



# CNU 데이터 분석 교육



7<sup>th</sup> lecture  
"Linear Analysis-1"

2022 - 11 - 07

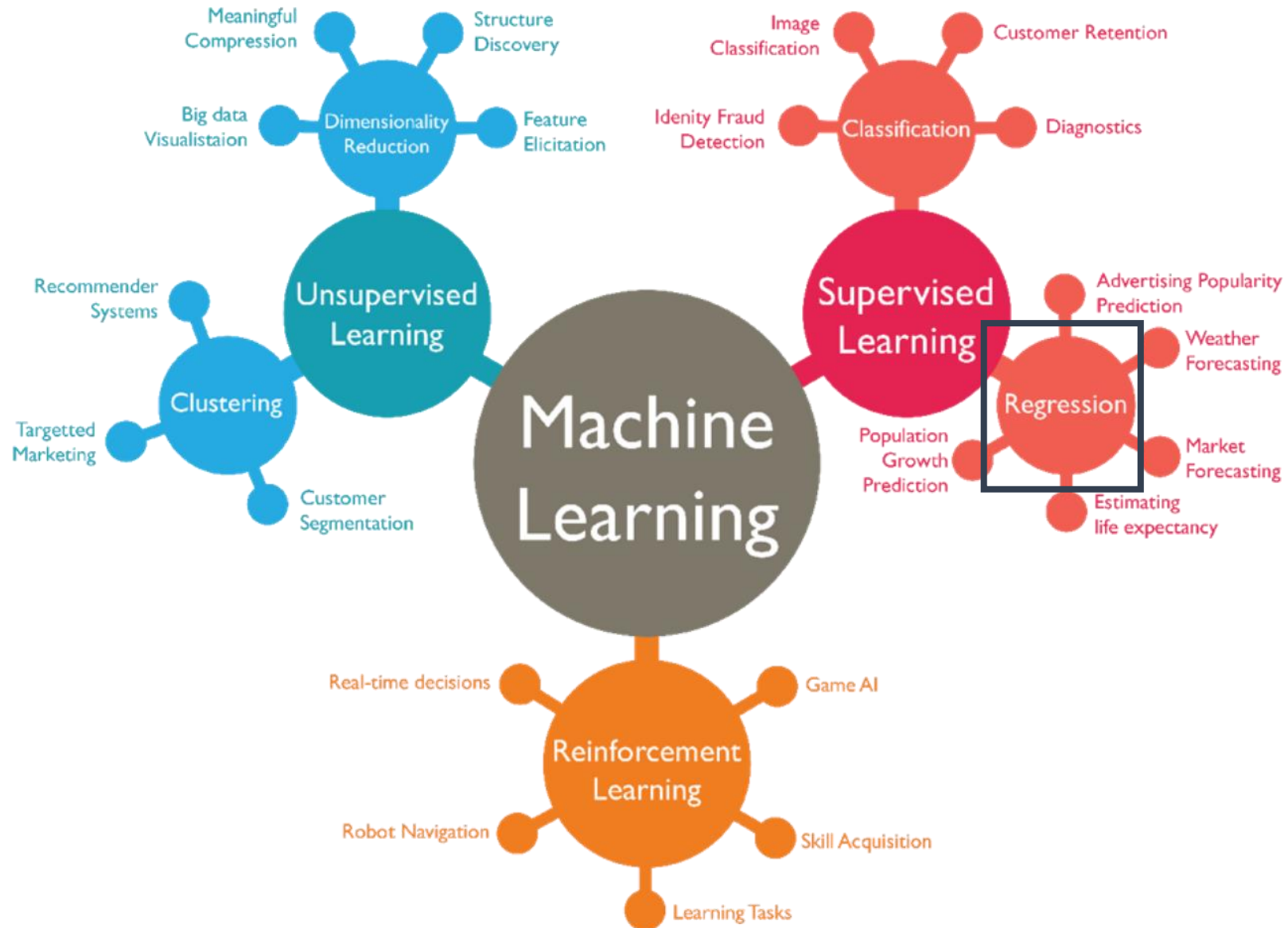
지난시간?

## 1. Matplotlib, Seaborn – Visualization

# 오늘은 무엇을?

## 1. 선형회귀(Linear regression) 모델 개발

- 두 변수의 관계를 가장 잘 설명할 수 있는 선형함수가 무엇일까?
- 그 선형함수를 찾는 데 어떤 수학적 방법이 이용될까?
- 수학적 아이디어를 코딩으로 나타내는 방법?

