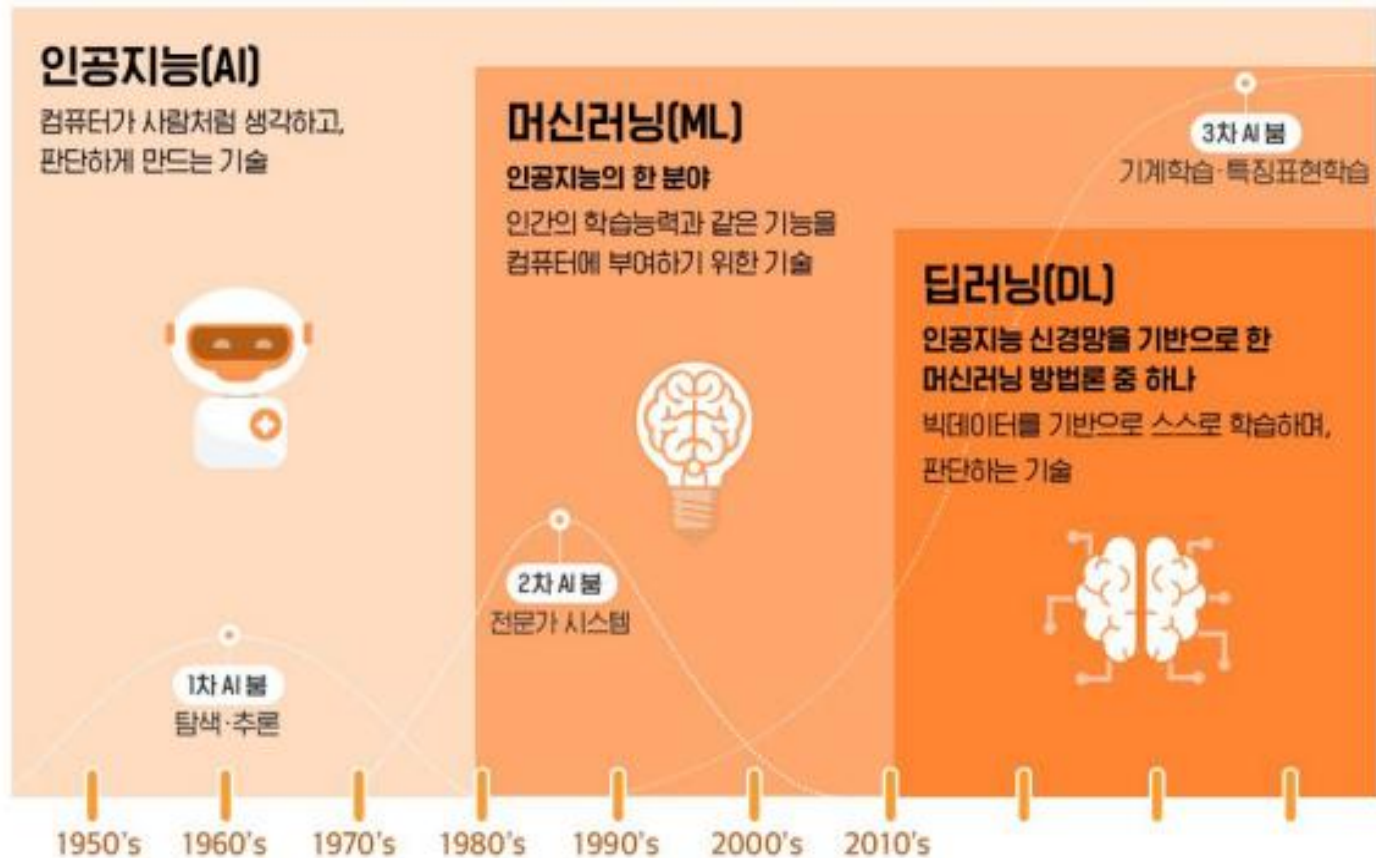


머신 러닝(Machine Learning)



1. 머신러닝

1959년, 아서 사무엘은 기계 학습을 "기계가 일일이 코드로 명시하지 않은 동작을 데이터로부터 학습하여 실행할 수 있도록 하는 알고리즘을 개발하는 연구 분야"라고 정의하였다. 기계 학습의 핵심은 표현 (representation) 과 일반화 (generalization) 에 있다. 표현이란 데이터의 평가이며, 일반화란 아직 알 수 없는 데이터에 대한 처리이다.

