

# 6장. 머신러닝(AI)



# 머신러닝(Machine Learning)

## ◆ 머신러닝이란?

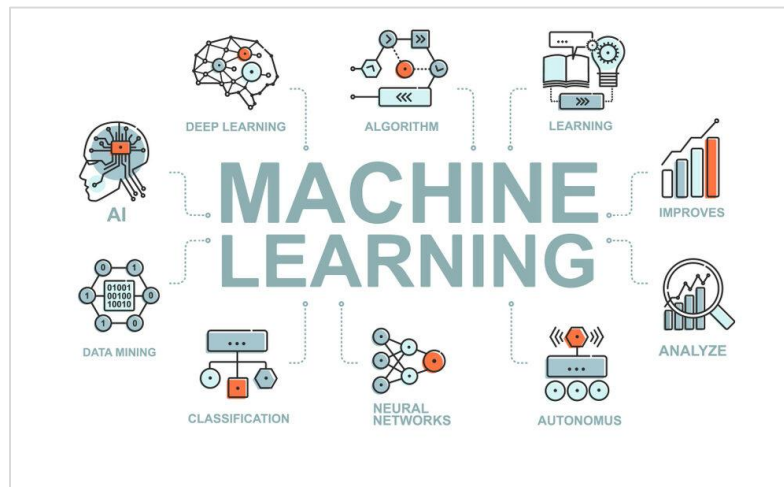
컴퓨터가 스스로 방대한 양의 데이터를 학습하고 지식을 습득하여 문제를 해결하는 기술이다.

컴퓨터가 학습할 때는 데이터와 함께 머신러닝 알고리즘이 필요하다.

머신러닝의 알고리즘의 종류로는 '지도 학습', '비지도학습', '강화 학습'이 있다.

**알고리즘(Algorithm)**

어떤 문제를 해결하기 위한 절차나 방법



# 머신러닝의 학습 방법

## ◆ 지도학습 / 비지도학습

	지도 학습(supervised)	비지도학습(unsupervised)
분석 모형 (알고리즘)	<ul style="list-style-type: none"><li>회귀</li><li>분류</li></ul>	<ul style="list-style-type: none"><li>군집</li><li>차원 축소</li></ul>
특징	정답(레이블)을 알고 있는 상태에서 학습	정답(레이블)이 없는 상태에서 서로 비슷한 데이터를 찾아서 그룹화(클러스터링)

# 지도학습(Supervised Learning)

## ➤ 분류(Classification)

주어진 데이터의 특성을 학습해 새로운 데이터가 어떤 범주(클래스)에 속하는지 예측함

- 이진분류 : 스팸/일반 메일 구분, 합격/불합격,  
CT사진에서 폐암이 보이는지 아닌지 예측
- 다중 분류 : 다양한 새의 종류, 지폐의 종류 등

## ➤ 회귀(Regression)

연속적인 값을 예측하는 것으로 다양한 값을 예측

- 서울의 집값 예측  
특징 데이터(방의 개수, 대중교통 접근성, 범죄율, 집값 등),  
정답 데이터는 집값

# 비지도학습(Unsupervised Learning)

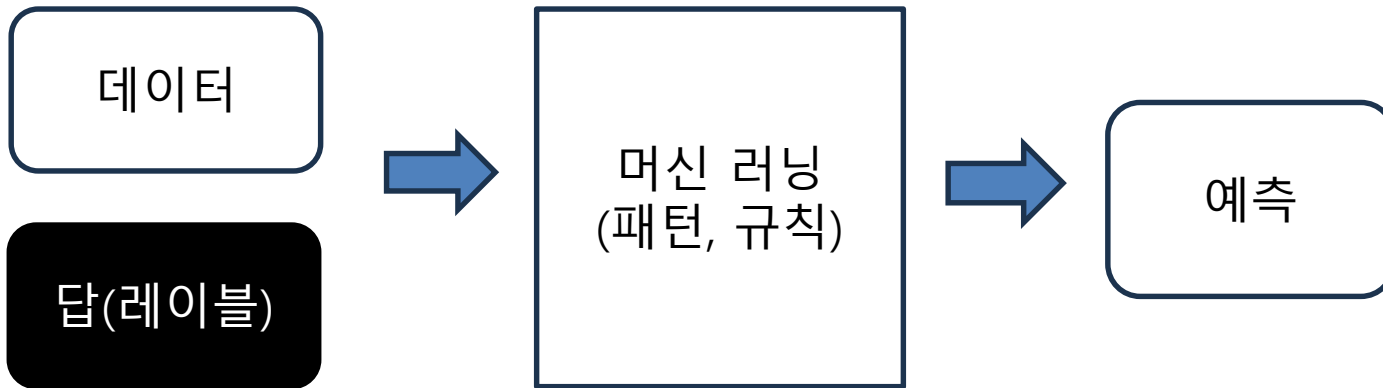
## ➤ 군집화(clustering)

