

# Ch. 2

## Statistical Learning

<https://www.youtube.com/watch?v=cyCJMVNhfnI&t=0s>

## ISLR 강의 by Hastie & Tibshirani

- <https://www.youtube.com/playlist?list=PLwyWqcwXkazRpdF7AUDb4GK3Nn8f9jUoK>

## 강의 by Hamed Hashemini

- <https://www.youtube.com/playlist?list=PL06ytJZ4Ak1rXmlvxTyAdOEfiVEzH00IK>

**보세요**

## 문제가 있나요?

### ■ 영어?

- 피할수 없어요
- 이 분야 대부분 자료, 지식이 영어로 되었어요
- 이런 자료를 이용 못하고 영어로 소통 못하면 3류에 그쳐요
- 다른 나라에 가 일하려면 필요해요

### ■ 이론은 필요없고 그냥 해법(cookbook) 만 알려줘요

- 이 정도의 이론도 모르는 사람이 한 분석을 나라면 "콩으로 메주를 쏜다 해도 안 믿어요"

### ■ 수학?

- 교재의 수학은 고등학교 1학년 수준을 넘지 않아요. 그나마도 개념만 알면 되요. 고1 수학책 봐요.

# Supervised Learning (지도 학습)

- 입력  $X$ 와 그에 따라 원하는 결과  $Y$ 가 훈련 데이터(Training Set)로 주어질 때,

$Y = f(X)$ :  $X, Y$  로 부터  $f$  를 잘~~ 추정함 이 목표

-  $(X, Y)$  는 관측된 데이터

- $X$  : 입력(input), 특징(feature/attribute), 독립변수(independent variable), predictor, covariate, 또는 그냥 “변수” 라고도 하며 보통 vector 형태
- $Y$  : 결과(response, output, outcome), 종속변수(dependent variable), target 등으로 불림
- 학습은  $X, Y$ 를 이용해 함수(알고리즘, 모델)  $f$  를 추정(estimate)하는 것

## Supervised Learning (지도학습) – cont.

- 학습을 통해 추정된  $\hat{f}$  는 완벽하지 않아  $\hat{f}(X)$ 가 실제로  $Y$ 와 똑같이 되기 어려움.  $\hat{f}(X)$ 를  $\hat{Y}$  (Y-hat) 으로 표기.  $\hat{Y} \rightarrow Y$  되도록 함이 학습의 목표
- $Y$ 가 숫자인 학습형태를 **regression** 이라 함
- $Y$ 가 카테고리 범주형이면 **classification** 이라 함
- $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$   $x_i$  : input,  $y_i$  : target response 쌍으로 구성된 training set (학습을 시키기 위한 훈련 자료)를 통해  $\hat{f}$ 를 구하는 과정이 학습
- Training set 학습으로 구한  $\hat{f}$ 에 테스트 데이터 (test set)를 적용하여  $\hat{f}$ 의 성능을 판단. 이를 Test Performance.
- Test Set: 모델을 훈련할 때 쓰지 않은 데이터(Training set에 없는 데이터)
- 개별  $(x_i, y_i)$ 를 instance, observation 또는 example이라 함

# Is it a Classification or a Regression?

- 모호할 때도 많다
- 많은 Classification이 사실은 내부적으로 Regression
  - 가령, target 클래스 레이블들이 확률분포...

## What are ;

Multiclass(multinomial) Classification

Multilabel Classification

Multioutput Regression

Multioutput Classification



- Is it a Dog or Cat : **Binary**

Ans : Dog

- Is it a Dog or Cat or Chicken : **Multiclass**

Ans : Dog



- Are there Dog, Cat, and Rabbit : **Multilabel**

Ans : Yes, Yes, No (wrong ans)

- **Multioutput** : for each output label  
do **Multiclass** operation

Ans : Dog, Cat, Rabbit

## 머신러닝을 나누는 또 다른 방법 - 훈련 데이터를 어떻게 활용할까?

- Batch Learning

- 갖고 있는 모든 training set 데이터를 한꺼번에 이용해 학습
- 대부분 offline 상태에서 학습
- takes time, and needs RAM to store the data