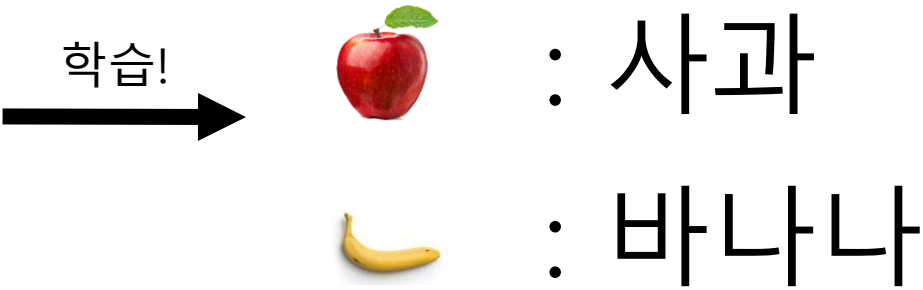


## 지도 학습 (Supervised learning)



# 지도 학습 (Supervised learning)

data			Label
			사과
			바나나

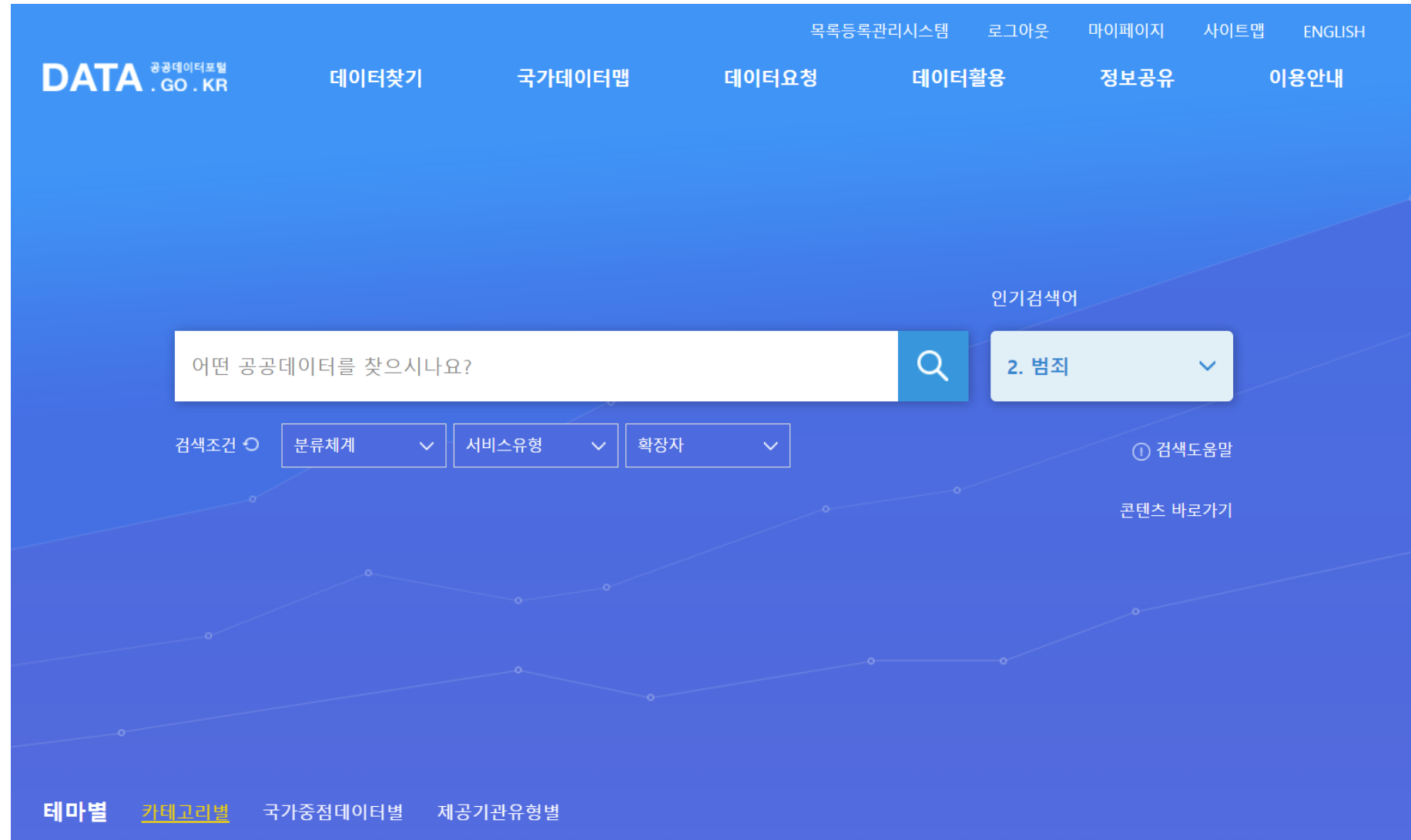


# 회귀분석(Regression)이란?

1. 독립변수  $X$ 와 종속변수  $Y$ 의 관계를 모델링하는 기법  
-너네 둘 무슨 사이니...?
2.  $X$ 는 일반적으로  $n$ 차원,  $Y$ 는 1차원
3. 둘 사이의 관계를 수치로 평가

