

KAIST 산학협동강좌 2024  
Feb. 19<sup>th</sup> (Mon), 2024

# Ryu Part 1/4: AI 작동 원리 요약 및 AI 활용 시 유의점

Seunghwa Ryu  
Professor of Mechanical Engineering

Korea Advanced Institute of Science and Technology (KAIST)

→ Feb 19<sup>th</sup>  
Afternoon  
Contents

01 강의1: AI작동원리 요약 및 AI 활용 시 주의점

02 강의2: 도메인 지식과 AI를 결합한 Data-efficient 설계 개괄

03 실습1: 전산유체(CFD) 모델링 데이터를 활용한 딥러닝 모델 학습

04 실습2: 학습된 딥러닝 모델을 활용한 유체장치 설계 실습

**AI 작동원리 요약**

**인공지능 활용 시 유의점?**

# 인공지능은 무엇을 하는가?

## 예측과 분류

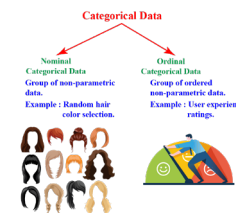
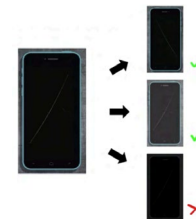
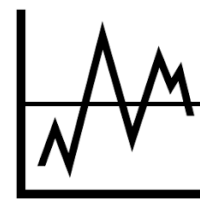
# 인공지능시스템 = AI모델 + Data

DNN, CNN, RNN,  
Support Vector  
Machine, Decision  
Tree,  
Random Forest, ...

→ PyCaret 오픈소스로  
20개 이상 알고리즘  
동시 학습 및 테스트 가능

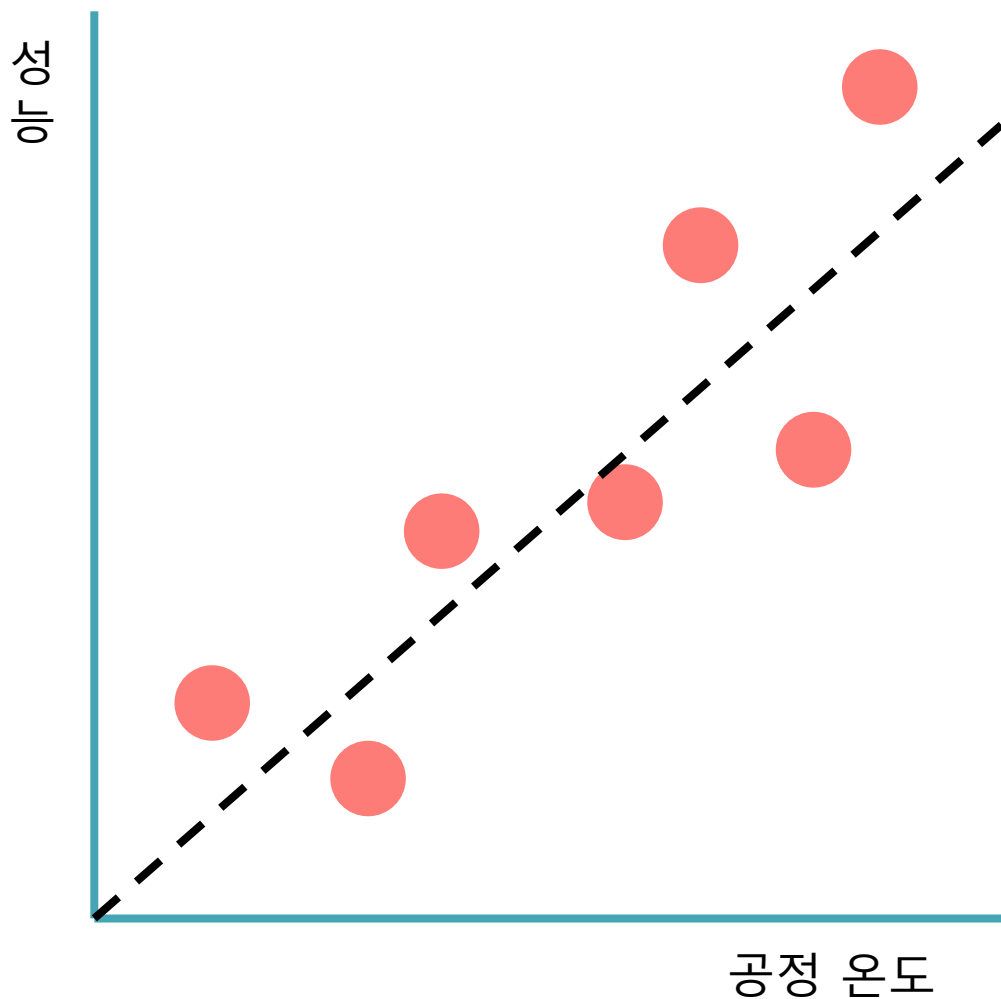
숫자데이터, 시계열데이터,  
그림데이터, 카테고리데이터

TEM	PRE	PER
24	3	40
46	4	35
43	6	26
33	1	50



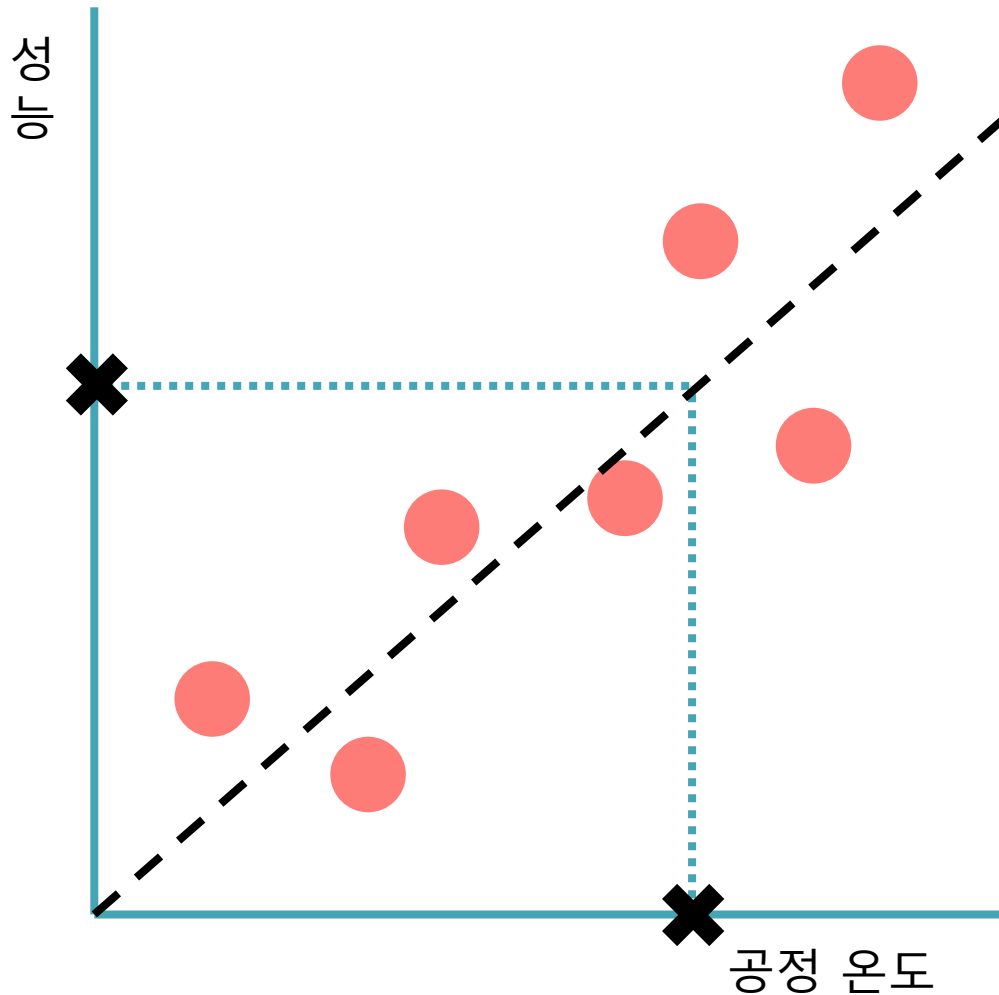
→ 데이터 오픈하는 기업 없음

# AI의 역할 I. 예측, 회귀분석



데이터 경향선 추출  
전반적인 제품 품질과  
공정 온도의 인과관계를 예측

# AI의 역할 I. 예측, 회귀분석



데이터 경향선 활용하여  
원래 데이터에 없던  
새로운 데이터에 대한  
“예측” 가능

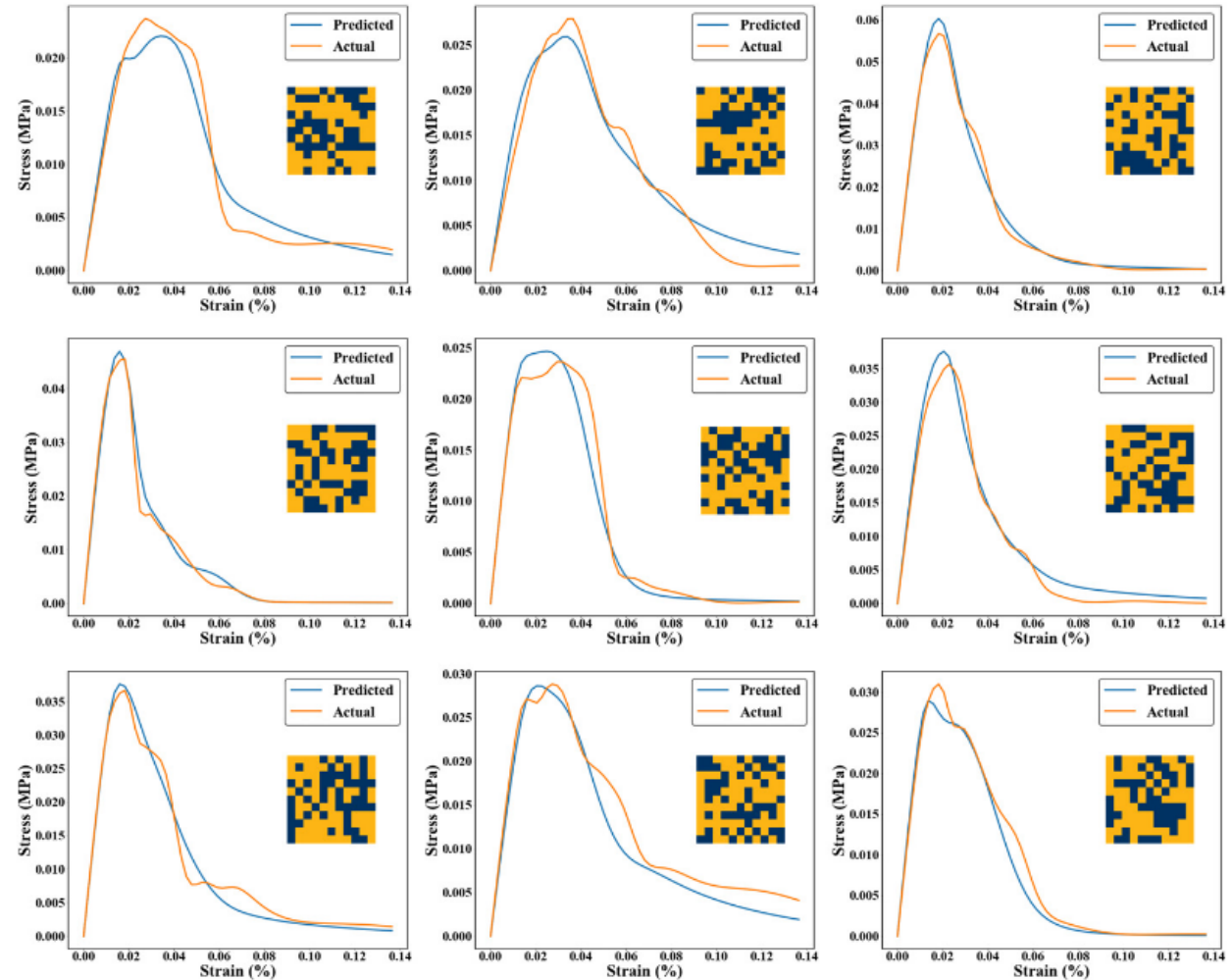
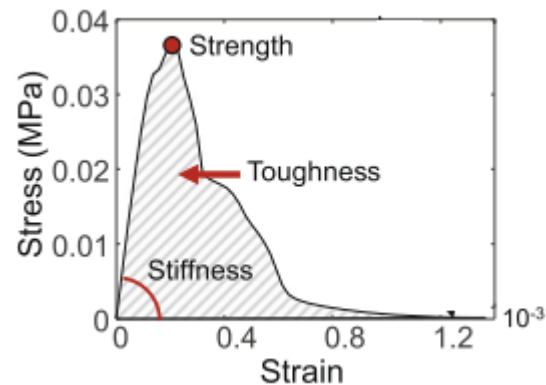
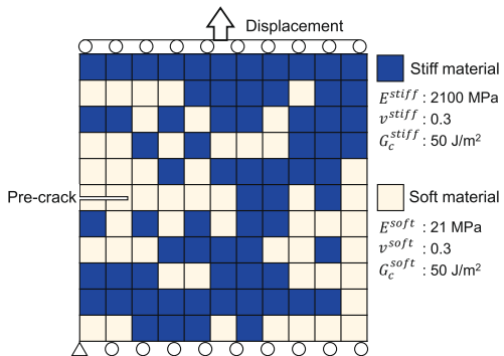
# AI의 역할 I. 예측, 회귀분석

Input & Output이

형상/패턴/이미지라도

딥러닝 회귀분석 가능

Input / Output





# AI의 역할 II. 분류

수 백 개의 제품 데이터

적정온도	불순물 검출	적정점도	품질기준
아니오	아니오	아니오	불합격
네	네	네	합격
네	네	아니오	불합격
네	아니오	아니오	합격
...	...	...	...

3가지 적정 공정조건으로  
품질기준 통과 예측 및 분류



# AI의 역할 II. 분류

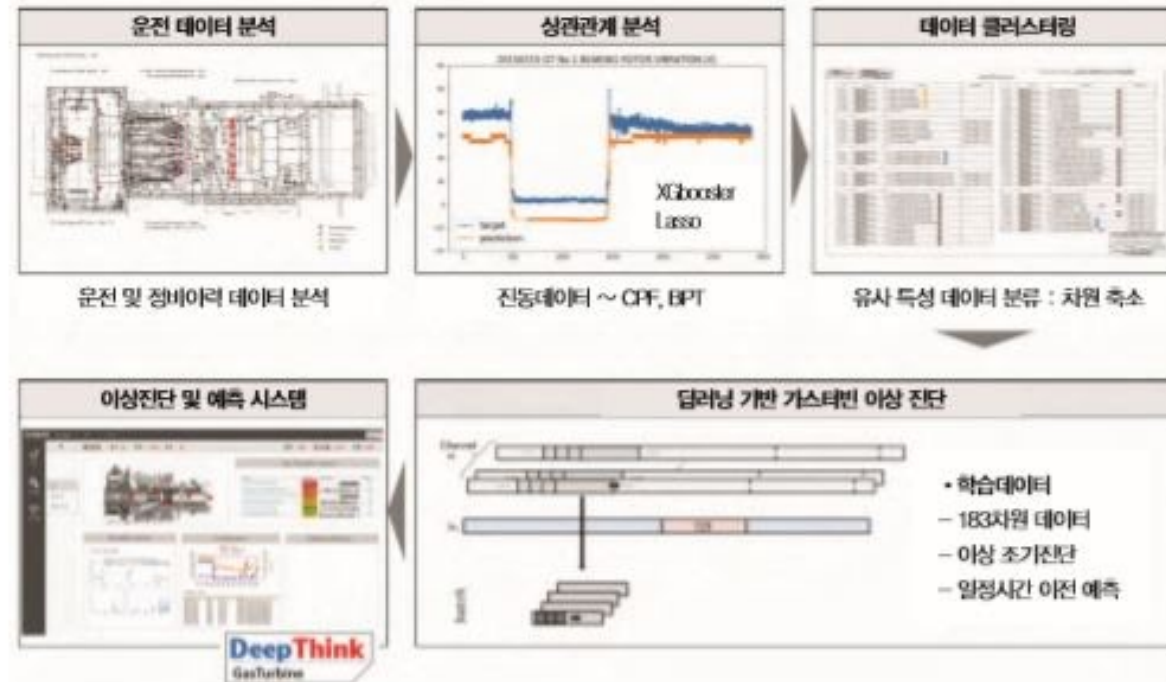


그림 31 답라닝 기반 가스터빈 이상진단 및 예측 개요도

<출처: 한전전력연구원, <http://www.keaj.kr/news/articleView.html?idxno=2789>>

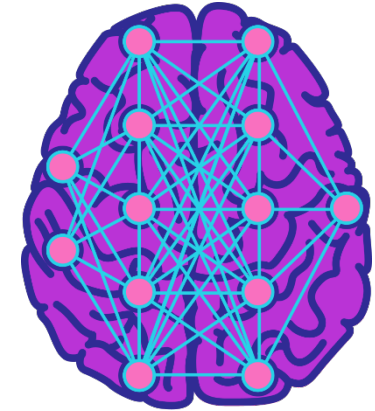
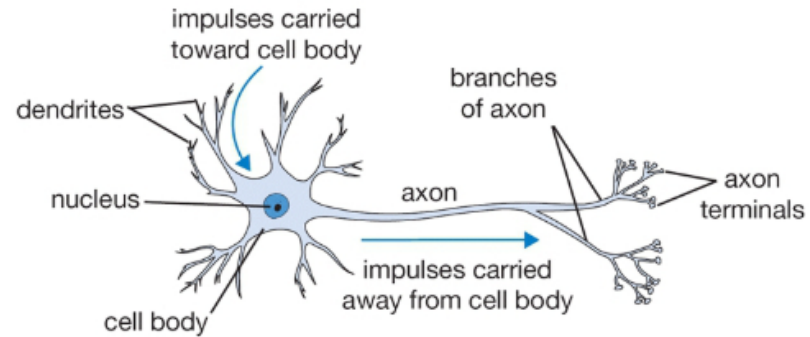
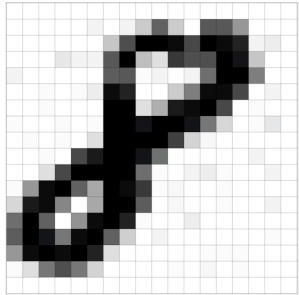
**Q. 압력, 진동, 온도 시계열 데이터에서 이상 진단?**

**Q. 시스템 수명에 무리를 안 주는 운영 조건?**

# 회귀분석 알고리즘 작동원리

심층신경망 – 유연한 비선형 회귀

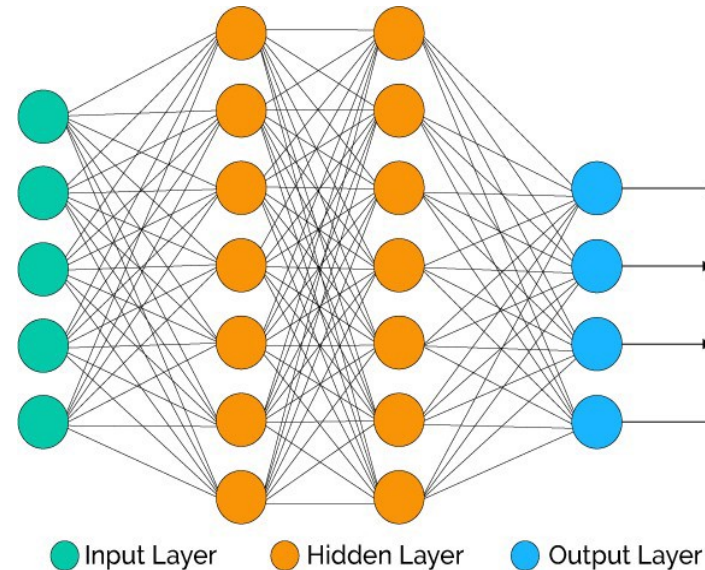
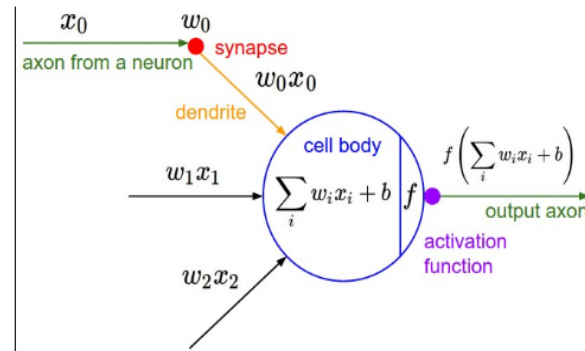
# 인공 신경망?!



