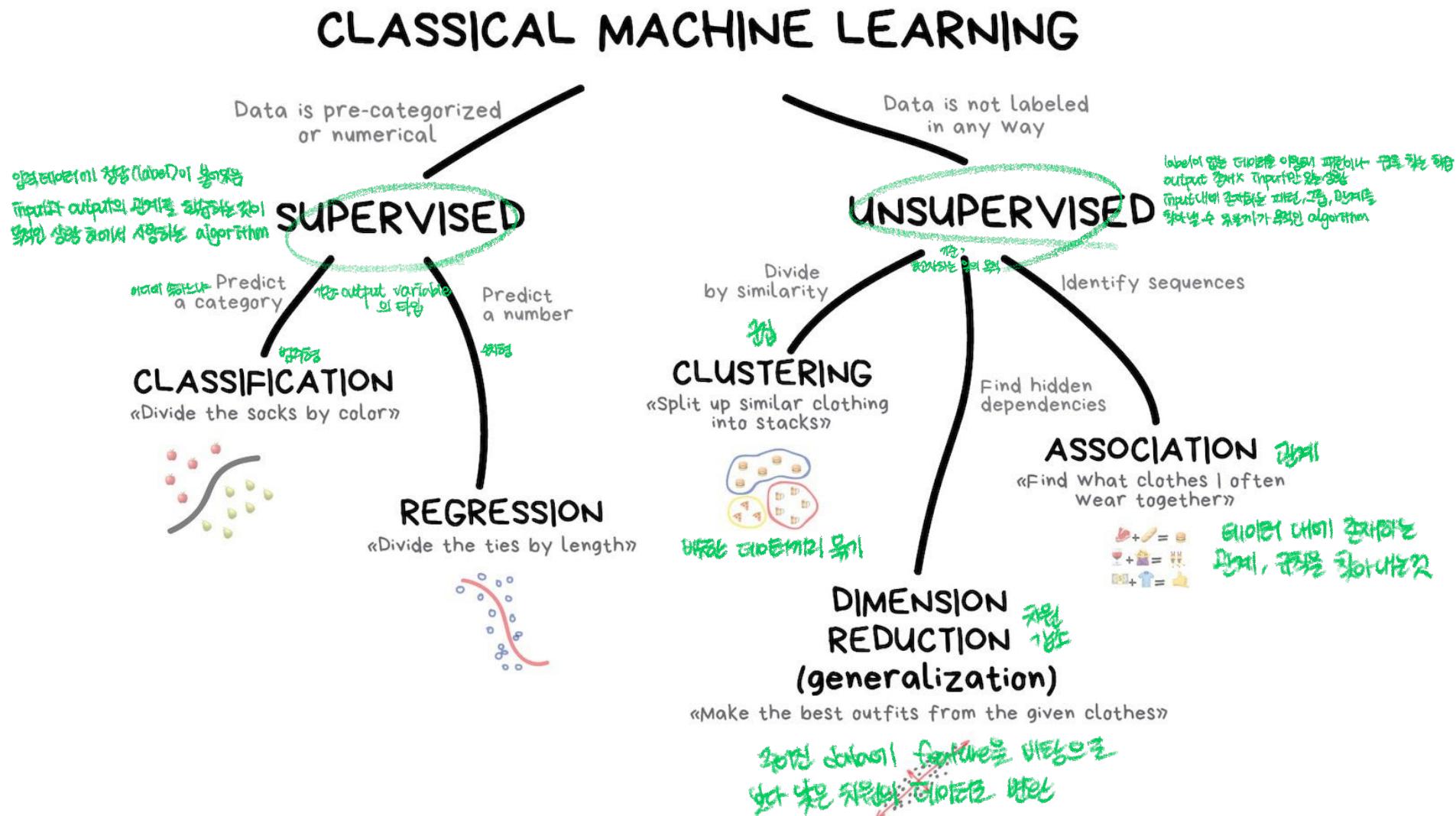
Optimization

- Why are optimization algorithms important for data analysis?
 - ▣ One of fundamental data analysis tasks is to seek a function that approximately maps \mathbf{x}_i to y_i for each observation, i
$$y = f(\mathbf{x})$$
 - ▣ The process of finding f based on data is called learning or training
 - During learning, optimization algorithms provide a tool to find the most appropriate f