$data$ 로 부터 원하는 값을 예측하고 싶을때 주로 사용.

# 공공데이터포털





파일데이터 (1,187건)

더보기 >

교육

국가행정기관

미리보기

CSV

JSON + XML

교육부\_학교건강검사 표본조사결과(건강조사)\_2017

교육부\_학교건강검사 표본조사결과(건강조사)\_2017

제공기관 교육부    수정일 2020-09-15    조회수 3029    다운로드 704    주기성 데이터 2    키워드 학교건강,건강검사,표본조사

다운로드

교육

국가행정기관

미리보기

CSV

JSON + XML

교육부\_학교건강검사 조사 결과 2016

전국 초·중·고등학생의 건강조사 표본조사 결과

제공기관 교육부    수정일 2020-09-03    조회수 2557    다운로드 426    키워드 학교건강,건강검사,표본조사

다운로드

교육

국가행정기관

미리보기

CSV

JSON + XML

교육부\_학생건강검사 키rawdata 2016

학생표본 신체(키) 검사 rawdata

제공기관 교육부    수정일 2020-09-22    조회수 3470    다운로드 846    키워드 학교건강,건강검사,표본조사

다운로드

교육

국가행정기관

미리보기

CSV

JSON + XML

교육부\_학생건강검사 결과분석 rawdata 서울 2015

2015년도 학생건강검사 결과(표본) 서울지역 학생 키, 몸무게, 혈당, 고지혈증, B형간염, 혈색소에 대한 raw data

제공기관 교육부    수정일 2020-09-03    조회수 4109    다운로드 1214    키워드 학교건강,건강검사,표본조사

의견수렴  
게시판

csv 교육부\_학생건강검사 결과분석 rawdata 서울 2015

다운로드

오류신고 및 담당자 문의

파일데이터 정보

메타데이터 다운로드

미리보기

※ 파일 데이터의 일부 내용을 제공하고 있으며, 전체 내용이 필요한 경우 해당 파일을 다운로드 받으시기 바랍니다.

파일데이터명	ID	최종가중치	학교ID	도시규모	도시규모별분석용	학년도	광역시도	시도별	학
분류체계	Aa011남10101	169.550665	Aa01	대도시/중소도시	특별/광역시	2015	서울	서울특별시교육청	1
관리부서명	Aa011남10102	169.550665	Aa01	대도시/중소도시	특별/광역시	2015	서울	서울특별시교육청	1
보유근거	Aa011남10103	169.550665	Aa01	대도시/중소도시	특별/광역시	2015	서울	서울특별시교육청	1
업데이트 주기	Aa011남10104	169.550665	Aa01	대도시/중소도시	특별/광역시	2015	서울	서울특별시교육청	1
매체유형	Aa011남10105	169.550665	Aa01	대도시/중소도시	특별/광역시	2015	서울	서울특별시교육청	1
확장자	Aa011남10106	169.550665	Aa01	대도시/중소도시	특별/광역시	2015	서울	서울특별시교육청	1
데이터 한계	Aa011남10107	169.550665	Aa01	대도시/중소도시	특별/광역시	2015	서울	서울특별시교육청	1
등록	Aa011남10108	169.550665	Aa01	대도시/중소도시	특별/광역시	2015	서울	서울특별시교육청	1
제공형태	Aa011남10109	169.550665	Aa01	대도시/중소도시	특별/광역시	2015	서울	서울특별시교육청	1
설명									
기타 유의사항									
비용부과유무									
이용허락범위									

GitHub 다운로드

<https://github.com/NumberL/NumberT>

# 데이터 업로드



🔍 드라이브에서 검색



새로 만들기



내 드라이브



컴퓨터



공유 문서함



최근 문서함

내 드라이브 > school > CNU\_2022\_data > data ▾

이름 ↑



학생건강검사 결과분석 rawdata\_서울\_2015\_20200114.csv

# 데이터 업로드

1. 드라이브 마운트 후 데이터를 올린 경로를 복사!
2. Encoding = 'cp949'로 설정 후 데이터가 잘 로딩이 되면 끝!