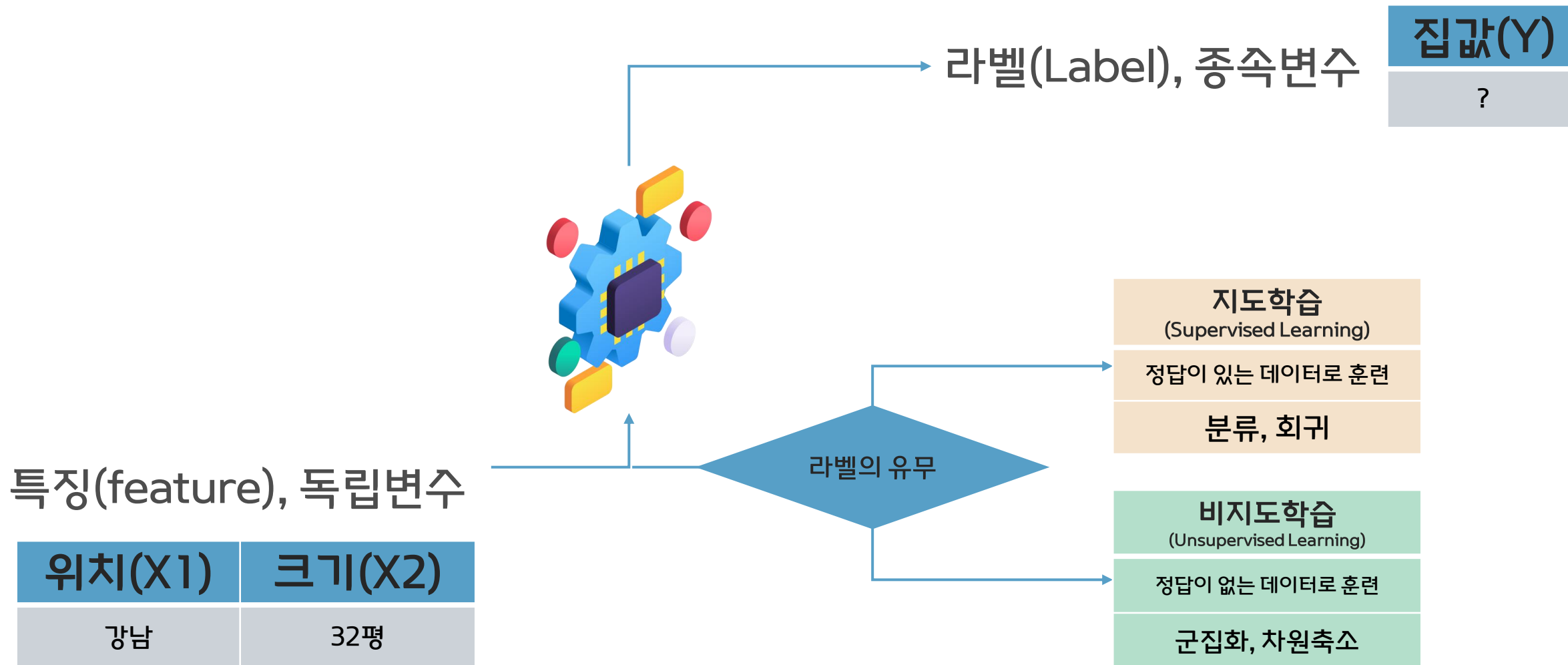


2. 지도학습과 비지도학습

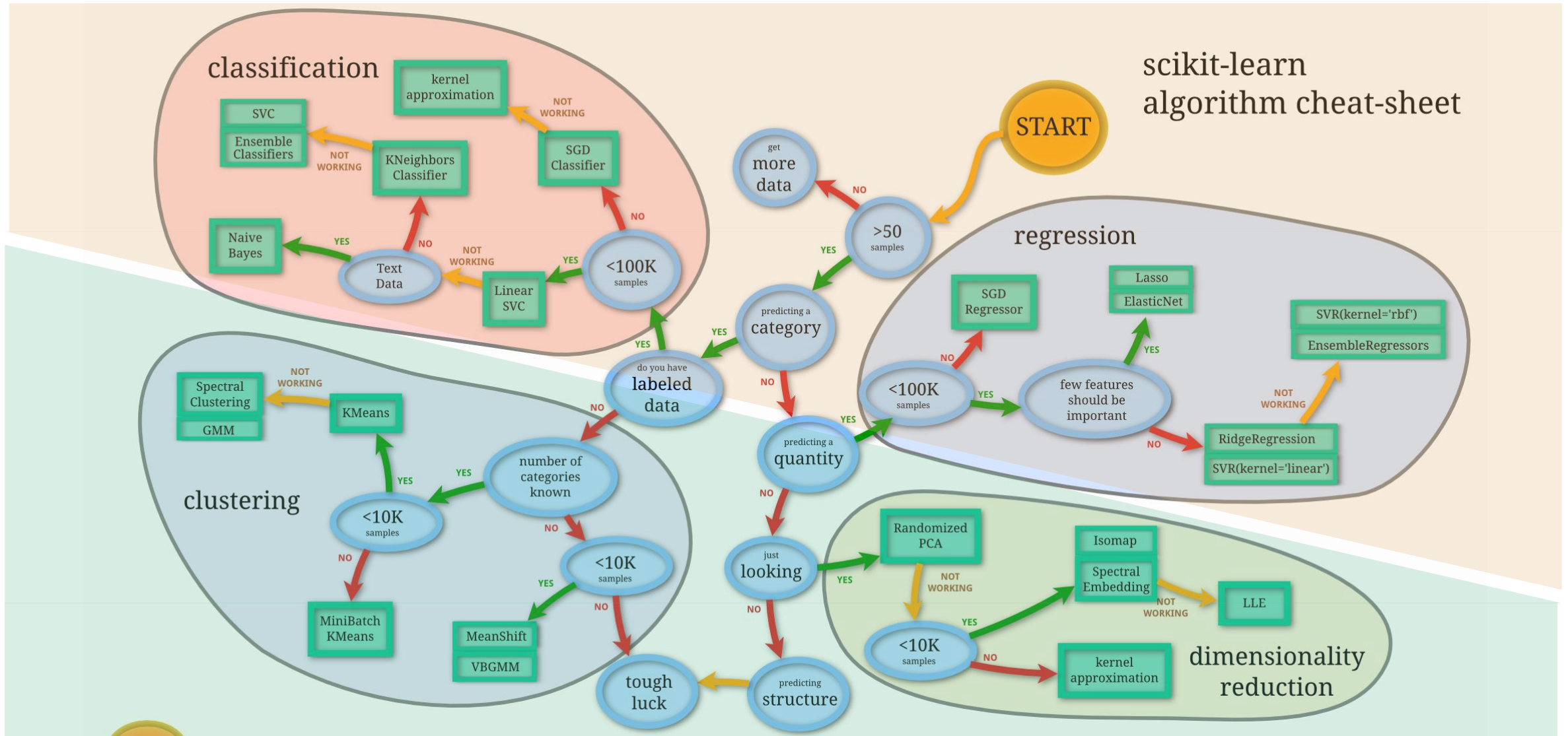


	지도학습 (Supervised Learning)	비지도학습 (Unsupervised Learning)	강화학습 (Reinforcement Learning)
훈련 방식	정답이 있는 데이터로 훈련	정답이 없는 데이터로 훈련	자신의 행동에 대한 보상을 받으며 목표를 달성하는 방향으로 학습
주요 알고리즘	분류, 회귀	군집화, 차원축소	로보틱스, 시뮬레이션
예시	<ul style="list-style-type: none">강아지와 고양이 사진 분류하기(분류)집 값 예측하기(회귀)	<ul