서로 비슷한 데이터를 찾아 그룹화 하는 과정

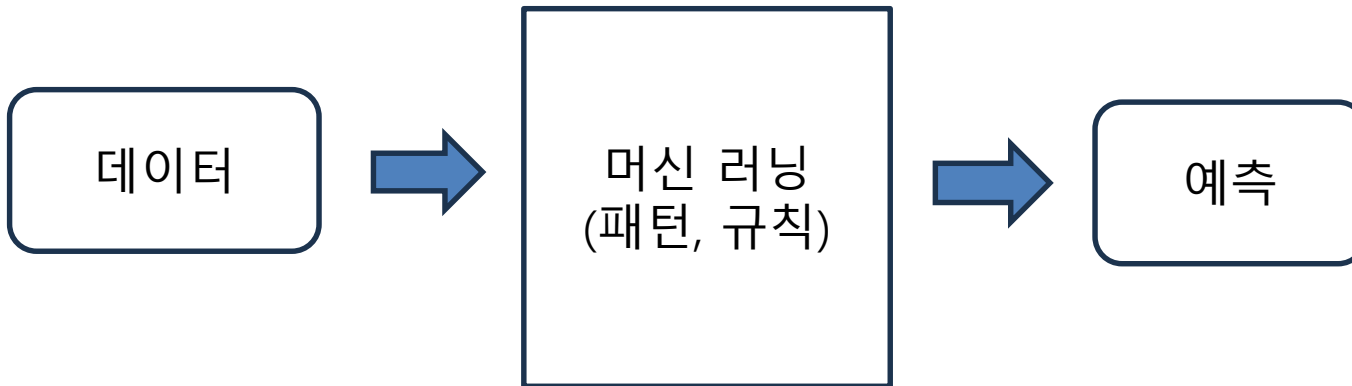
- 유튜브 : 추천 영상

- 네이버, 카카오 : 사람이 구매한 내역을 보고 그 사람들을 여러 그룹으로 나눔. 알맞은 콘텐츠 제공

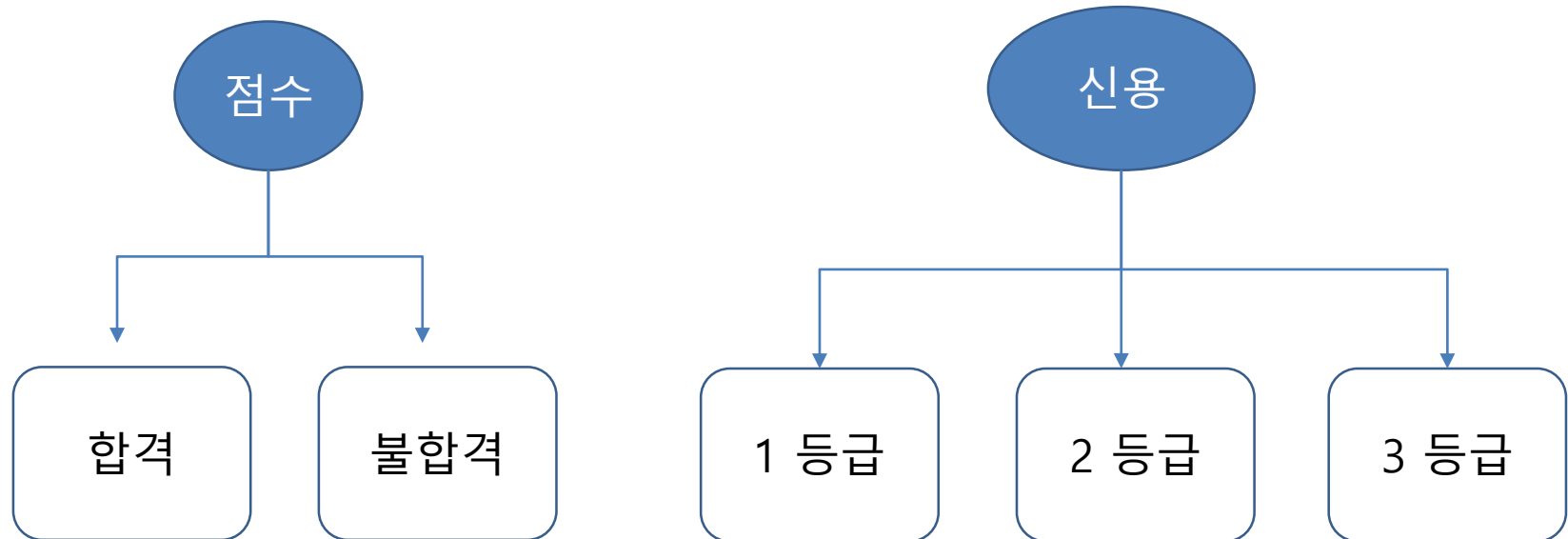
# 지도학습(Supervised Learning)



# 비지도학습(Unsupervised Learning)

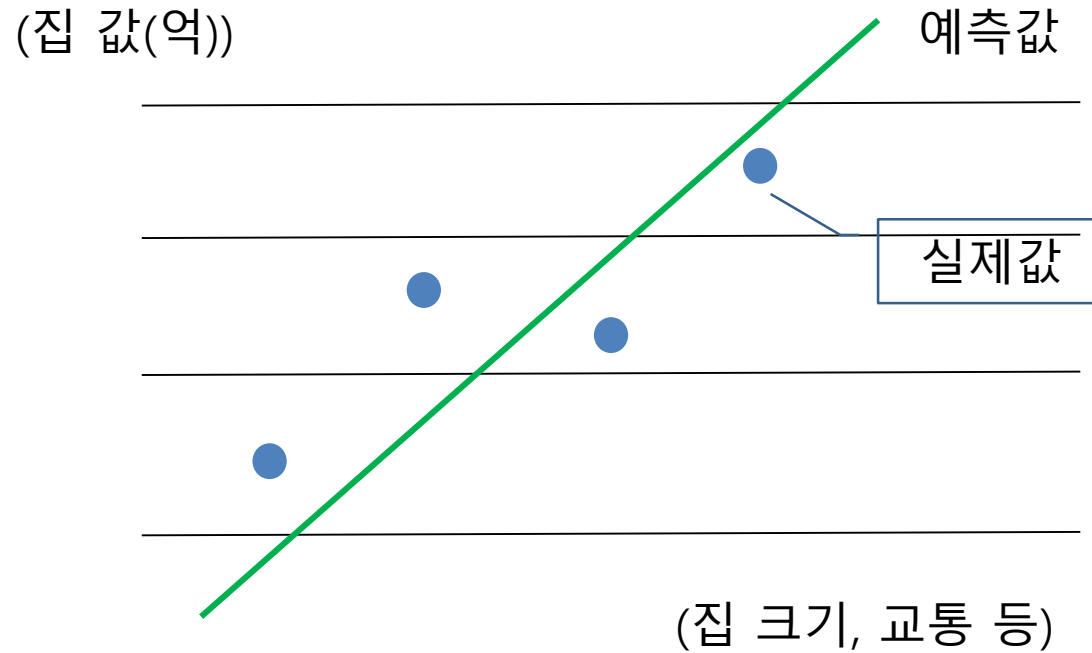


# 분류(Classification)





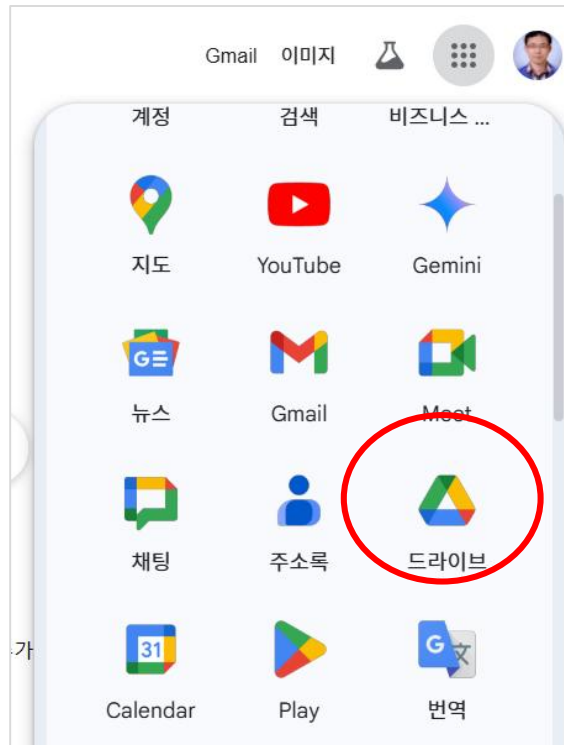
# 회귀(Regression)



# 구글 코랩(Colab)

## ❖ 구글 코랩(Colab)

구글에서 제공하는 클라우드 기반의 Jupyter 노트북 환경이다.