The screenshot shows the Google Colab interface. On the left, the file explorer shows a folder named 'drive' which is mounted as 'MyDrive'. A red box highlights the 'drive' folder, and a red arrow points from it to the first instruction in the list above. Below the 'drive' folder, a large black box with white text says '개인정보ㅎㅎ' (Personal information hehe). The main code area shows three code cells. The first cell imports matplotlib and pandas. The second cell mounts the drive. The third cell sets the figure size and loads a CSV file from the drive. The file path is: `"/content/drive/MyDrive/school/CNU_2022_data/data/학생건강검사 결과분석 rawdata_서울_2015_20200114.csv"`. The encoding is set to 'cp949'. The output of the third cell shows the first few rows of the data as a table.

```
import matplotlib.pyplot as plt
import pandas as pd

[7] # from google.colab import drive
# drive.mount('/content/drive')

[8] # setup the figsize
plt.rcParams["figure.figsize"] = (12,8)

# data load
# 한글이 들어있는 데이터는 cp949, 혹은 utf8로 인코딩하자. 안그러면 깨짐!
data_raw = pd.read_csv("/content/drive/MyDrive/school/CNU_2022_data/data/학생건강검사 결과분석 rawdata_서울_2015_20200114.csv", encoding='cp949')

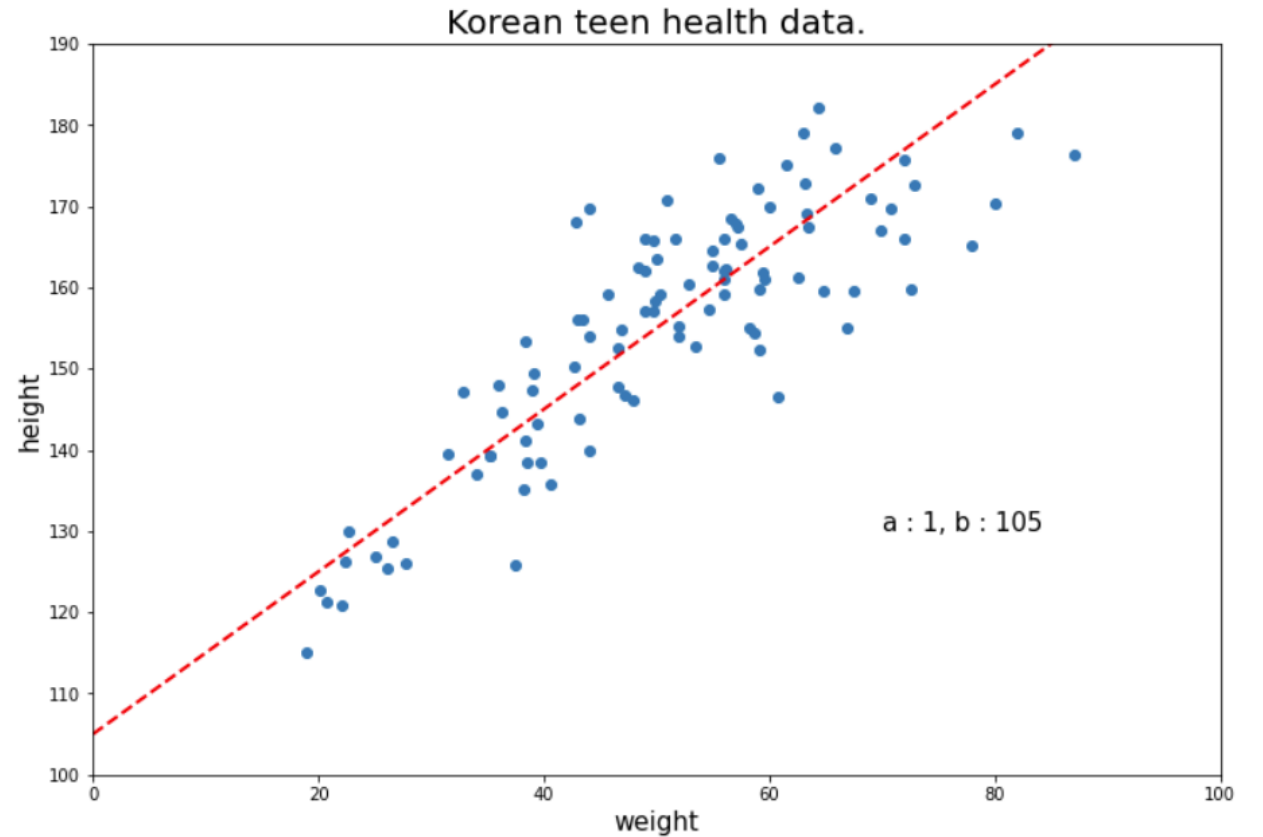
# data_raw = pd.DataFrame(data_raw)
data_raw.head()

# type(data_raw)
```

