

Machine Learning

Chapter 1 머신러닝 개요(intro)

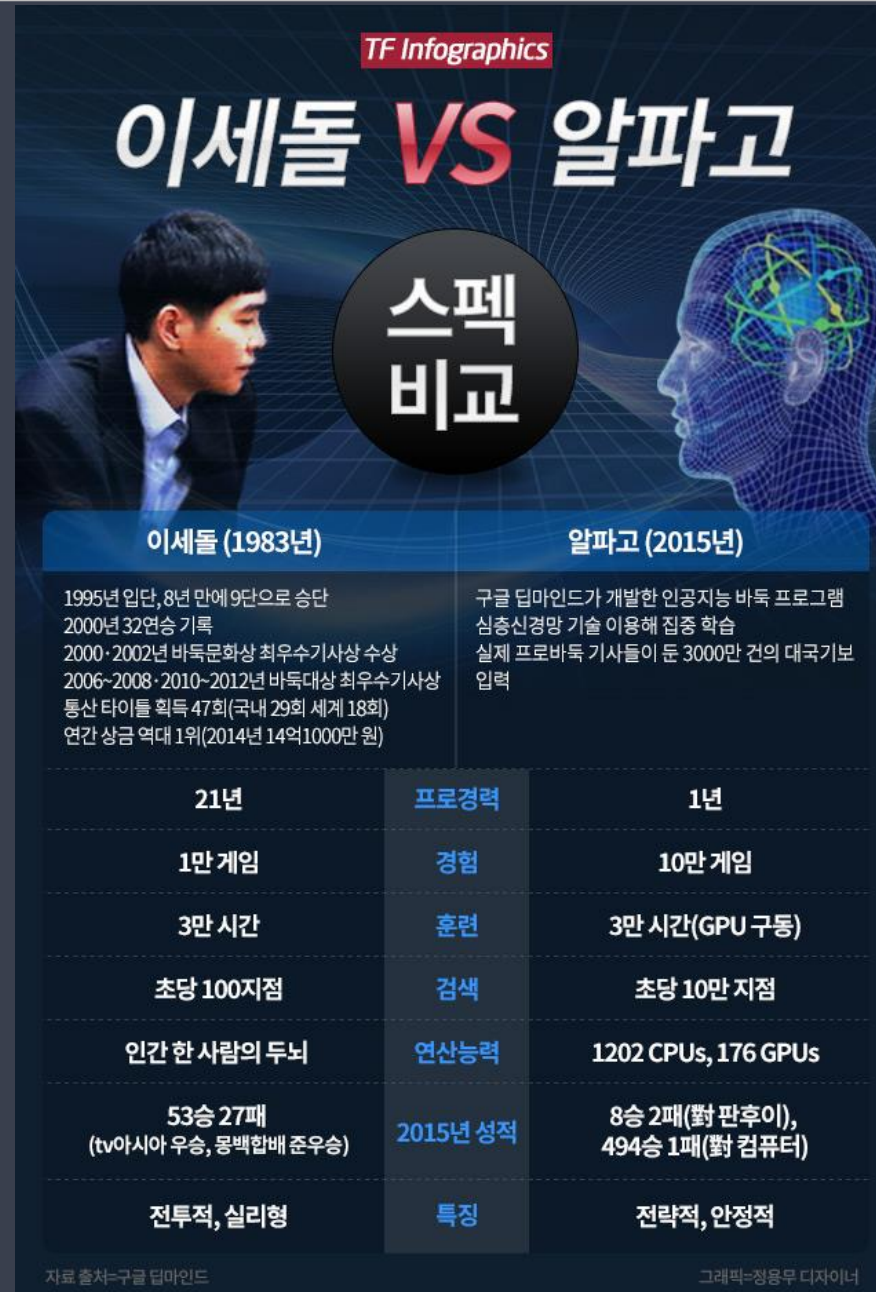


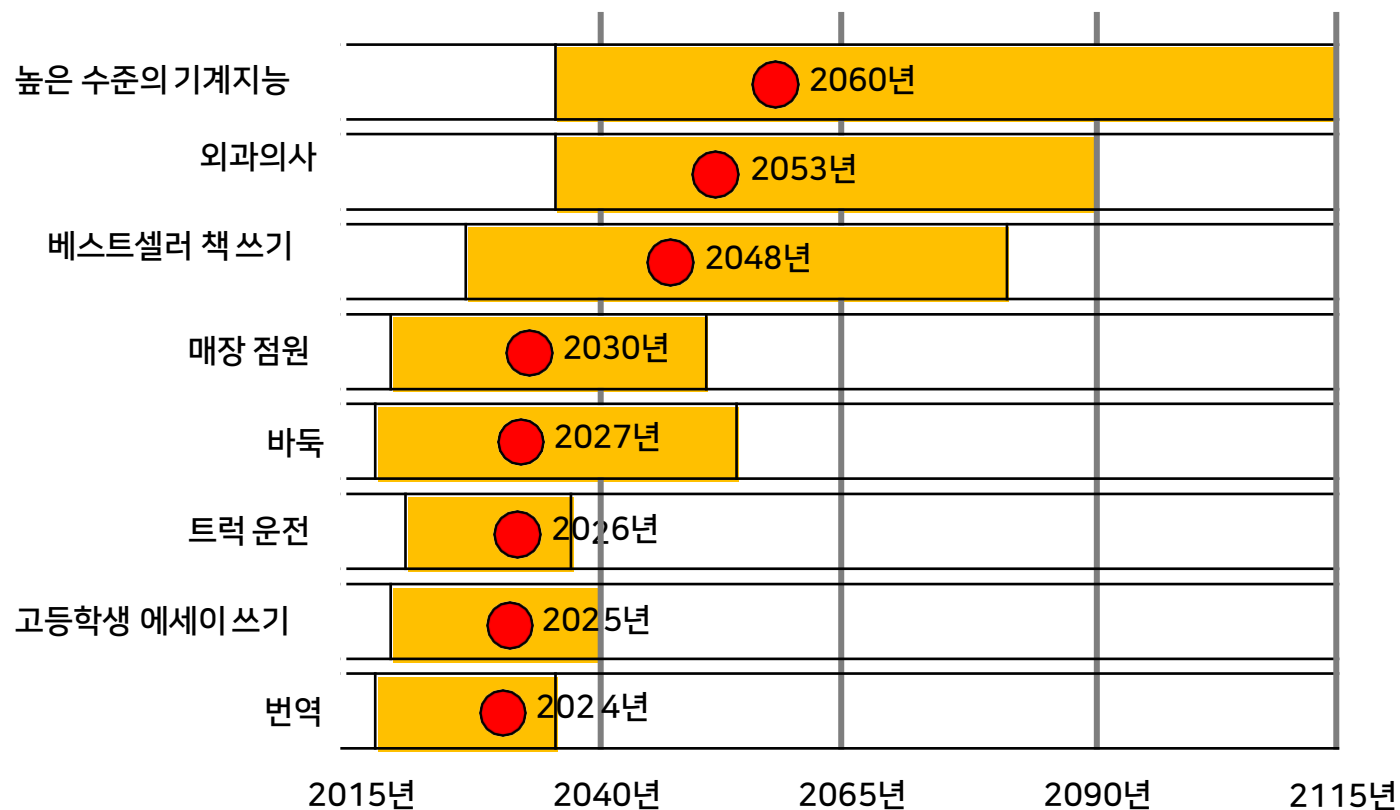
START

- Machine Learning 개념을 이해 할 수 있다.