style="list-style-type: none">라벨이 없는 데이터를 n개의 집단으로 구분(군집화)변수의 여러가지 특징을 보다 작은 수로 축소(차원 축소)	<ul style="list-style-type: none">자율주행 차알파고 등

2. 지도학습과 비지도학습

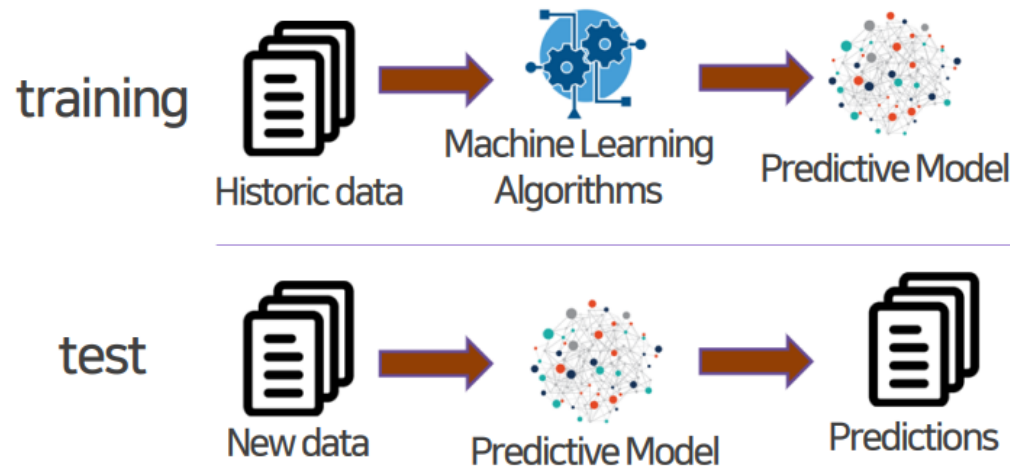


2. 지도학습과 비지도학습



2. 지도학습과 비지도학습 – 지도학습의 경우

- 본 적이 없는(학습데이터에 없었던) test data의 output을 예측(prediction)
- 어떤 input이 output에 어떻게 영향을 미쳤는지 이해하고 분석(inference)
- 모델을 평가하고, 다시 훈련하는 반복과정을 거쳐 성능을 향상시킴