	ID	최종가중치	학교 ID	도시 규모	도시 규모 별분 석용	학년 도	광역 시도	시도 별	학교 급 별	학교 명	...	키	몸무 게	혈당식 전 mgdI	총콜레 스테롤 mgdI	ASTUL	ALTUL	혈색 소 gdI	간염 검사	수...
0	Aa011	남	169.550665	Aa01	대도 시/중 소도	특별/ 광역	2015	서울	서울 특별 시	서울 대도 초...	...	125.8	27.3	NaN	NaN	NaN	NaN	NaN	NaN	77.

# 문제 1

- 주어진 scatter plot 위에 직선을 그려주는 함수를 def로 선언한 후 plt.plot 기능을 이용하여 다음과 같은 그래프를 그려라.



# 문제 1

- 정답

```
# linear function 생성  
# y = ax + b
```

```
def height_pred(x, a, b):  
    return a*x + b
```

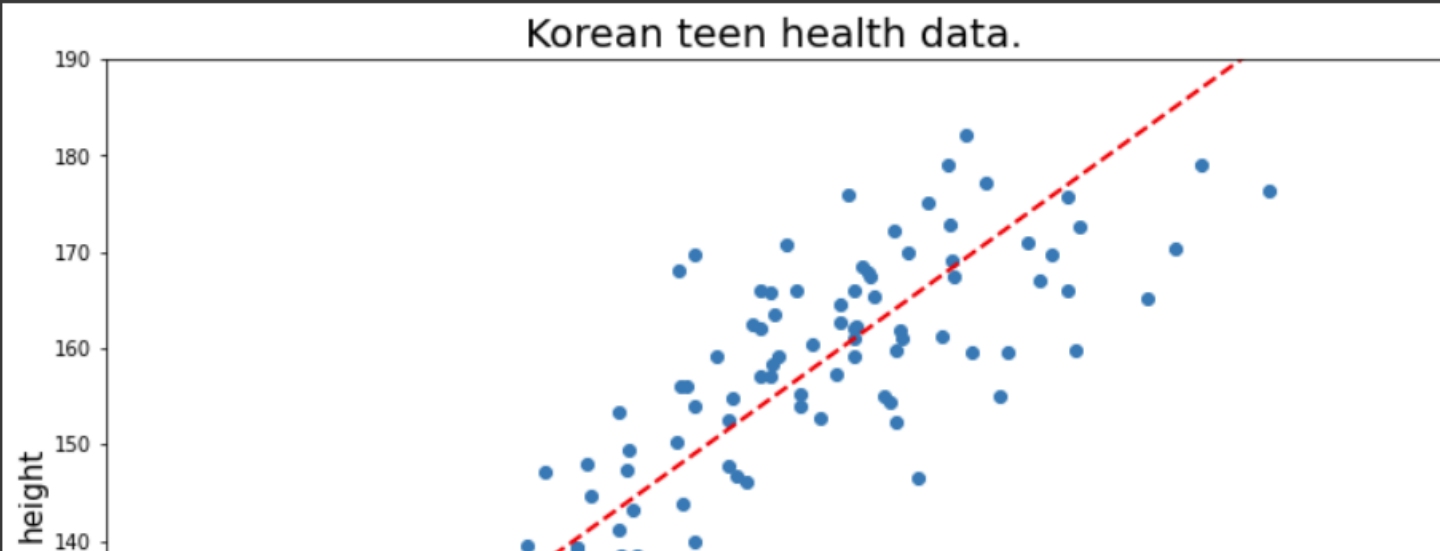
▶ # Line 그리기

```
a = 1  
b = 105
```