이를 모사하여 **신경망(퍼셉트론)**이 활성화되면서 학습하는 모델!

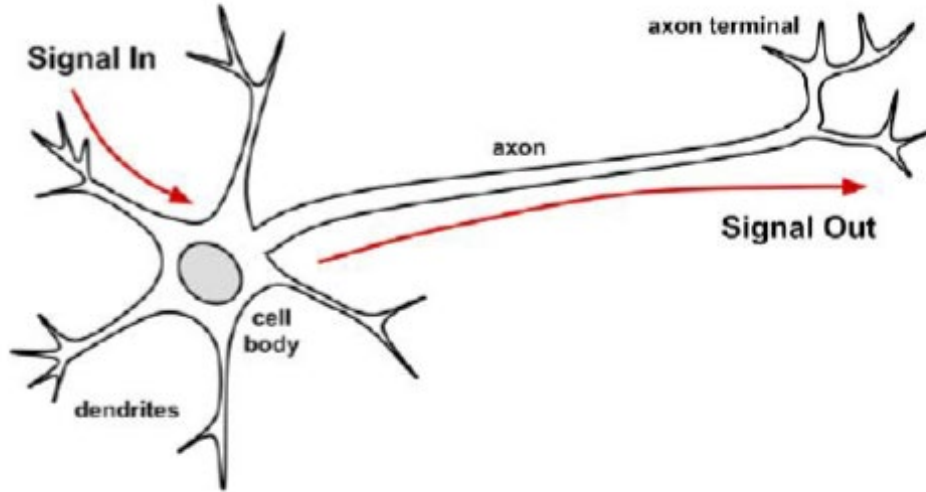
```

0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 1 12 0 11 39 137 37 0 152 147 84 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 1 0 0 0 0 41 168 258 255 235 162 255 238 206 11 13 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 16 9 9 158 251 45 21 184 159 154 255 233 40 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
10 0 0 0 0 0 0 145 146 3 10 0 11 124 253 255 187 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 3 0 4 15 236 216 0 0 38 189 247 240 169 0 11 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 0 2 0 0 0 253 253 23 62 224 241 255 164 0 5 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
6 0 0 4 0 3 252 258 228 255 255 234 112 28 0 2 17 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 2 1 4 0 21 255 253 251 255 172 31 8 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 4 0 163 225 251 255 229 120 0 0 0 0 0 11 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 21 162 255 255 254 255 126 6 0 10 14 6 0 0 9 0 0 0 0 0 0 0 0 0 0 0 0 0 0
3 79 242 255 141 66 255 245 189 7 8 0 0 0 5 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
26 221 237 98 0 67 251 255 144 0 8 0 0 7 0 0 11 0 0 0 0 0 0 0 0 0 0 0 0 0 0
125 255 141 0 87 244 255 208 3 0 0 13 0 1 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
145 248 228 116 235 255 141 34 0 11 0 1 0 0 0 1 3 0 0 0 0 0 0 0 0 0 0 0 0 0 0
85 237 253 246 255 218 21 1 0 1 0 0 6 2 4 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
6 23 112 157 114 32 0 0 0 0 0 2 0 8 0 7 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
    
```

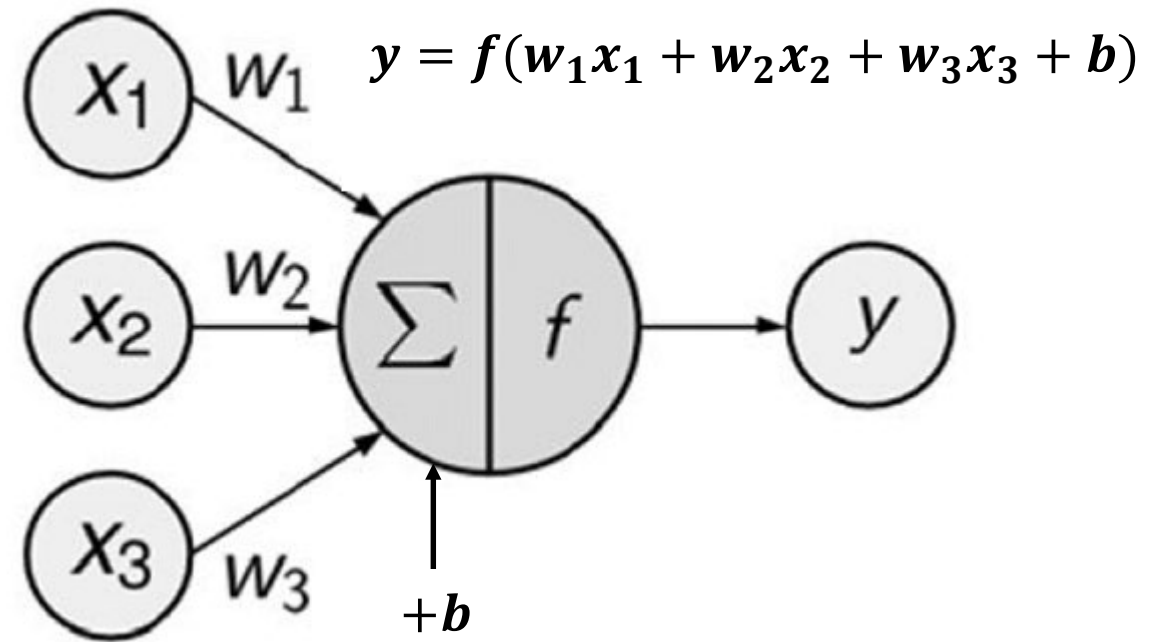


# 신경망이란?

## 사람의 뉴런



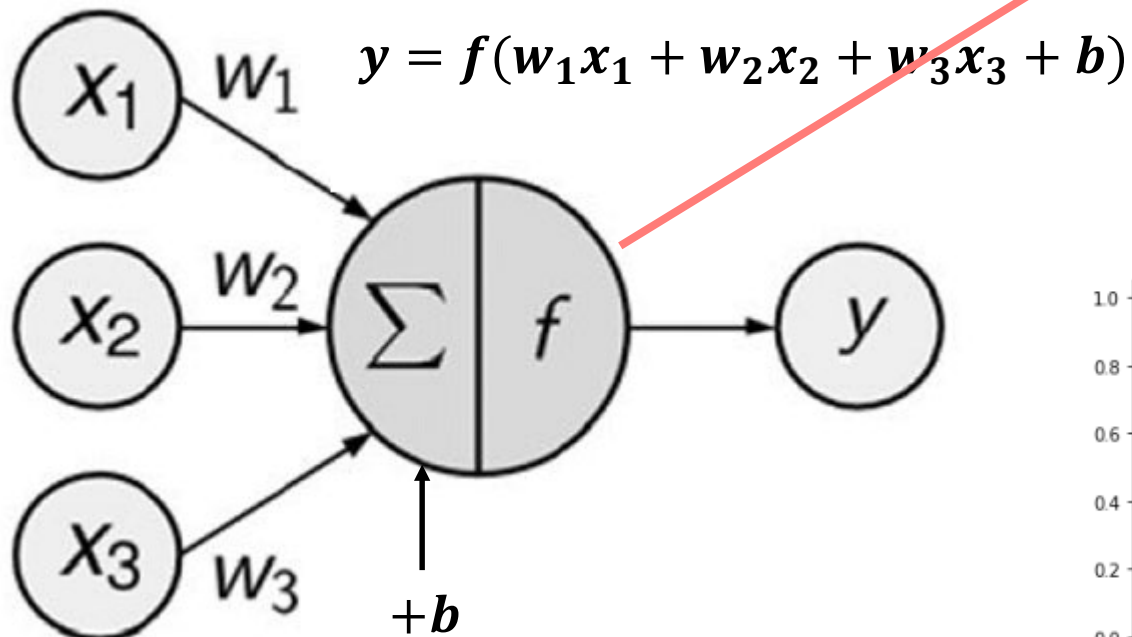
## 인공 신경망의 퍼셉트론 (Perceptron)



입력값  $x_1, x_2, x_3$  에 대한 출력값  $y$

# 신경망이란?

## 인공지능 퍼셉트론

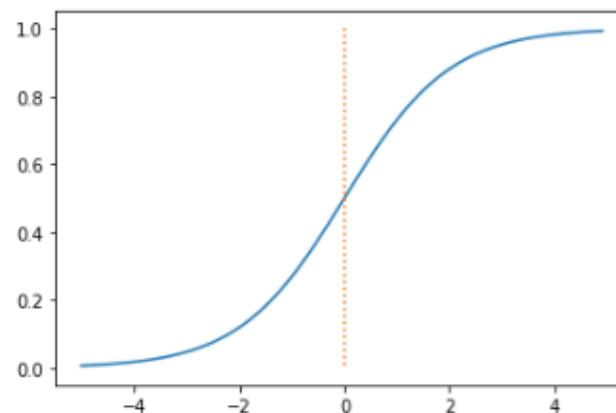


$f$  : **활성 함수** (activation function)

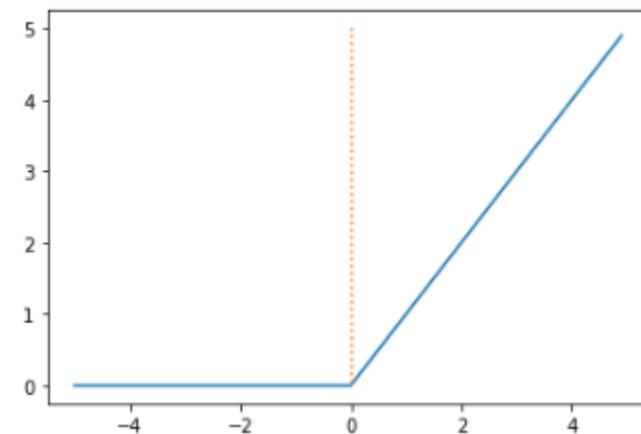
Sigmoid / Hyperbolic tangent / ReLu ...

→ 활성함수가 선형이면 선형회귀와 동일해짐.

Sigmoid 함수

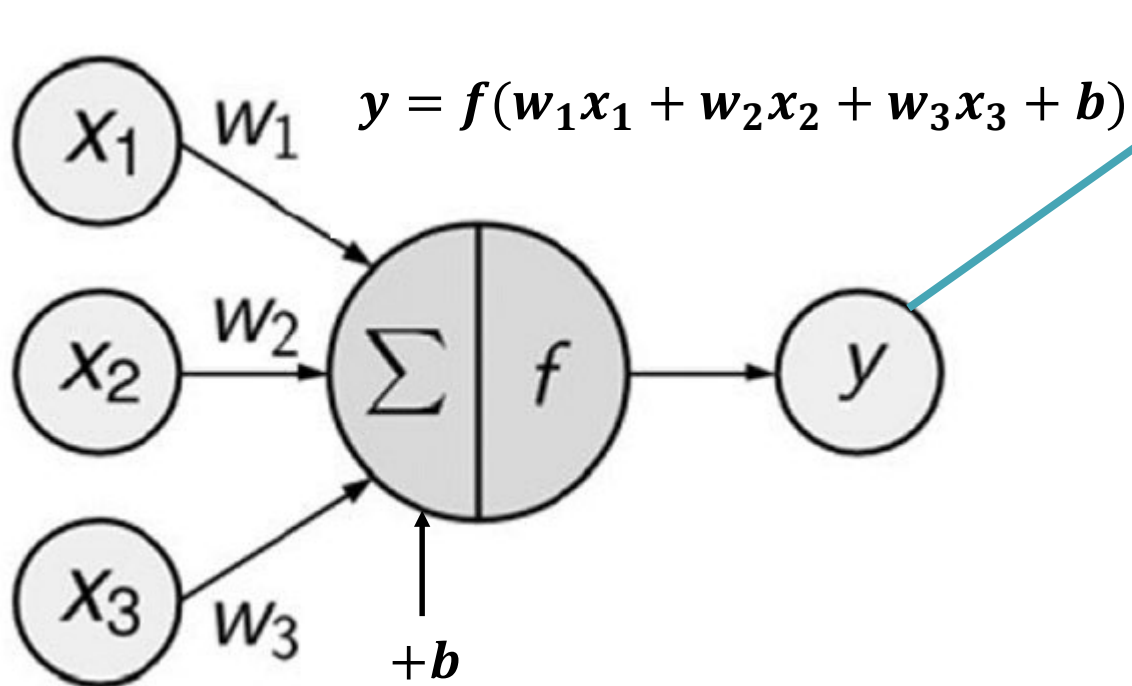


ReLu 함수



# 신경망이란?

## 인공지능 퍼셉트론



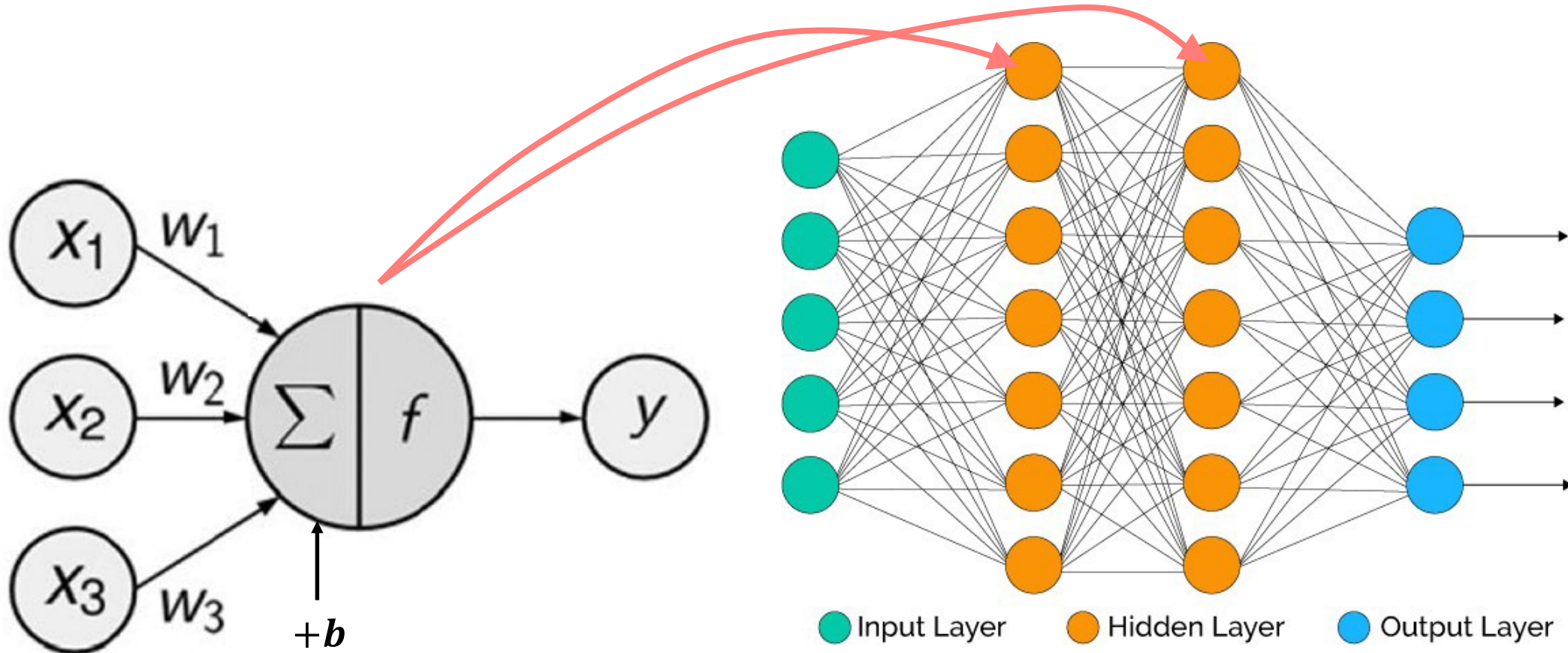
실제  $y$  와  $f(wx + b)$ 의 차이  
= 신경망의 오차 = 손실 함수(loss)

→ 신경망을 훈련한다는 것은,  
 $y$ 값을 잘 예측하는  $w$ ,  $b$  를 찾는다는 것!

→ 결국 손실 함수의 값을 줄이는 것이 신경망  
학습의 목표 방향!!