데이터 준비 (데이터-라벨 쌍)

데이터 정제

데이터 훈련

성능 평가

최종 결론

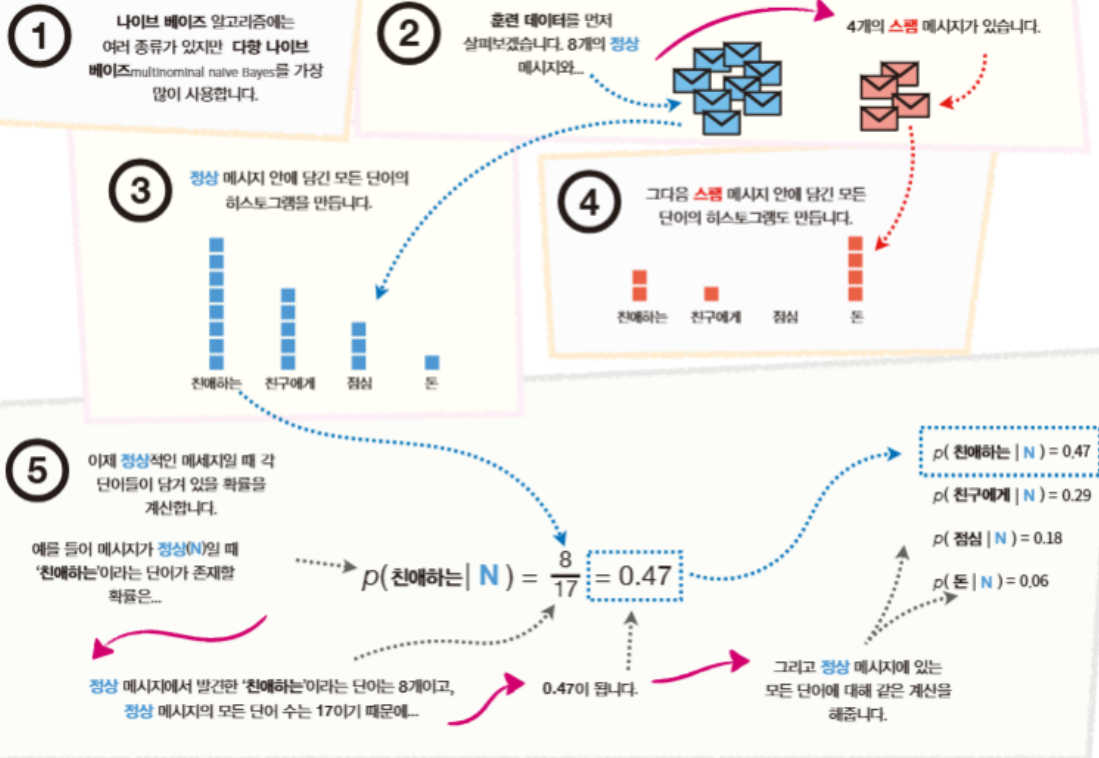
지도학습 알고리즘

	알고리즘 명	분류(classification)	회귀(regression)
1	나이브 베이즈 분류(Naïve Bayes Classification)	O	X
2	로지스틱 회귀(Logistic Regression)	O	X
3	서포트 벡터 머신(SVM)	O	O
4	K-최근접 이웃(K-nearest neighbors)	O	O
5	결정 트리(Decision Tree)	O	O
6	선형 회귀(Linear Regression)	X	O
7	Lidge, Lasso	X	O
8	ElasticNet	X	O