- Machine Learning의 종류 및 과정을 알 수 있다.
- 기계학습과 관련된 기본 용어를 알 수 있다.



머신러닝이란?

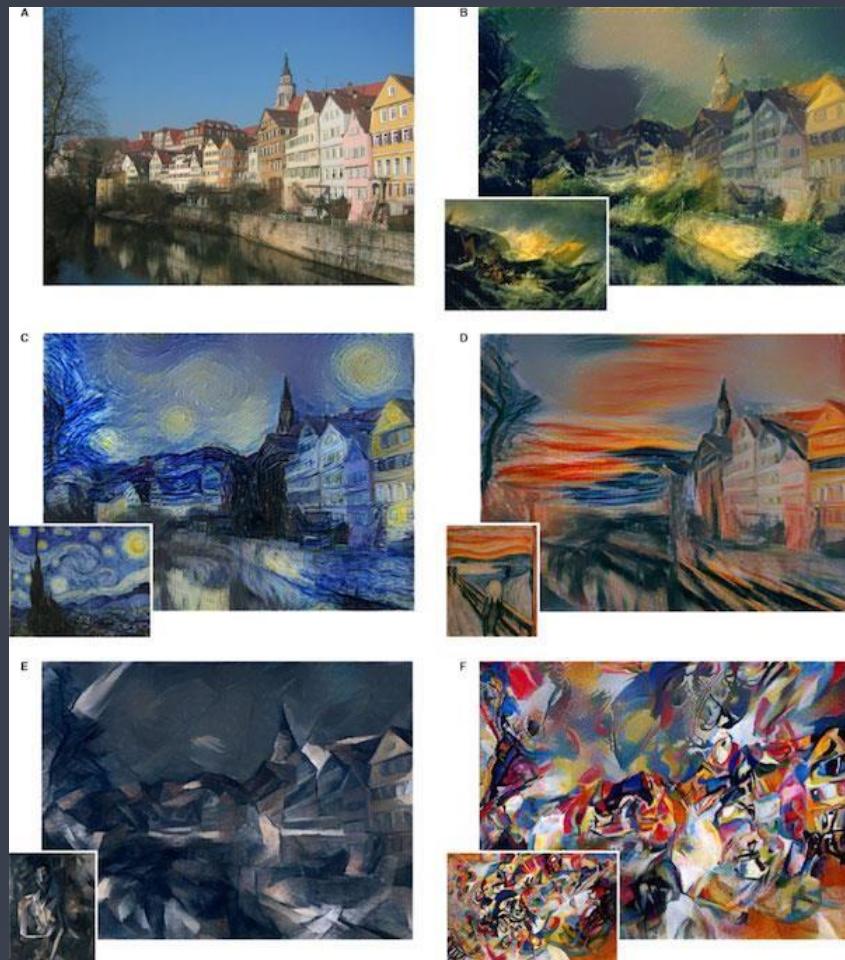




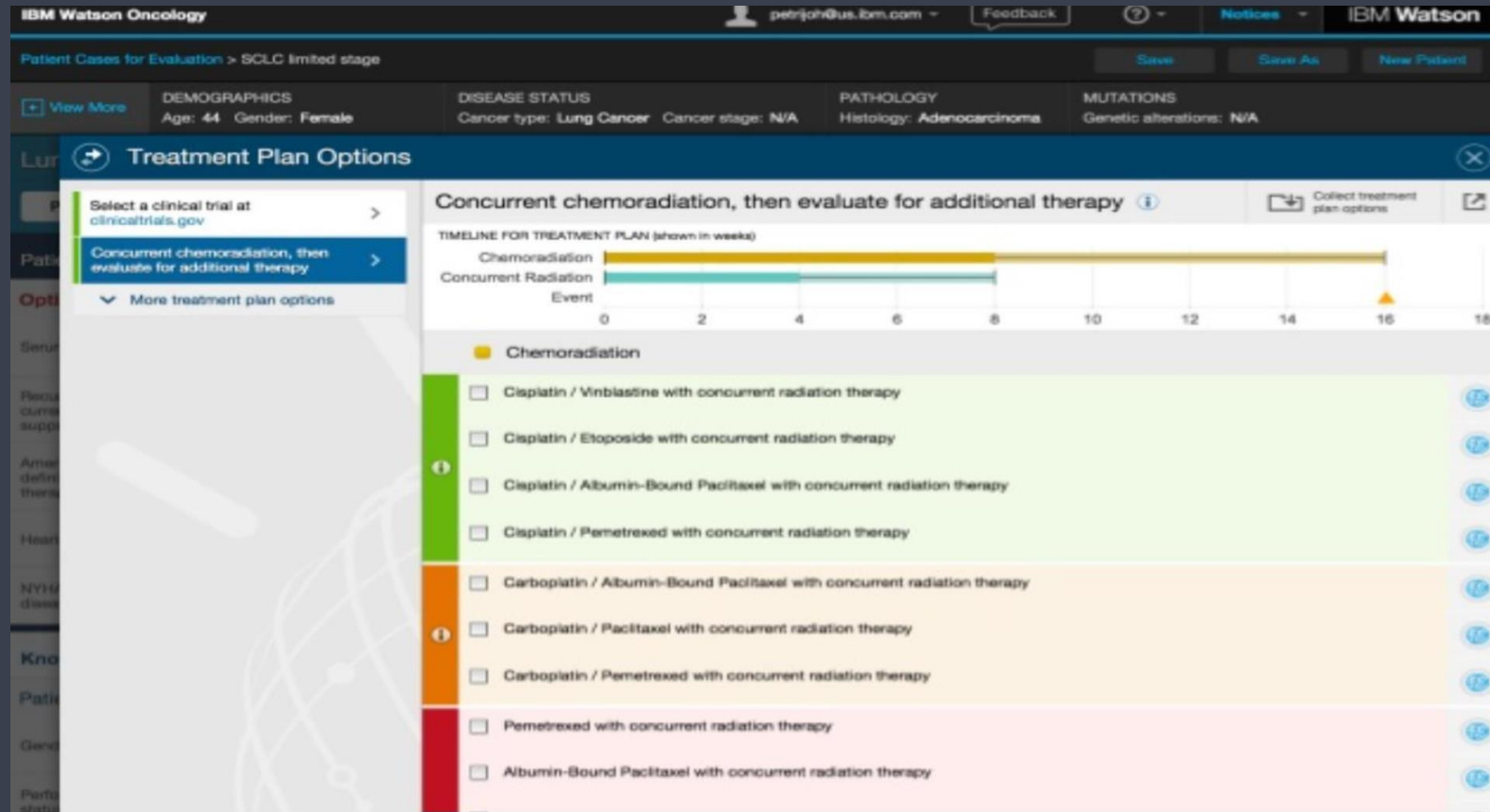
- AI 전문가를 대상으로 AI의 성능이 언제 사람을 능가하는가 (352명)