# 다층 신경망(Multilayer Perceptron)



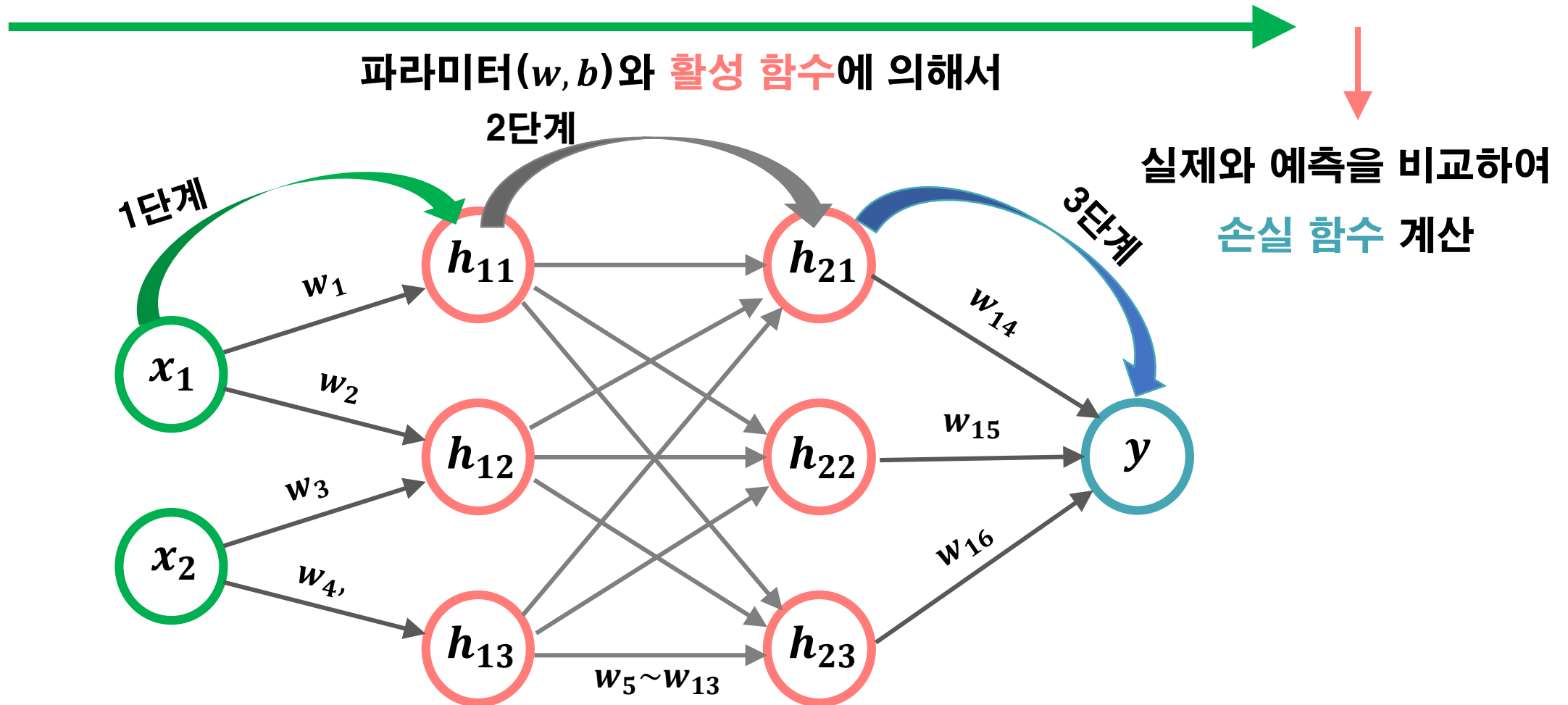
여러 층의 신경망으로 이루어진 다층 신경망

→ 심층신경망, 딥 러닝

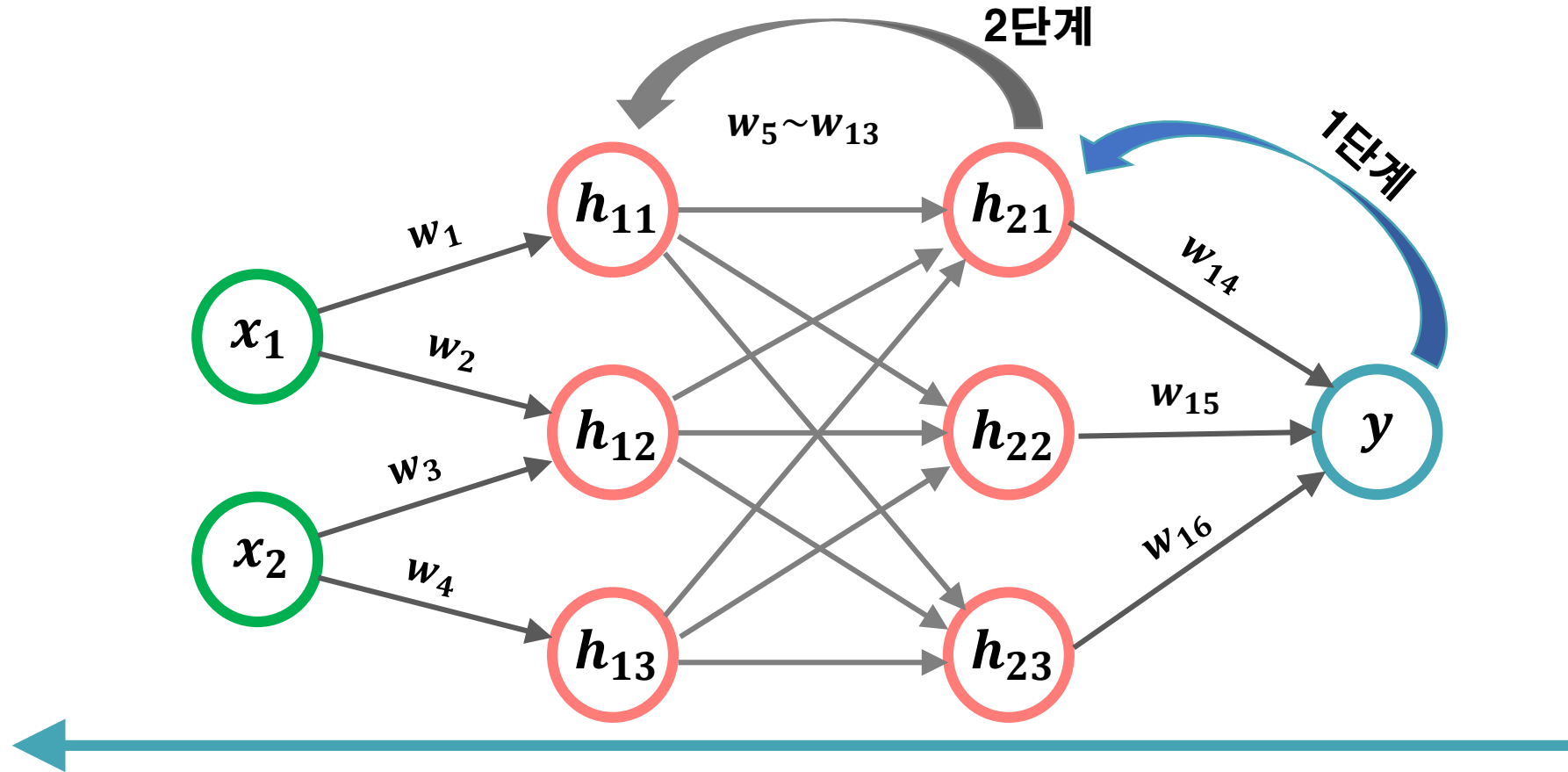
# 순전파(Forward Propagation)

입력값  $x_1, x_2, \dots$  을 넣으면

예측되는 출력  $y$ 가 계산



# 역전파(Back Propagation)



손실 함수 최소화를 목표로 출력 값에 가까운 순서로 파라미터( $w, b$ ) 조정

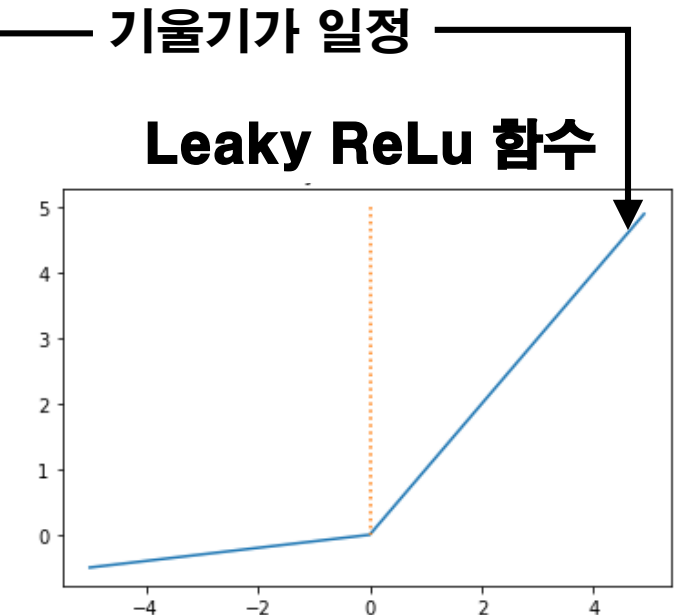
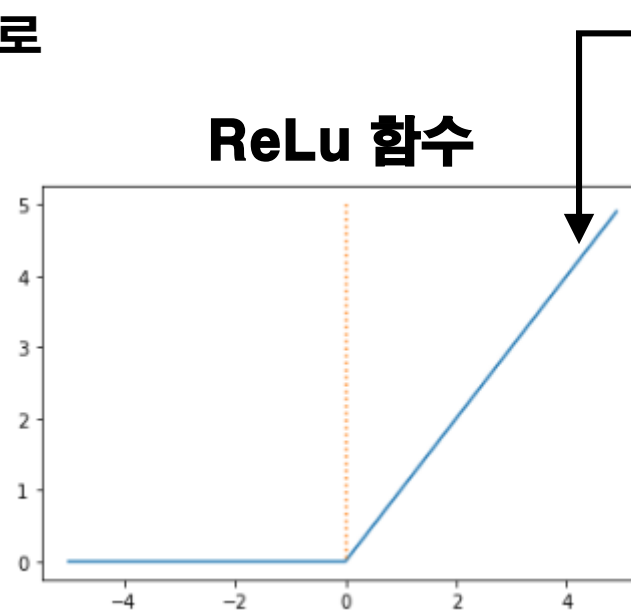
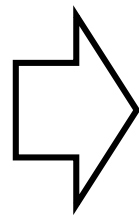
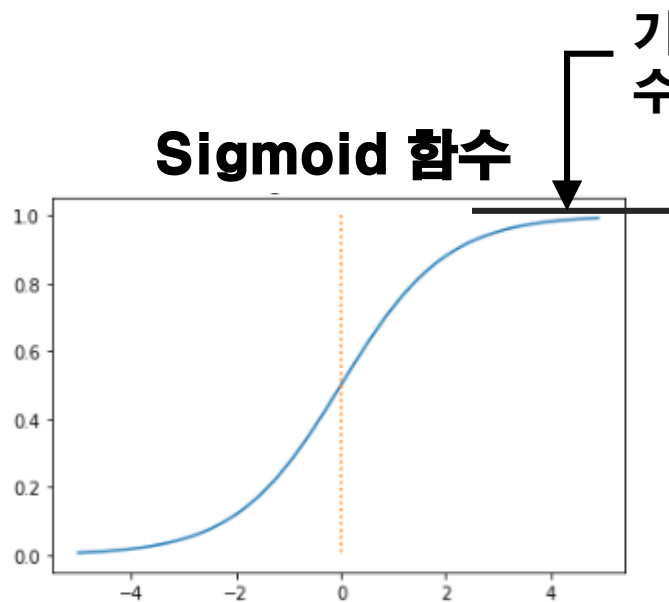
**How?** → 경사 하강법(Gradient descent):  $\frac{\partial Loss}{\partial w}$  를 계산하여 퍼셉트론 학습

# 활성 함수(activation function)

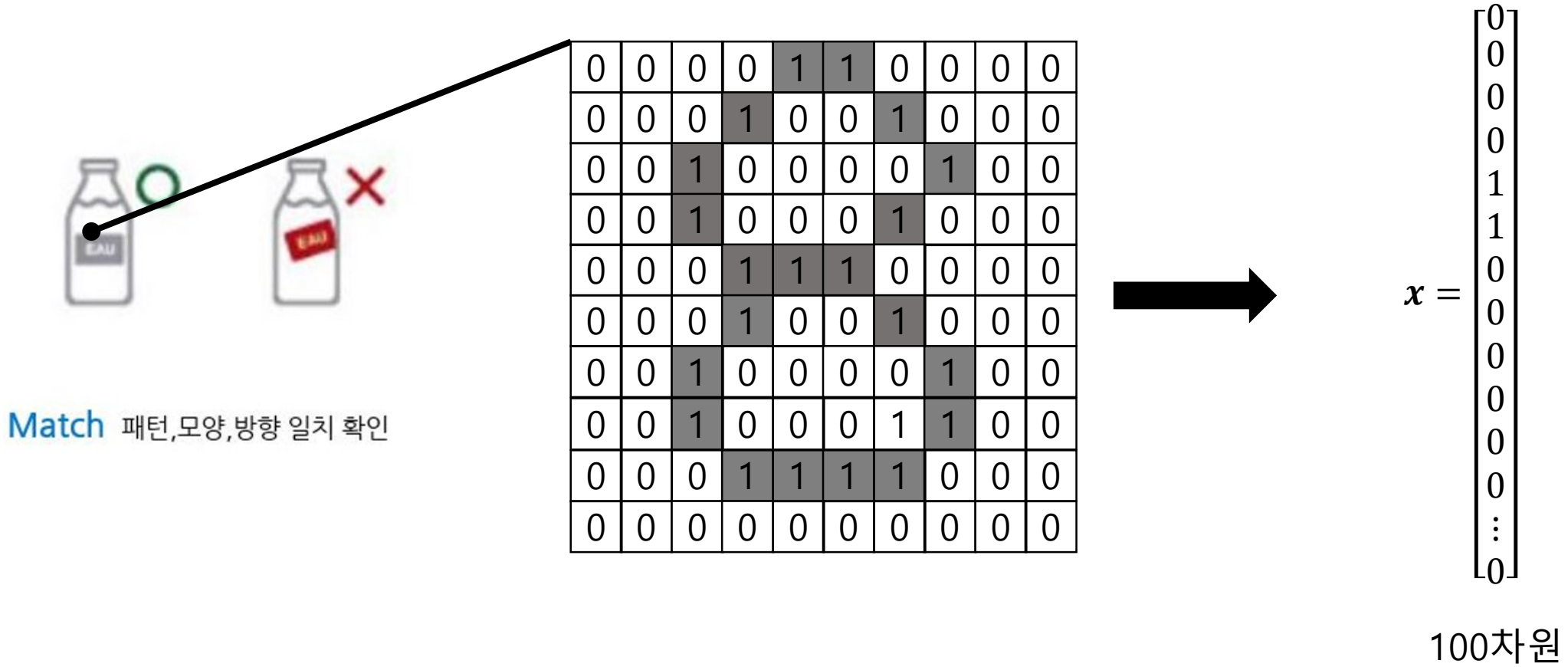
→ 은닉층(hidden layer)이 두꺼울수록 역전파 과정에서 기울기가 잘 전달되지 않음

$$\frac{df}{dx_1} = \frac{df}{df_1} \frac{df_1}{df_2} \frac{df_2}{df_3} \frac{df_3}{df_4} \frac{df_4}{df_5} \dots : \text{하나만 0에 가까워져도 곱은 0에 수렴함.}$$

→ ReLu 함수 기반의 **활성 함수** 사용

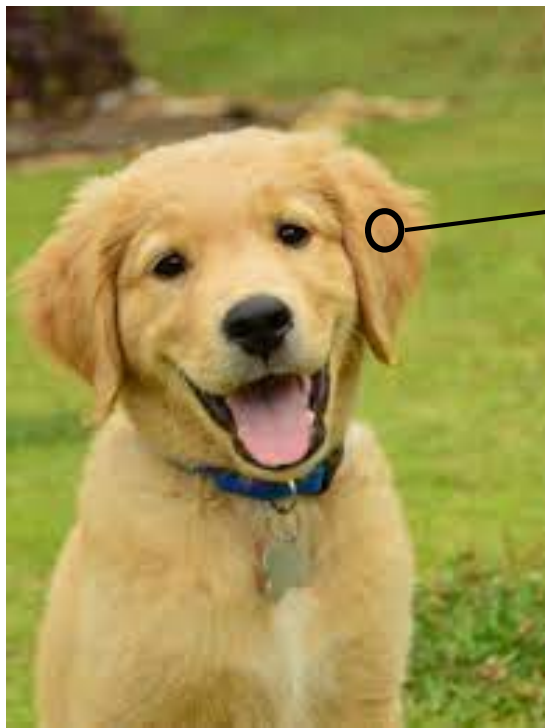


# 이미지 제조데이터의 분석



→ **합성곱 신경망을 활용하면, 단순한 숫자 나열 이상의 분석이 가능함.**

# 합성곱 신경망(CNN)



R:198

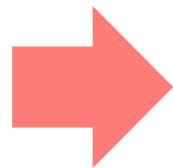
G:148

B:85

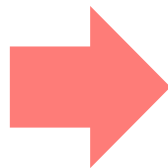


→ 단순히 각 점의 절대값만 아는 것으로는 판별할 수가 없다.