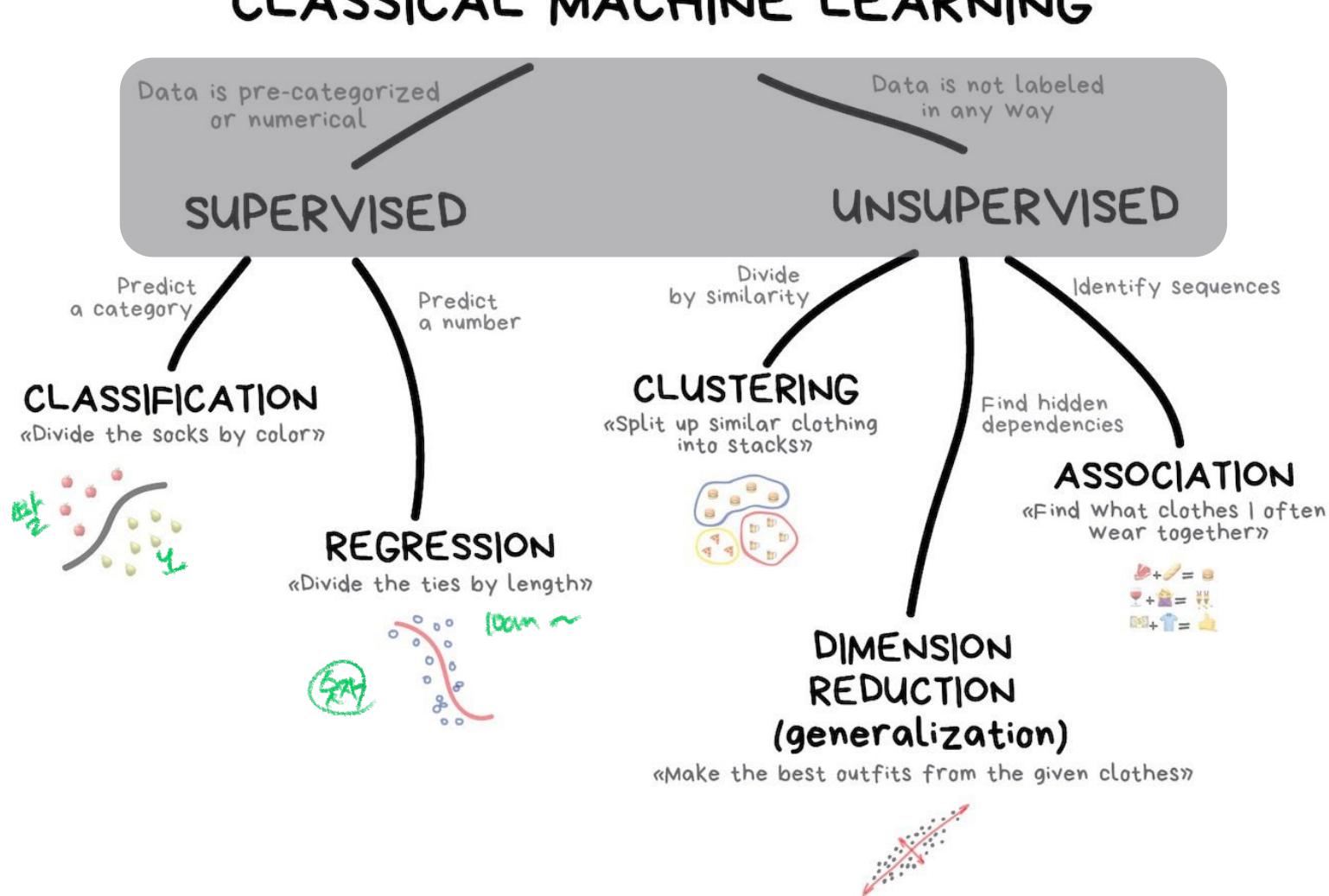
widley used to final model predict sth

Basic Terminologies

Topics Covered in This Class

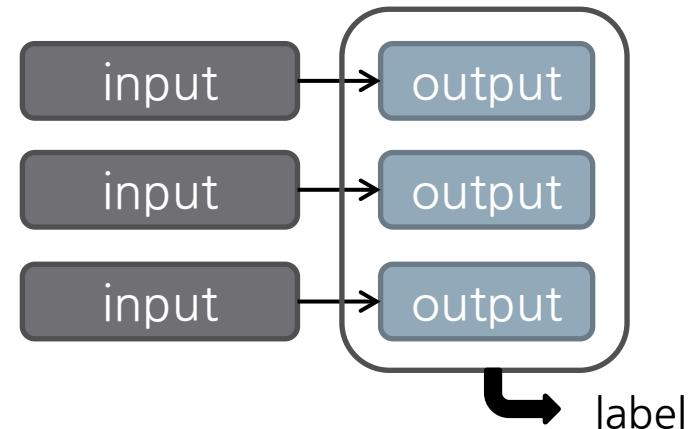


Topics Covered in This Class

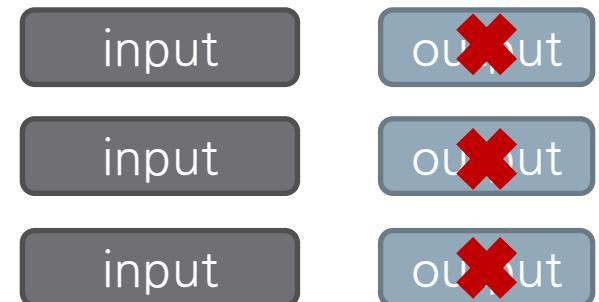


Types of Learning

- Supervised learning
 - ▣ We have knowledge of output
 - We call such data labeled
 - We know answer
 - ▣ Goal
 - Estimate output for unlabeled input

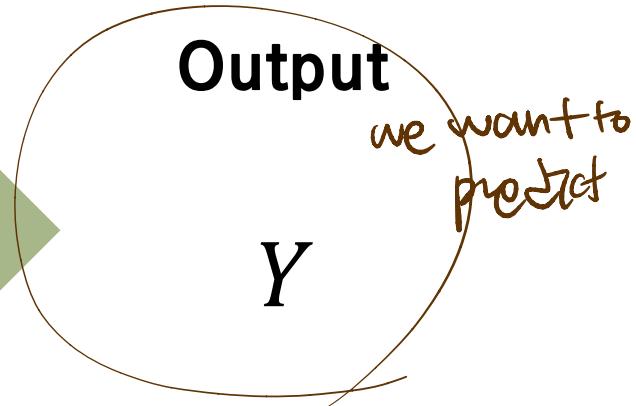
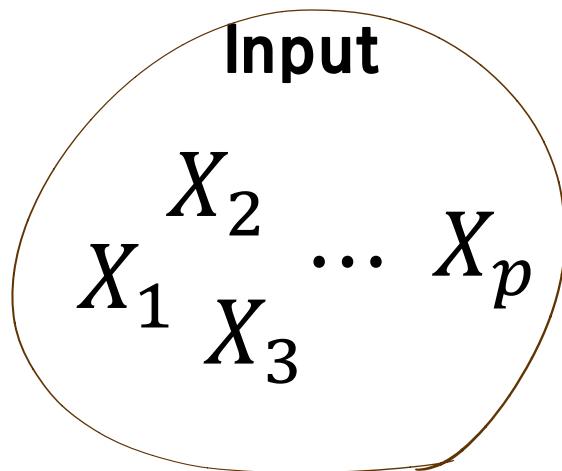


- Unsupervised learning
 - ▣ No output
 - We call such data unlabeled
 - ▣ Goal
 - Find patterns, groups, or relation



Supervised learning

ex)



