

사이킷런(sklearn)

1. 사이킷런(sklearn) 개요
2. 회귀분석(boston 주택가격 예측하기)
3. 분류 분석(iris 붓꽃분류)

1. 사이킷런(sklearn) 개요

❖ 사이킷런

- 파이썬으로 머신러닝을 수행하기 위한 쉽고 효율적인 개발 라이브러리를 제공
- 다양한 [분류](#), [회귀](#), [서포트 벡터 머신](#), [랜덤 포레스트](#), [그라디언트 부스팅](#), [k-평균](#), [DBSCAN](#)을 포함한 [클러스터링](#) 알고리즘을 지원
- 파이썬의 수치 및 과학 라이브러리 [NumPy](#) 및 [SciPy](#)와 함께 운용되도록 설계
- 다양한 머신러닝 분석용 데이터셋 을 제공

- load_boston: 보스턴 집값 데이터
- load_iris: 아이리스 붓꽃 데이터
- load_diabetes: 당뇨병 환자 데이터
- load_digits: 손글씨 데이터
- load_linnerud: multi-output regression 용 데이터
- load_wine: 와인 데이터
- load_breast_cancer: 위스콘신 유방암 환자 데이터

1. 사이킷런(sklearn) 개요

❖ 분석 평가 지표

- 회귀 분석 결과에 대한 평가 지표는 예측값과 실제값의 차이인 오류의 크기가 됨
- 정확한 평가를 위해 오류의 절대값 평균이나 제곱의 평균, 제곱 평균의 제곱근 또는 분산 비율을 사용

평가 지표	수식	사이킷런 라이브러리
MAE: Mean Absolute Error	$\frac{1}{n} \sum_{i=1}^n Y_i - \hat{Y}_i $	<code>metrics.mean_absolute_error()</code>
MSE: Mean Squared Error	$\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$	<code>metrics.mean_squared_error()</code>
RMSE: Root Mean Squared Error	$\sqrt{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}$	없음
R^2 : Variance score, 결정 계수coefficient of determination	$\frac{\text{예측값의 분산}}{\text{실제값의 분산}}$	<code>metrics.r2_score()</code>

2. 회귀분석(boston 주택가격 예측하기)

❖ 데이터 수집, 준비 및 탐색

- 주택가격분석'으로 노트북 페이지를 추가하고 입력

In [1]:	!pip install sklearn # 사이킷런 설치
In [2]:	<pre>import numpy as np import pandas as pd from sklearn.datasets import load_boston boston = load_boston() print(boston)</pre>

In [2]: 사이킷런에서 제공하는 데이터셋sklearn.datasets 중에서 보스턴 주택 가격 데이터셋을 사용하기 위해 load_boston을 임포트하고, 데이터셋을 로드하여load_boston() 객체boston를 생성

2. 회귀분석(boston 주택가격 예측하기)