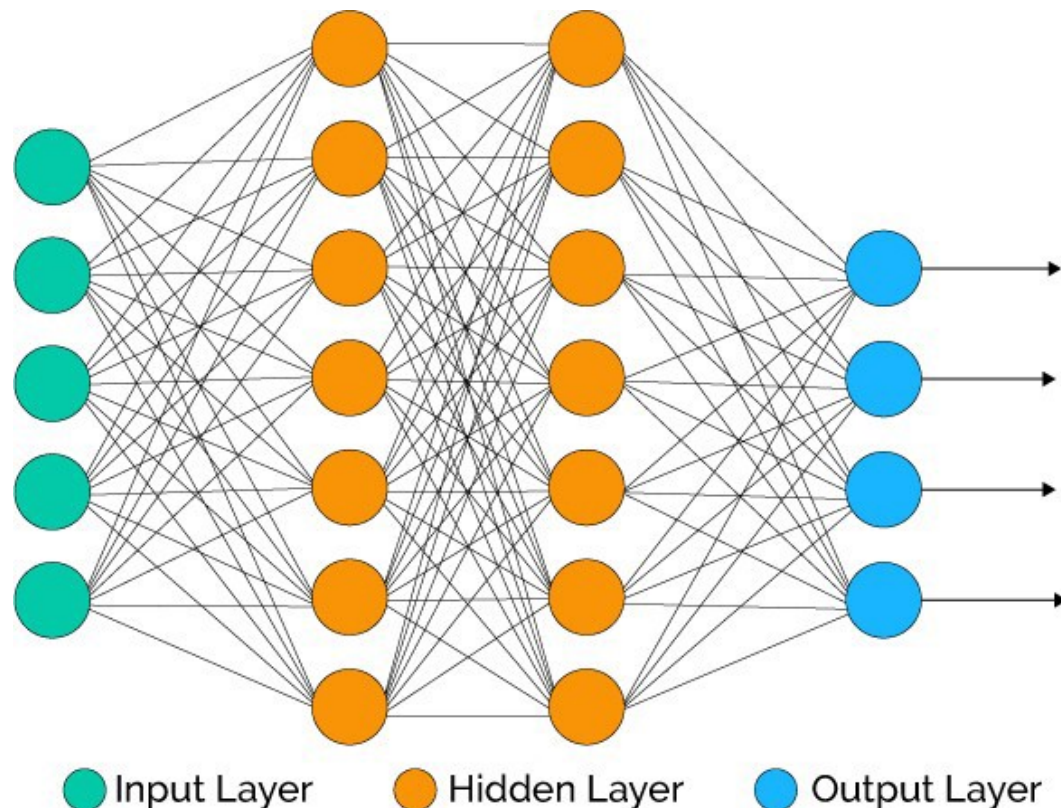
# 합성곱 신경망(CNN)



2차원 이미지를  
1차원으로 변경



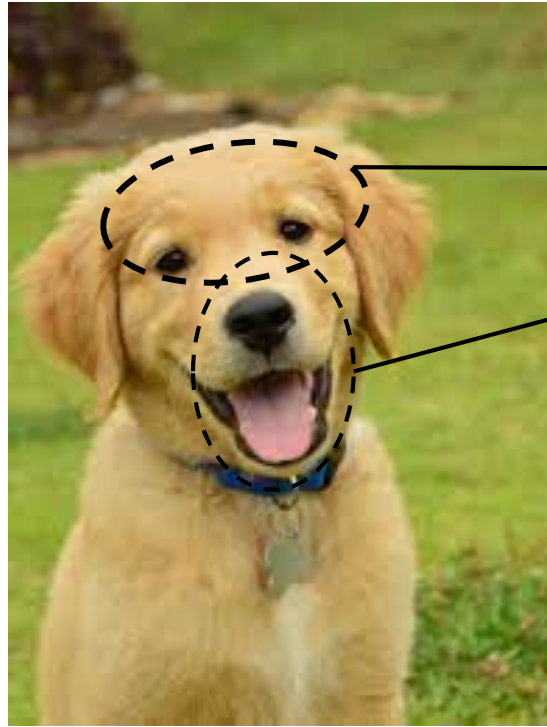
1차원 행렬을  
입력치로 활용



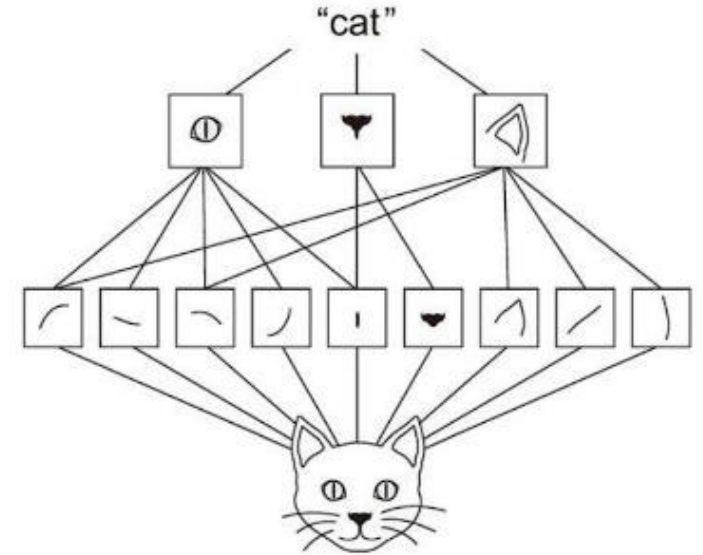
일반적인 신경망은 주위 픽셀의 정보를 고려하기 어려움



# 합성곱 신경망(CNN)



또다른 예시



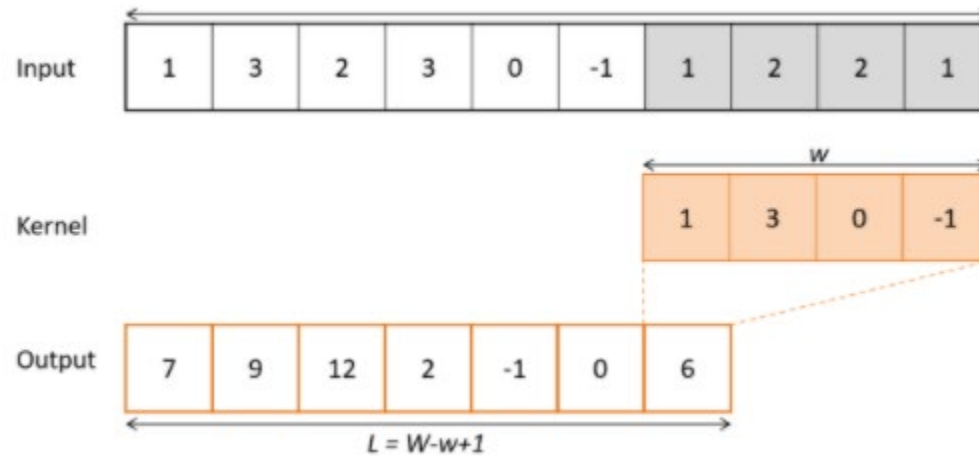
→ 여러 주위 점들 사이의 관계를 보고 이미지를 판별.

→ 인공지능 신경망도 이렇게 주변과의 **관계를 학습하면**  
더 좋지 않을까?



# 합성곱

## 1차원 합성곱의 개념



## 1차원 합성곱의 예



Input: 매일의 주식 가격

Kernel → 5일 이동평균: 5칸에 모두 1/5

→ 20일 이동평균: 20칸에 모두 1/20

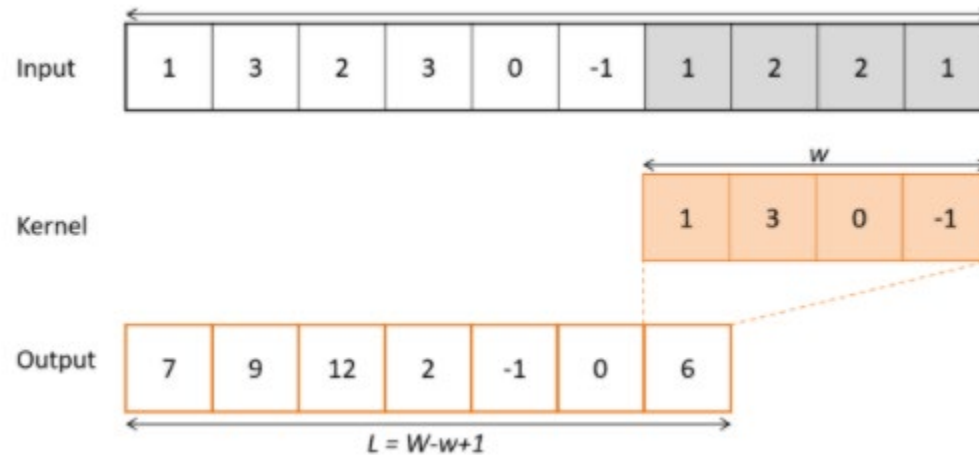
→ 60일 이동평균: 60칸에 모두 1/60

Output: 5일, 20일, 60일 이동평균

각 Kernel 합성곱 후 단순한 숫자 나열 이상의 정보를 표시

# 합성곱

## 1차원 합성곱의 개념



## 1차원 합성곱의 예

Input: 매일의 코로나 확진자

Kernel → 7일 이동평균: 7칸에 모두 1/7

Output: 7일씩 이동평균데이터

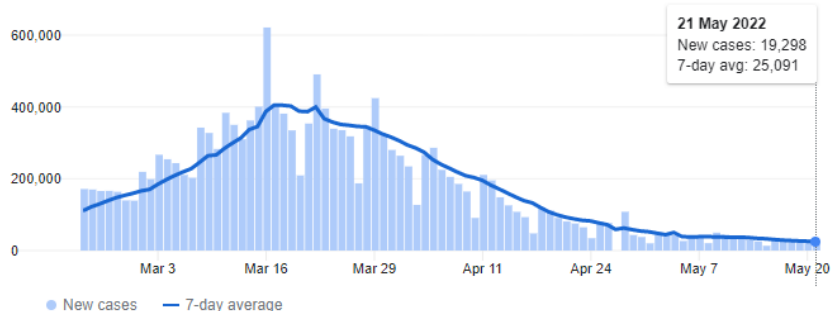
요일 별로 검사받는 사람수가 다른 걸 보정하여 패턴 파악

### Statistics

New cases Deaths Vaccinations Tests

From JHU CSSE COVID-19 Data · Last updated: 7 hours ago

South Korea 3 months



Each day shows new cases reported since the previous day · [About this data](#)

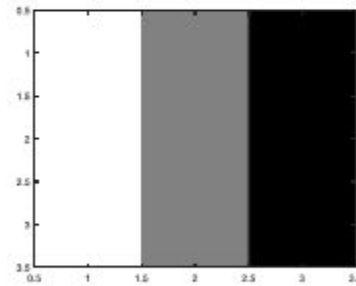
[Feedback](#)

# 합성곱



Image

$$\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$$



Kernel



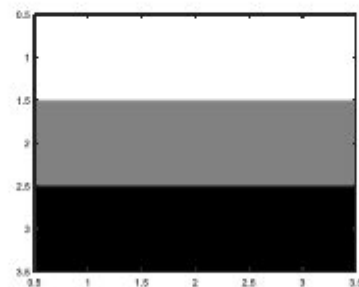
Output

# 합성곱



Image

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$$

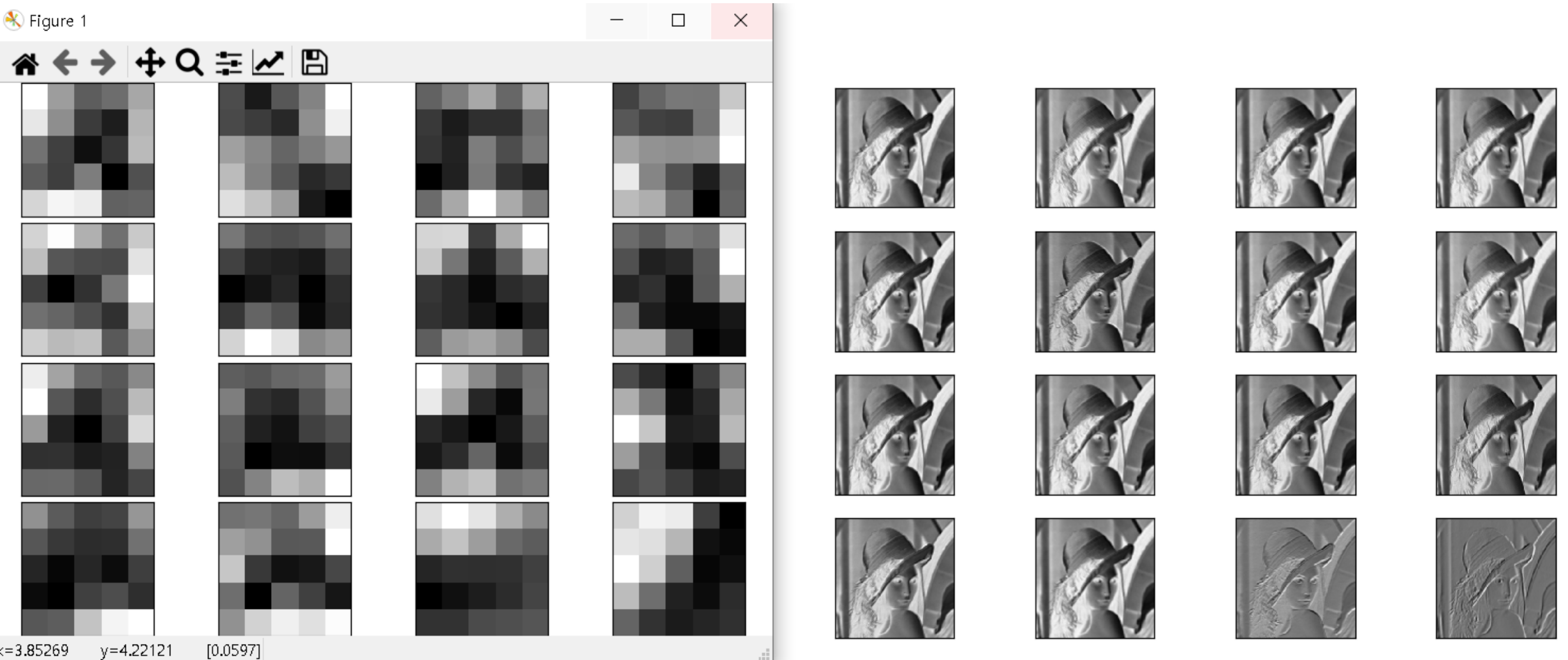


Kernel



Output

# 합성곱



**CNN으로 학습된 커널(좌)과 합성곱 처리된 이미지(우)**

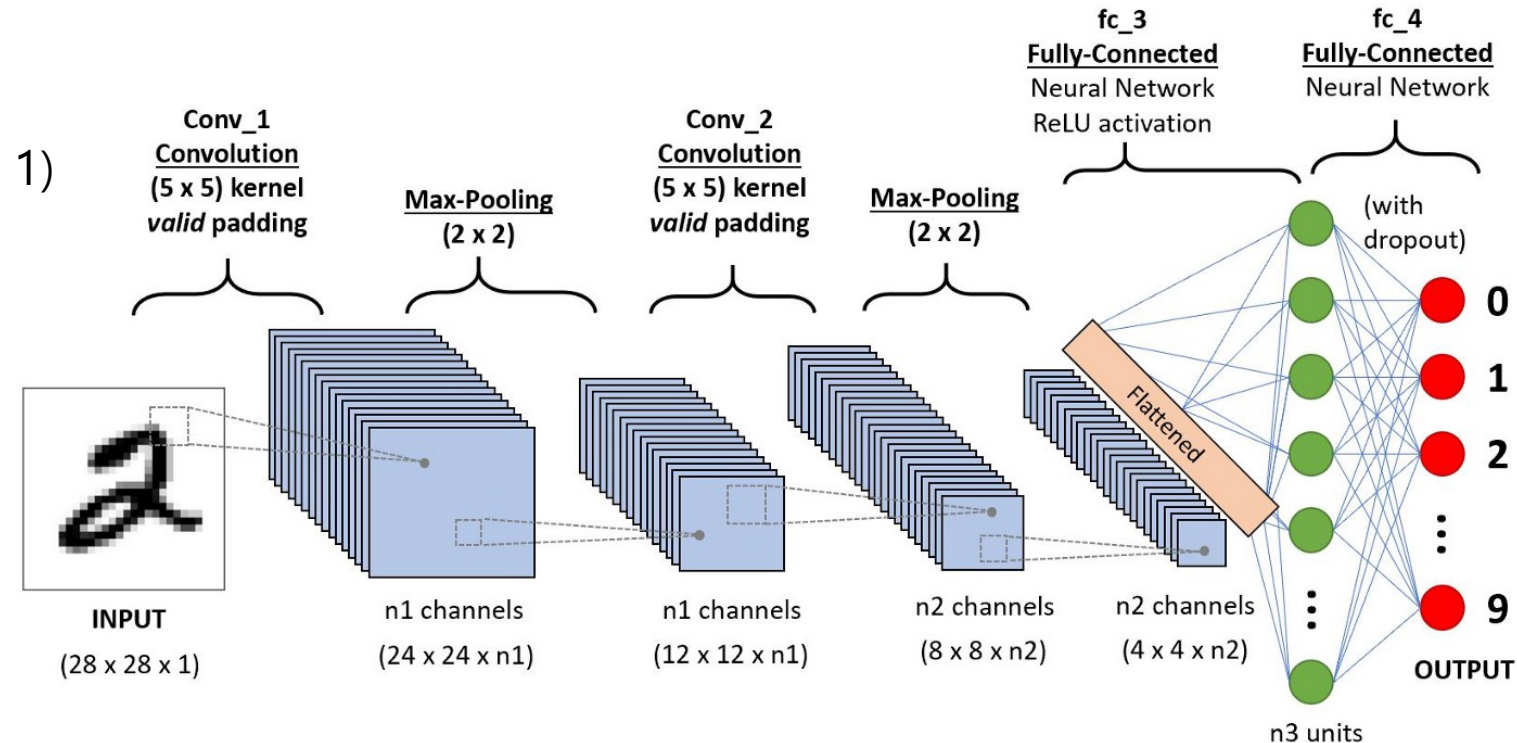
출처: <https://huangdi.tistory.com/36>

# 합성곱 신경망(CNN)

입력 값(Input feature)들 사이의 관계가 있을 때

→ 합성곱 신경망 사용.

→ 특히 이미지 학습에 적합함.





# 합성곱 신경망 활용 예시



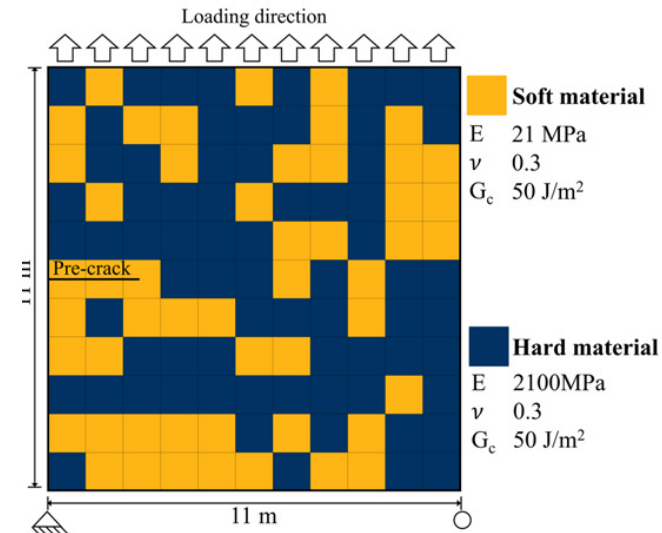
이미지 → RGB값으로 분리



Periodic table

Image

화합물의 조성비



복합재의 패턴

# 합성곱 신경망 활용 예시

## 1. 수치 기반 DNN

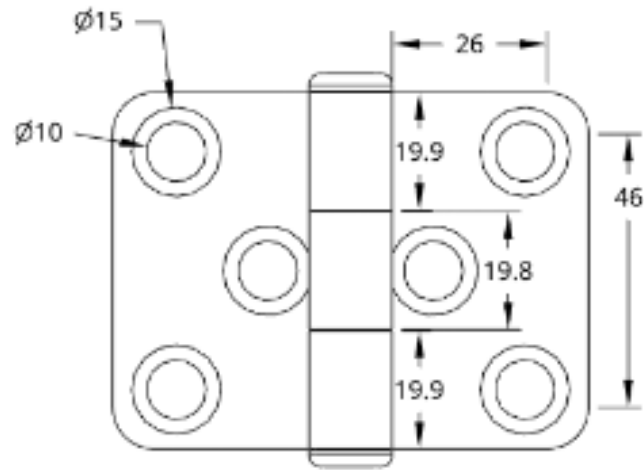
- 입력:

$L_1, L_2, L_3, \dots$

$\phi_1, \phi_2, \phi_3, \dots$

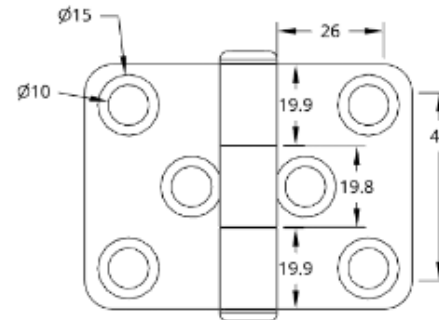
- 출력:

제품 성능



## 2. 이미지 기반 CNN

입력:



출력:

제품성능



# 분류 알고리즘 작동원리

## 의사 결정 나무

# 정보 불순도(Impurity)



항아리 1.