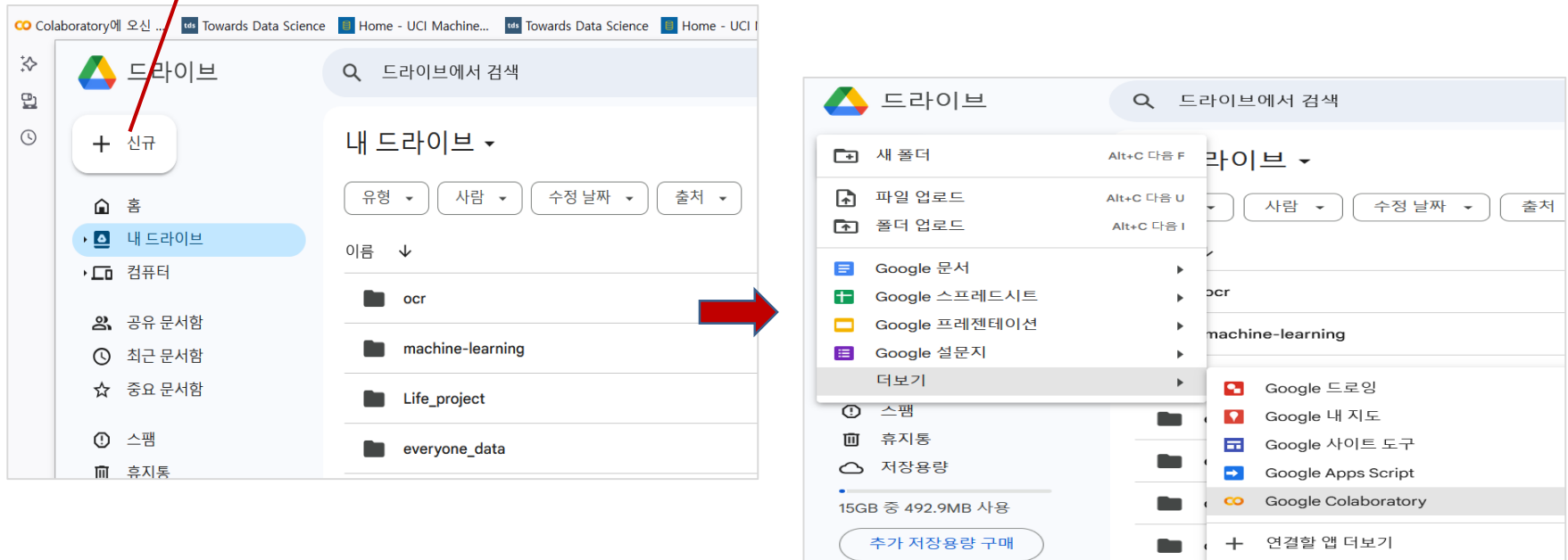
1. 드라이브로 연결하기

크롬 등 브라우저 > 구글 드라이브

# 구글 코랩

## ❖ 구글 코랩(Colaboratory)

신규 > 더보기 > Google Colaboratory



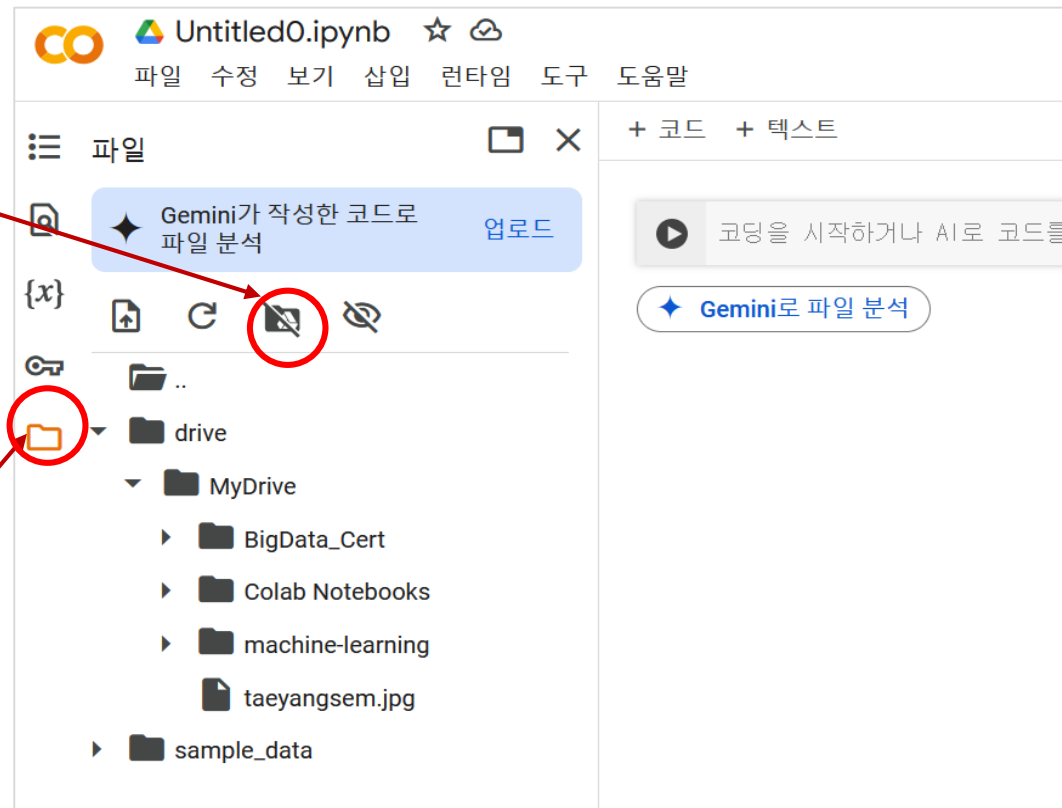
# 구글 코랩

## ❖ 구글 코랩(Colaboratory)

➤ Colab 웹으로 접속하기 : 새 노트 클릭

② 드라이브 연결

① 폴더 열기

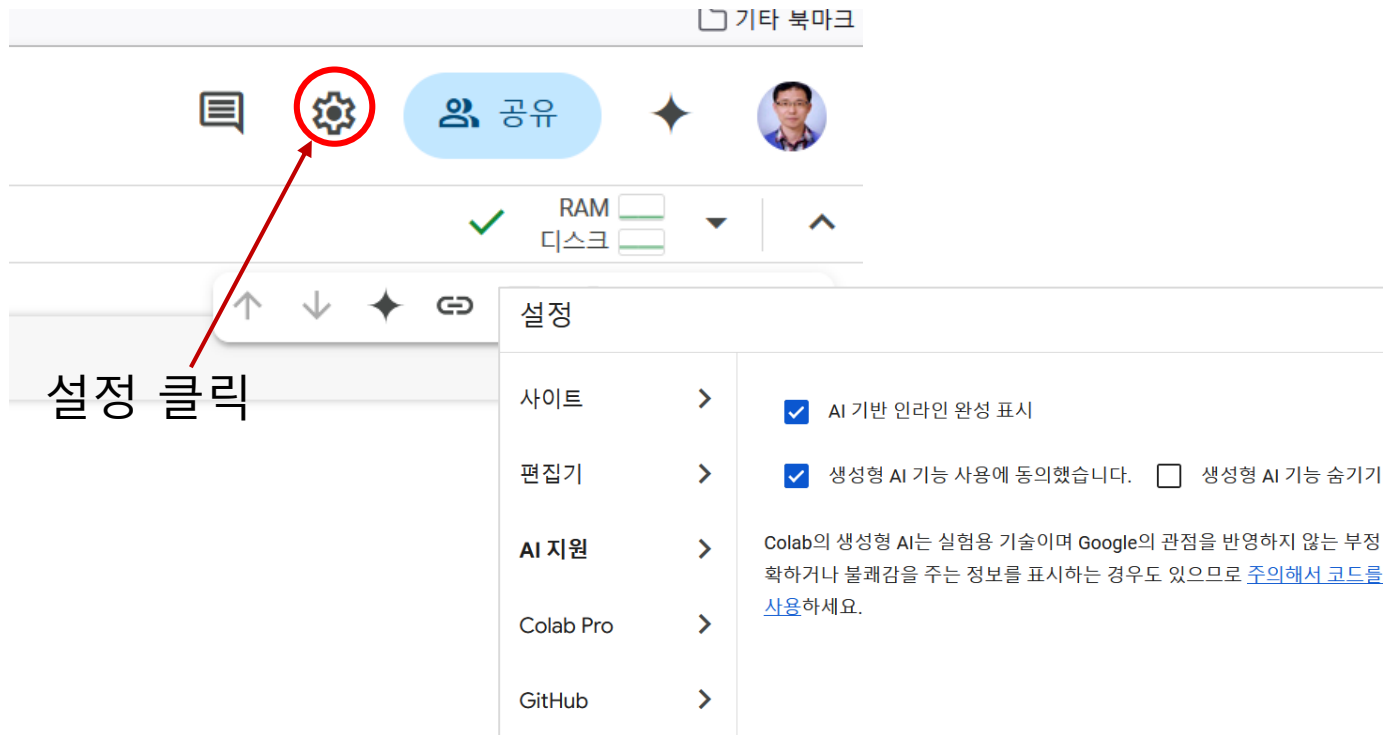


# 구글 코랩

## ❖ 구글 코랩(Colaboratory)

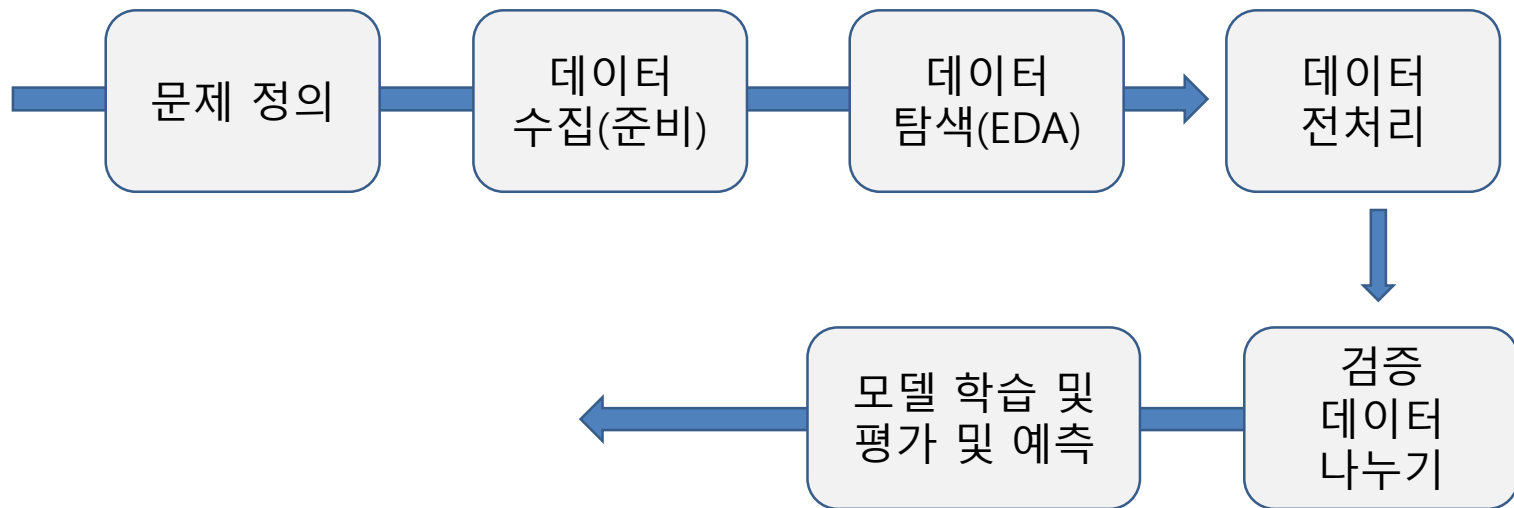
➤ Colab 웹으로 접속하기

AI(Gemini) 지원 : 설정 > AI 지원 > 동의