\downarrow Supervised Learning

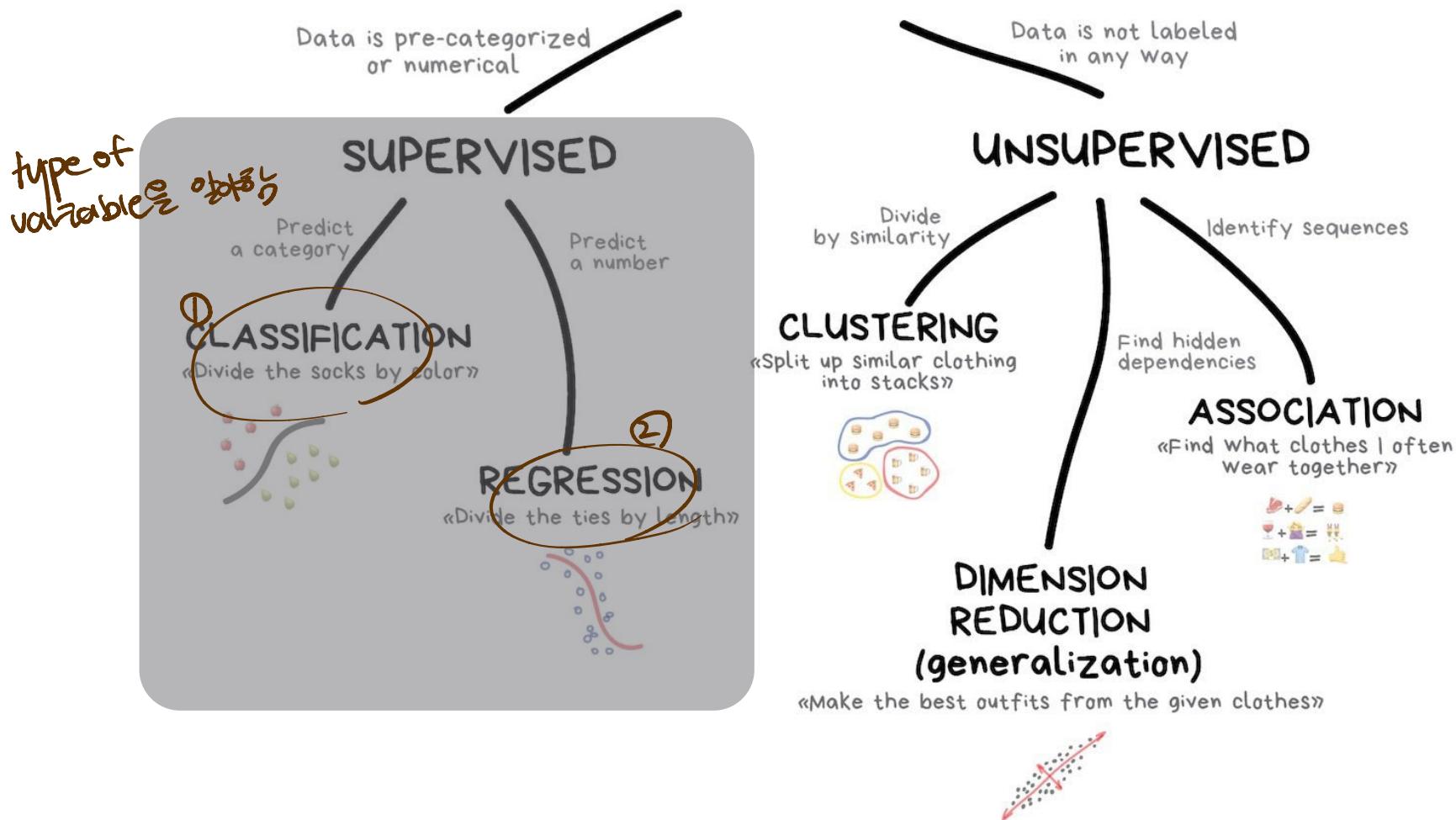
$$Y = f(X_1, X_2, \dots, X_p)$$

function

check every functional form

Topics Covered in This Class

CLASSICAL MACHINE LEARNING



Data for Data Mining: Structured Data

↖ dataset can be express as table .

- Example of data set
 - ▣ The input data set is usually expressed as a set of independent instances

instance,
sample,
example

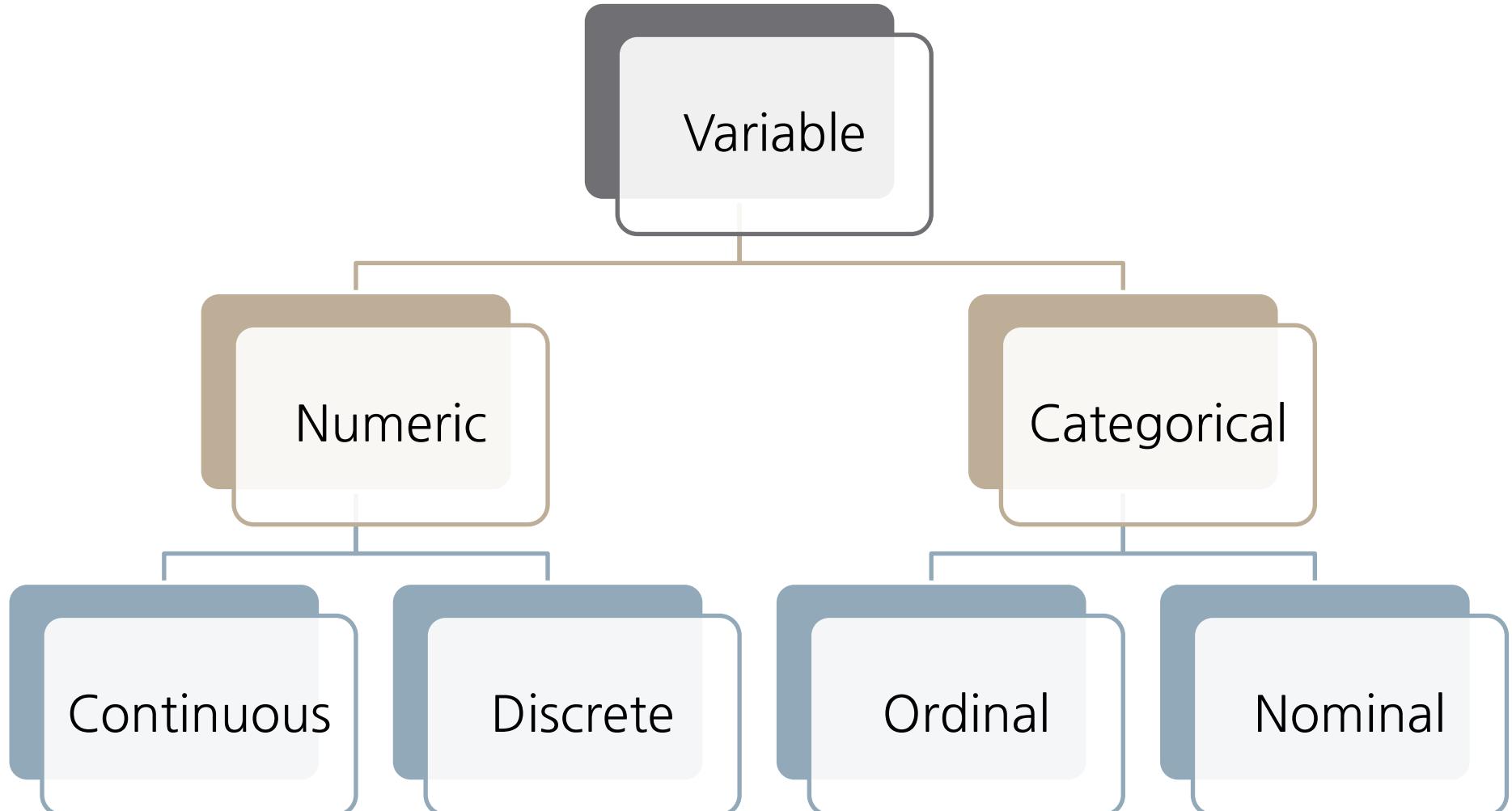
Outlook	Temperature(°F)	Humidity(%)	Windy	Play Time(min)
Sunny	85	85	false	5
Sunny	80	90	true	0
Rainy	70	96	false	40
Rainy	68	80	false	65
Sunny	72	95	false	0
Sunny	69	70	false	70
Rainy	75	80	true	45

variable,
attribute,
feature

Types of Data

- Structured
 - ▣ Values of variable reside in a fixed field
 - ▣ Examples
 - Numeric
 - Date
 - Restricted terms: (male, female), (Mr., Ms., Mrs.)
 - Address
- Unstructured
 - ▣ Values of variable do not reside in a fixed field
 - ▣ Examples
 - Documents
 - Webpages
 - Images
 - Videos