항아리 2.



항아리 3.

항아리1 & 항아리 3 → 순도 100%

항아리2 → 불순도 높다.

Q. 정보 불순도를 숫자로 측정할 수 있을까? 정보 엔트로피, 지니 지수

Q. 항아리2의 불순도를 낮추는 좋은 기준? 색깔, 크기, 모양

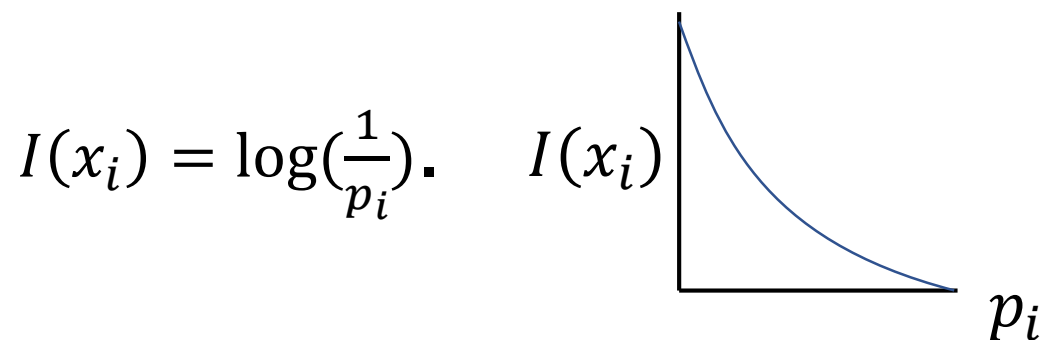
# 정보 엔트로피?

정보 엔트로피(무질서도) =  $\sum_{i=1}^c p_i \times I(x_i)$

$c$ : 사건의 개수

$p_i$ :  $x_i$ 라는 사건이 발생할 확률

$I(x_i)$ :  $x_i$ 라는 사건의 정보량



확률이 낮은 사건일수록 정보량이 높다. 혹은, 흔한 사건은 정보량이 없다.

# 정보 엔트로피?



항아리 2.

붉은 색 뽑을 확률=1/2    붉은 색 뽑는 사건 정보량 =  $\log(2)$

파란 색 뽑을 확률=1/2    파란 색 뽑는 사건 정보량 =  $\log(2)$

정보엔트로피=  $1/2 \times \log(2) + 1/2 \times \log(2) = \log(2)$



항아리 1.

붉은 색 뽑을 확률=0

파란 색 뽑을 확률=1    파란 색 뽑는 사건 정보량 =  $\log(1)=0$

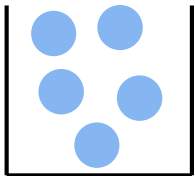
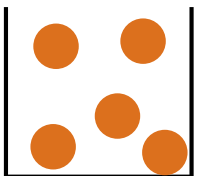
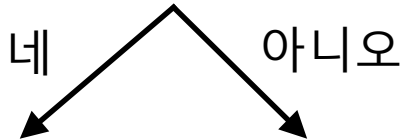
정보엔트로피=  $0 + 1 \times \log(1) = 0$

# 정보 엔트로피?



항아리 2.

붉은색?



붉은 색 뽑을 확률=1/2   붉은 색 뽑는 사건 정보량 =  $\log(2)$

파란 색 뽑을 확률=1/2   파란 색 뽑는 사건 정보량 =  $\log(2)$

정보엔트로피=  $1/2 \times \log(2) + 1/2 \times \log(2) = \log(2)$

정보엔트로피(좌)=  $1 \times \log(1)=0$

정보엔트로피(우)=  $1 \times \log(1)=0$

정보엔트로피(좌+우)=  $0+0=0 \rightarrow$  무질서도를 낮추는 기준

# 정보 엔트로피?

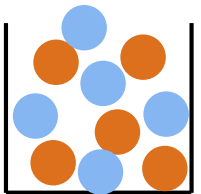


항아리 2.

원형?

네

아니오



붉은 색 뽑을 확률=1/2    붉은 색 뽑는 사건 정보량 =  $\log(2)$

파란 색 뽑을 확률=1/2    파란 색 뽑는 사건 정보량 =  $\log(2)$

정보엔트로피=1/2 x  $\log(2)$  + 1/2x $\log(2)$  =  $\log(2)$

정보엔트로피(좌)= $\log(2)$

정보엔트로피(우)=0

정보엔트로피(좌+우)= $\log(2)$ +0= $\log(2)$  → 무질서도 변화없음

# 지니지수?

$$\text{지니 지수} = 1 - \sum_{i=1}^c p_i^2$$

$c$ : 사건의 개수

$p_i$ :  $x_i$ 라는 사건이 발생할 확률

지니지수가 높을수록 무질서도가 크다.

# 지니지수?



항아리 2.

붉은 색 뽑을 확률=1/2

파란 색 뽑을 확률=1/2

$$\text{지니지수} = 1 - (1/2)^2 - (1/2)^2 = 1 - 0.25 - 0.25 = 0.5$$



항아리 1.

붉은 색 뽑을 확률=0

파란 색 뽑을 확률=1

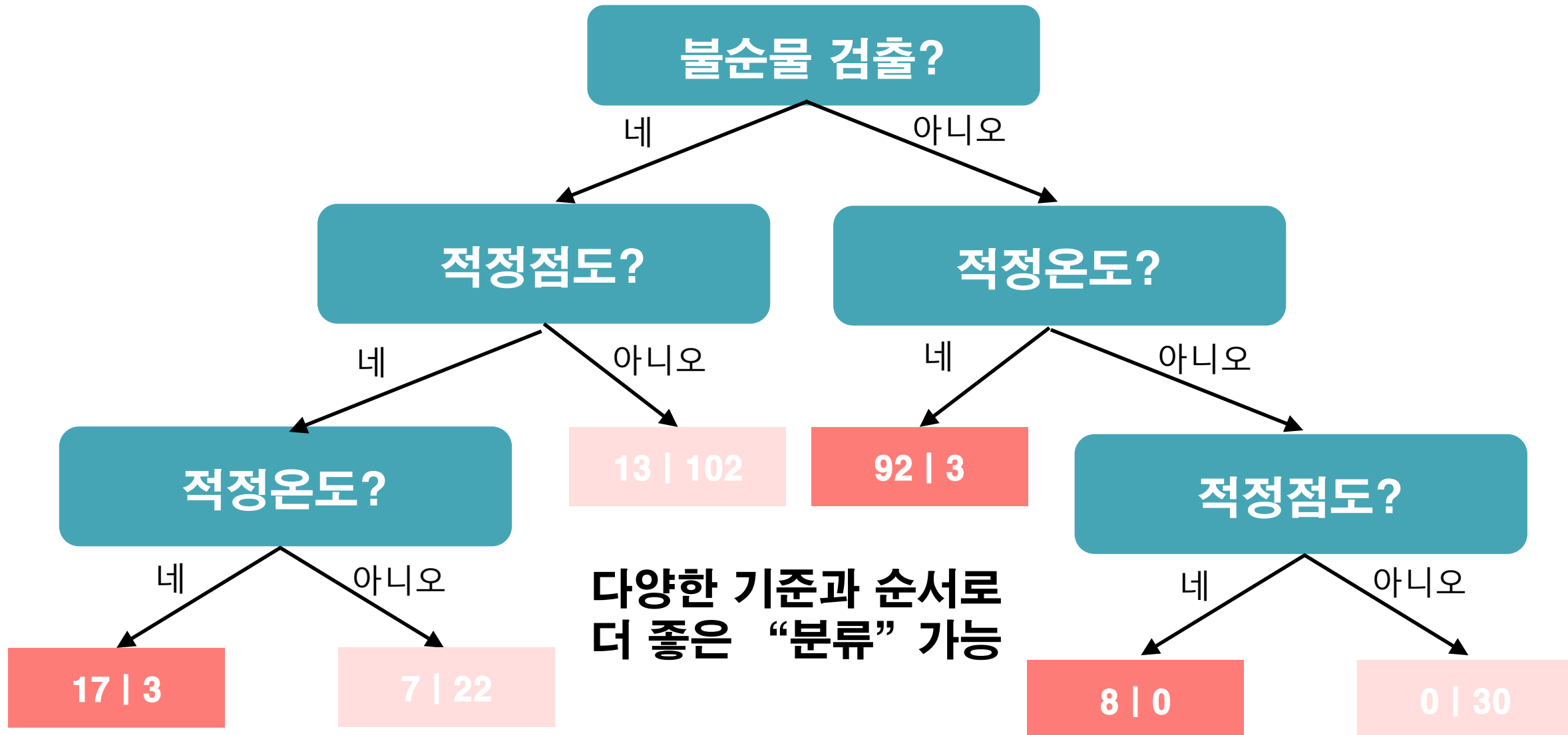
$$\text{지니지수} = 1 - 0^2 - 1^2 = 1 - 0 - 1 = 0$$

정보엔트로피, 지니지수 모두 무질서도가 크면 커지는 수치 !



# 의사 결정 나무란?

- 의사 결정 나무(decision tree) : 특정 기준(질문)에 따라 데이터를 구분하는 모델



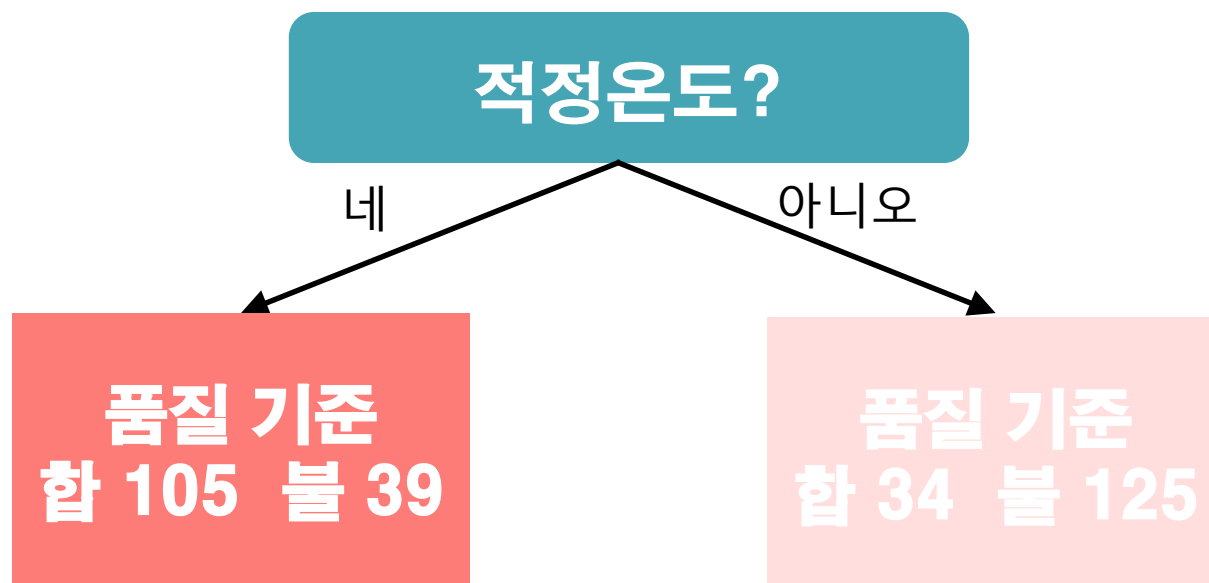
# 의사 결정 나무란?

공정 조건에 따른 제품 품질 기준 통과 여부 예시

적정온도	불순물 검출	적정점도	품질기준
아니오	아니오	아니오	불합격
네	네	네	합격
네	네	아니오	불합격
네	아니오	?	합격
...	...	...	...

# 의사 결정 나무란?

적정온도에 따른 제품 품질 기준 통과 여부

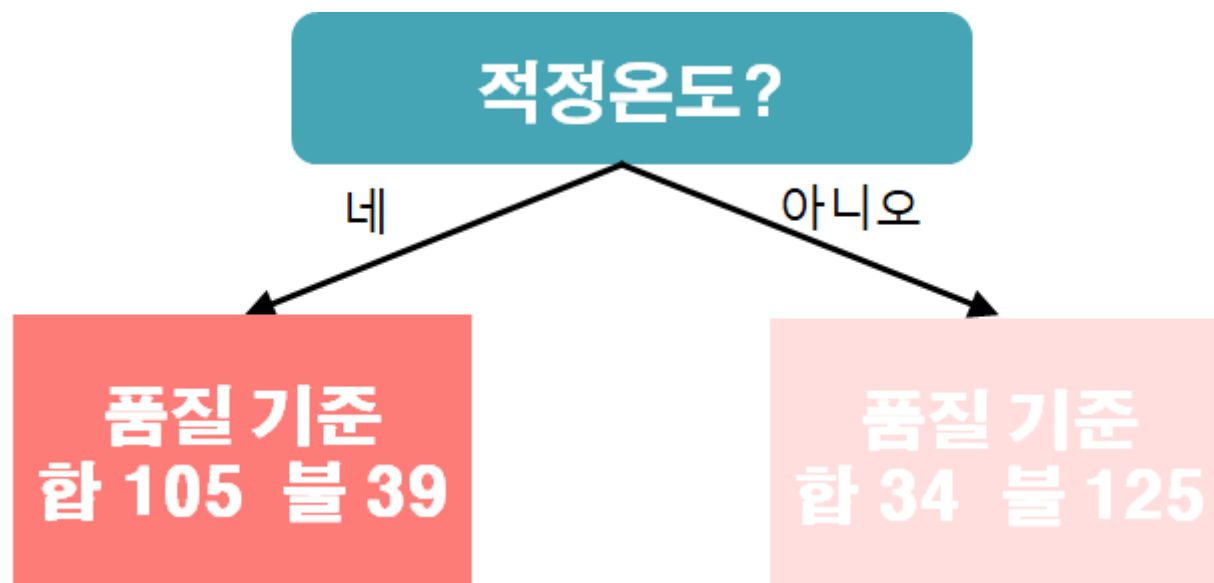


적정온도가 품질 조건을 만족시키는데 얼마나 중요한 요인인지 궁금함.

정보 엔트로피 혹은 지니 지수(Gini index)로 불순도를 최소화하는 기준 탐색!

👉 본 강의는 지니 지수 기준으로 구성하였으나, 정보 엔트로피를 사용해도 유사한 결과 나옴.

# 의사 결정 나무란?



지니 지수(Gini index) 계산법

지니 지수(Gini index) =  $1 - (\text{Yes의 확률})^2 - (\text{No의 확률})^2$

지니 지수 = 그룹1의 비율 × 그룹1의 지니 지수 + 그룹2의 비율 × 그룹2의 지니 지수

→ 지니 지수가 낮을수록 불순도를 낮추는 중요한 구분 기준!

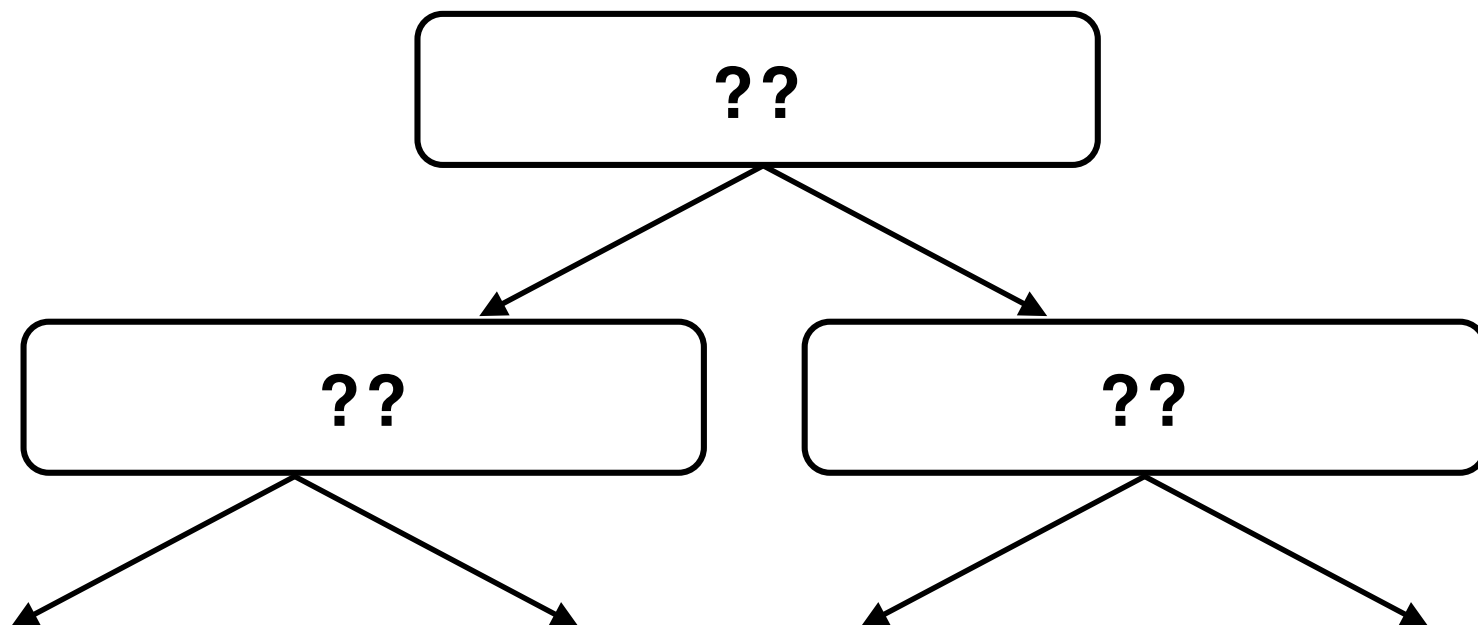
# 의사 결정 나무란?