인공지능이 인간을 대체할 수 있는가 ?

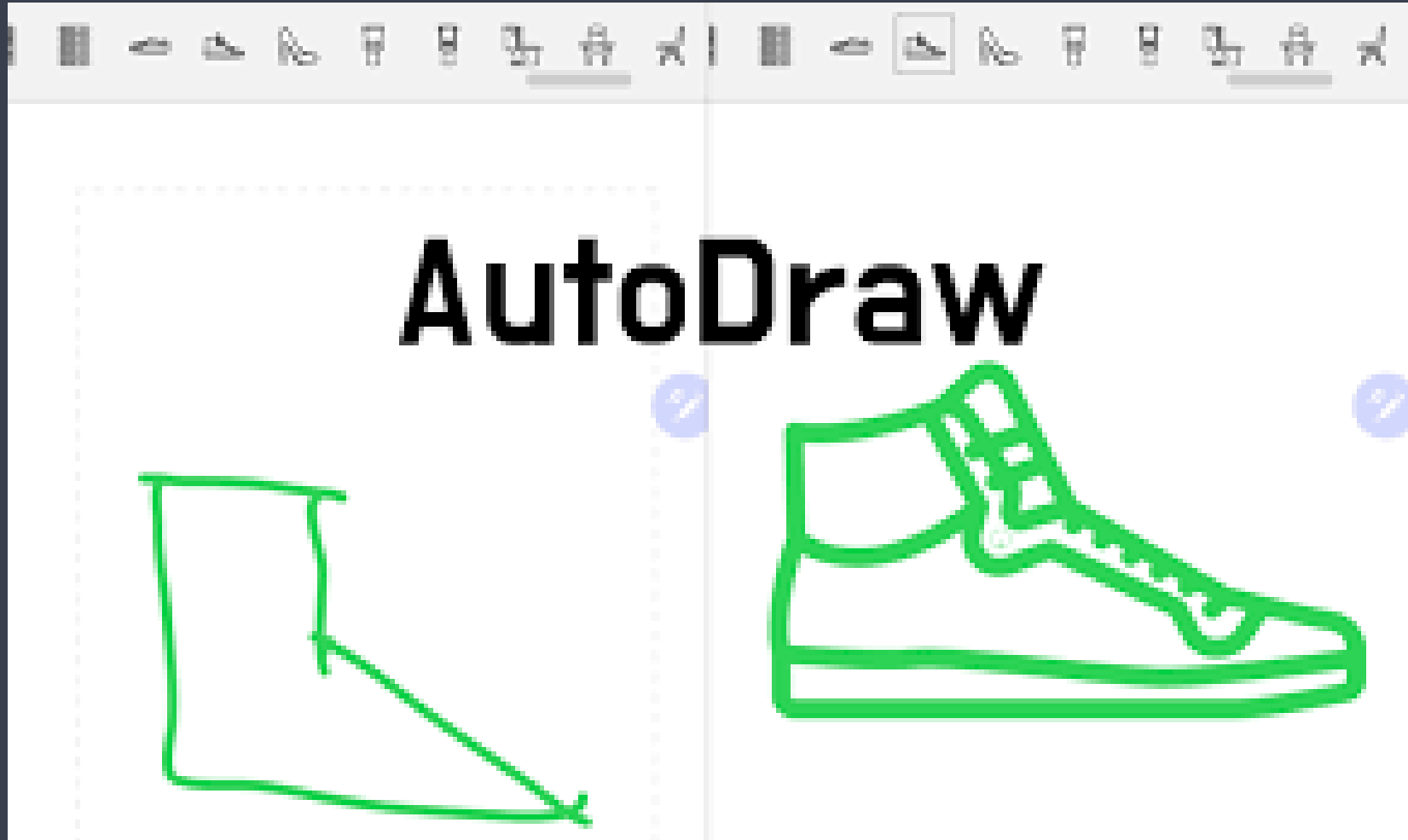


복잡한 의료 데이터 분석 및 insight 도출 (Expert System) - IBM Watson



영상 의료/병리 데이터의 분석 및 판독 (Deep Learning) - 영상의학과 전문의





Artificial Intelligence

인공지능

사고나 학습 등 인간이 가진
지적 능력을 컴퓨터를 통해
구현하는 기술



Machine Learning

머신러닝

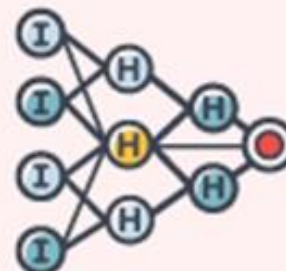
컴퓨터가 스스로 학습하여
인공지능의 성능을