❖ 데이터 수집, 준비 및 탐색

- 데이터가 이미 정리된 상태이므로 데이터셋 구성을 확인

In [3]:	print(boston.DESCR)																																																																																										
In [4]:	boston_df = pd.DataFrame(boston.data, columns = boston.feature_names) boston_df.head()																																																																																										
Out[4]:	<table><tr><th></th><th>CRIM</th><th>ZN</th><th>INDUS</th><th>CHAS</th><th>NOX</th><th>RM</th><th>AGE</th><th>DIS</th><th>RAD</th><th>TAX</th><th>PTRATIO</th><th>B</th><th>LSTAT</th></tr><tr><td>0</td><td>0.00632</td><td>18.0</td><td>2.31</td><td>0.0</td><td>0.538</td><td>6.575</td><td>65.2</td><td>4.0900</td><td>1.0</td><td>296.0</td><td>15.3</td><td>396.90</td><td>4.98</td></tr><tr><td>1</td><td>0.02731</td><td>0.0</td><td>7.07</td><td>0.0</td><td>0.469</td><td>6.421</td><td>78.9</td><td>4.9671</td><td>2.0</td><td>242.0</td><td>17.8</td><td>396.90</td><td>9.14</td></tr><tr><td>2</td><td>0.02729</td><td>0.0</td><td>7.07</td><td>0.0</td><td>0.469</td><td>7.185</td><td>61.1</td><td>4.9671</td><td>2.0</td><td>242.0</td><td>17.8</td><td>392.83</td><td>4.03</td></tr><tr><td>3</td><td>0.03237</td><td>0.0</td><td>2.18</td><td>0.0</td><td>0.458</td><td>6.998</td><td>45.8</td><td>6.0622</td><td>3.0</td><td>222.0</td><td>18.7</td><td>394.63</td><td>2.94</td></tr><tr><td>4</td><td>0.06905</td><td>0.0</td><td>2.18</td><td>0.0</td><td>0.458</td><td>7.147</td><td>54.2</td><td>6.0622</td><td>3.0</td><td>222.0</td><td>18.7</td><td>396.90</td><td>5.33</td></tr></table>		CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98	1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14	2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03	3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94	4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33						
	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT																																																																														
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98																																																																														
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14																																																																														
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03																																																																														
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94																																																																														
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33																																																																														
In [5]:	boston_df['PRICE'] = boston.target boston_df.head()																																																																																										
Out[5]:	<table><tr><th></th><th>CRIM</th><th>ZN</th><th>INDUS</th><th>CHAS</th><th>NOX</th><th>RM</th><th>AGE</th><th>DIS</th><th>RAD</th><th>TAX</th><th>PTRATIO</th><th>B</th><th>LSTAT</th><th>PRICE</th></tr><tr><td>0</td><td>0.00632</td><td>18.0</td><td>2.31</td><td>0.0</td><td>0.538</td><td>6.575</td><td>65.2</td><td>4.0900</td><td>1.0</td><td>296.0</td><td>15.3</td><td>396.90</td><td>4.98</td><td>24.0</td></tr><tr><td>1</td><td>0.02731</td><td>0.0</td><td>7.07</td><td>0.0</td><td>0.469</td><td>6.421</td><td>78.9</td><td>4.9671</td><td>2.0</td><td>242.0</td><td>17.8</td><td>396.90</td><td>9.14</td><td>21.6</td></tr><tr><td>2</td><td>0.02729</td><td>0.0</td><td>7.07</td><td>0.0</td><td>0.469</td><td>7.185</td><td>61.1</td><td>4.9671</td><td>2.0</td><td>242.0</td><td>17.8</td><td>392.83</td><td>4.03</td><td>34.7</td></tr><tr><td>3</td><td>0.03237</td><td>0.0</td><td>2.18</td><td>0.0</td><td>0.458</td><td>6.998</td><td>45.8</td><td>6.0622</td><td>3.0</td><td>222.0</td><td>18.7</td><td>394.63</td><td>2.94</td><td>33.4</td></tr><tr><td>4</td><td>0.06905</td><td>0.0</td><td>2.18</td><td>0.0</td><td>0.458</td><td>7.147</td><td>54.2</td><td>6.0622</td><td>3.0</td><td>222.0</td><td>18.7</td><td>396.90</td><td>5.33</td><td>36.2</td></tr></table>		CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	PRICE	0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98	24.0	1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14	21.6	2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03	34.7	3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94	33.4	4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33	36.2
	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	PRICE																																																																													
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98	24.0																																																																													
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14	21.6																																																																													
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03	34.7																																																																													
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94	33.4																																																																													
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33	36.2																																																																													

In [3]: 데이터셋에 대한 설명 `boston.DESCR` 을 확인

In [4]: 데이터셋 객체의 data 배열 `boston.data`, 즉 독립 변수 X가 되는 피쳐들을 DataFrame 자료형으로 변환하여 `boston_df`를 생성, `boston_df`의 데이터 5개를 확인 `boston_df.head()`

In [5]: 데이터셋 객체의 target 배열 `boston.target`, 즉 종속 변수인 주택 가격('PRICE') 컬럼을 `boston_df`에 추가 `boston_df`의 데이터 5개를 확인 `boston_df.head()`

2. 회귀분석(boston 주택가격 예측하기)