## 지니 지수(Gini index) 계산 예시

적정온도 지니 지수=0.364

불순물 검출 지니 지수=0.360

적정점도 지니 지수=0.381



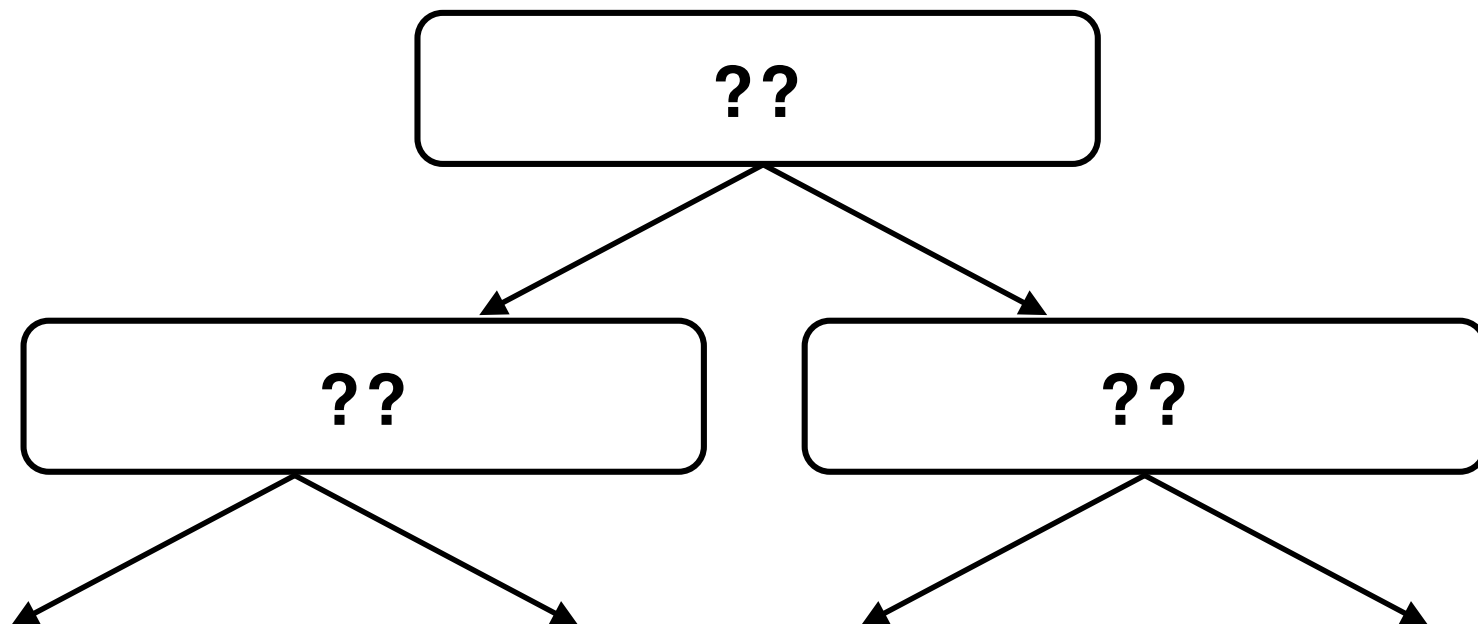
# 의사 결정 나무란?

## 지니 지수(Gini index) 계산 예시

적정온도 지니 지수=0.364

불순물 검출 지니 지수=0.360 → 지니 지수가 가장 낮은 불순물 검출을 의사 결정 나무 상단에 배치

적정점도 지니 지수=0.381



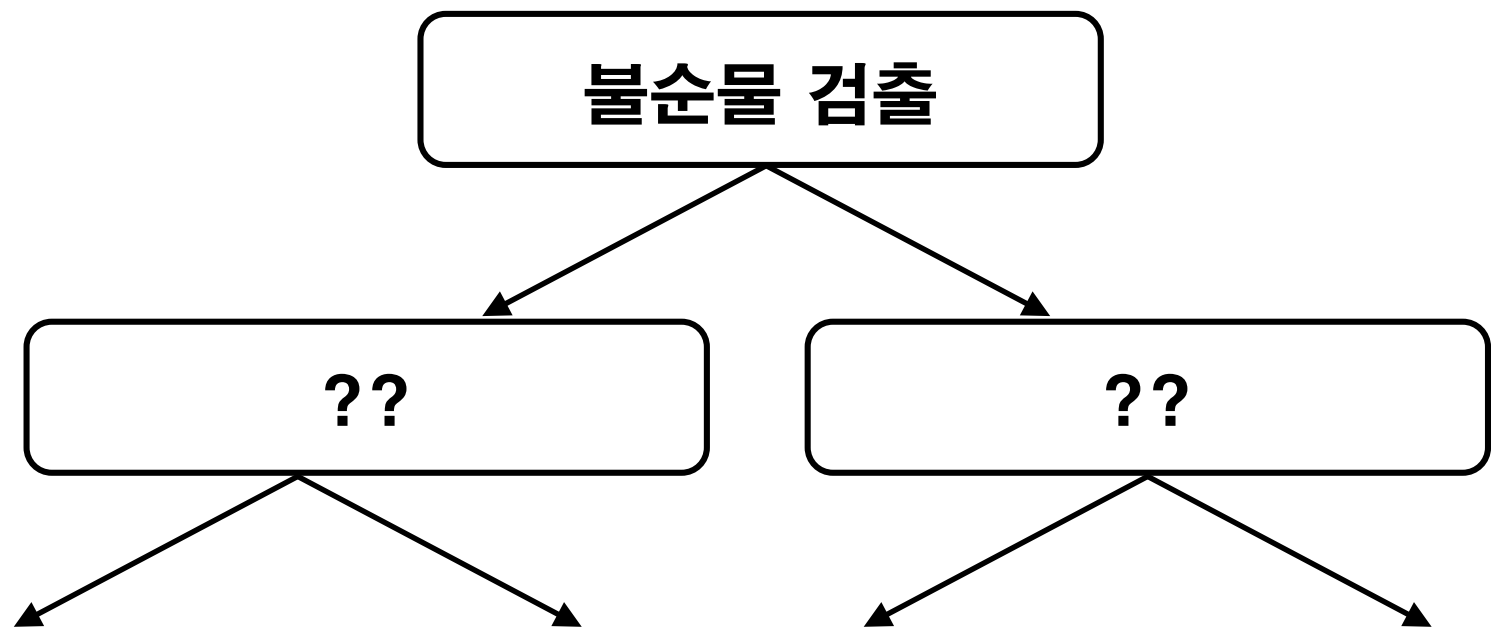
# 의사 결정 나무란?

## 지니 지수(Gini index) 계산 예시

적정온도 지니 지수=0.364

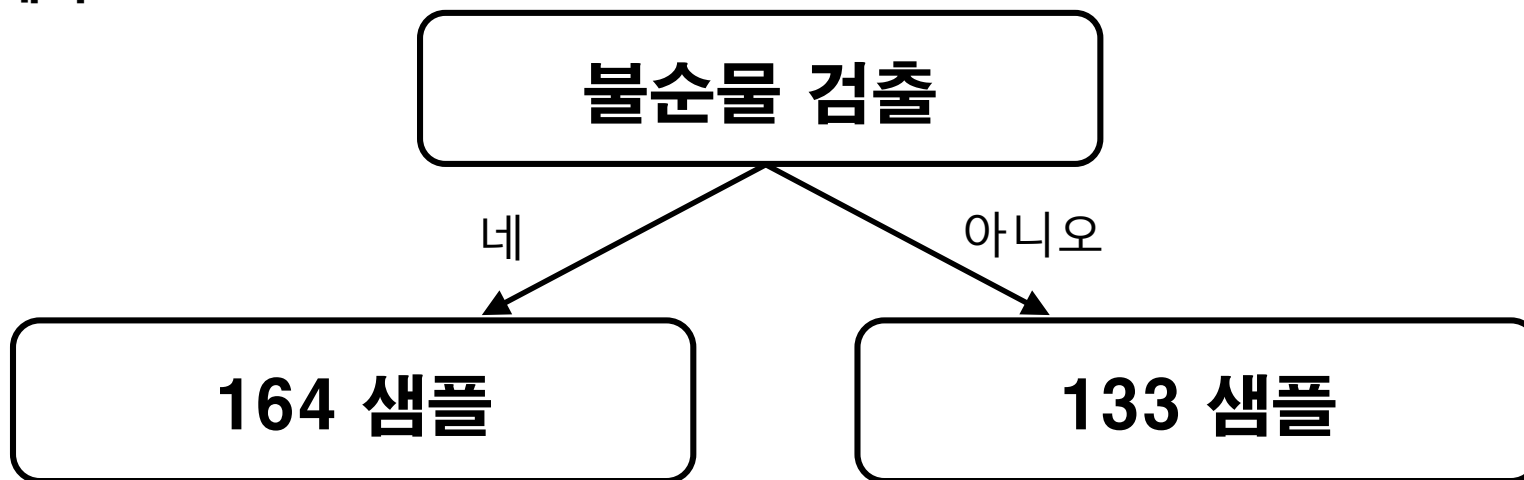
불순물 검출 지니 지수=0.360 → 지니 지수가 가장 낮은 불순물 검출을 의사 결정 나무 상단에 배치

적정점도 지니 지수=0.381



# 의사 결정 나무란?

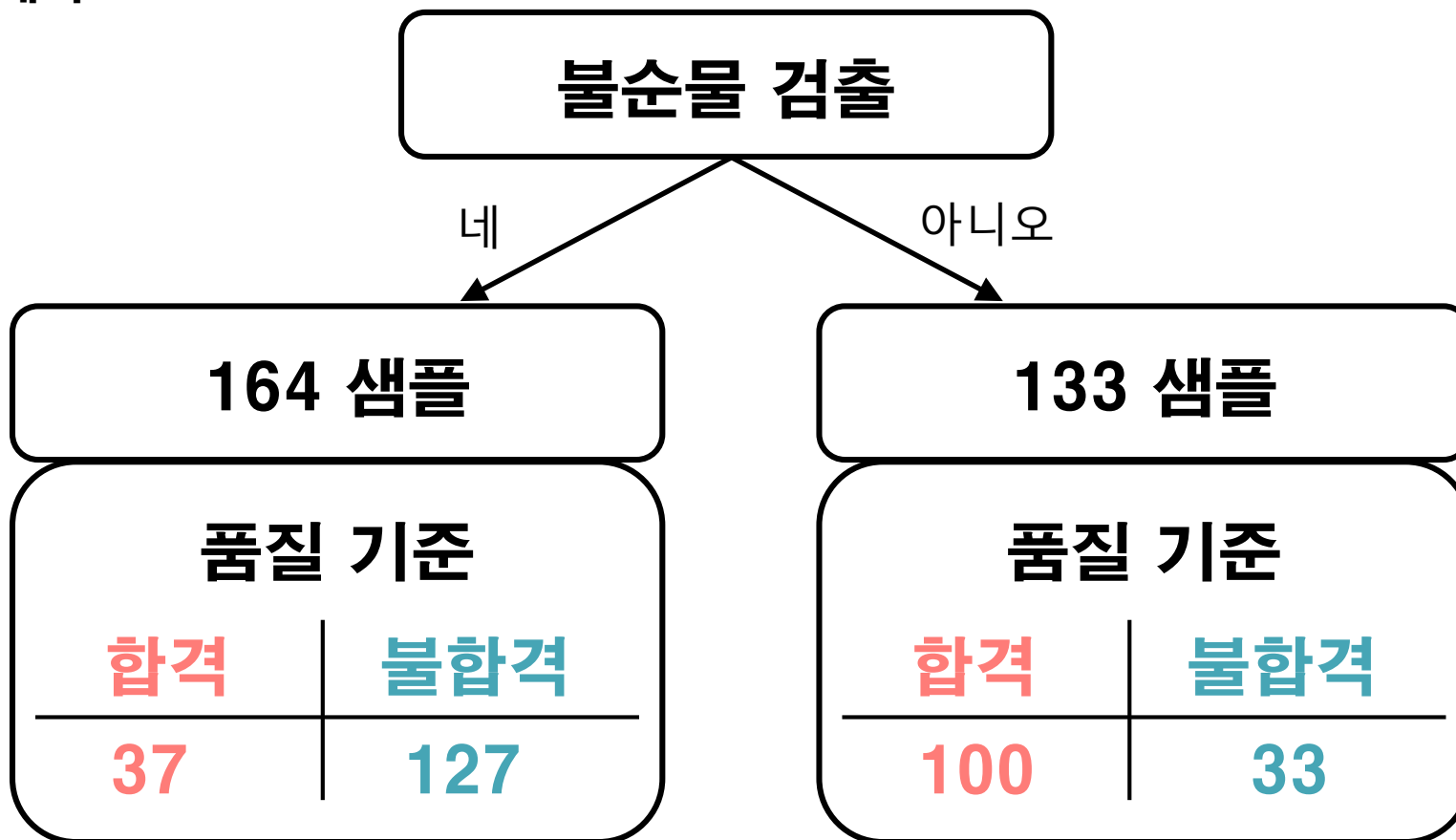
의사 결정 나무 예시





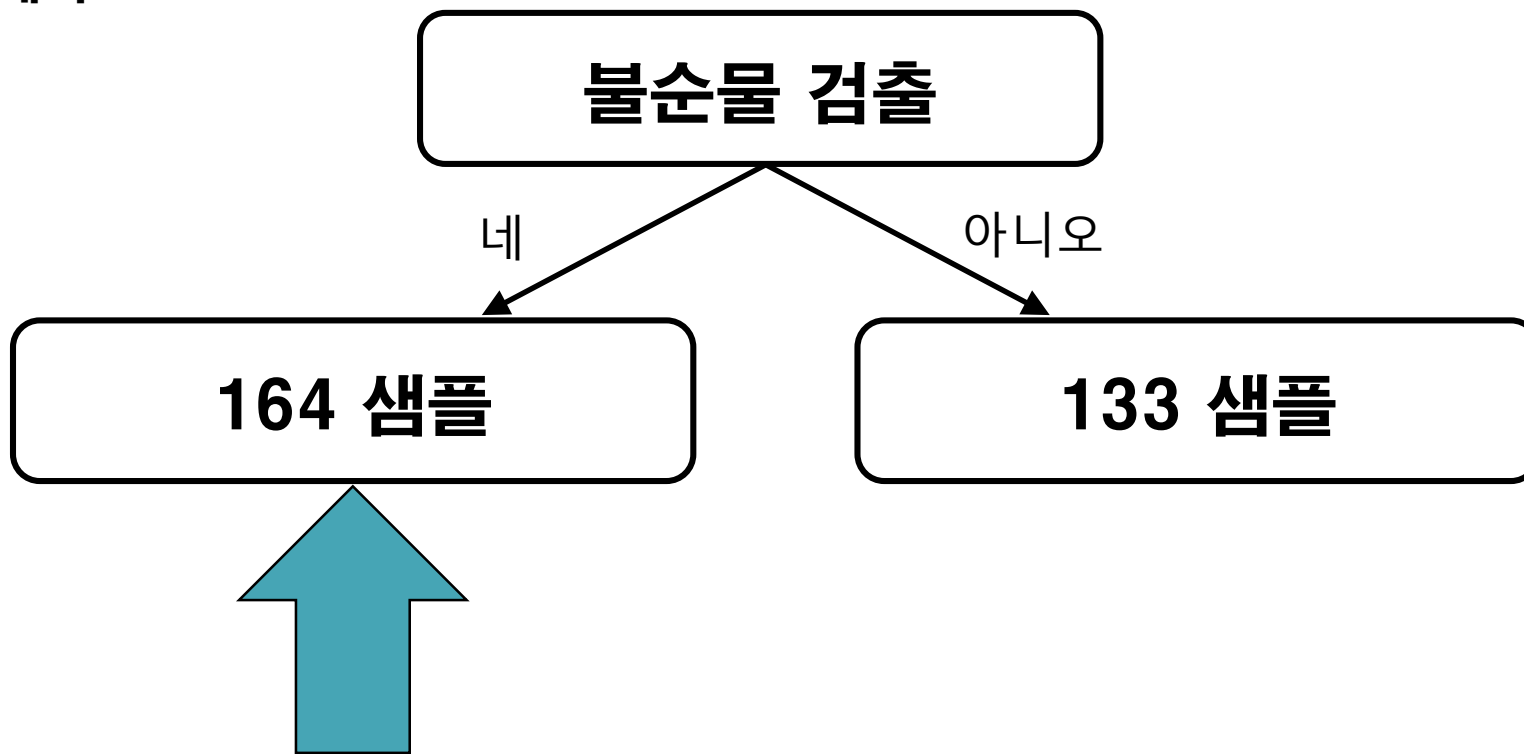
# 의사 결정 나무란?