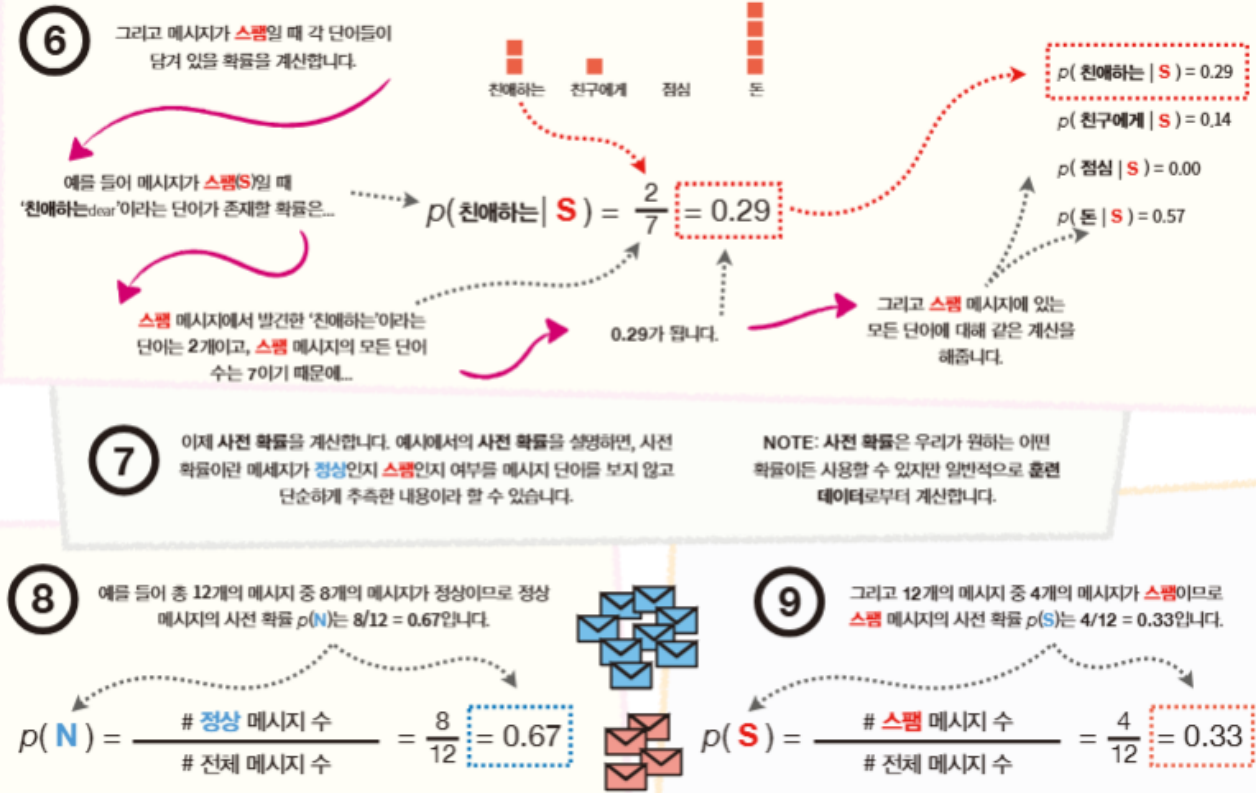
3. 지도학습 - 분류

나이브 베이즈 분류기

다항 나이브 베이즈: 자세히 살펴보기 Part 1



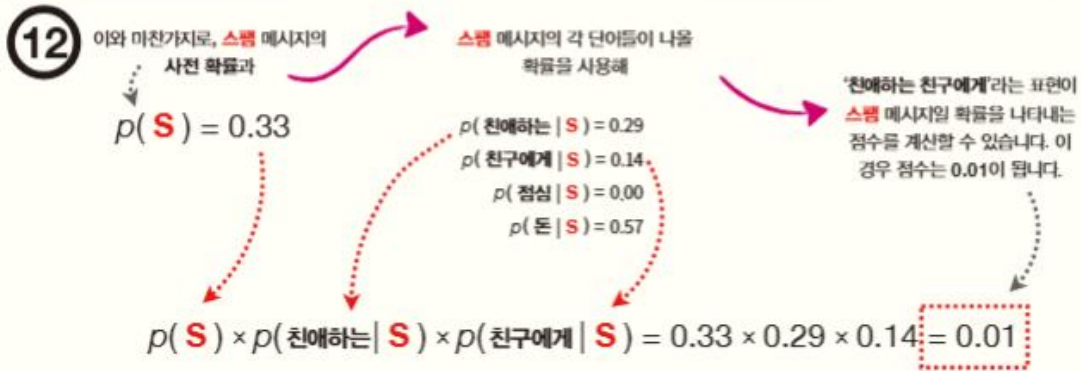
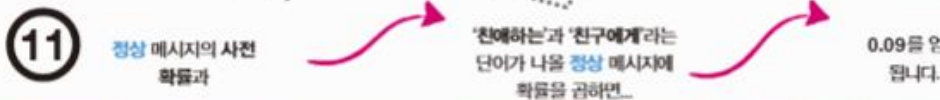
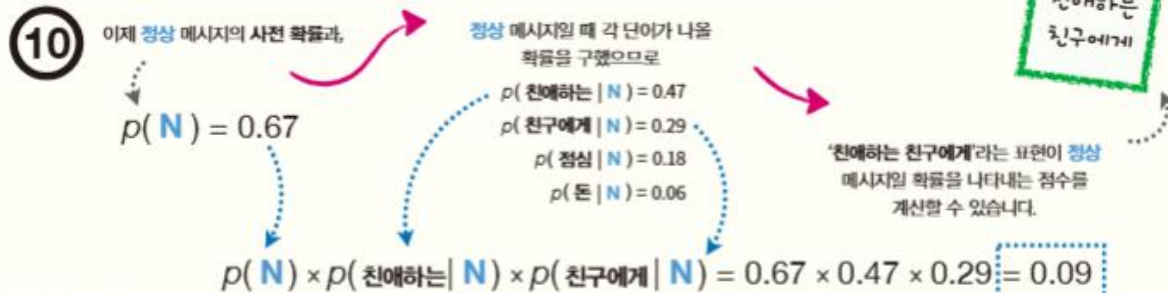
다항 나이브 베이즈: 자세히 살펴보기 Part 2