향상 시키는 기술 방법



Deep Learning

딥러닝

인간의 뉴런과 비슷한
인공신경망 방식으로
정보를 처리



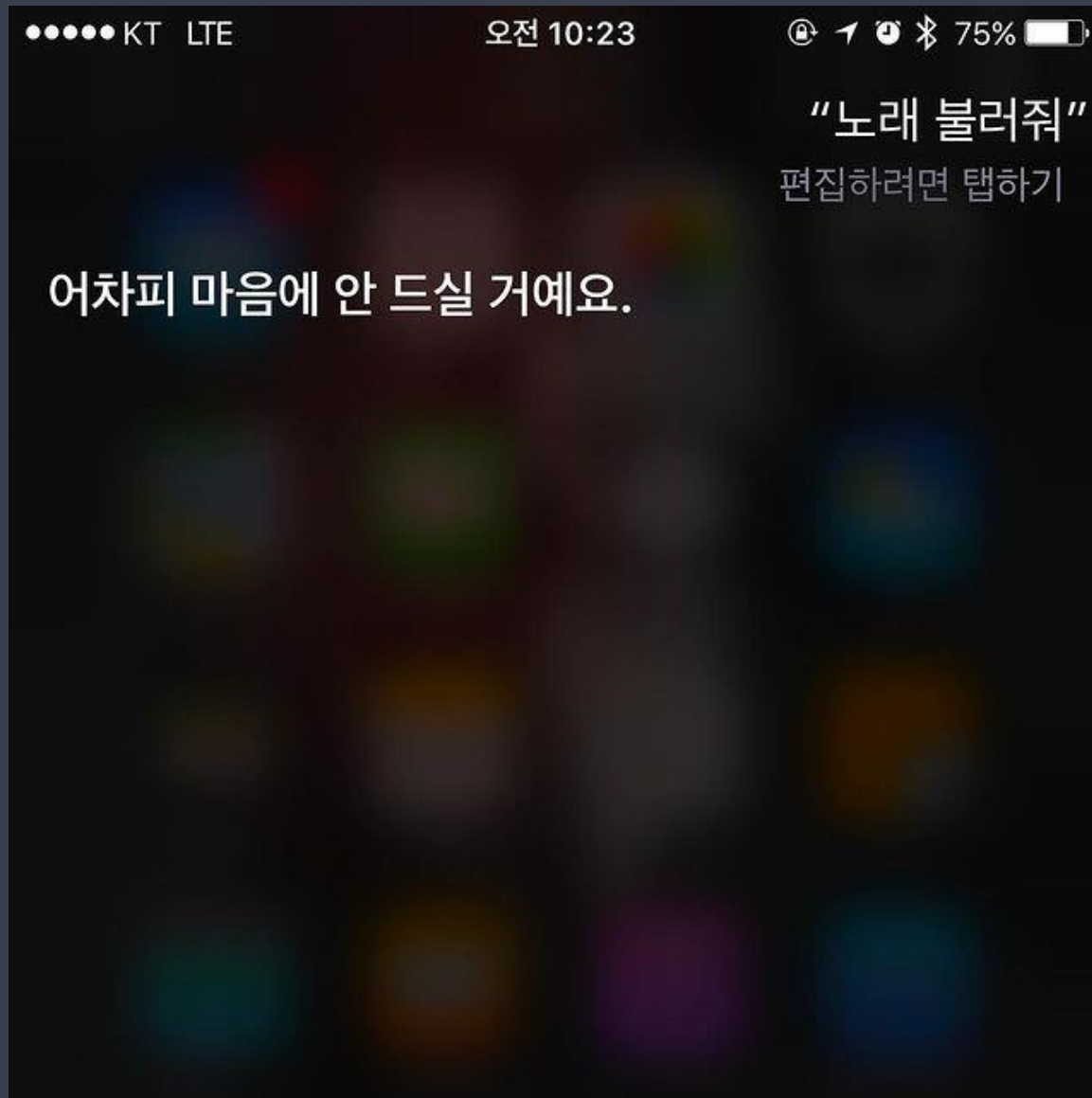
머신러닝(Machine Learning)이란?

- 데이터를 기반으로 학습을 시켜서 예측하게 만드는 기법
- 인공지능의 한 분야로 컴퓨터가 학습할 수 있도록 하는 알고리즘과 기술을 개발하는 분야
- 통계학, 데이터 마이닝, 컴퓨터 과학이 어우러진 분야

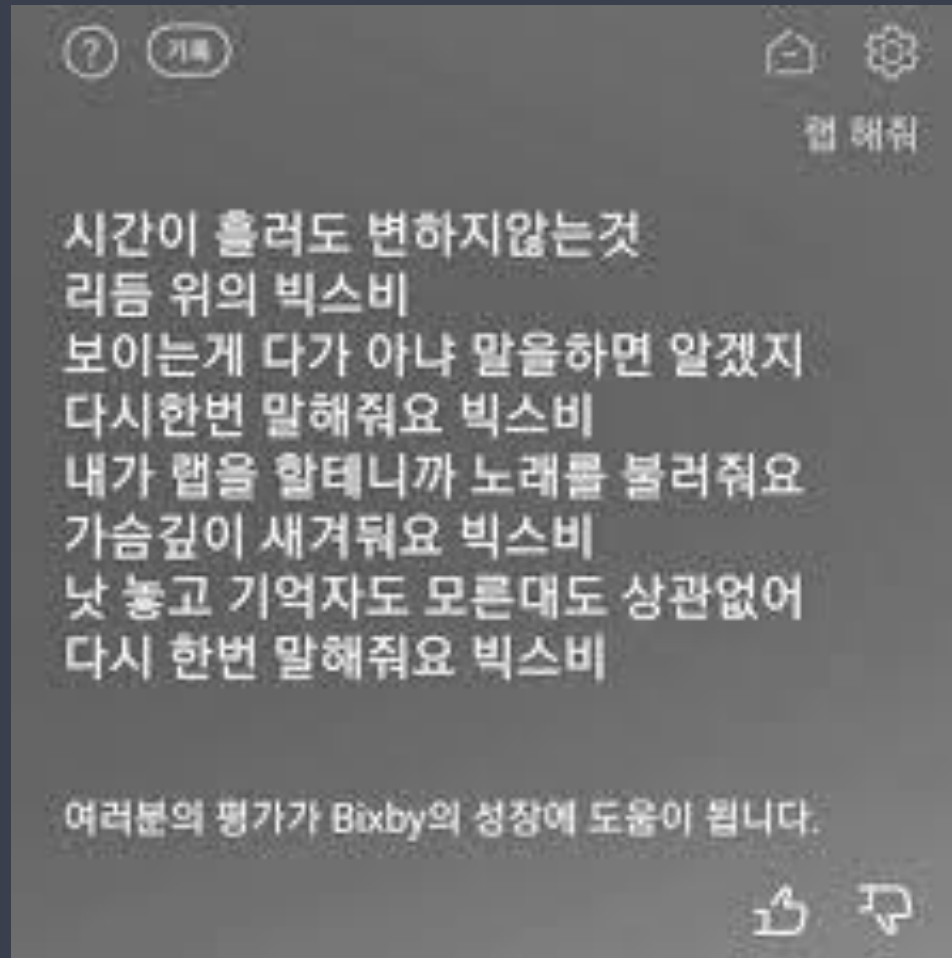
Rule-based expert system (규칙 기반 전문가 시스템)