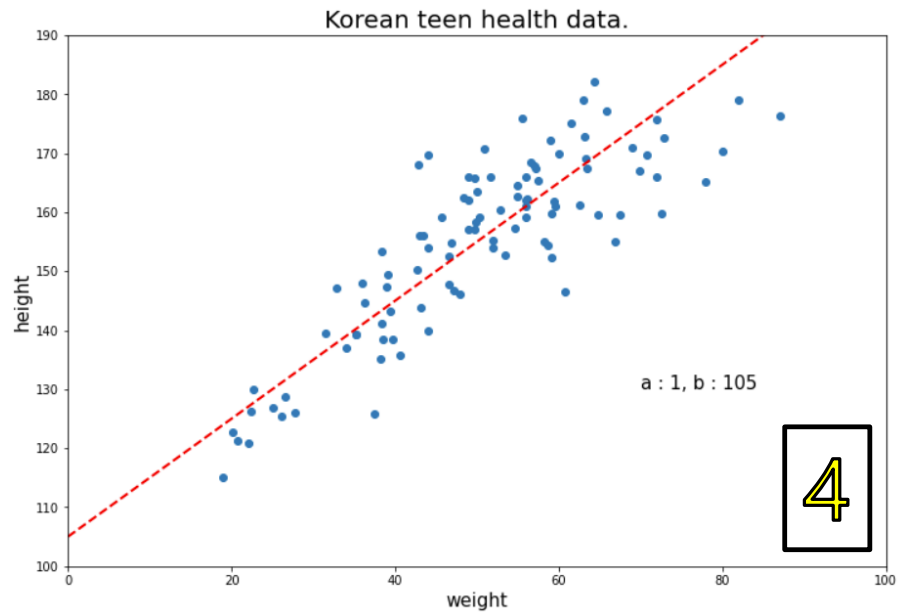
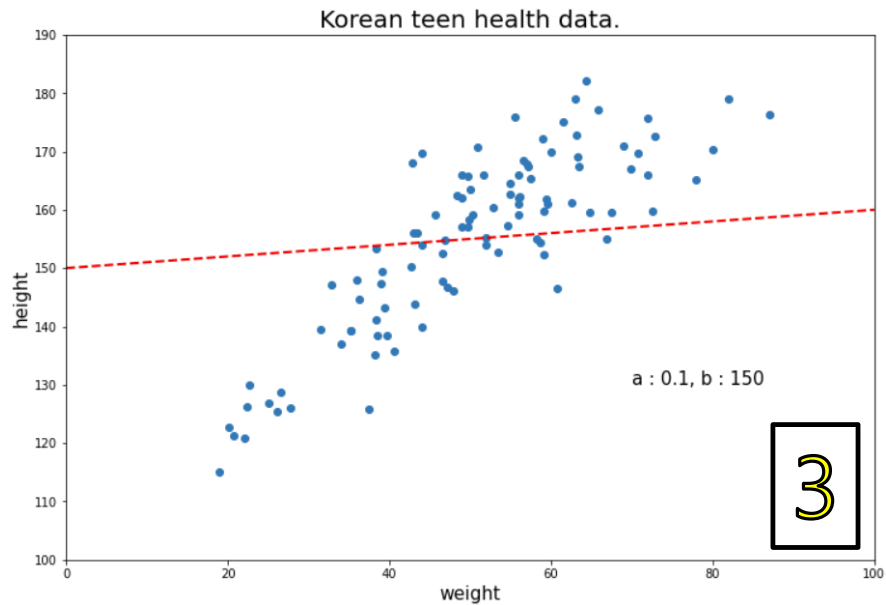
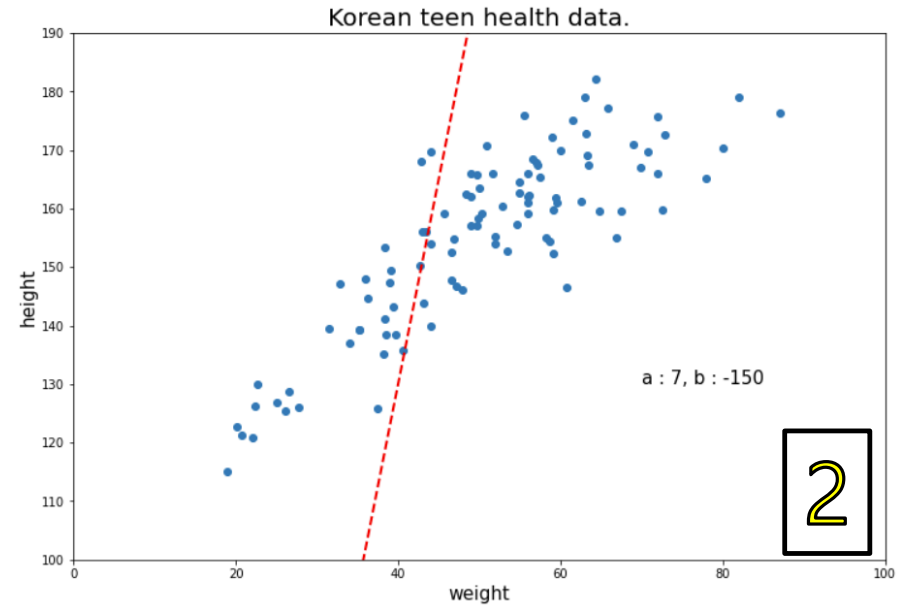
```
x = np.linspace(0,100,10001)  
y = height_pred(x, a, b)
```

```
plt.scatter(weight_100, height_100)  
plt.plot(x,y, linestyle = '--', linewidth = 2, color = 'red')  
plt.title("Korean teen health data.", fontsize = 20)  
plt.xlabel("weight", fontsize = 15)  
plt.ylabel("height", fontsize = 15)  
plt.text(70, 130, "a : {}, b : {}".format(a,b), fontsize = 15)
```

```
plt.xlim([0,100])  
plt.ylim([100,190])  
plt.show()
```