3. 지도학습 - 분류

나이브 베이즈 분류기

다항 나이브 베이즈: 자세히 살펴보기 Part 3



다항 나이브 베이즈: 자세히 살펴보기 Part 4

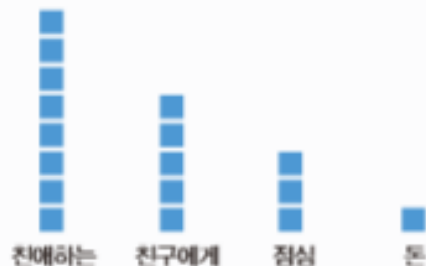
13

다시 복습해봅시다. 우리의 목표는 '친애하는 친구에게'라는 메시지가 **정상 메시지**인지 **스팸**인지 여부를 구별하는 것입니다.

8개의 **정상** 메시지와 4개의 **스팸** 메시지가 있는 훈련 데이터를 사용했습니다.



그다음 메시지에 담긴 각 단어의 히스토그램을 만들었습니다.



$$\begin{aligned} p(\text{친애하는} | N) &= 0.47 \\ p(\text{친구에게} | N) &= 0.29 \\ p(\text{점심} | N) &= 0.18 \\ p(\text{돈} | N) &= 0.06 \end{aligned}$$

그리고 히스토그램을 사용해 확률을 계산했습니다.



$$\begin{aligned} p(\text{친애하는} | S) &= 0.29 \\ p(\text{친구에게} | S) &= 0.14 \\ p(\text{점심} | S) &= 0.00 \\ p(\text{돈} | S) &= 0.57 \end{aligned}$$

그다음 메시지가 **정상** 혹은 **스팸**일 조건에 따라 사전 확률과 각 단어가 나올 확률로 '친애하는 친구에게'라는 메시지의 점수를 계산합니다.

$$p(N) \times p(\text{친애하는} | N) \times p(\text{친구에게} | N) = 0.67 \times 0.47 \times 0.29 = 0.09$$

$$p(S) \times p(\text{친애하는} | S) \times p(\text{친구에게} | S) = 0.33 \times 0.29 \times 0.14 = 0.01$$

$$p(N) = \frac{\text{\#정상 메시지 수}}{\text{\# 전체 메시지 수}} = 0.67$$

$$p(S) = \frac{\text{\#스팸 메시지 수}}{\text{\# 전체 메시지 수}} = 0.33$$

먼저 **정상** 혹은 **스팸**일 수 있는 메시지에 담긴 내용을 보지 않고 추측한 확률인 사전 확률을 계산했습니다.

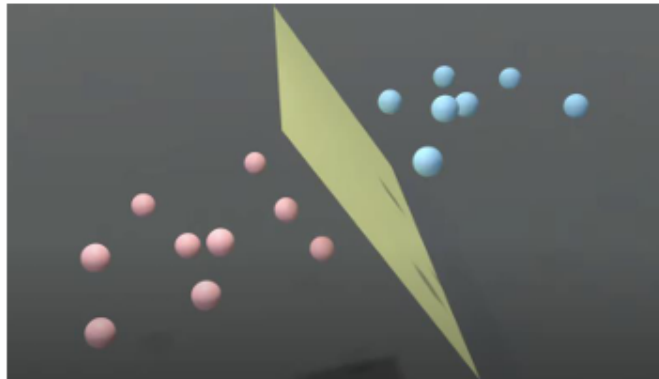
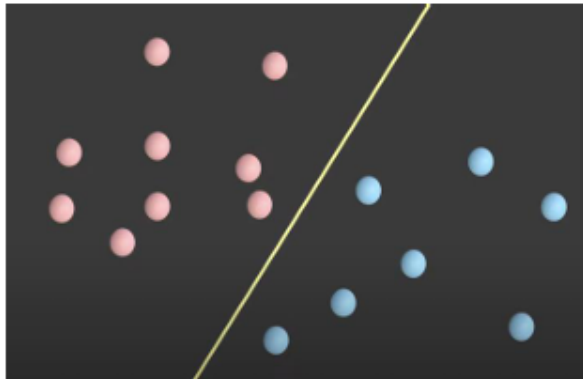
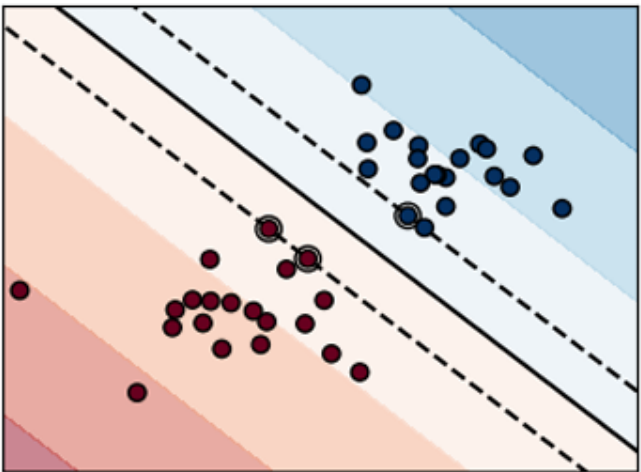
이제 '친애하는 친구에게'라는 메시지를 분류할 수 있게 되었습니다. 왜냐하면 **정상** 메시지일 확률, 즉 점수는 0.09로 해당 메시지가 **스팸**일 점수(0.01)보다 크기 때문에 우리는 이 메시지를 **정상** 메시지로 구분합니다.