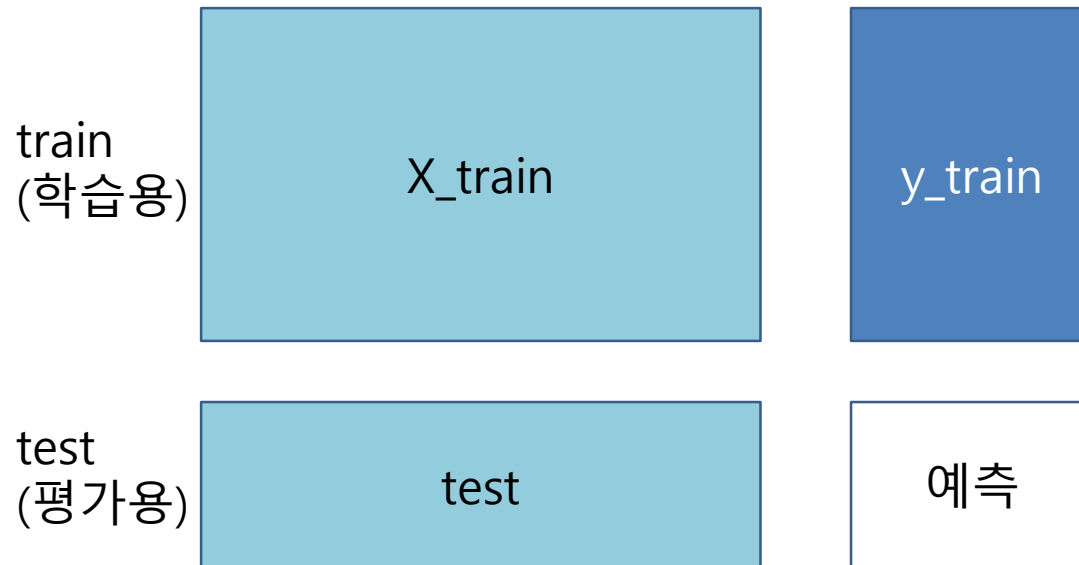
# 머신러닝 프로세스

- 머신러닝 모델 예측 프로세싱



# 데이터 수집(준비)

- 데이터는 학습용과 평가용으로 준비한다.



# 데이터 전처리

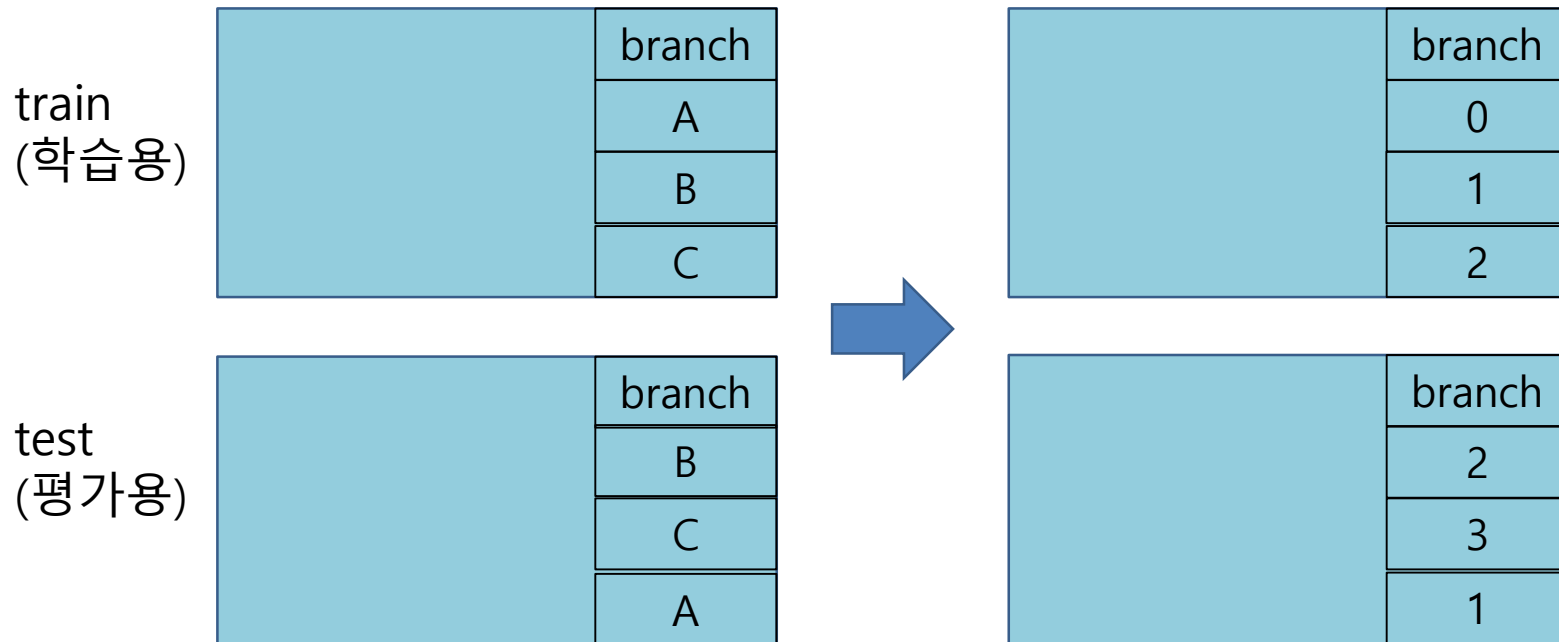
- 데이터 처리를 위한 모듈

	알고리즘	필수 여부
레이블 인코딩	sklearn.preprocessing(모듈) LabelEncoder	필수 (원-핫 인코딩중 선택)
스케일링	sklearn. preprocessing(모듈) StandardScaler	선택
데이터 분할	sklearn. model_selection(모듈) train_test_split	필수



# 데이터 전처리

- 인코딩(Encoding) – 문자형을 숫자형으로 변환
  - ✓ 레이블 인코딩(Label Encoding) : 각 칼럼을 숫자로 매핑하는 인코딩 방법이다.