- Online (mini-batch) Learning

- 데이터를 조금씩 순차적으로 이용해 학습
- each learning step is fast & cheap and can be done on the fly
- 초기 학습은 보통 offline 상태에서 (즉, not on the live system)



## Parametric Model (Linear models, LDA, Naïve Bayes, 단순한 신경망)

- $f$ (모델 형태)를 가정하고, 모델의 특성을 규정하는 parameter (매개변수,  $\beta_i$ ) 들을 학습을 통해 추정하는 형태

$$Y = f(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

- 위의 경우, 모델 원형  $f$ 를 feature들  $X_1, \dots, X_p$ 로 된 **선형식**으로 가정
- linear parametric 모델은 parameter  $\beta_i$  ( $i=0, \dots, p$ )로 규정됨
- 위의 linear parametric 모델에서  $\hat{f}$ 를 추정함은  $\beta_0, \dots, \beta_p$ 들의 추정치  $\hat{\beta}_0, \dots, \hat{\beta}_p$ 들을 구하는 것
- 따라서, parametric model의 학습은 training set를 이용해  $\hat{\beta}_i$ 들을 잘 구해  $\hat{f}(X) = \hat{Y}$ 의 값이 실제  $Y$ 에 최대한 가깝도록 함이 목표

## Non-Parametric Model

- $f$ 의 형태(구조)가 어떨 것이라는 것을 미리 가정하지 않음
- (More) Data Driven
- Training set의 크기가 작을 때 문제가 될 수 있음
- KNN, Decision Tree, Boosting/Random Forest, Support Vector Machine, 복잡한 신경망

## Remember,

All the models are wrong!

우리 모델은 그 무슨 수학 원리를, 우주 법칙을 말하는 것이 아니고

Most models perform lousy if you got insufficient data

\* 모델의 강력함(flexibility/capacity)에 따른 prediction - inference, 일반화(generalize)의 중요함, training 성능 - test 성능, bias - variance, underfit - overfit, training 시간 등을 2장에서 얘기합니다. 머신러닝의 중요한 개념이라 정확히 알아야 합니다. 하지만, 이 개념을 **교조적으로 이해하면 안됩니다.**

## Vector-Matrix 방정식 복습 – ISLR 2.1.2 에서의 표현식

$n$  : 입력 data points (instance)들의 개수. 이것들을 이용해 모델  $f$  를 훈련

$X$  : 위 입력 데이터들로 이루어진 predictor 집합.

$x_i$  :  $X$  내의 개별 입력 데이터 포인트 중  $i$ 번째 data point 값. ( $i=1, \dots, n$ )

$x_{ij}$  :  $i$ 번째 입력 데이터 포인트의  $j$ 번째 predictor 값.  $X$ 의 data point가 여러 개의 feature 들로 되어 있을 때 (보통  $X$ 는 여러 feature들로 됨). 이 책에서는 feature들의 개수를  $p$ 라 일반적으로 놓음 (즉,  $j = 1, \dots, p$ )

\* ISLR에서  $i$ 번째 data point(instance)를 가르킬때는  $x_i$ , ( $i=1, \dots, n$ )

$i$ 번째 feature set 를 가르킬때는  $X_j$  로 표시함, ( $j=1, \dots, p$ )

\* 어떤 책에서는  $i$ 번째 data point(instance)를 가르킬때  $x^{(i)}$  으로 표기. 이 경우 3번째 instance의 8번째 feature를  $x^{(3)}_8$  로 표현

$Y$  :  $X$ 에 대응되는 response들의 집합

$y_i$  :  $Y$  내의 개별 입력 데이터 포인트 중  $i$ 번째 data point 값. ( $i=1, \dots, n$ ). 즉, 입력이  $x_i$  일 때의 response.