“if”와 “else”로 하드 코딩된 명령을 사용하는 시스템

Rule-based expert system



Rule-based expert system



Rule-based expert system 문제점

- 스팸 메일 필터
- 얼굴 인식 시스템

많은 상황에 대한 규칙들을 모두 만들어 낼 수 없다

- 제작한 로직이 특정 작업에만 국한된다. 작업이 조금만 변경되더라도 전체 시스템을 다시 만들어야 할 수 있다.
- 규칙을 설계하려면 해당 분야에 대해서 잘 알고 있어야 한다.

Data → Model
(알고리즘)

학습을 통해 기계가 스스로 규칙을 만들어낸다.

머신러닝(Machine Learning)

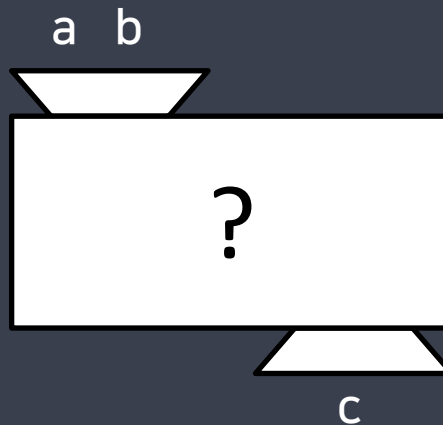
Data



Model
(알고리즘)

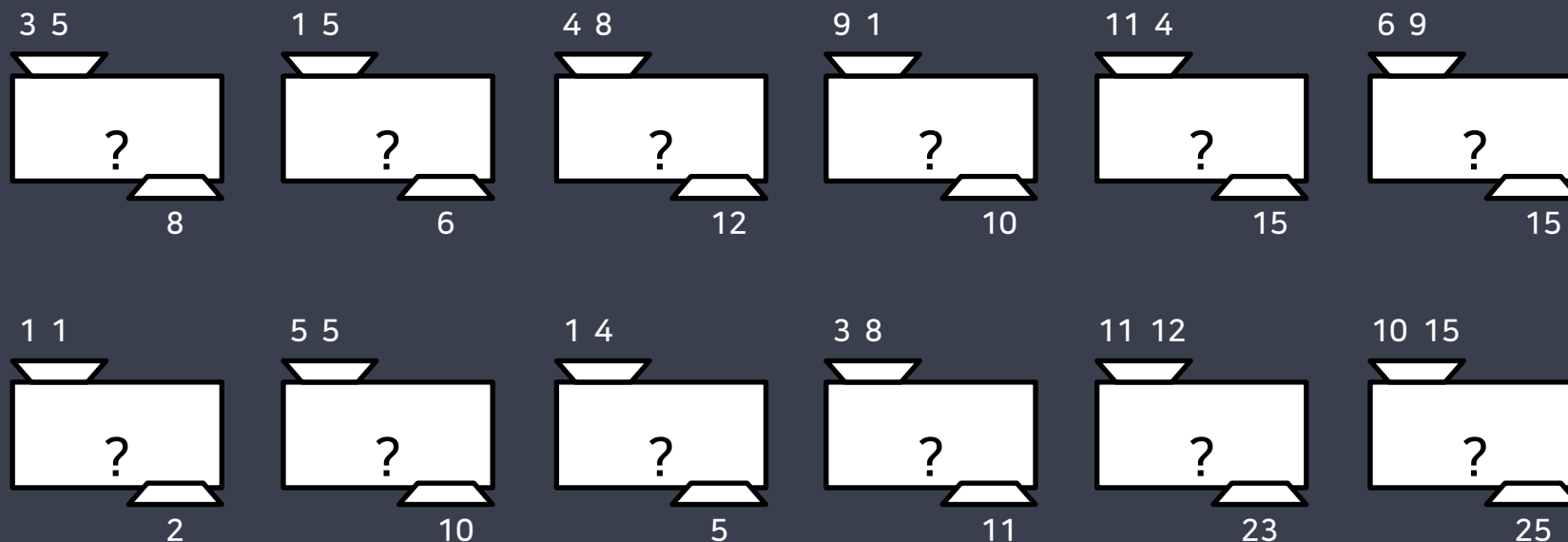
데이터를 이용하여 **특성과 패턴을 학습**하고,
그 결과를 바탕으로 미지의 데이터에 대한 **미래결과**
(값, 분포)를 예측하는 것

요구 사항 분석 단계



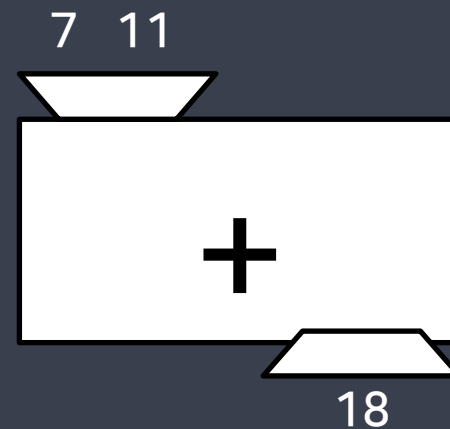
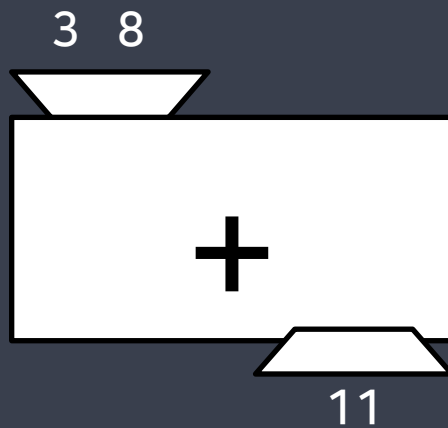
a와 b를 입력하면 안에서 무슨 일이 일어나서 c가 나올까요 ?