의사 결정 나무 예시



# 의사 결정 나무란?

의사 결정 나무 예시



위의 샘플을 대상으로 다시 걱정온도와 걱정점도의 지니 지수를 계산

예시

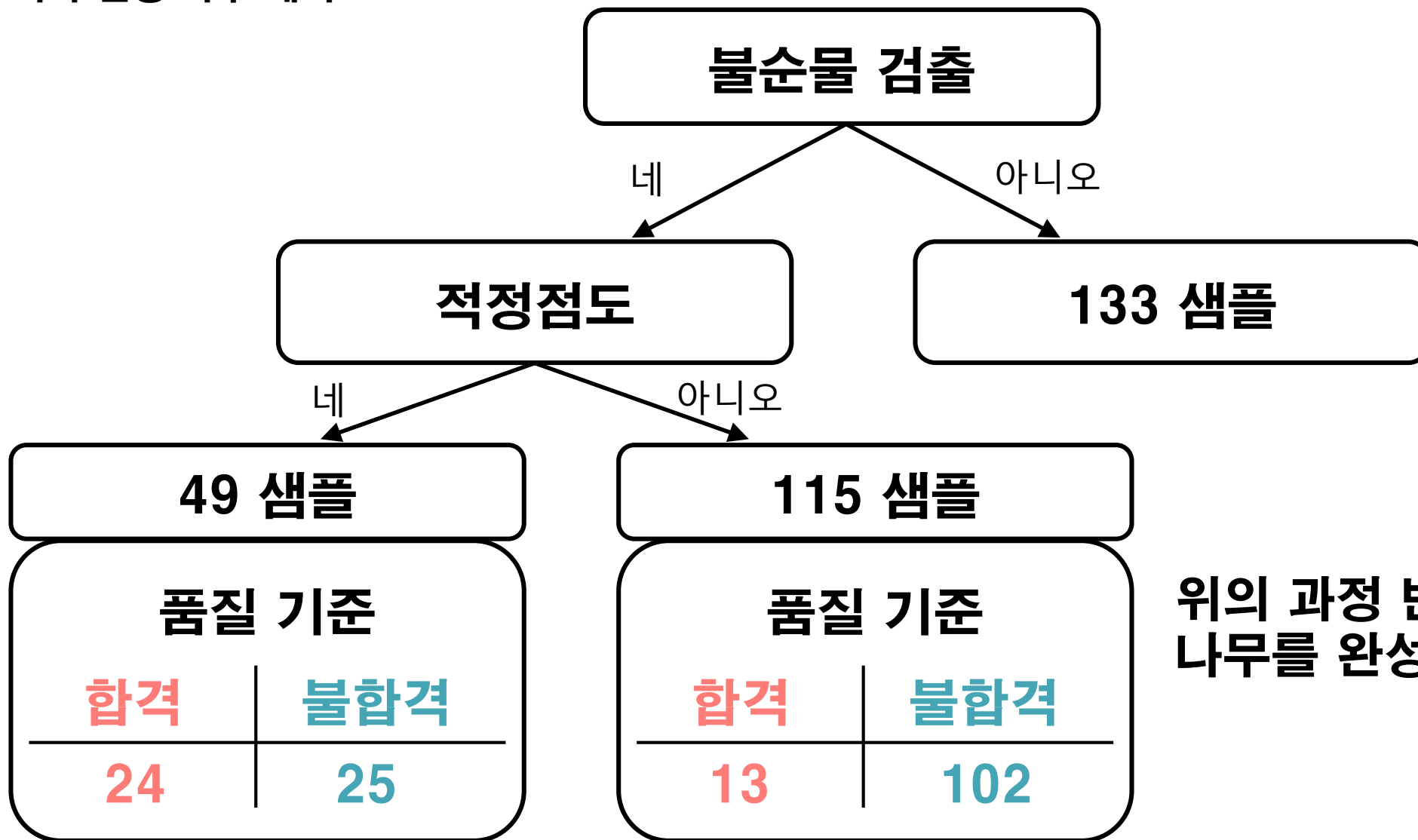
걱정온도 지니 지수=0.3

걱정점도 지니 지수=0.29

← 지니 지수가 낮은 항목 선택

# 의사 결정 나무란?

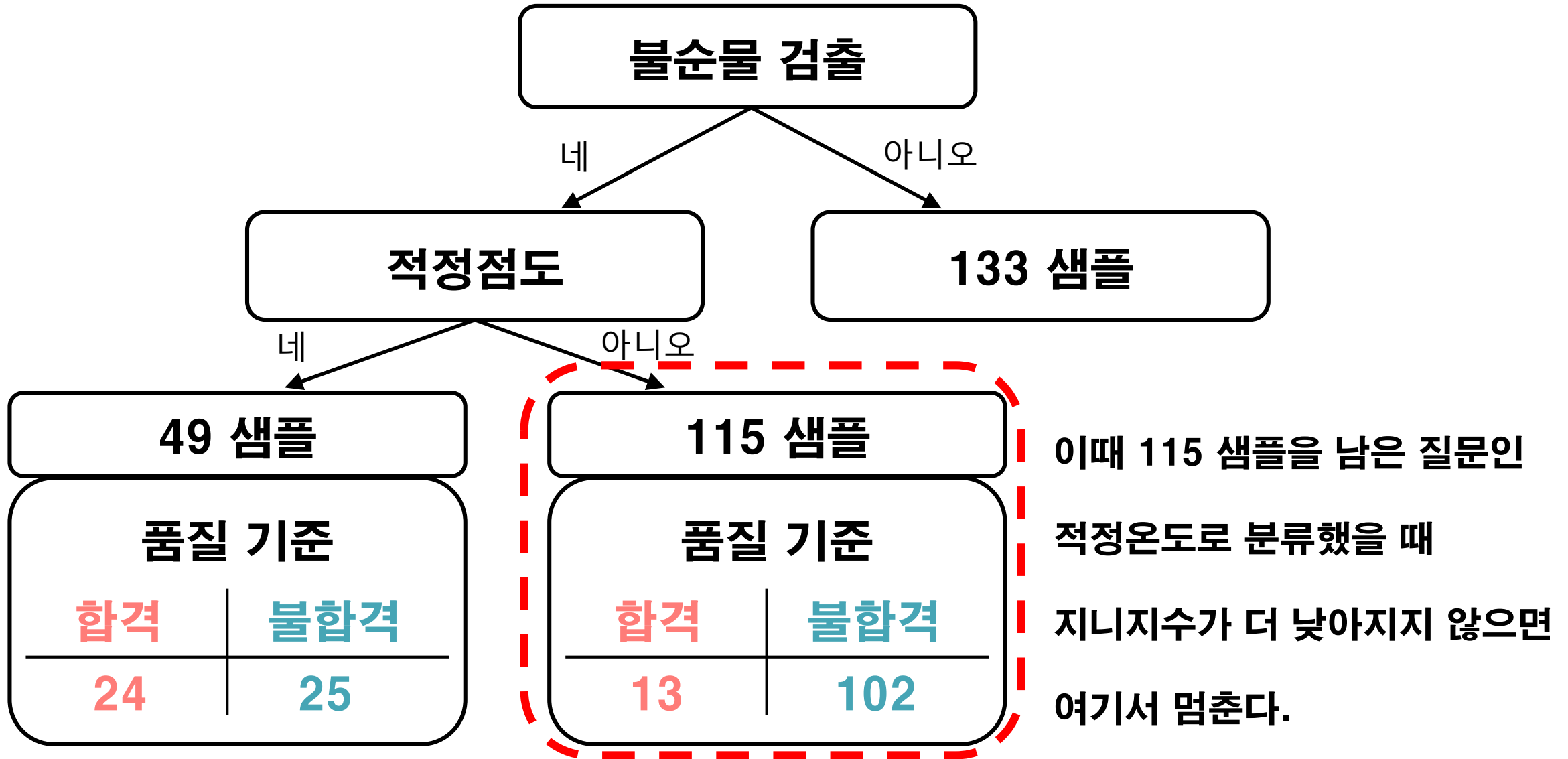
의사 결정 나무 예시



위의 과정 반복하여 의사 결정 나무를 완성

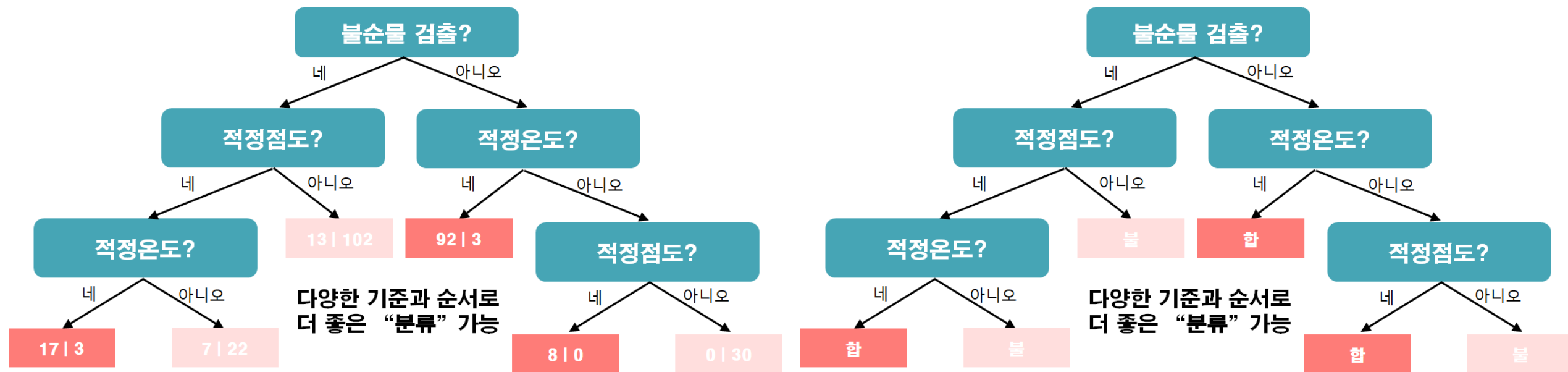
# 의사 결정 나무란?