# 어느 것이 적당해 보임?



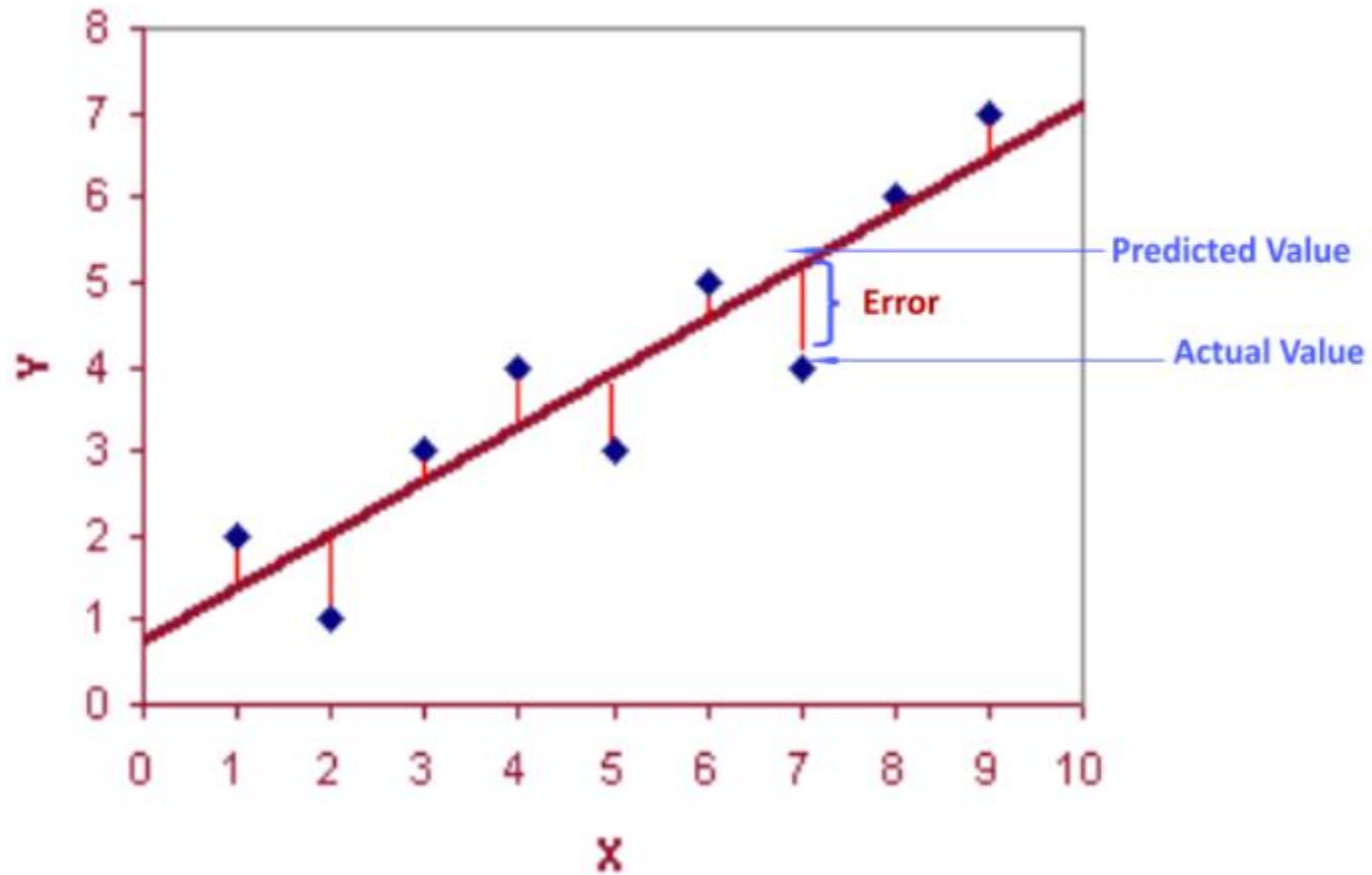


어느 것이 적당해 보임?