Structured Data: Types of Variables



Structured Data: Types of Variables

- Numeric (Quantitative)
 - ▣ A broad category that includes any variable that can be counted, or has a numerical
- Continuous
 - ▣ A variable with infinite number of values
 - ▣ Example
 - Many numeric variables: temperature, weight, height, pressure and etc.
- Discrete *> count*
 - ▣ A variable that can only take on a certain number of values or have a countable number of values between any two values
 - ▣ Example
 - The number of cars in a parking lot
 - the number of flaws or defects

Structured Data: Types of Variables

- Categorical
 - A variable that contains a finite number of categories or distinct groups
- Nominal
 - A Variable that has two or more categories, but there is no intrinsic ordering to the categories.
 - Example
 - (Male, Female), (Class 1, Class 2, Class 3), (Red, Yellow, Green)
- Ordinal
 - Similar to a nominal variable, but the difference between the two is that there is a clear ordering of the variables.
 - Example
 - Score: A+, A, A-, B+, B, B-, C+, C, C-, D, F
 - Size: S, M, L, XL, XXL

Example: The Input to a Data Mining

- Example of data set

can't be non-integer

num-of-doors	body-style	wheel-base	length	make
2	convertible	88.6	168.8	Audi
2	convertible	88.6	168.8	BMW
2	hatchback	94.5	171.2	Chevrolet
4	sedan	99.8	176.6	BMW
4	sedan	99.4	176.6	Audi
2	sedan	99.8	177.3	Audi
4	wagon	105.8	192.7	Chevrolet

Types:

Discrete

Nominal

Continuous

Continuous

Nominal

can't be compare

non-negative.

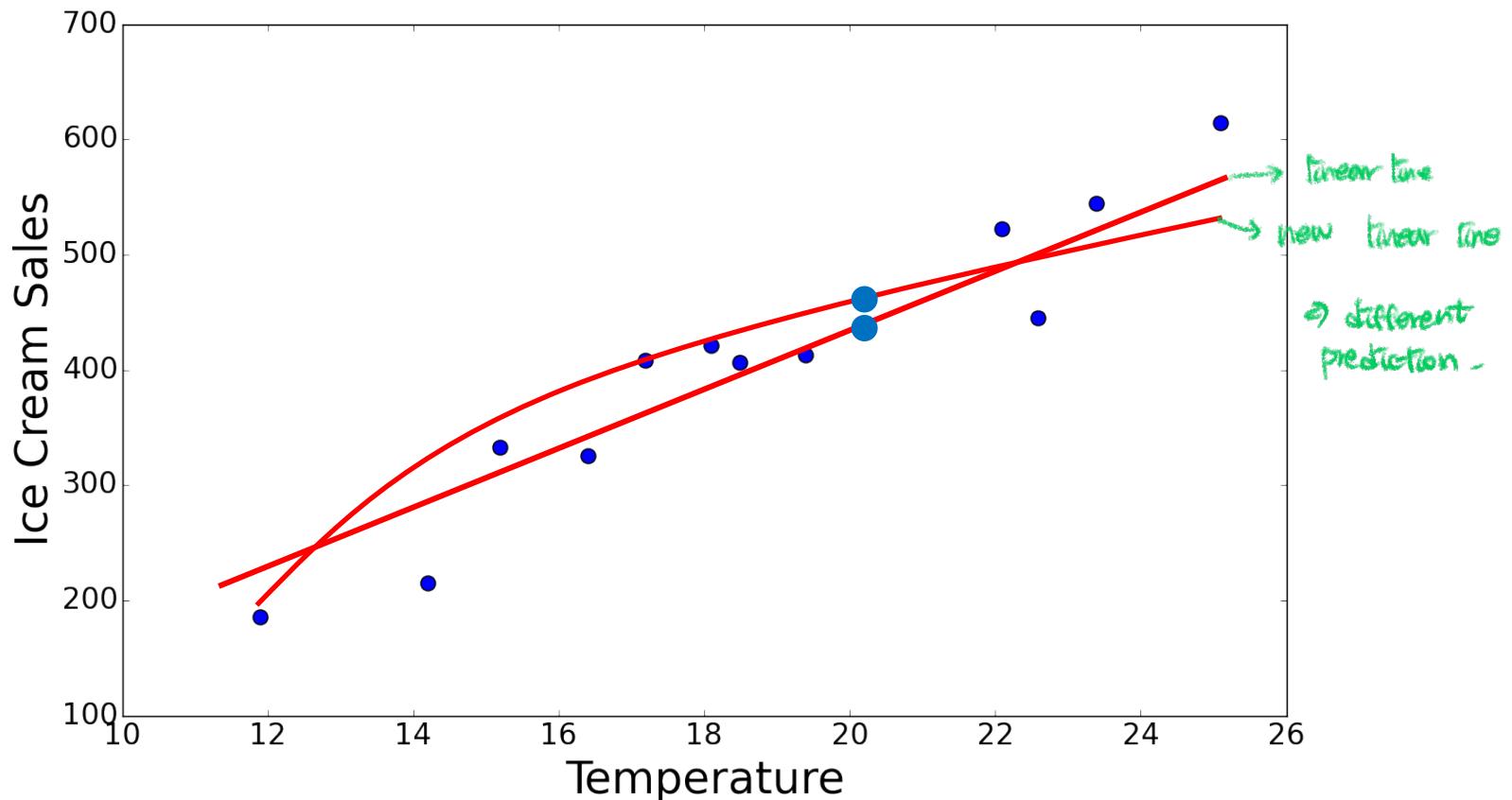
Supervised Learning: Regression

→ numeric \rightarrow types of output

- Temperature vs. Ice Cream Sales
 - How about 21°C?

find the relation

→ predict Icecream sales for one day
with new average temperature.