# 머신러닝 프로세스

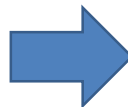
- 인코딩(Encoding) – 문자형을 숫자형으로 변환
  - ✓ 원-핫 인코딩(One-Hot Encoding) : 각 칼럼에 대해 새로운 칼럼을 만들고 새로운 칼럼에 0 또는 1이 저장된다.

train  
(학습용)

branch
A
B
C

test  
(평가용)

branch
B
C
A



branch_A	branch_B	branch_C
1	0	0
0	1	0
0	0	1

branch_A	branch_B	branch_C
0	1	0
0	0	1
1	0	0

# 머신러닝 프로세스

- 스케일링(Scaling)

설명(독립) 변수 열들이 갖는 데이터의 상대적 크기 차이를 없애기 위한 정규화 과정이다.

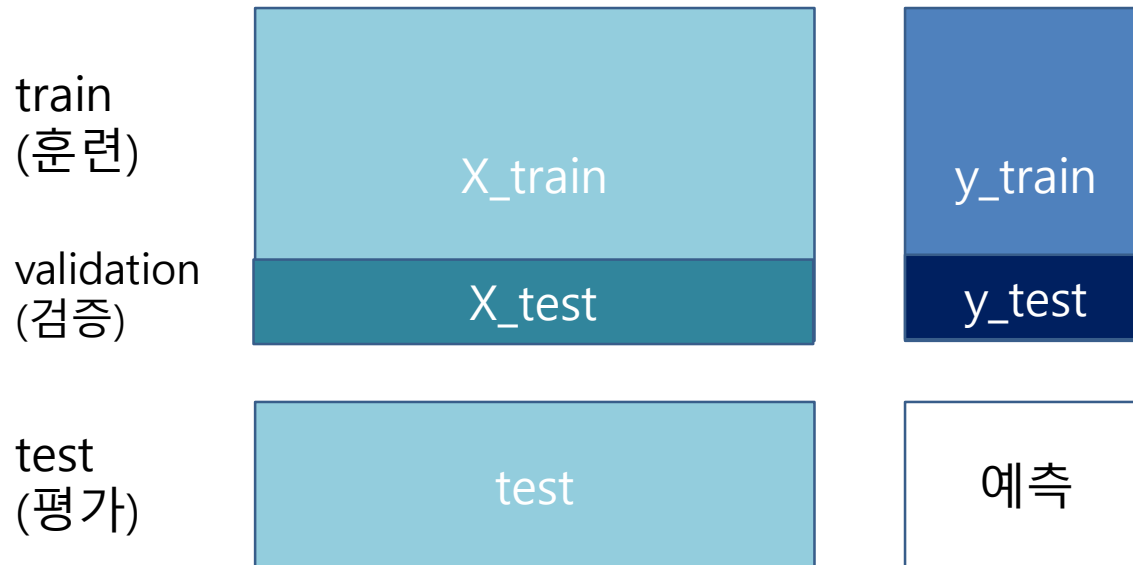
- 0 ~ 10000인 경우 이를 0~1 사이의 값으로 변경함
- A 칼럼 수치데이터가 0 ~ 10000이고, B 칼럼의 데이터가 0~10인 경우 두 칼럼간 데이터 차이가 커서 B 칼럼의 경우 유효하지 않은 데이터가 될 수 있음.
- 스케일링을 통해 같은 범위(0~1)로 맞추으로써 모델 성능을 향상시킬 수 있음

A 칼럼    2000, ,999, 10, 10000    ➡    0.222, 0.09, 0.001, 1.000

B 칼럼    2, 9, 1, 4    ➡    0.022, 0.09, 0.001, 0.004

# 학습용 데이터 분할

- 학습용(훈련/검증) 데이터 분할
  - 훈련용 데이터 중 일부를 검증용 데이터로 분리한다.  
(8:2 정도비율)



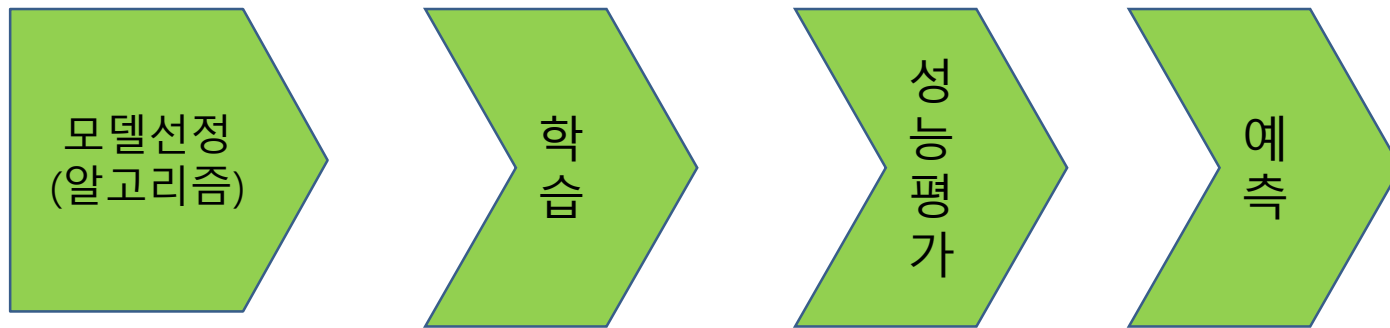
# 모델 학습 및 평가

- 모델 학습 및 평가

지도 학습	알고리즘	평가 지표
분류	sklearn.ensemble(모듈) RandomForestClassifier	sklearn.metrics(모듈) roc_auc_score, f1_score
회귀	sklearn.linear_model(모듈) LinearRegression  sklearn.ensemble(모듈) RandomForestRegressor	sklearn.metrics(모듈) mean_squared_error mean_absolute_error

# 모델 학습 및 평가

- 모델 학습, 평가 및 예측



# 선형 회귀 분석

- 단순 선형 회귀

- $y = Wx + b$
- $W$  를 가중치(Weight),  $b$ 를 편향(bias)라고 부른다
- 그래프의 형태는 **직선**

- 다중 선형 회귀

- $y = W_1x_1 + W_2x_2 + \dots + W_nx_n + b$
- 만약 2개의 독립 변수면 그래프는 **곡선**으로 나타남

# 선형 회귀 분석

- 공부한 시간의 차이에 따른 시험 성적 예측하기
  - 성적을 변하게 하는 요소를  $x$ 라 하고 이  $x$ 값에 따라 변하는 '성적'을  $y$ 라 하자.
  - 독립 변수( $x$ ) : 공부한 시간, 사교육비 지출, 학생의 지능지수 등
  - 종속 변수( $y$ ) : 성적
  - 어떤 변수가 다른 변수에 영향을 준다면 두 변수 사이에 선형관계가 있다.