학습단계

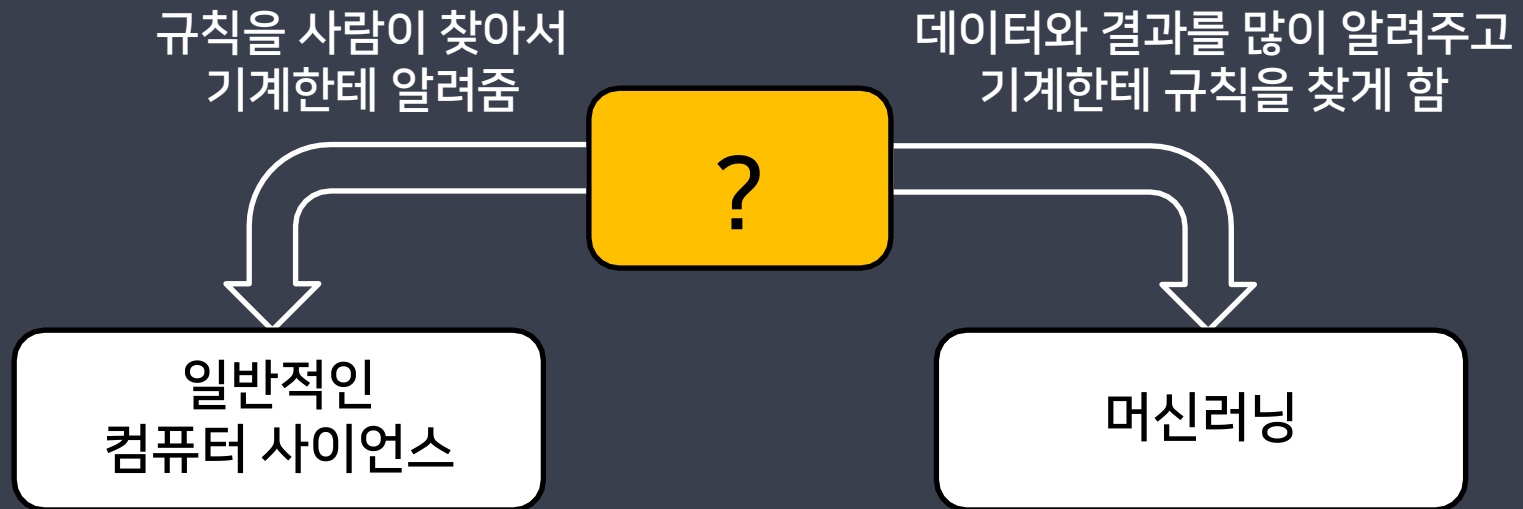


?는 무엇일까요 ?

예측단계



머신러닝(Machine Learning)





scikit-learn

- 파이썬에서 쉽게 사용할 수 있는 머신러닝 프레임워크, 라이브러리
- 회귀, 분류, 군집, 차원 축소, 특성 공학, 전처리, 교차 검증, 파이프라인 등 머신러닝에 필요한 기능을 갖추
- 학습을 위한 샘플 데이터도 제공

미스터리 사인 학습하기

문제[Data]

N1	N2	Result
1	1	1
1	2	3
1	3	5
1	5	9

답[Label]

미스터리 사인 학습하기

Model
(알고리즘)

RandomForest 회귀 모델을 사용

미스터리 사인 학습하기

```
random_forest = RandomForestRegressor()
```

```
random_forest.fit(문제 , 답)
```

```
random_forest.predict(값을 얻고 싶은데이터 )
```

미스터리 사인을 학습 해보자

지도학습 (Supervised Learning)

비지도학습 (Unsupervised Learning)

강화학습 (Reinforcement Learning)

지도 학습 (Supervised Learning)

- 데이터에 대한 Label(명시적인 답)이 주어진 상태에서 컴퓨터를 학습시키는 방법.
- 사람이 직접 개입하기 때문에 정확도가 높은 데이터 사용가능
- 사람이 직접 Label을 달아야 하므로 구할 수 있는 데이터의 한계가 있음.
- 분류(Classification)와 회귀(Regression)로 나뉘어진다.

Kaggle Titanic 데이터

Feature

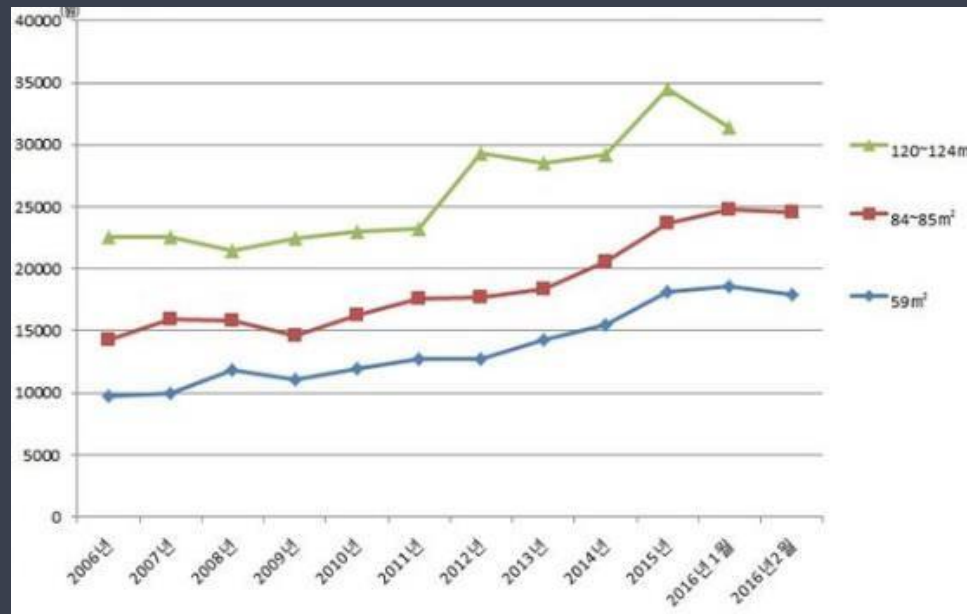
Class
Label

Survived	Pclass	Name	Sex	Age	SibSp
0	3	Braund, M	male	22	1
1	1	Cumings, M	female	38	1
1	3	Heikkinen, female		26	0
1	1	Futrelle, M	female	35	1
0	3	Allen, Mr.	male	35	0
0	3	Moran, Mr	male		0
0	1	McCarthy, male		54	0
0	3	Palsson, M	male	2	3
1	3	Johnson, M	female	27	0
1	2	Nasser, Mr	female	14	1
1	3	Sandstrom	female	4	1

지도 학습 (Supervised Learning)



스팸 메일 분류



집 가격 예측

분류 (Classification)

- 미리 정의된 여러 클래스 레이블 중 하나를 예측하는 것.
- 속성 값을 입력, 클래스 값을 출력으로 하는 모델
- 붓꽃(iris)의 세 품종 중 하나로 분류, 암 분류 등.
- 이진분류, 다중 분류 등이 있다.