$y_{ij}$  :  $i$ 번째 response data point의  $j$ 번째 response값 (Response가 행렬로 표현). 이 경우는 개별 response data point가 여러 개의 response 값들로 이루어진 vector 형태일 때 적용된다 (Multioutput). **ISLR에서는 이 경우를 다루지 않음**

## Note

ISLR (Introduction to Statistical Learning with R) 과  
ESL (The Elements of Statistical Learning) 책에서

$$\hat{Y} = X^T \hat{\beta}, \quad (2.2)$$

$$\hat{y} = \mathbf{X} \hat{\beta}$$

두 형태의 표현식이 보인다. 이 경우 위 식의  $X^T$ 는  $X^T = (X_1, X_2, \dots, X_p)$ 로써 p개의 predictor로 구성된 하나의 observation을 뜻하고 row vector임을 나타냄. 이  $X^T$ 가 만약 i번째 observation이라면 ISLR/ESL은  $x_i^T$ 로도 표현.

- **Exception** : 어떤 경우  $X_j$ 가 n개 observation들의 j번째 predictor들로 표현되기도 한다. 이 경우  $X_j$ 는 vector (of length n)가 되겠고  $\mathbf{X}$ 는 matrix가 됨. Predictor  $X_j$ 가 scalar인지 또는 n observation 개수 만큼의 vector\_of\_length n 인지는 문맥을 통해 알 수 있다. 또 경우에 따라서 대문자  $X$ 는 vector를, 소문자  $x$ 는 scalar를 나타내기도 함.

반면 아랫 식의 bold  $\mathbf{X}$ 는 matrix(행렬)를 나타냄. 만약 predictor data set이 n개의 observation, 각 observation이 p개의 predictor들로 되어있으면  $\mathbf{X}$ 는 "n x p" matrix.

기본적으로 ISLR/ESL은 벡터가 Column Vector라 가정. 따라서 어떤 벡터  $x$ 가 row vector일 시 그것을  $x^T$ 로 표시.

## Linear System의 Vector-Matrix 방정식 의미:

$$f(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p. \quad (2.4) \quad \text{Linear model로 가정하고,}$$

↓ fit to given (X, Y) training set

$$Y \approx \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p. \quad \text{ie. if } Y \text{ can be approximated as a linear combination of } X_1, X_2, \dots, X_p$$

↓ expanding...

$$y_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \beta_3 x_{13} + \dots + \beta_p x_{1p}$$

$$y_2 = \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \beta_3 x_{23} + \dots + \beta_p x_{2p}$$

$$y_3 = \beta_0 + \beta_1 x_{31} + \beta_2 x_{32} + \beta_3 x_{33} + \dots + \beta_p x_{3p}$$

⋮

$$y_n = \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \beta_3 x_{n3} + \dots + \beta_p x_{np}$$

\*  $\beta_i$  : 우리가 구하려는 패러미터



$$\begin{pmatrix} Y \\ y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{pmatrix} = \beta_0 + \beta_1 \begin{pmatrix} X_1 \\ x_{11} \\ x_{21} \\ x_{31} \\ \vdots \\ x_{n1} \end{pmatrix} + \beta_2 \begin{pmatrix} X_2 \\ x_{12} \\ x_{22} \\ x_{32} \\ \vdots \\ x_{n2} \end{pmatrix} + \beta_3 \begin{pmatrix} X_3 \\ x_{13} \\ x_{23} \\ x_{33} \\ \vdots \\ x_{n3} \end{pmatrix} + \dots + \beta_p \begin{pmatrix} X_p \\ x_{1p} \\ x_{2p} \\ x_{3p} \\ \vdots \\ x_{np} \end{pmatrix}$$

$$\begin{matrix} Y \\ Y - \beta_0 \end{matrix} = \begin{pmatrix} \begin{matrix} X \\ x_{11} & x_{12} & x_{13} & \dots & x_{1p} \\ x_{21} & x_{22} & x_{23} & \dots & x_{2p} \\ x_{31} & x_{32} & x_{33} & \dots & x_{3p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & x_{n3} & \dots & x_{np} \end{matrix} \\ \begin{matrix} \beta \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \vdots \\ \beta_p \end{matrix} \end{pmatrix}$$

①  $(X, Y)$ 는 Training set 데이터

② RSS를 최소화하는  $\hat{\beta}$ 들을 구하자  
(linear regression의 parameter를 구하자)

③ 그리고는 새로운  $X$ 가 주어지면  $\hat{y} = X\hat{\beta}$ 하여  $\hat{Y}$  추정