# 선형 회귀 분석

- 공부한 시간의 차이에 따른 시험 성적 예측하기

공부한 시간	2시간	4시간	6시간	8시간
성적	81점	93점	91점	97점

$$x = \{2, 4, 6, 8\}$$

$$y = \{81, 93, 91, 97\}$$

직선의 방정식을 구한다.

$$y = ax + b(\text{일차함수})$$

기울기  $a$ 와 절편  $b$ 를 알아야 함

## 최소 제곱법

- 공부한 시간의 차이에 따른 시험 성적 예측하기
  - 기울기  $a$  구하기

$$a = \frac{(x - x_{\text{평균}})(y - y_{\text{평균}}) \text{의 합}}{(x - x_{\text{평균}})^2 \text{의 합}}$$

공부한 시간( $x$ ) 평균 :  $(2 + 4 + 6 + 8) \div 4 = 5$

성적( $y$ ) 평균:  $(81 + 93 + 91 + 97) \div 4 = 90.5$

기울기  $a = 2.3$

$x$ 의 편차(각 값과 평균과의 차이)를 제공해서 합한 값을 분모로 놓고,  
 $x$ 와  $y$ 의 편차를 곱해서 합한 값을 분자로 놓으면 기울기가 나옴

# 최소 제곱법

- 공부한 시간의 차이에 따른 시험 성적 예측하기
  - 편차(deviation)  
관측값에서 평균을 뺀 값
  - 분산 (variance)  
편차(관측값-평균)를 제곱하고 그것을 모두 더한 후  
전체 개수로 나눈것. 즉 편차의 제곱의 평균을 말한다.
  - 표준편차(standard deviation)  
분산의 제곱근, 즉 분산 값에 루트를 씌운것을 말하며  
stdev라고 함.

## 최소 제곱법

- 공부한 시간의 차이에 따른 시험 성적 예측하기
  - y절편 b값 구하기

$$b = y \text{의 평균} - (\text{기울기 } a \times (x \text{의 평균}))$$

$$b = 90.5 - (2.3 \times 5) = 79$$

최적의 직선(방정식)

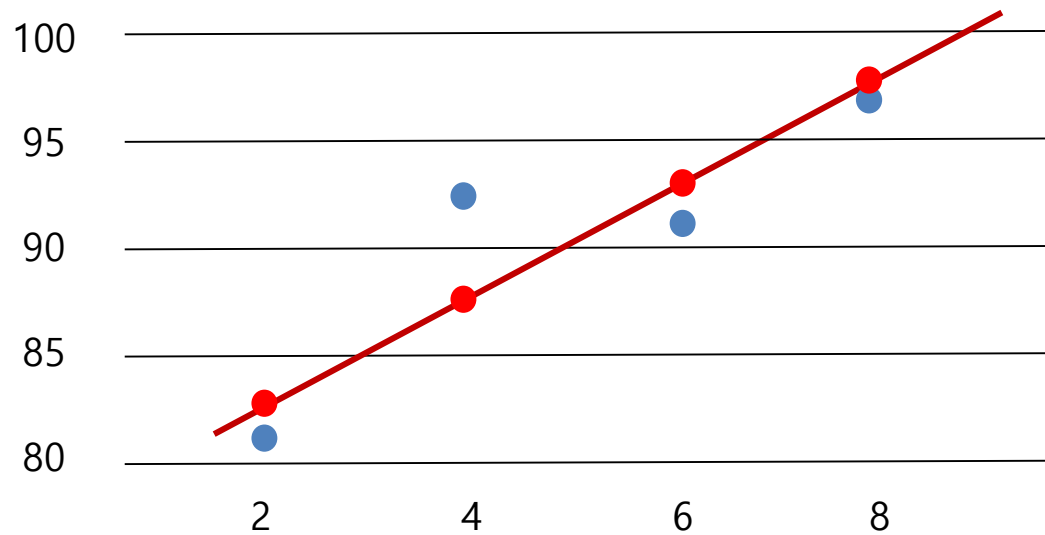
$$y = 2.3x + 79$$

공부한 시간	2	4	6	8
성적	81	93	91	97
예측 값	83.6	88.2	92.8	97.4

# 최소 제곱법

- 공부한 시간의 차이에 따른 시험 성적 예측하기

공부한 시간	2	4	6	8
성적	81	93	91	97
예측 값	83.6	88.2	92.8	97.4



# 최소 제공법

- 코드로 구현하기

```
import numpy as np

x = np.array([2, 4, 6, 8]) #공부한 시간
y = np.array([81, 93, 91, 97])

print(x, y)

mx = np.mean(x) # x의 평균값
my = np.mean(y) # y의 평균값

print(mx, my)
```

```
[2 4 6 8] [81 93 91 97]
5.0 90.5
```

# 최소 제곱법

- 코드로 구현하기

```
# 분자 계산 함수
# x와 y의 편차값을 곱하여 누적해서 합산함
def top(x, mx, y, my):
    d = 0
    for i in range(len(x)):
        d = d + (x[i] - mx) * (y[i] - my)
    return d

#x의 편차를 제공한 값
for i in x:
    print((i - mx)**2)

# 분모 - x의 편차를 제공해서 합한 값
dividend = sum([(i - mx)**2 for i in x])

#분자
divisor = top(x, mx, y, my)
# print(divisor, dividend)
```

$$d = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

# 최소 제곱법

- 코드로 구현하기

```
# 기울기
a = divisor / dividend

# y 절편
b = my - (mx * a)

print("기울기 a = ", a) #2.3
print("y절편 b = ", b) #79

# 최적의 직선
# y = 2.3x + 79.0
```

```
9.0
1.0
1.0
9.0
46.0 20.0
기울기 a = 2.3
y절편 b = 79.0
```

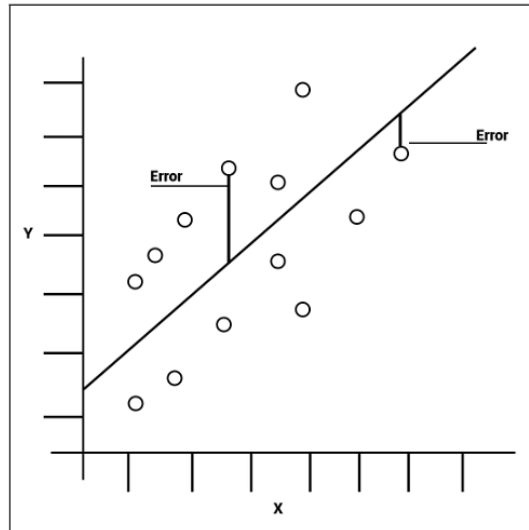


# 평균 제곱 오차

## ✓ 평균 제곱 오차

가설을 하나 세운 후(먼저 선을 긋고) 이 값이 주어진 요건을 충족하는지 판단해서 조금씩 변화를 주고 이 변화가 긍정적이면 오차가 최소가 될 때까지 이 과정을 계속 반복하는 방법.

선형 회귀는 임의의 직선을 그어 이에 대한 평균 제곱 오차를 구하고, 이 값을 가장 작게 만들어주는 기울기(a)와 절편(b)을 찾아가는 작업이다.