회귀 (Regression)

- 연속적인 숫자를 예측하는 것.
- 속성 값을 입력, 연속적인 실수 값을 출력으로 하는 모델
- 어떤 사람의 교육수준, 나이, 주거지를 바탕으로 연간 소득 예측.
- 예측 값의 미묘한 차이가 크게 중요하지 않다.

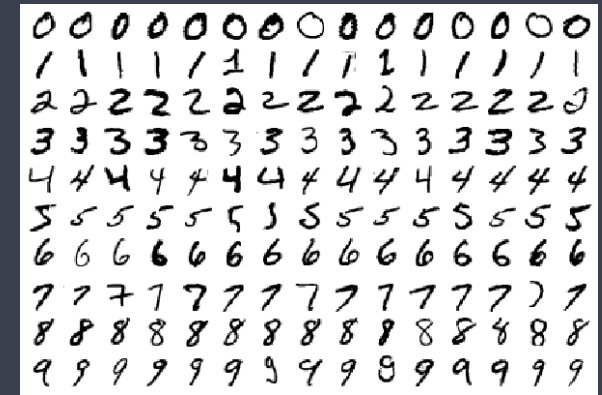
비지도 학습 (Unsupervised Learning)

- 데이터에 대한 Label(명시적인 답)이 없는 상태에서 컴퓨터를 학습시키는 방법.
- 데이터의 숨겨진 특징, 구조, 패턴을 파악하는데 사용.
- 데이터를 비슷한 특성끼리 묶는 클러스터링(Clustering)과 차원축소(Dimensionality Reduction)등이 있다.

비지도 학습 (Unsupervised Learning)

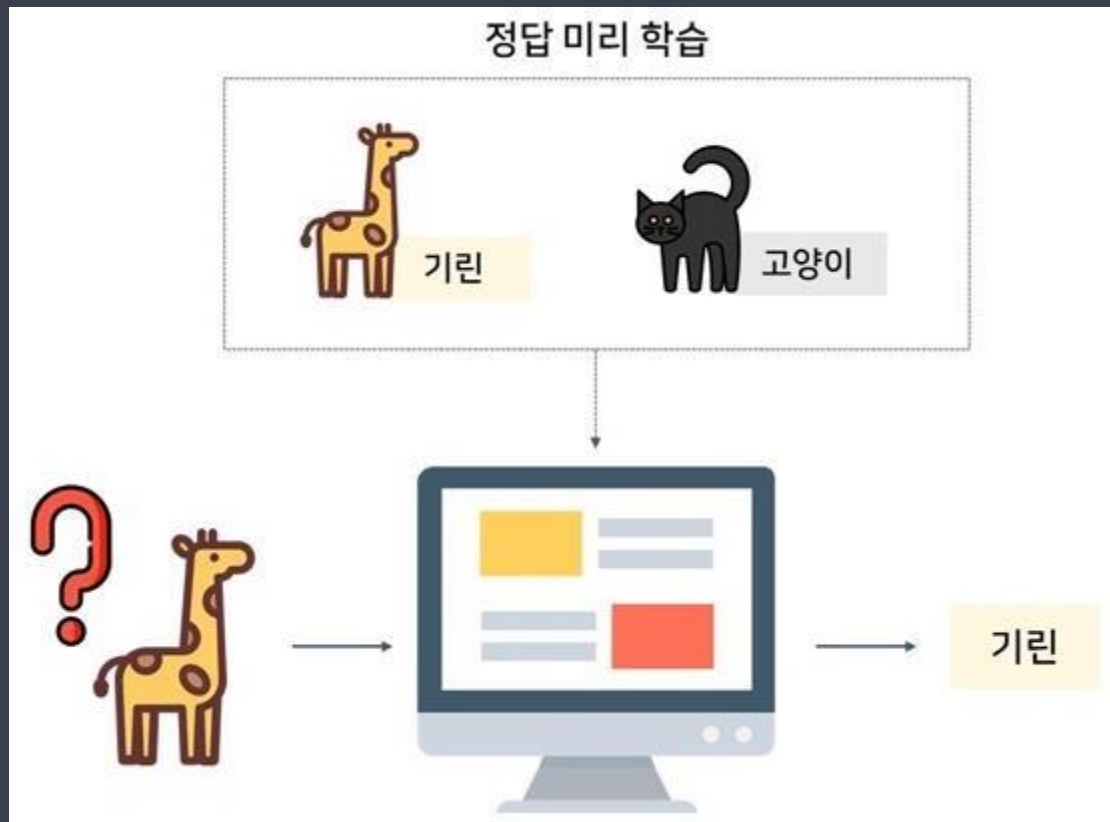


이미지 감색 처리

소비자 그룹 발견을
통한 마케팅

손 글씨 숫자 인식

지도 학습 (Supervised Learning)



비지도 학습 (Unsupervised Learning)



강화 학습 (Reinforcement Learning)

- 지도학습과 비슷하지만 완전한 답(Label)을 제공하지 않는 특징이 있다.
- 기계는 더 많은 보상을 얻을 수 있는 방향으로 행동을 학습
- 주로 게임이나 로봇을 학습시키는데 많이 사용

머신러닝(Machine Learning)이 유용한 분야

- 기존 솔루션으로는 많은 수동 조정과 규칙이 필요한 문제
- 전통적인 방식으로는 전혀 해결 방법이 없는 복잡한 문제
- 새로운 데이터에 적응해야하는 유동적인 환경
- 대량의 데이터에서 통찰을 얻어야 하는 문제

1. Problem Identification(문제정의)
2. Data Collect(데이터 수집)
3. Data Preprocessing(데이터 전처리)
4. EDA(탐색적 데이터분석)
5. Model 선택, Hyper Parameter 조정
6. Training(학습)
7. Evaluation(평가)

1. Problem Identification(문제정의)

- 비즈니스 목적 정의
모델을 어떻게 사용해 이익을 얻을까?
- 현재 솔루션의 구성 파악
- 지도 vs 비지도 vs 강화
- 분류 vs 회귀

2. Data Collect(데이터 수집)

- File (CSV, XML, JSON)
- Database
- Web Crawler (뉴스, SNS, 블로그)
- IoT 센서를 통한 수집
- Survey