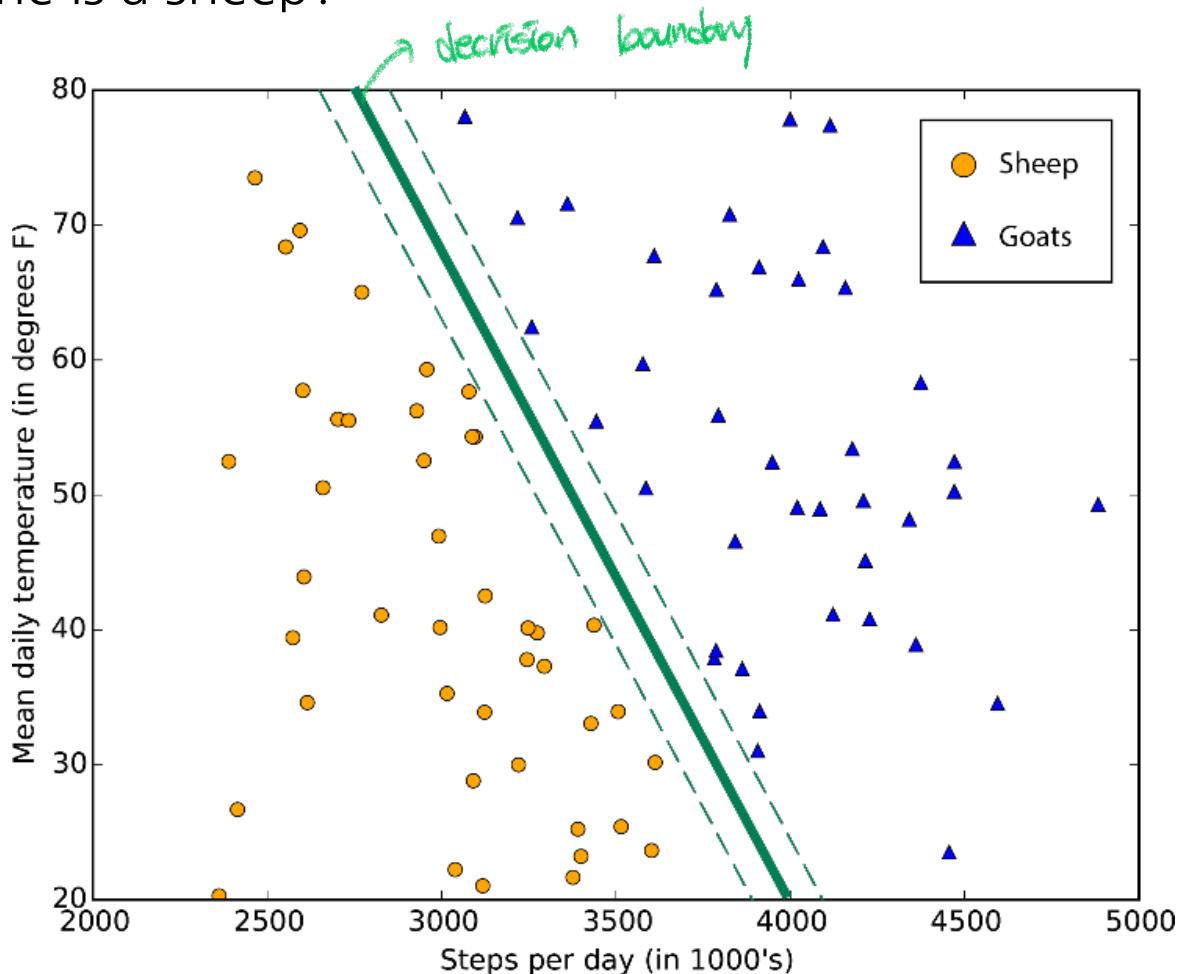


이것 (회귀선)을 따라 예상된 출력(결과값)을 예측하는 것

Supervised Learning: Classification

find the function to discriminate two classes using this two factors

- Which one is a sheep?



Question

- Suppose you are working on weather prediction, and you would like to predict whether or not it will be raining at 5pm tomorrow. You want to use a learning algorithm for this. Would you treat this as a classification or a regression problem?
① Regression
 ② Classification

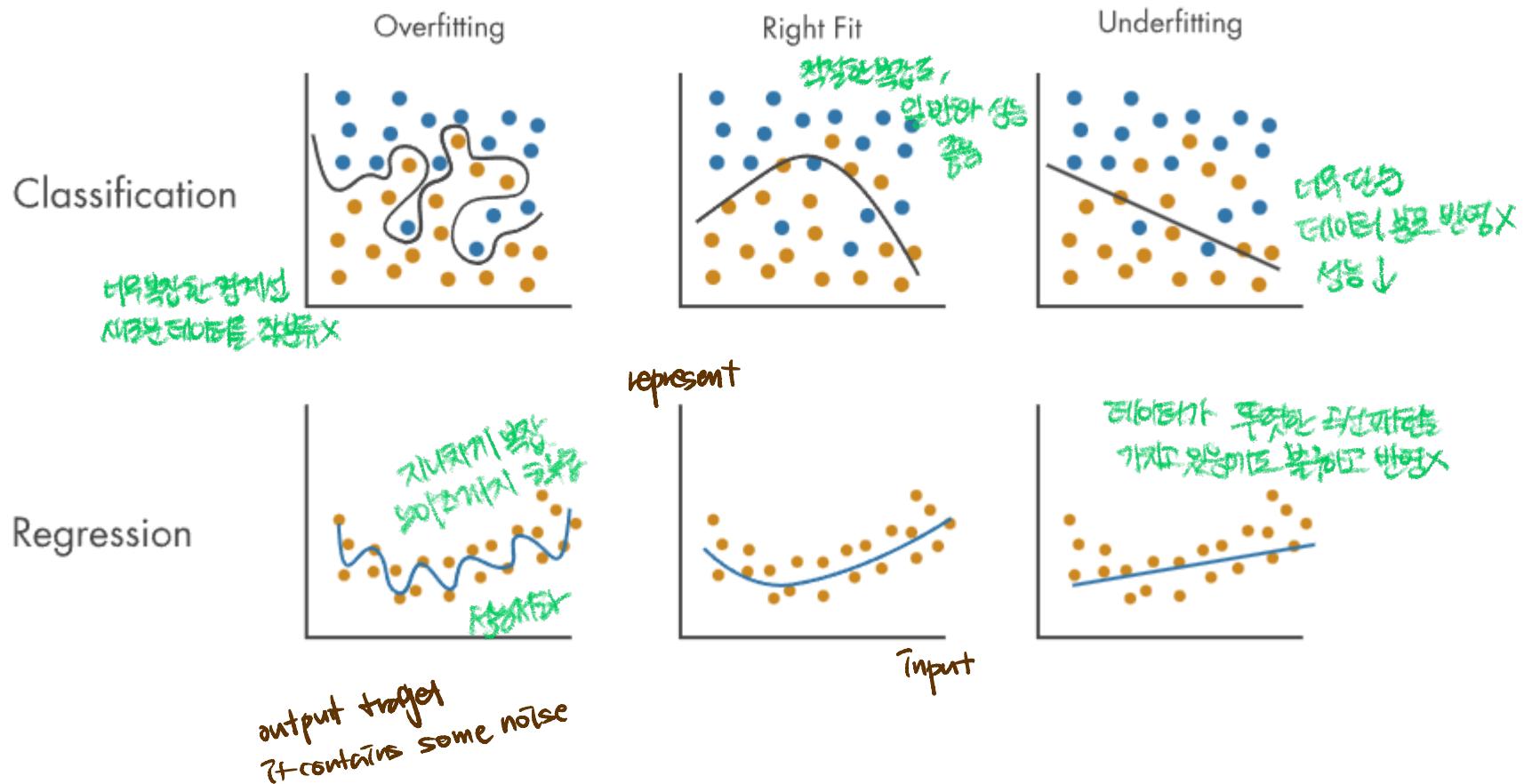
- Suppose you are working on stock market prediction, and you would like to predict the price of the specific stock tomorrow (measured in dollars). You want to use a learning algorithm for this. Would you treat this as a classification or a regression problem?
 ① Regression
② Classification