# 평균 제곱 오차

✓ 평균 제곱 오차 / 평균 절대값 오차

- MSE(Mean Squared Error)

- 평균 제곱오차
- 오차의 제곱을 평균으로 나눈 것

$$MSE = \frac{\sum_{i=1}^n (y - \hat{y})^2}{n}$$

- MAE(Mean Absolute Error)

- 평균 절대 오차
- 오차의 차이를 절대값으로 변환한 뒤 합산

$$MAE = \frac{\sum |y - \hat{y}|}{n}$$

## 평균 제곱 오차

✓ 오차 = 실제값 - 예측값

공부한 시간	2	4	6	8
성적	81	93	91	97
예측 값	82	88	94	100
오차	1	-5	3	3

분모 :  $1 + 25 + 9 + 9 = 44$

$MSE = 44 / 4 = 11$

$MAE = 1 + 5 + 3 + 3 = 12$

# 평균 제곱 오차

## ✓ 평균 제곱 오차 계산

```
import numpy as np

# 가상의 기울기 a와 절편 b를 정함
fake_a = 3
fake_b = 76

# 공부시간 x와 점수 y의 넘파이 배열을 만듦
x = np.array([2, 4, 6, 8])
y = np.array([81, 93, 91, 97])

#  $y = ax + b$ 에 가상의 a 값과 b 값을 대입한 결과를 출력하는 함수
def predict(x):
    return fake_a * x + fake_b
print("예측점수:", predict(2))
```

# 평균 제곱 오차

## ✓ 평균 제곱 오차 계산

```
# 예측값 배열(리스트)
predict_result = []

# 예측값 구하는 함수
for i in range(len(x)):
    predict_result.append(predict(x[i]))
    print(f"공부시간={x[i]}, 실제점수={y[i]}, 예측점수={predict(x[i])}")
# print(predict_result)

# mse(mean square error) - 평균 제곱 오차 함수
# 분산 - 실제 점수와 예측 점수를 전달받아 제공하고 그합을 전체 개수로 나누어 반환함
n = len(x)
def mse(y, y_pred):
    print(y - y_pred)
    return (1/n) * sum((y - y_pred) ** 2)
```

## 평균 제곱 오차

### ✓ 평균 제곱 오차 계산

```
# 평균 제곱 오차값  
# print(mse(y, predict_result))  
print(f"평균 제곱 오차: {mse(y, predict_result)}")
```

```
예측점수: 82  
공부시간=2, 실제점수=81, 예측점수=82  
공부시간=4, 실제점수=93, 예측점수=88  
공부시간=6, 실제점수=91, 예측점수=94  
공부시간=8, 실제점수=97, 예측점수=100  
[-1  5 -3 -3]  
평균 제곱 오차: 11.0
```

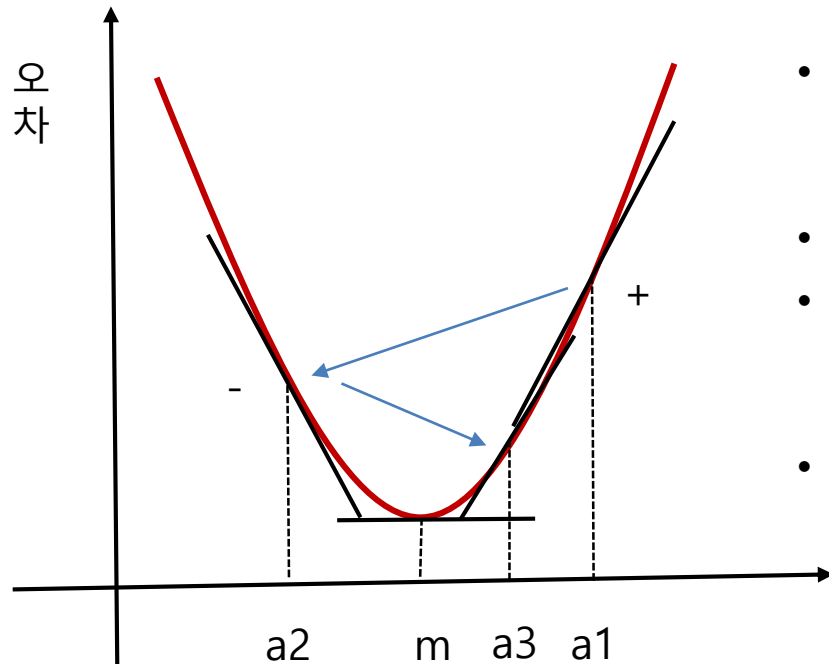
가정한  $a = 3$ ,  $b = 76$ 은 오차가 약 11.0이다.

이제 할 일은 이 오차를 줄이면서 새로운 선을 긋는 것이다.

이에 경사하강법을 사용하여 오차를 줄여갈 것이다.

# 경사 하강법

- 경사 하강법(기울기 하강법)



- 기울기를 크게 잡으면 오차가 커지고, 작게 잡아도 오차가 커진다. 기울기와 오차 사이에는 상관관계가 있다.
- 기울기와 오차 사이에는 이차 함수의 관계가 있다.
- 함수의 기울기(경사)를 구하고 기울기의 반대 방향으로 계속 이동시켜서 최소값에 이를 때까지 반복시킨다.
- 결국  $m$ 에 이르게 되면 최적의 기울기를 찾는 것인데 이러한 방법을 **경사하강법**이라 한다.

미분값이 0인값을 찾을  
이때 가중치는 0이 됨

# 경사 하강법

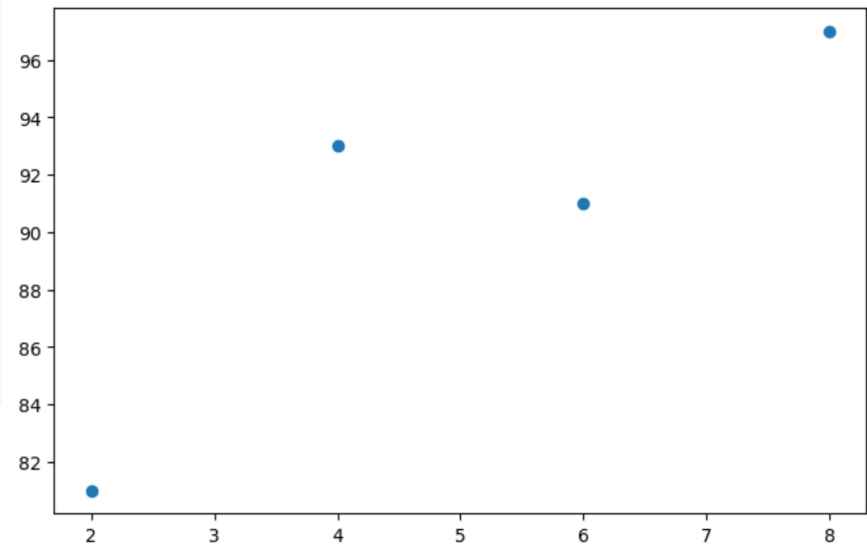
- 경사 하강법

```
import numpy as np
import matplotlib.pyplot as plt

# 공부시간 X와 성적 Y의 리스트 만들기
x = np.array([2, 4, 6, 8])
y = np.array([81, 93, 91, 97])

print(x, y)

# 산점도 그래프로 나타내기
plt.figure(figsize=(8,5))
plt.scatter(x, y)
plt.show()
```





# 경사 하강법

- 경사 하강법

```
a = 0 #기울기
b = 0 #절편

lr = 0.03 #학습률(learning rate)

epochs = 2000 #반복 횟수

# 경사하강법 실행 함수
n = len(x)
for i in range(epochs):
    y_pred = a * x + b #예측값
    error = y - y_pred #오차 = 실제값 - 예측값

    # 오차 제곱합을 a, b로 각각 편미분한 값
    a_diff = -(1/n) * sum(x * (error))
    b_diff = -(1/n) * sum(error)

    # 기울기 하강법(Gradient Descent) 업데이트
    a = a - lr * a_diff #학습률을 곱해 기존의 a값을 업데이트함
    b = b - lr * b_diff
```

$$\frac{\partial L}{\partial a} = -\frac{2}{n} \sum x_i(y_i - (ax_i + b))$$

$$\frac{\partial L}{\partial b} = -\frac{2}{n} \sum (y_i - (ax_i + b))$$

## 편미분 공식

여기서는 2를 생략했는데, 어차피 학습률이 조절해주기 때문에 결과에는 큰 문제 없음.