친애하는 친구에게

BAM!!!

3. 지도학습 - 분류

서포트 벡터(Support Vector)머신



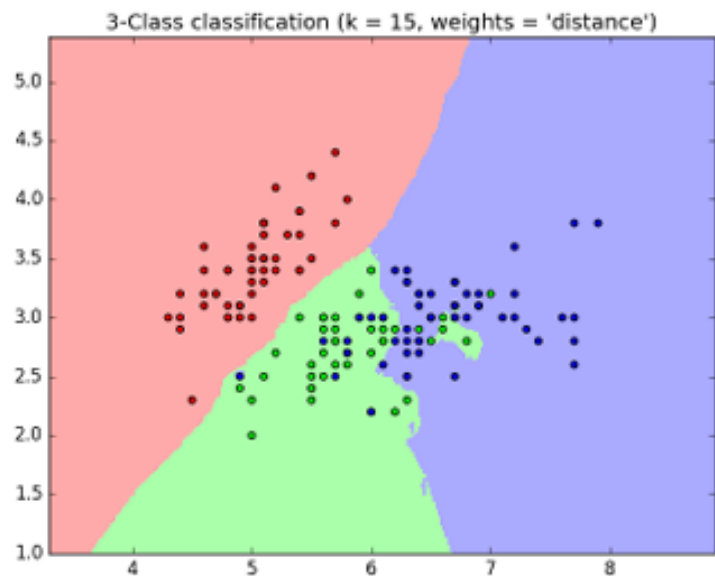
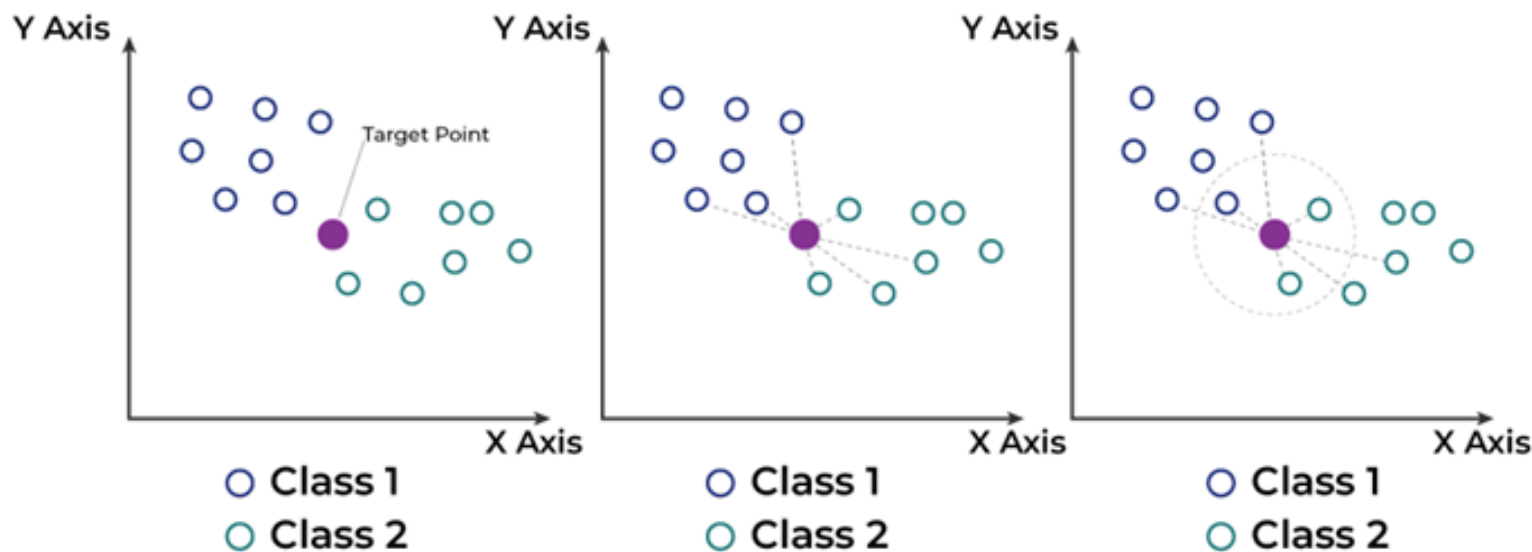
데이터의 특징들(x, features) 중 가능한 한 먼 결정 경계를 찾아내는 것