❖ 데이터 수집, 준비 및 탐색

- 데이터가 이미 정리된 상태이므로 데이터셋 구성을 확인

In [6]:	print('보스톤 주택 가격 데이터셋 크기: ', boston_df.shape)
Out[6]:	보스톤 주택 가격 데이터셋 크기: (506, 14)
In [7]:	boston_df.info()
Out[7]:	<div><div><class 'pandas.core.frame.DataFrame'> RangeIndex: 506 entries, 0 to 505 Data columns (total 14 columns): CRIM 506 non-null float64 ZN 506 non-null float64 INDUS 506 non-null float64 CHAS 506 non-null float64 NOX 506 non-null float64 RM 506 non-null float64 AGE 506 non-null float64 DIS 506 non-null float64 RAD 506 non-null float64 TAX 506 non-null float64 PTRATIO 506 non-null float64 B 506 non-null float64 LSTAT 506 non-null float64 PRICE 506 non-null float64 dtypes: float64(14) memory usage: 55.4KB</div><div>In [6]: 데이터셋의 형태 boston_df.shape, 즉 행의 개수(데이터 개수)와 열의 개수(변수 개수)를 확인, 행의 개수가 506이므로 데이터가 506개 있으며, 열의 개수가 14이므로 변수가 14개 있음 변수 중에서 13개는 독립 변수 x가 되고, 마지막 변수 'PRICE'는 종속 변수 y가 됨 In [7]: boston_df에 대한 정보를 확인</div></div>

■ 14개의 독립 변수(피처)의 의미

- CRIM: 지역별 범죄 발생률
- ZN: 25,000평방피트를 초과하는 거주 지역 비율
- INDUS: 비상업 지역의 넓이 비율
- CHAS: 찰스강의 더미변수(1은 강에 인접, 0은 강에 인접 아님)
- NOX: 일산화질소 농도
- RM: 거주할 수 있는 방 개수
- AGE: 1940년 이전에 건축된 주택 비율
- DIS: 5개 주요 고용센터까지 가중 거리
- RAD: 고속도로 접근 용이도
- TAX: 10,000달러당 재산세 비율
- PTRATIO: 지역의 교사와 학생 수 비율
- B: 지역의 흑인 거주 비율
- LSTAT: 하위 계층의 비율
- PRICE(MEDV): 본인 소유 주택 가격의 중앙값

2. 회귀분석(boston 주택가격 예측하기)

❖ 분석 모델 구축, 결과 분석 및 시각화

- 선형 회귀를 이용해 분석 모델 구축하기

- 사이킷런의 선형 분석 모델 패키지 `sklearn.linear_model`에서 선형 회귀 `LinearRegression`를 이용하여 분석 모델을 구축

In [8]:	<pre>from sklearn.linear_model import LinearRegression from sklearn.model_selection import train_test_split from sklearn.metrics import mean_squared_error, r2_score</pre>
In [9]:	<pre>#X, Y 분할하기 Y = boston_df['PRICE'] X = boston_df.drop(['PRICE'], axis = 1, inplace = False)</pre>
In [10]:	<pre><i>#훈련용 데이터와 평가용 데이터 분할하기</i> X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.3, random_state = 156)</pre>

In [8]: 사이킷런을 사용하여 머신러닝 회귀 분석을 하기 위한 `LinearRegression`과 데이터셋 분리 작업을 위한 `train_test_split`, 성능 측정을 위한 평가 지표인 `mean_squared_error, r2_score`를 임포트

In [9]: PRICE 피처를 회귀식의 종속 변수 Y로 설정하고 PRICE를 제외 `drop()`한 나머지 피처를 독립 변수 X로 설정

In [10]: X와 Y 데이터 506개를 학습 데이터와 평가 데이터로 7:3 비율로 분할 `test_size=0.3`

2. 회귀분석(boston 주택가격 예측하기)

❖ 분석 모델 구축, 결과 분석 및 시각화

- 선형 회귀를 이용해 분석 모델 구축하기
 - 사이킷런의 선형 분석 모델 패키지 `sklearn.linear_model`에서 선형 회귀 `LinearRegression`를 이용하여 분석 모델을 구축

In [11]:	#선형 회귀 분석 : 모델 생성 <code>lr = LinearRegression()</code>
In [12]:	#선형 회귀 분석 : 모델 훈련 <code>lr.fit(X_train, Y_train)</code>
Out[12]:	<code>LinearRegression()</code>
In [13]:	#선형 회귀 분석 : 평가 데이터에 대한 예측 수행 -> 예측 결과 <code>Y_predict</code> 구하기 <code>Y_predict = lr.predict(X_test)</code>

In [11]: 선형 회귀 분석 모델 객체 `lr`을 생성

In [12]: 학습 데이터 `XX_train`와 `YY_train`를 가지고 학습을 수행 `fit()`.