Overfitting vs. Underfitting

- Overfitting *모델이 훈련 데이터에 지나치게 맞춰진 경우*
 - ▣ Overfitting is a machine learning problem that occurs when a model is too closely aligned to training data, causing it to perform poorly on new data.
 - ▣ How it happens
 - The model is too complex *복잡*
 - The training data is too small or contains irrelevant information *작고 헛된 정보*
 - The model memorizes subtle patterns in the training data *훈련 데이터를 완벽히 배운다*
 - ▣ Why it's a problem
 - An overfit model can't generalize well to new data *일반화 불가*
 - It can give inaccurate predictions *예측 정확성 X*
 - It can't perform well for all types of new data *다양한 예제에 대한 일반화는 성공하지 X*
- Underfitting *B*
 - ▣ Underfitting occurs when a machine learning model is too simple to capture the underlying patterns in data. *모델이 너무 단순해서 데이터를 잘 찾지 못해 X*

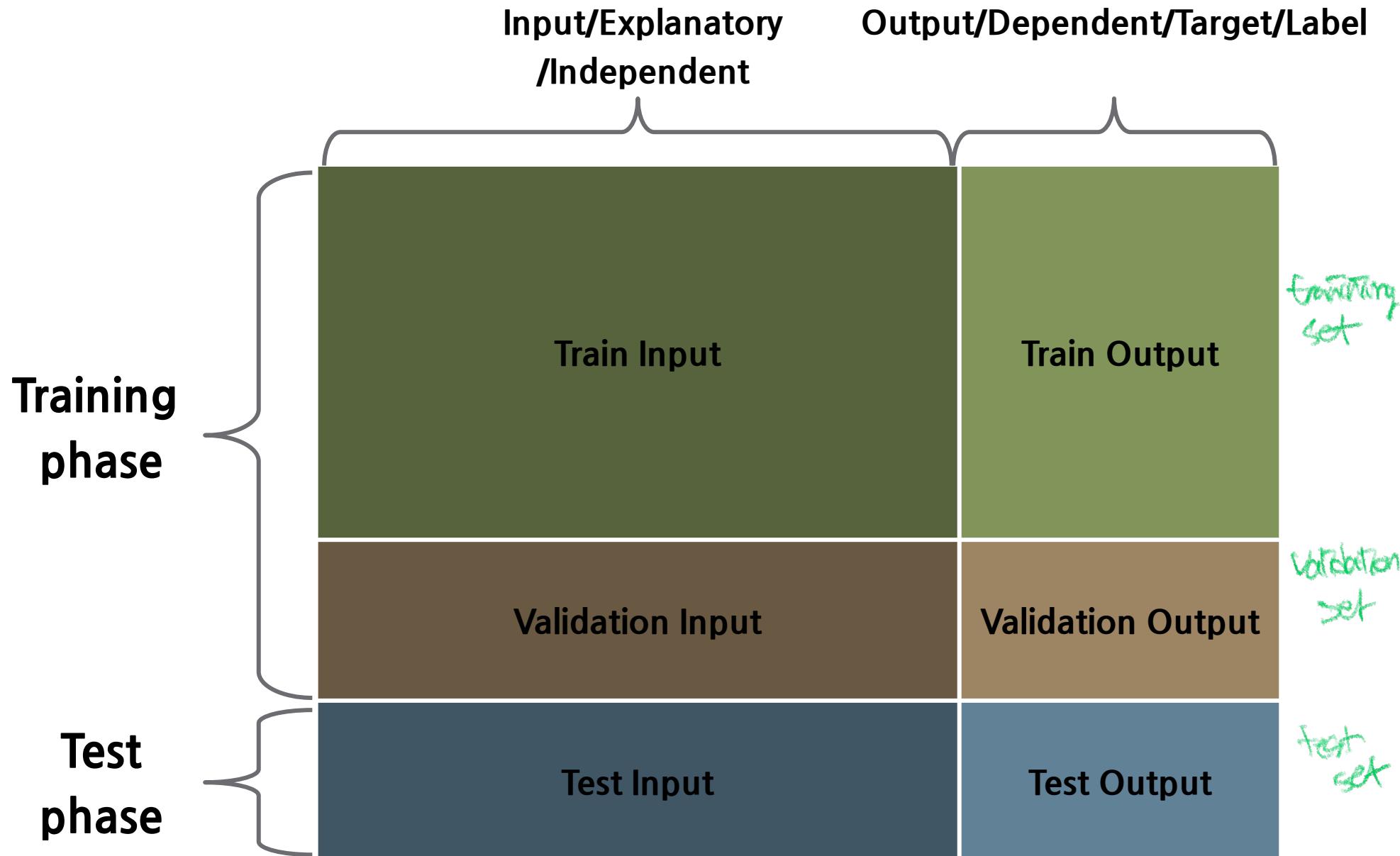
*일반화는 잘되나
성능↓ 예측률↓*

Overfitting vs. Underfitting



Data Partition