의사 결정 나무 예시



# 의사 결정 나무란?

## 의사 결정 나무 예시

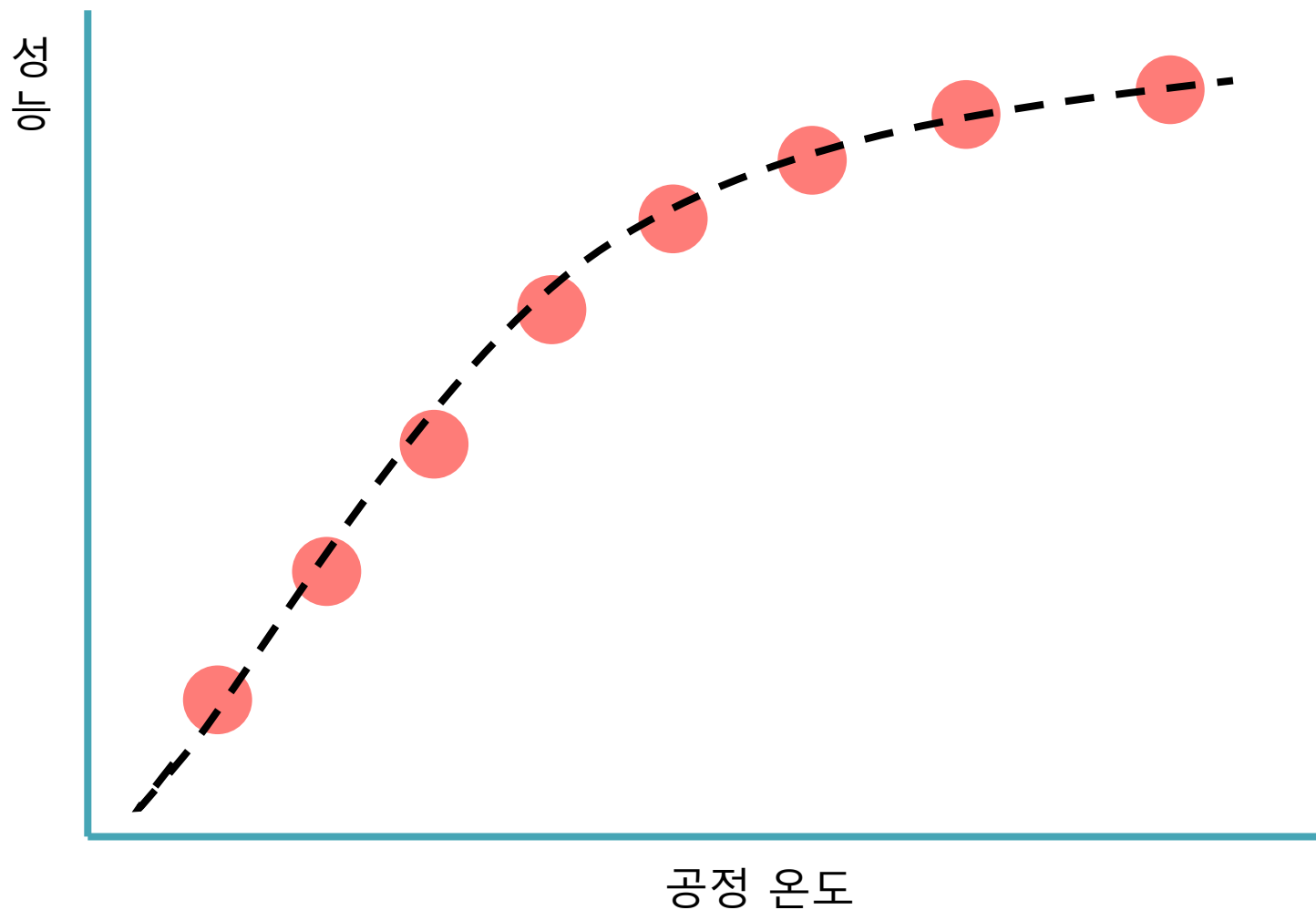


# 인공지능 활용 시 유의점?

적은 오차 데이터 확보

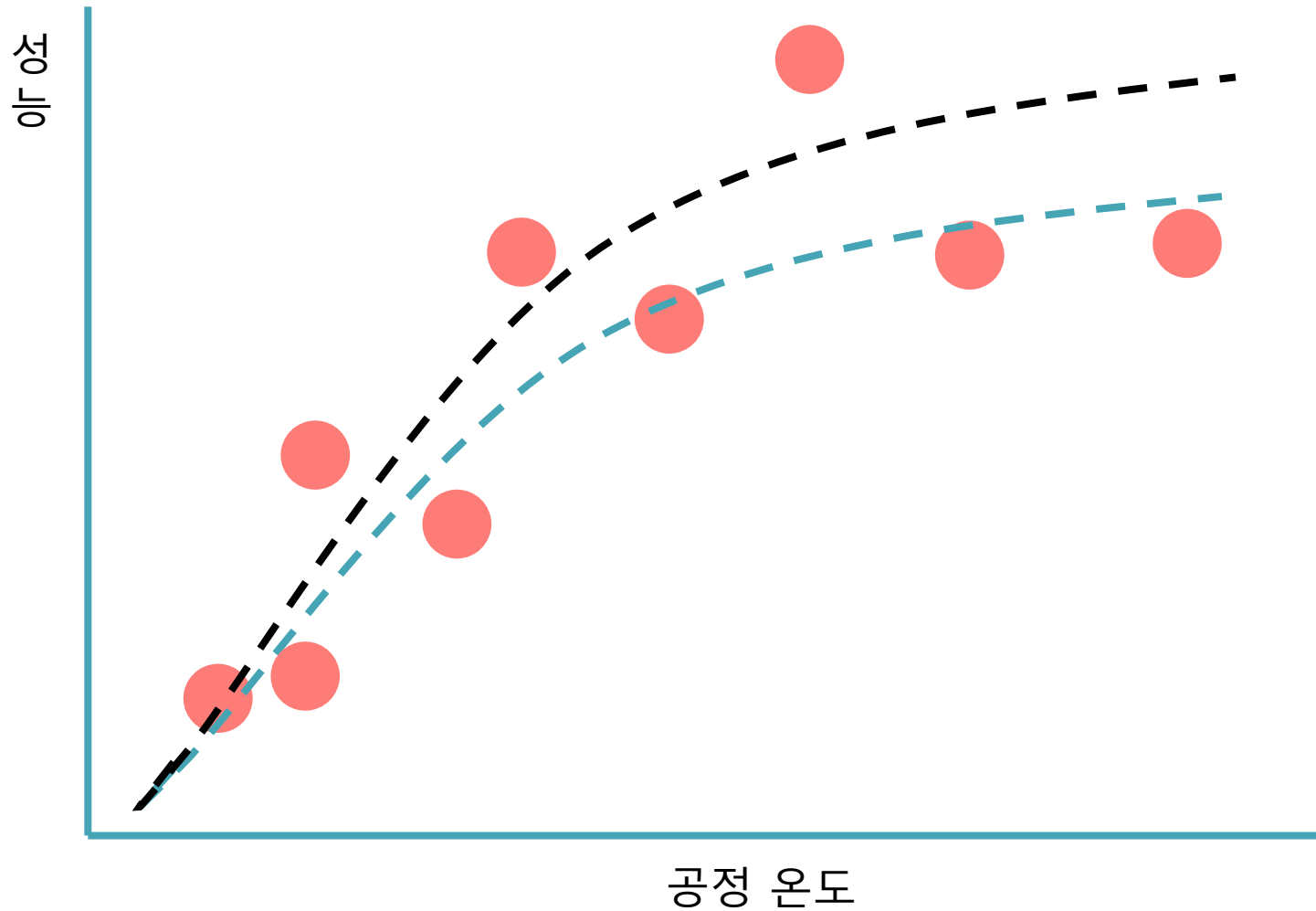
데이터 불균형 유의

# 적은 오차 데이터의 중요성



노이즈 없는  
좋은 데이터:  
학습된 모델의  
예측력 우수

# 적은 오차 데이터의 중요성



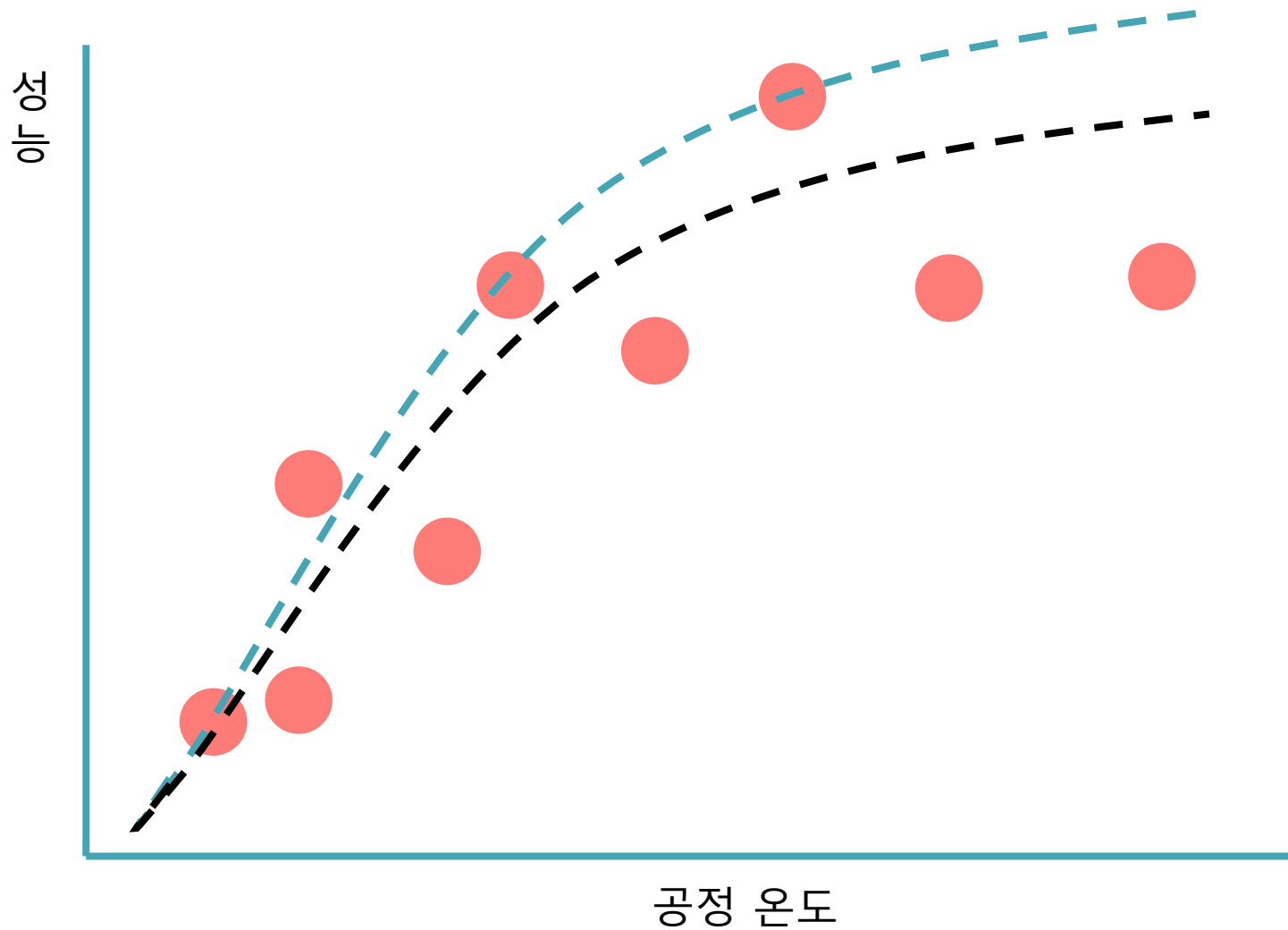
노이즈 많은 경우

학습 어려움

예측 어려움



# 적은 오차 데이터의 중요성

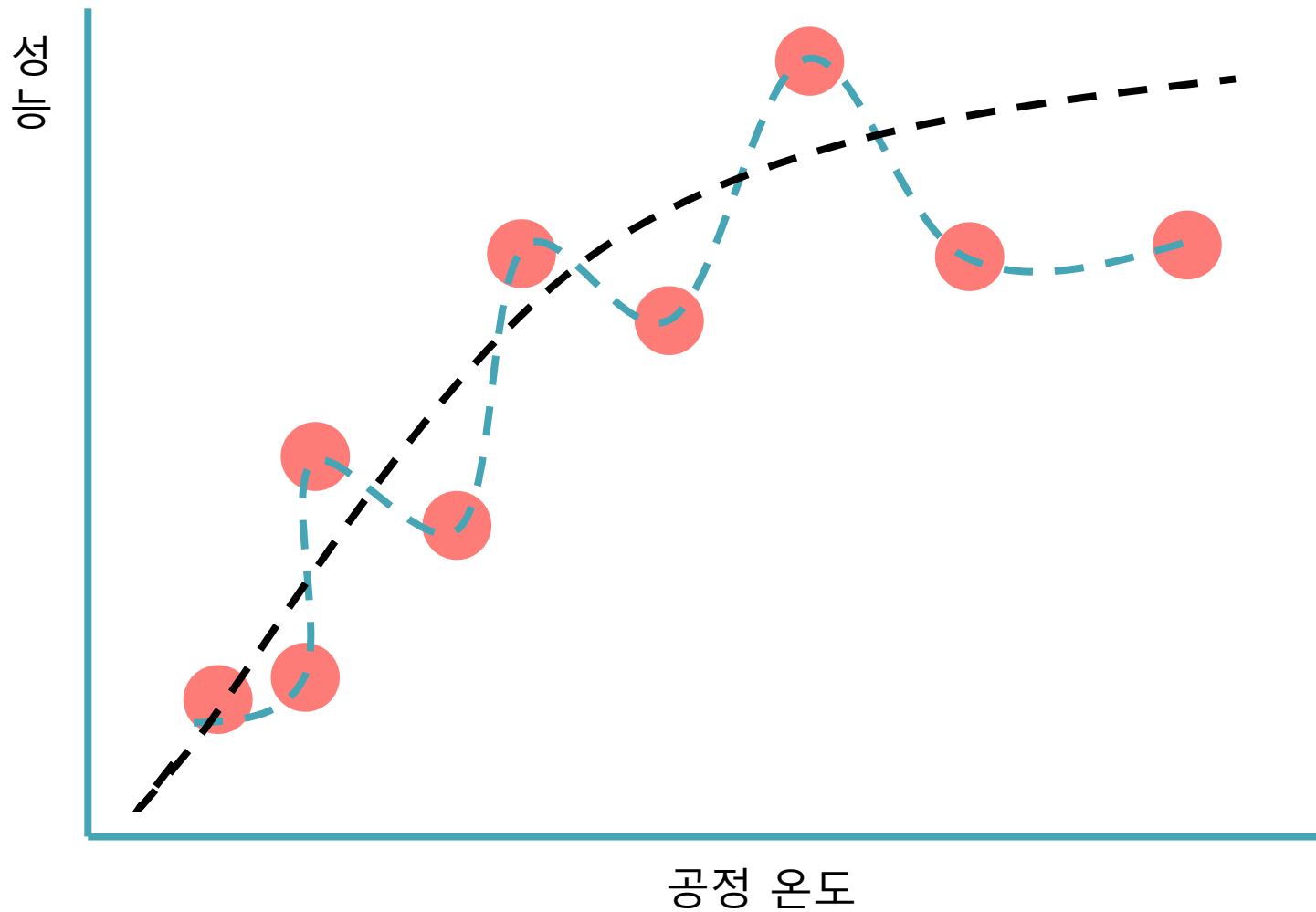


노이즈 많은 경우

학습 어려움

예측 어려움

# 적은 오차 데이터의 중요성



노이즈 많은 경우

학습 어려움

예측 어려움

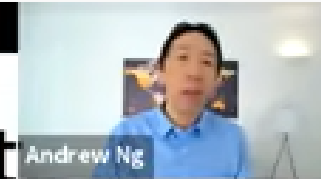
# Andrew Ng 교수님 Comments

## Improving the code vs. the data

	Steel defect detection	Solar panel	Surface inspection
Baseline	76.2%	75.68%	85.05%
Model-centric	+0% (76.2%)	+0.04% (75.72%)	+0.00% (85.05%)
Data-centric	+16.9% (93.1%)	+3.06% (78.74%)	+0.4% (85.45%)

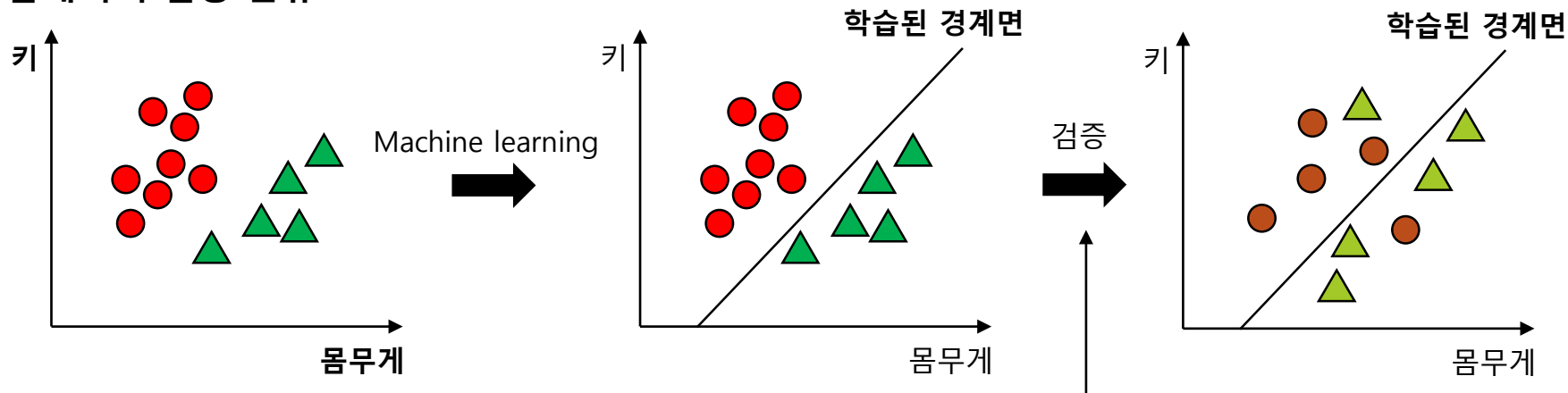
<https://www.youtube.com/watch?v=06-AZXmwHjo&t=1536s>

Google에서 Andrew Ng MLOPS 검색



# 데이터 균형의 중요성

## 학습데이터 활용 분류



● : 건강(Train)  
▲ : 고혈압(Train)

- 학습되지 않은 검증 데이터

Data #	키	몸무게	상태
1	163cm	89kg	고혈압 ▲
⋮	⋮	⋮	⋮
100	182cm	67kg	건강 ●

● : 건강(Test)  
▲ : 고혈압(Test)

## 균형 잡힌 검증 데이터셋

건강 : 50명 → Test { 건강 : 40, 고혈압 : 10 }  
고혈압 : 50명 → Test { 건강 : 10, 고혈압 : 40 }

$$\text{Accuracy} = \frac{40+40}{40+10+40+10} = 80\%$$

## 불균형 검증 데이터셋

건강 : 90명 → Test { 건강 : 90, 고혈압 : 0 }  
고혈압 : 10명 → Test { 건강 : 10, 고혈압 : 0 }

$$\text{Accuracy} = \frac{90}{90+10} = 90\% \text{ Better ?}$$

<

# 분류 성능 판단

## - 혼동행렬

True Positive(TP) = 1을 1로 잘 예측함

False Negative(FN) = 1을 0으로 잘못 예측함

False Positive(FP) = 0을 1로 잘못 예측함

True Negative(TN) = 0을 0으로 잘 예측함

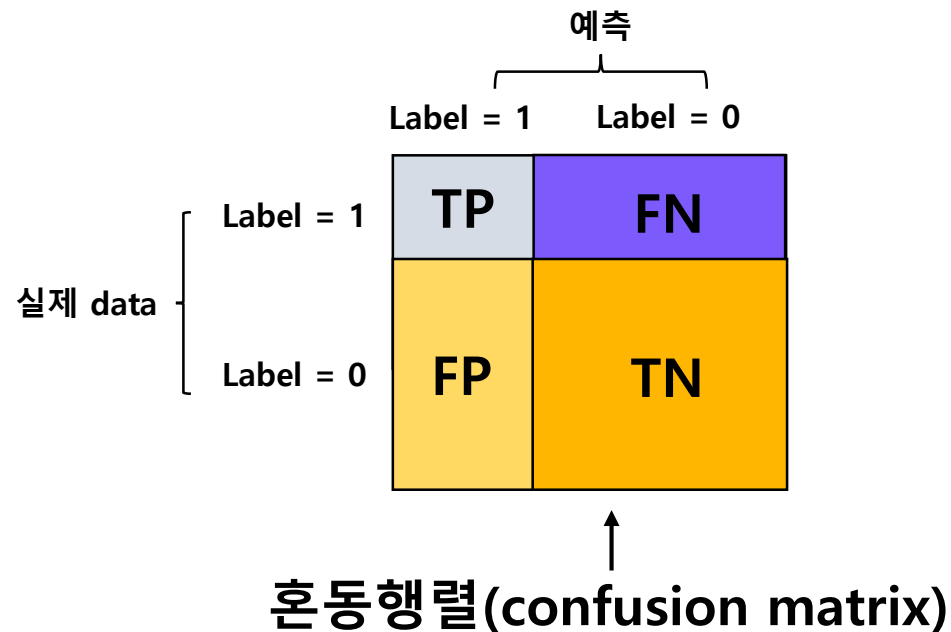
1 : Positive

예측을 뭐로 했나?

0 : Negative

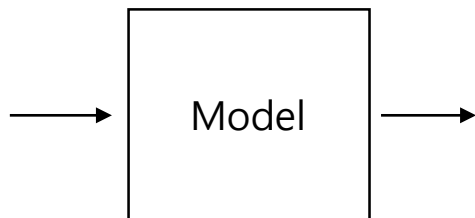
False Positive

실제와 예측이 같나?



Ex) 100개의 공정데이터

다수  
90 : 정상  
소수  
10 : 불량



정상 90개 중 80개 맞게 판단  
불량 10개 중 3개 맞게 판단

혼동행렬

	예측 불량	예측 정상
실제 불량	3	7
실제 정상	10	80

# 분류 성능 판단

## - 혼동행렬

		예측	
		Label = 1	Label = 0
실제 data	Label = 1	TP	FN
	Label = 0	FP	TN

예측을 뭐로 했나?

False Positive

실제와 예측이 같나?

인공지능 모델 1번	예측 불량	예측 정상
실제 불량	3	7
실제 정상	10	80

$$\text{정확도} = \frac{83}{100} = 83\%$$

$$\text{민감도} = \frac{3}{10} = 30\%$$

$$\text{특이도} = \frac{80}{90} = 88.9\%$$

인공지능 모델 2번	예측 불량	예측 정상
실제 불량	10	0
실제 정상	17	73

$$\text{정확도} = \frac{83}{100} = 83\%$$

$$\text{민감도} = \frac{10}{10} = 100\%$$

$$\text{특이도} = \frac{73}{90} = 81.1\%$$

인공지능 모델 3번	예측 불량	예측 정상
실제 불량	1	9
실제 정상	0	90

$$\text{정확도} = \frac{91}{100} = 91\%$$

$$\text{민감도} = \frac{1}{10} = 10\%$$

$$\text{특이도} = \frac{90}{90} = 100\%$$

## - 민감도(Sensitivity) or 재현율(Recall)

$$\text{Sensitivity} = \frac{TP}{TP+FN} : \text{실제 1중에 1이라 예측한 비율}$$

## - 특이도(specificity)

$$\text{Specificity} = \frac{TN}{FP+TN} : \text{실제 0중에 0이라 예측한 비율}$$

## - 정확도 : weighted average of recall & Specificity

$$\begin{aligned} &0 \text{의 개수} : 1 \text{의 개수} \\ &= 9 : 1 \end{aligned} \longrightarrow p = \frac{1}{9+1} = 0.1$$

$$\text{정확도} = p * \text{민감도} + (1-p) * \text{특이도}$$

# 분류 성능 판단

## - 혼동행렬

		예측	
		Label = 1	Label = 0
실제 data	Label = 1	TP	FN
	Label = 0	FP	TN

예측을 뭐로 했나?

False Positive

실제와 예측이 같나?

인공지능 모델 1번	예측 불량	예측 정상
실제 불량	3	7
실제 정상	10	80

$$\text{정확도} = \frac{83}{100} = 83\%$$

$$\text{민감도} = \frac{3}{10} = 30\%$$

$$\text{특이도} = \frac{80}{90} = 88.9\%$$

$$\text{정밀도} = \frac{3}{13} = 23.1\%$$

F1 점수 = 26 %

인공지능 모델 2번	예측 불량	예측 정상
실제 불량	10	0
실제 정상	17	73

$$\text{정확도} = \frac{83}{100} = 83\%$$

$$\text{민감도} = \frac{10}{10} = 100\%$$

$$\text{특이도} = \frac{73}{90} = 81.1\%$$

$$\text{정밀도} = \frac{10}{27} = 37.0\%$$

F1 점수 = 54 %

인공지능 모델 3번	예측 불량	예측 정상
실제 불량	1	9
실제 정상	0	90

$$\text{정확도} = \frac{91}{100} = 91\%$$

$$\text{민감도} = \frac{1}{10} = 10\%$$

$$\text{특이도} = \frac{90}{90} = 100\%$$

$$\text{정밀도} = \frac{1}{1} = 100\%$$

F1 점수 = 18 %

## - 민감도(Sensitivity) or 재현율(Recall)

$$\text{Sensitivity} = \frac{TP}{TP+FN} : \text{실제 1중에 1이라 예측한 비율}$$

## - 정밀도(precision)

$$\text{precision} = \frac{TP}{TP+FP} : \text{예측으로 1중에 실제 1인 비율}$$

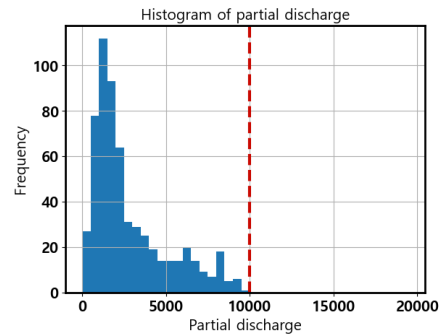
$$\text{F1 점수} = 2 \frac{\text{정밀도} * \text{민감도}}{\text{정밀도} + \text{민감도}}$$

# 회귀 분석 시 데이터 불균형 보완 사례

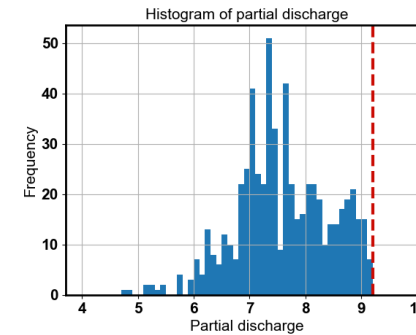
$$\ln(y) \rightarrow \text{Prediction } \hat{y} \rightarrow \exp(\hat{y})$$



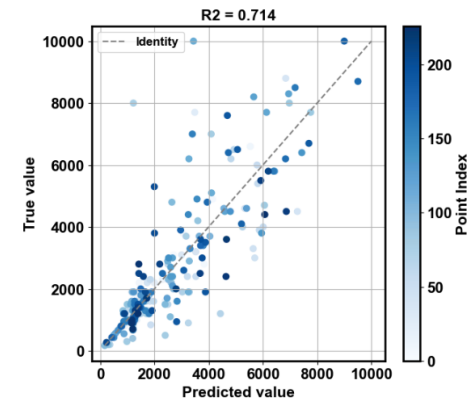
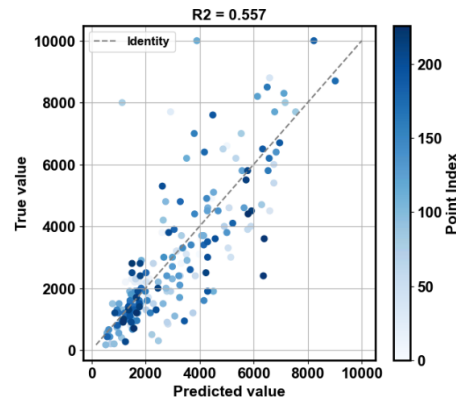
학습 데이터셋의 histogram



학습 데이터셋의 histogram



간단한 log 처리만으로도  
데이터 분포가 더 균일해짐.



학습 성능  $R^2$  상승!



좋은 데이터셋 구성을 위한 요소

도메인 전문가와 주요 입력변수 선정

표준화된 포맷

적은 노이즈

성공/실패 사례 모두 포함한 균형 데이터