당신들이 선택한것의 공통점이 무엇인가?

1. 데이터와 직선간의 거리가 가깝다  
->점과 직선사이의거리...?

# 어느 것이 적당해 보임?



# 손실함수 (loss function, cost function)

Means Square Error

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

오차의 정도를 '양수'로 표현해준다.

Object!

Root Mean Square Error

$$RMSE = \sqrt{\sum_{i=1}^n \frac{(\hat{y}_i - y_i)^2}{n}}$$

**이것을 최소화 시키자!**

\*손실함수를 목적함수(Object function)이라고도 한다.

손실함수 (loss function)

HOW??????

가능한 모든  $a, b$ 의 쌍을 입력해서 그중 최소의 RMSE를 선택!

## 문제2

- $a$ (기울기)와  $b$ ( $y$ 절편)의 범위를 설정 한 후 각각의 대한 RMSE 를 DataFrame으로 작성하라.

	a	b	RMSE
0	0.9	90.0	21.278450
1	0.9	100.0	12.617387
2	0.9	110.0	8.100892
3	1.0	90.0	16.810529
4	1.0	100.0	9.479130
5	1.0	110.0	9.854639
6	1.1	90.0	12.943688
7	1.1	100.0	8.696726
8	1.1	110.0	13.554595

## 문제 2

- 정답

```
# 문제 2
```

```
a_range = np.linspace(0.9, 1.1, 3)
```

```
b_range = np.linspace(90, 110, 3)
```

```
a_lst = []
```

```
b_lst = []
```

```
rmse_lst = []
```

```
x = np.linspace(0, 100, 10001)
```

```
y = height_pred(x, a, b)
```

```
df = pd.concat([weight_100, height_100], axis = 1)
```

```
for a in a_range:
```

```
    for b in b_range:
```

```
        df["예상 키"] = height_pred(df["몸무게"], a, b)
```

```
        diff = df["키"] - df["예상 키"]
```

```
        RMSE = np.sqrt(sum(diff**2)/len(diff))
```

```
        a_lst.append(a)
```

```
        b_lst.append(b)
```

```
        rmse_lst.append(RMSE)
```

```
data_rmse = pd.DataFrame([a_lst, b_lst, rmse_lst]).T
```

```
data_rmse.columns = ["a", "b", "RMSE"]
```

```
data_rmse
```

# 손실함수(loss function)

## HOW??????

가능한 모든  $a, b$ 의 쌍을 입력해서 그중 최소의 RMSE를 선택!

그러면 우리가 찾은 **parameter**  $a$ 와  $b$ 는  
한국청소년들의 키와 몸무게를 설명하는 선형함수 모델에서  
가장 좋은 성능을 나타내는 것이다.

이러한 방법을 “

요즘|이ㄱ? 😊

그거슨 다음시간에 확인

## 정리 요약

1. 선형회귀(Linear regression)란 데이터를 가장 잘 설명해줄 수 있는 선형함수를 찾는 것이다
2. 데이터를 가장 잘 설명해준다는 것은 정의된 "loss function" 을 최소화 하는 것이다.
3. 우리는 데이터의 수( $n$ )과 parameter의 범위( $a\_range, b\_range$ )를 **입력(input)**하여 loss function을 최소화하는  $a, b$ 를 **출력(output)**는 함수(SLR)를 만들었다.
4. 데이터로 부터 가장 좋은 값을 얻을 수 있는  $a, b$ 를 학습(learn)하였다.



# 회귀분석(Regression test)이란?

1. 독립변수  $X$ 와 종속변수  $Y$ 의 관계를 모델링하는 기법  
-너네 둘 무슨 사이니...?
2.  $X$ 는 일반적으로  $n$ 차원,  $Y$ 는 1차원
3. Error의 척도를 가장 낮추는 parameter를 구하는 방법

# 공지

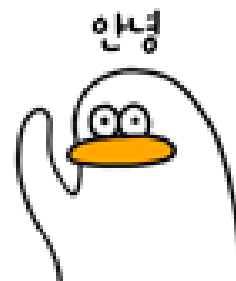
## 1. 어려웠?



누워서 떡먹기



끼  
트



담에뵈시당