# 경사 하강법

- 경사 하강법

```
if i % 100 == 0: #100번 반복할때마다 출력
    print(f"epoch = {i}, 기울기 a = {a}, y절편 b = {b}")

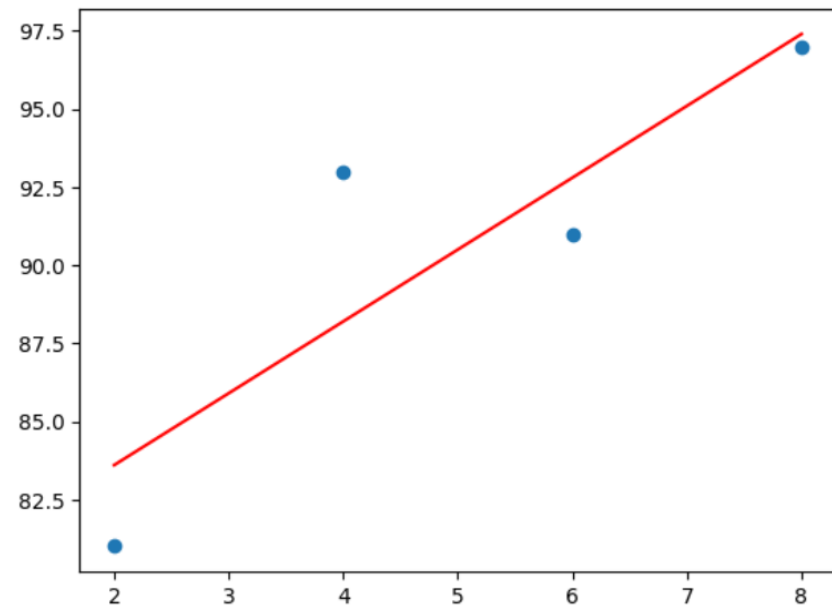
y_pred = a * x + b
print("실제 점수:", y)
print("예측 점수:", y_pred)
print("기울기:", a)
print("절편:", b)

# 그래프 그리기
plt.figure(figsize=(8, 5))
plt.scatter(x, y)
plt.plot(x, y_pred, 'r') #학습된 회귀 직선
plt.show()
```

# 경사 하강법

## ● 경사 하강법

epoch=0, 기울기=27.8400, 절편=5.4300  
epoch=100, 기울기=7.0739, 절편=50.5117  
epoch=200, 기울기=4.0960, 절편=68.2822  
epoch=300, 기울기=2.9757, 절편=74.9678  
epoch=400, 기울기=2.5542, 절편=77.4830  
epoch=500, 기울기=2.3956, 절편=78.4293  
epoch=600, 기울기=2.3360, 절편=78.7853  
epoch=700, 기울기=2.3135, 절편=78.9192  
epoch=800, 기울기=2.3051, 절편=78.9696  
epoch=900, 기울기=2.3019, 절편=78.9886  
epoch=1000, 기울기=2.3007, 절편=78.9957  
epoch=1100, 기울기=2.3003, 절편=78.9984  
epoch=1200, 기울기=2.3001, 절편=78.9994  
epoch=1300, 기울기=2.3000, 절편=78.9998  
epoch=1400, 기울기=2.3000, 절편=78.9999  
epoch=1500, 기울기=2.3000, 절편=79.0000  
epoch=1600, 기울기=2.3000, 절편=79.0000  
epoch=1700, 기울기=2.3000, 절편=79.0000  
epoch=1800, 기울기=2.3000, 절편=79.0000  
epoch=1900, 기울기=2.3000, 절편=79.0000  
epoch=2000, 기울기=2.3000, 절편=79.0000  
실제 점수: [81 93 91 97]  
예측 점수: [83.59999984 88.19999992 92.8 97.40000008]  
기울기: 2.3000000409418653  
절편: 78.99999975567644



# 실습 문제 - 회귀

## ● 실습 문제 - 회귀

문제) 통신사에서 고객에서 청구될 총 금액을 예측하시오.

- 제공된 데이터 목록: churn\_train.csv, churn\_test.csv
- 예측할 칼럼: TotalCharges(총 청구액)

학습용 데이터(churn\_train.csv)를 이용하여 총 판매액을 예측하는 모델을 만든 후 이를 평가용 데이터(churn\_test.csv)에 적용해 얻은 예측값을 다음과 같은 형식의 CSV 파일로 생성하시오.

제출 파일은 다음 1개의 칼럼만 포함해야 한다.

- pred: 예측된 총 청구액
- 제출 파일명: 'result.csv'

제출한 모델의 성능은 MAE(Mean Absolute Error) 평가지표에 따라 채점한다.

- 제출 csv 파일명 및 형태: result.csv(수치형 데이터)

```
pred
3510
800
3970
950
```

# 실습 문제 - 회귀

- 실습 문제 - 회귀

```
import pandas as pd

train = pd.read_csv("/content/drive/MyDrive/BigData_Cert/회귀/churn_train.csv")
test = pd.read_csv("/content/drive/MyDrive/BigData_Cert/회귀/churn_test.csv")

# 데이터 탐색
# print(train.head())
train.info()

# print(train.describe(include="object"))
# print(test.describe(include="object"))

# customerID - 삭제
train = train.drop("customerID", axis=1)
test = test.drop("customerID", axis=1)

# object 칼럼 인코딩
cols = train.select_dtypes(include='object').columns
# print(cols)

# 원-핫 인코딩
train = pd.get_dummies(train)
test = pd.get_dummies(test)
```

## 실습 문제 - 회귀

- 실습 문제 - 회귀

```
# 레이블 인코딩
# from sklearn.preprocessing import LabelEncoder, StandardScaler
# encoder = LabelEncoder()
# for col in cols:
#     train[col] = encoder.fit_transform(train[col])
#     test[col] = encoder.fit_transform(test[col])

# train.info()

# 학습용(훈련/검증) 데이터 분할
from sklearn.model_selection import train_test_split
# 변수 선택
X = train.drop("TotalCharges", axis=1) # 독립(설명) 변수
y = train["TotalCharges"] # 종속(예측) 변수

X_train, X_test, y_train, y_test = train_test_split(X,
                                                    y, test_size=0.2, random_state=30)
print(X_train.shape, X_test.shape, y_train.shape, y_test.shape)
```

## 실습 문제 - 회귀

- 실습 문제 - 회귀

```
# 모델 학습 및 예측, 평가
from sklearn.ensemble import RandomForestRegressor
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_absolute_error, mean_squared_error
# 랜덤 포레스트
# model = RandomForestRegressor(random_state=0)
# 선형 회귀
model = LinearRegression()

model.fit(X_train, y_train)
pred = model.predict(X_test)
# print(pred)

# 평균 절대값 오차
mae = mean_absolute_error(y_test, pred)
print("MAE:", mae)

# 평균 제곱 오차
mse = mean_squared_error(y_test, pred)
print("MSE:", mse)
```

## 실습 문제 - 회귀

- 실습 문제 - 회귀

```
# 평균 제곱근 오차
import numpy as np
rmse = np.sqrt(mse)
print("RMSE:", rmse)
```

```
# 평가용 데이터 - 예측 및 파일 제출
pred2 = model.predict(test)
print(pred2)
```

```
result = pd.DataFrame({'pred': pred2})
result.to_csv("/content/drive/MyDrive/BigData_Cert/회귀/result.csv", index=False)

pd.read_csv("/content/drive/MyDrive/BigData_Cert/회귀/result.csv")
```

	pred
0	3365.262956
1	901.756573
2	4030.186814
3	793.861079
4	1385.789030
...	...
1759	268.579922
1760	4140.701446