정제데이터를
기반으로 사용,
원본은 사용



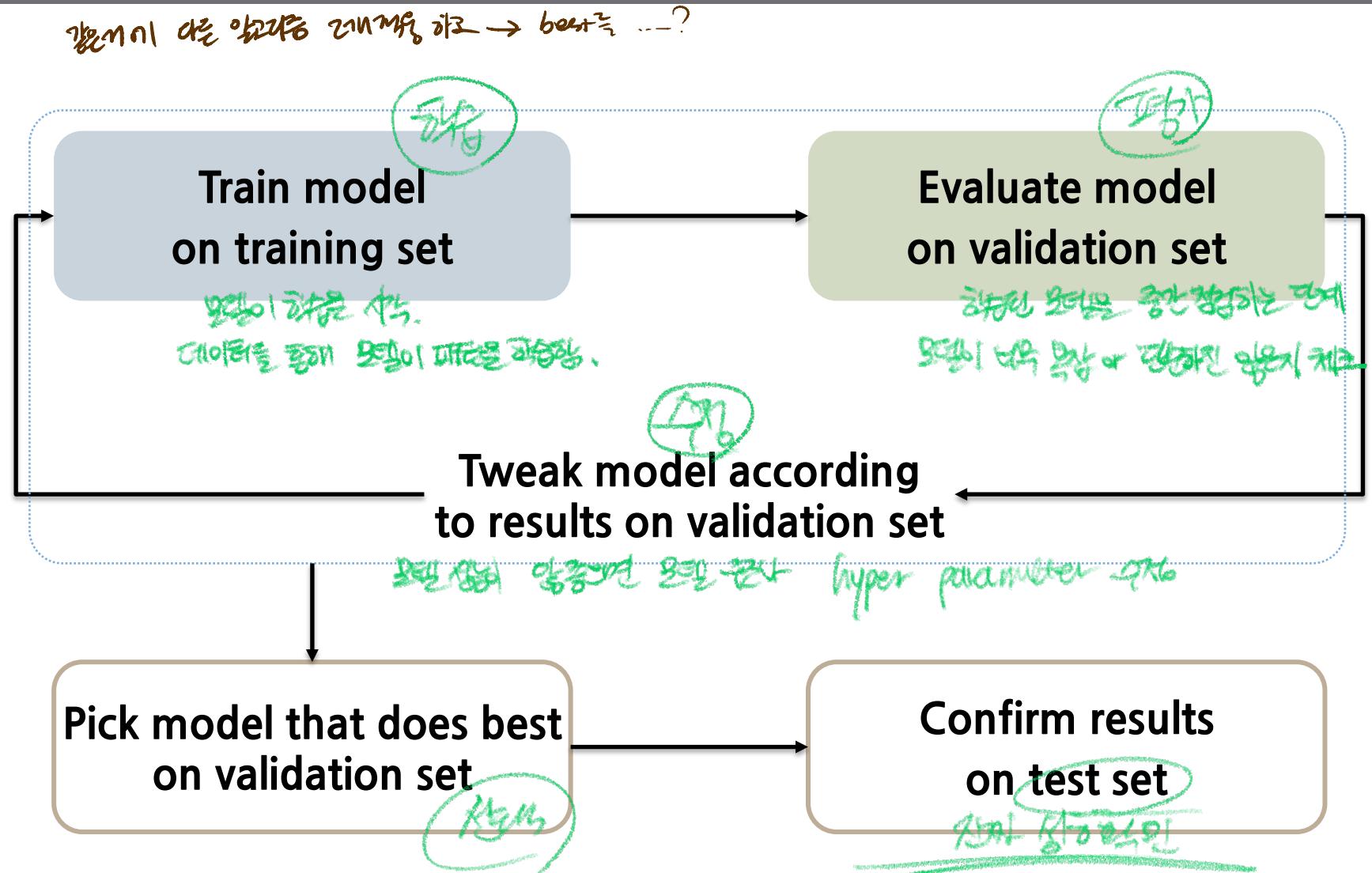
Data Partition

- Training set *학습 데이터로 모델의 파라미터를 찾는다*
 - Purpose: The training set is used to train the model. It contains the labeled examples the model will learn from. During the training phase, the model's parameters are adjusted based on the data in the training set.

- Validation set *학습한 데이터로 모델 성능을 평가하기 위해서*
 - Purpose: The validation set is used to tune hyperparameters and evaluate the model during training. It helps in selecting the best version of the model.
여기서 ↗ 학습하고 초기의 모델을 평가

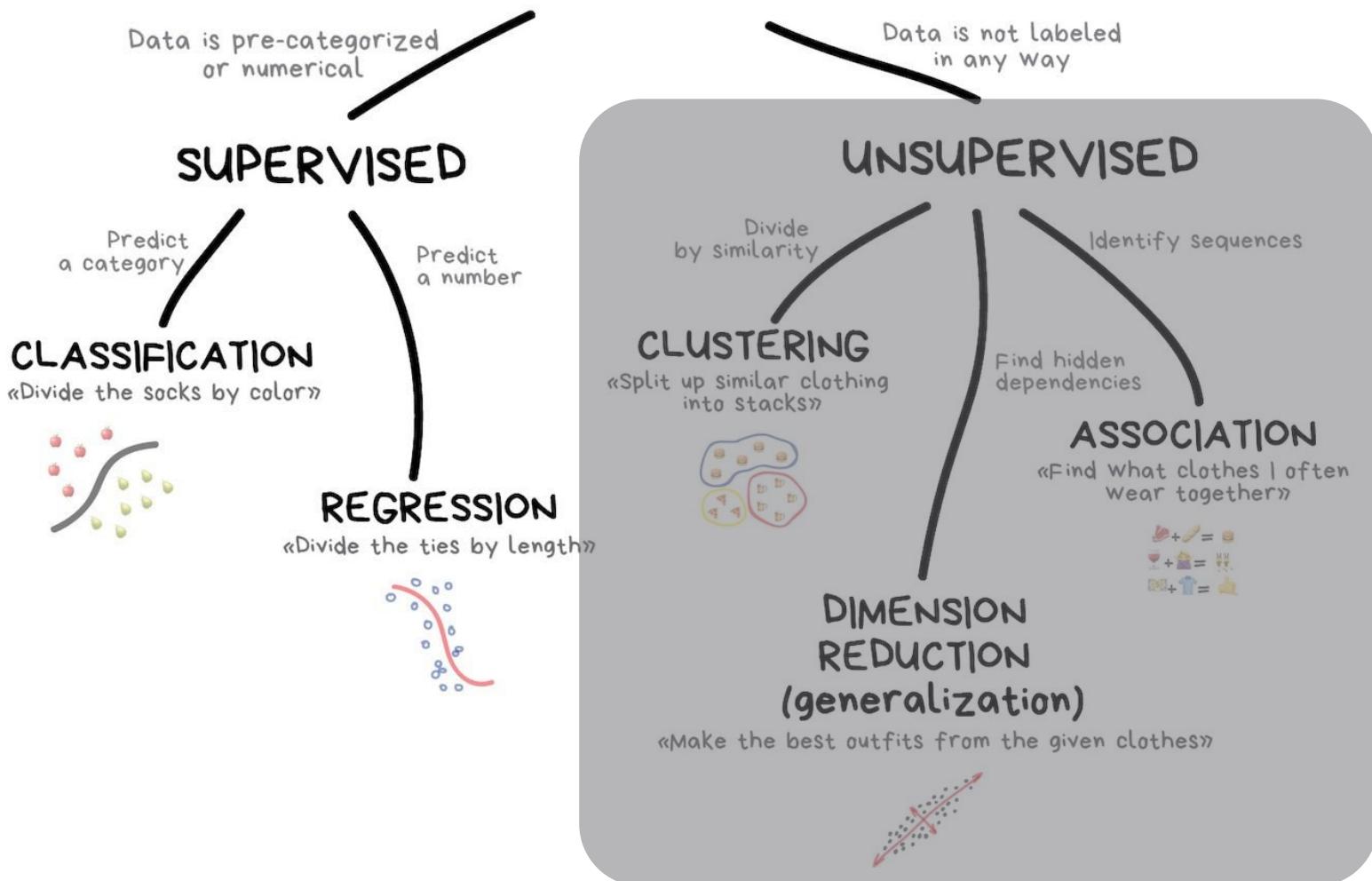
- Test set *학습한 데이터로 모델의 최종 성능을 평가하기 위해서*
 - Purpose: The test set is used to evaluate the model's final performance after training and validation. This set simulates new, unseen data, giving an unbiased estimate of how the model will perform in a real-world scenario.
모든 것들에 대해서 모델을 시험해보는 그룹 . *얼마나 잘 돌아갈지*