3. Data Preprocessing(데이터 전처리)

- 결측치, 이상치 처리
- Feature Engineering (특성공학)
 - Scaling (단위 변환),
 - Transform (새로운 속성 추출),
 - Encoding (범주형 -> 수치형),
 - Binning (수치형 -> 범주형)

4. EDA(탐색적 데이터분석)

- 기술통계, 변수간 상관관계
- 시각화
`pandas, matplotlib, seaborn`
- Feature Selection (사용할 특성 선택)

5. Model 선택, Hyper Parameter 조정

- 목적에 맞는 적절한 모델 선택
- KNN, SVM, Linear Regression, Ridge, Lasso, Decision Tree, Random forest, CNN, RNN ...
- Hyper Parameter
model의 성능을 개선하기위해 사람이 직접 넣는 parameter

6. Model Training(학습)

- `model.fit(X_train,y_train)`
train 데이터와 test 데이터를 7:3 정도로 나눔
- `model.predict (X_test)`

머신러닝(Machine Learning) 과정



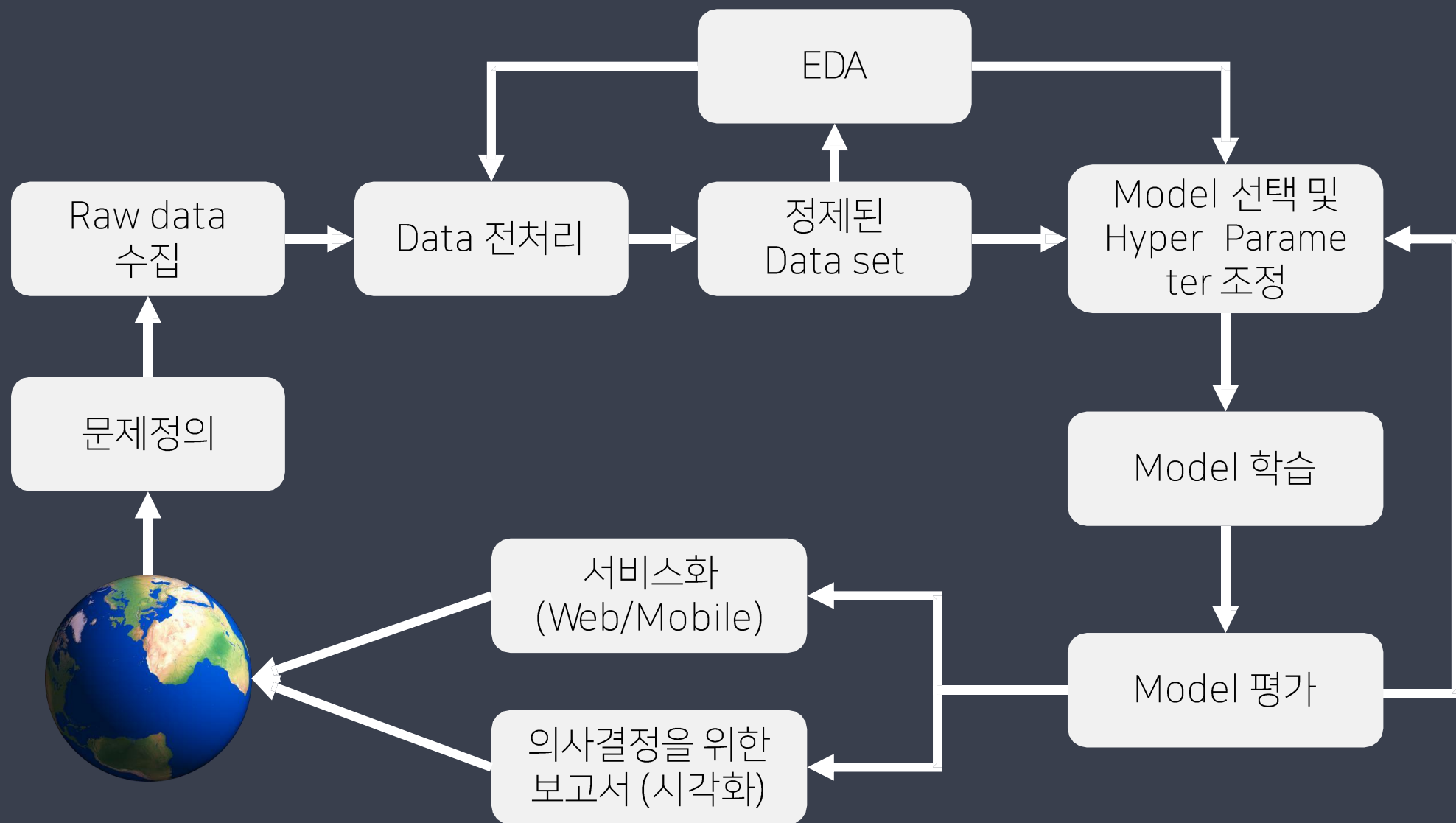
7. Evaluation(평가)

- accuracy(정확도)
- recall(재현율)
- precision(정밀도)
- f1 score

환자 100명이 있다고 가정했을 때 그 중에 5명이 암이다.

만약 100명 전부 암 환자가 아니라고 예측한다면
Accuracy(정확도)는 95%

머신러닝(Machine Learning) 과정



다음 학기 성적 점수를 예측하려면
어떤 특성이 필요할까?

직전 학기 성적, 알바 진행 여부, 연애,
잠자는 시간, 출석률, 학교와 집 사이 거리

비만도 데이터를 이용해 학습 해보자.

	A	B	C	D
1	Gender	Height	Weight	Label
2	Male	174	96	Obesity
3	Male	189	87	Normal
4	Female	185	110	Obesity
5	Female	195	104	Overweight
6	Male	149	61	Overweight
7	Male	189	104	Overweight
8	Male	147	92	Extreme Obesity
9	Male	154	111	Extreme Obesity
10	Male	174	90	Overweight
11	Female	169	103	Obesity
12	Male	195	81	Normal
13	Female	159	80	Obesity
14	Female	192	101	Overweight
15	Male	155	51	Normal
16	Male	191	79	Normal
17	Female	153	107	Extreme Obesity
18	Female	157	110	Extreme Obesity
19	Male	140	129	Extreme Obesity
20	Male	144	145	Extreme Obesity

회귀 문제?
분류 문제?

시각화를 통해 분류가 가능한
문제인지 확인해보자.