마진(결정 경계에 가장 가까운 점과 결정 경계 사이의 거리)의 최대화

- 최적의 초평면 찾기 -> 초평면에 가장 가까운 데이터를 '서포트 벡터'로 식별 -> 서포트벡터를 기준으로 결정경계 구축
- 장점 : 일반화 능력(다양한 케이스에 적용 가능)과 높은 성능을 보여줌
- 단점 : 파라미터를 어떻게 세팅하느냐에 따라 결과의 변동성이 큼, 모델 해석이 어려움
- 주의할 점 : 소프트마진(데이터를 완전히 분류하기 어려울 때 일부 데이터가 마진 안에 포함되는 경우)가 발생할 수 있음

3. 지도학습 - 분류

K-최근접 이웃



주변 K개 데이터 포인트와의 거리를 계산하여, 다수결로 더 가까운 데이터로 분류/회귀

주어진 데이터에 대해 주변 k개 데이터와의 유클리드 거리를 계산

- 장점 : 구현이 쉽고, 데이터 포인트가 비교적 섞여 있는 경우에도 안정적으로 분류 가능
- 단점 : 리소스의 소모가 크고 과적합에 취약

3. 지도학습 - 분류

결정 트리



조건 분기에 따라 학습 데이터를 나누어 분류

불순도를 최소화하도록 데이터를 나누며 분류 수행

- 장점 : 직관적인 형태로 모델 이해 가능, 스케일링 불필요, 비선형 데이터에 대해 강건
- 단점 : 과적합 가능성이 높음, 데이터 변화에 민감
- 분류는 지니 계수를 기준으로, 회귀의 경우 MSE를 기준으로 노드를 나누어 분류/회귀 모두 사용
- 불순도 지표 : 지니 계수

2. 지도학습과 비지도학습 – 일반적인 머신 러닝의 순서

- 1 데이터 불러오기
- 2 데이터 확인하기(통계적 특징, 데이터의 크기 등)
- 3 데이터 전처리(결측치 및 이상값 정리, 스케일링 등)
- 4 train_test_split 및 x, y 데이터의 정의
- 5 머신러닝 모델 정의 및 훈련, 검증

2. 지도학습과 비지도학습 – 비지도학습의 경우

- 데이터의 특징을 활용하여 군집화/차원 축소
- 분석가는 군집화/차원 축소의 결과물을 활용하여 데이터를 분석하거나, 지도학습을 수행
- 목표와 일치하지 않는 군집화/차원 축소인 경우, 다른 모델을 쓰거나 재군집화/차원 축소

A selection from the 64-dimensional digits dataset

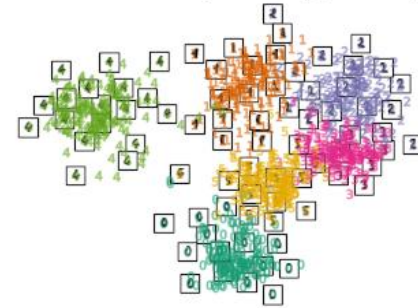
0	1	2	3	4	5	0	1	2	3
4	5	0	1	2	3	4	5	0	5
5	5	0	4	1	3	5	1	0	0
2	2	2	0	1	2	3	3	3	3
4	4	1	5	0	5	2	2	0	0
1	3	2	1	4	3	1	3	1	4
3	1	4	0	5	3	1	5	4	4
2	2	2	5	5	4	4	0	0	1
2	3	4	5	0	1	2	3	4	5
0	1	2	3	4	5	0	5	5	5

비지도
학습

Truncated SVD embedding (time 0.003s)



Linear Discriminant Analysis embedding (time 0.006s)



데이터 준비

데이터 정제

데이터 훈련

성능 평가

최종 결론