Process of Supervised Learning with Partitioned Data



Topics Covered in This Class

CLASSICAL MACHINE LEARNING

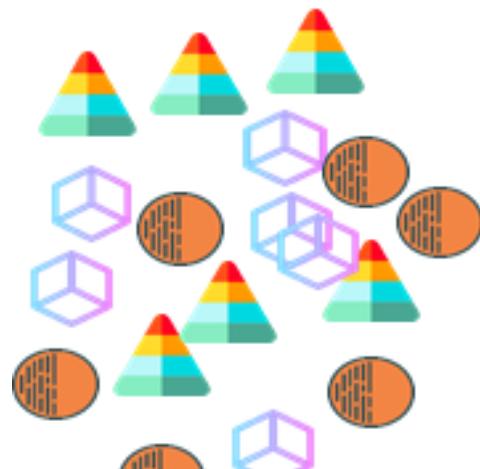


Unsupervised Learning: Clustering

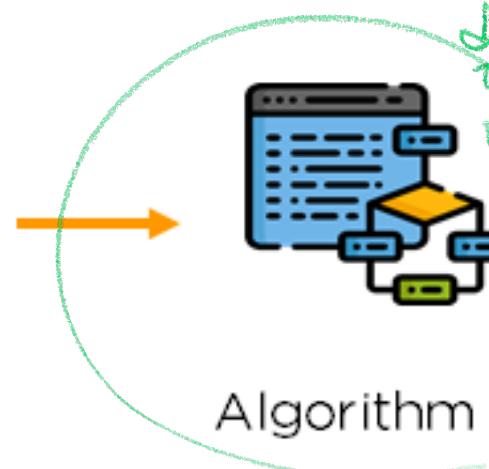
- Grouping data points

- How to determine which group does each data belongs to?

make ~~one~~ group ~~한 그룹 만들기~~ ~~한 그룹~~

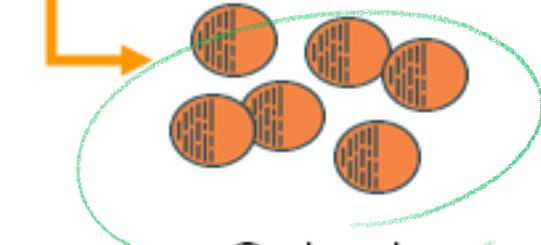
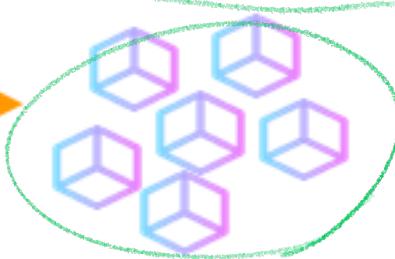


Raw Data
label X



원하는 결과
제공하는 모델
방법

divide
into
multiple
groups

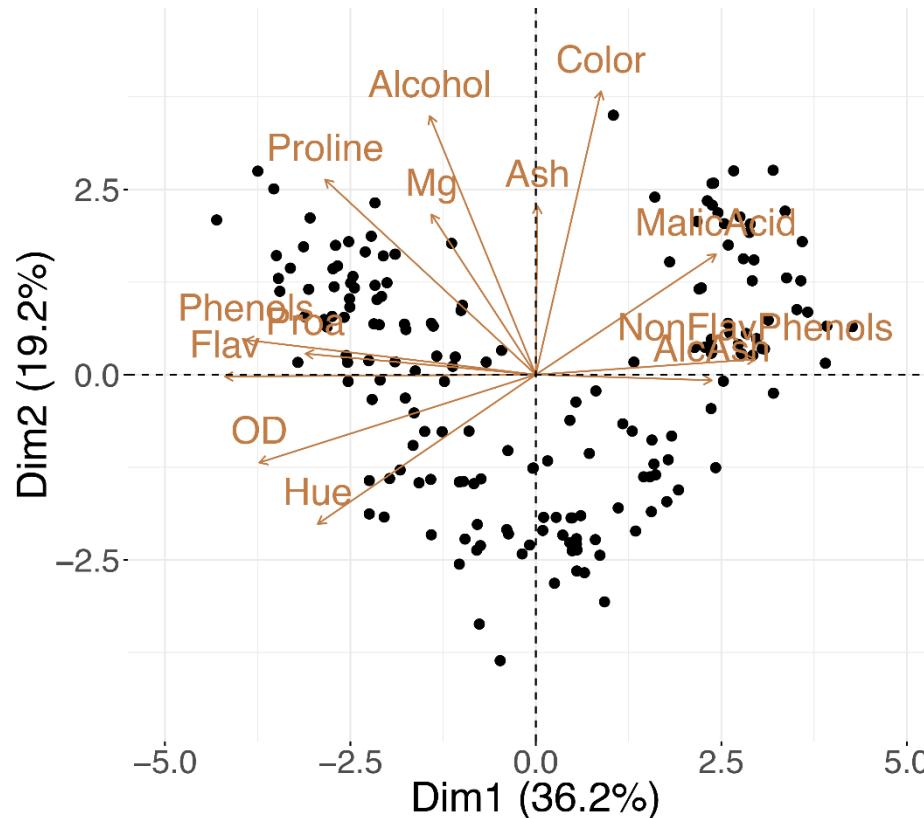


Output
groupings

Unsupervised Learning: Dimensionality Reduction

차원削減

- Dimensionality reduction
 - ▣ The process of reducing the number of random variables under consideration by obtaining a set of principal variables
 - ▣ High dimension → Low dimension



Unsupervised Learning: Association Rule Mining

- Find useful information from transactions

Datetime	Customer	Items
2015-07-15 14:03	1	orange juice, banana
2015-07-15 16:20	2	orange juice, milk
2015-07-16 10:14	3	detergent, banana, orange juice
2015-07-25 19:34	2	milk, bread, soda
2015-07-29 09:41	4	detergent, window cleaner
2015-08-01 20:55	1	bread, milk

- One of useful information is information like “If item A then item B”
 - This information is called association rule
- Find pair of items that are more likely to be purchased together based on transactions

Question

- Of the following examples, which would you address using an unsupervised learning algorithm? (Find all that apply.)
 - ① Given email labeled as spam/not spam, learn a spam filter.
 - ② Given a set of news articles found on the web, group them into set of articles about the same story.
 - ③ Given a database of customer data, automatically discover market segments and group customers into different market segments.
 - ④ Given a dataset of patients diagnosed as either having diabetes or not, learn to classify new patients as having diabetes or not.

Overall Description