In [13]: 평가 데이터 `XX_test`를 가지고 예측을 수행하여 `predict()` 예측값 `YY_predict`를 구함

2. 회귀분석(boston 주택가격 예측하기)

❖ 분석 모델 구축, 결과 분석 및 시각화

- 선형 회귀를 이용해 분석 모델 구축하기
 - 선형 회귀 분석 모델을 평가 지표를 통해 평가하고 회귀 계수를 확인하여 피처의 영향을 분석

In [14]:	<pre>mse = mean_squared_error(Y_test, Y_predict) rmse = np.sqrt(mse) print('MSE : {0:.3f}, RMSE : {1:.3f}'.format(mse, rmse)) print('R^2(Variance score) : {0:.3f}'.format(r2_score(Y_test, Y_predict)))</pre>
Out[14]:	<pre>MSE : 17.297, RMSE : 4.159 R^2(Variance score) : 0.757</pre>
In [15]:	<pre>print('Y 절편 값: ', lr.intercept_) print('회귀 계수 값: ', np.round(lr.coef_, 1))</pre>
Out[15]:	<pre>Y 절편 값: 40.995595172164336 회귀 계수 값: [-0.1 0.1 0. 3. -19.8 3.4 0. -1.7 0.4 -0. -0.9 0. -0.6]</pre>

In [14]: 회귀 분석은 지도 학습이므로 평가 데이터 x에 대한 결과값 \hat{y} 를 이미 알고 있는 상태에서 평가 데이터 y 와 In [13]에서 구한 예측 결과 \hat{y} 의 오차를 계산하여 모델을 평가. 평가 지표 MSE를 구하고 `mean_squared_error()` 구한 값의 제곱근을 계산하여 `np.sqrt(mse)` 평가 지표 RMSE를 구함 그리고 평가 지표 R2 을 구함 `r2_score()`

In [15]: 선형 회귀의 Y절편 `lr.intercept_`과 각 피처의 회귀 계수 `lr.coef_`를 확인

2. 회귀분석(boston 주택가격 예측하기)

❖ 분석 모델 구축, 결과 분석 및 시각화

- 선형 회귀를 이용해 분석 모델 구축하기
 - 선형 회귀 분석 모델을 평가 지표를 통해 평가하고 회귀 계수를 확인하여 피처의 영향을 분석

```
In [16]: coef = pd.Series(data = np.round(lr.coef_, 2), index = X.columns)
coef.sort_values(ascending = False)
```

```
Out[16]: RM      3.35
CHAS      3.05
RAD       0.36
ZN        0.07
INDUS     0.03
B         0.01
AGE       0.01
TAX      -0.01
CRIM     -0.11
LSTAT   -0.57
PTRATIO -0.92
DIS      -1.74
NOX     -19.80
dtype: float64
```

In [16]: 회귀 모델에서 구한 회귀 계수 값 `lr.coef_` 과 피처 이름 `X.columns` 을 묶어서 Series 자료 형으로 만들고, 회귀 계수 값을 기준으로 내림차순으로 정렬하여 `ascending=False` 확인 `sort_values()`

- 회귀 모델 결과를 토대로 보스톤 주택 가격에 대한 회귀식

$$\text{YPRICE} = -0.11\text{XCRIM} + 0.07\text{XZN} + 0.03\text{XINDUS} + 3.05\text{XCHAS} - 19.80\text{XNOX} + 3.35\text{XRM} + 0.01\text{XAGE} - 1.74\text{XDIS} + 0.36\text{XRAD} - 0.01\text{XTAX} - 0.92\text{XPTRATIO} + 0.01\text{XB} - 0.57\text{XLSTAT} + 41.00$$

2. 회귀분석(boston 주택가격 예측하기)

❖ 회귀 분석 결과를 산점도 + 선형 회귀 그래프로 시각화하기

- 선형 회귀를 이용해 분석 모델 구축하기

In [17]:	<pre>import matplotlib.pyplot as plt import seaborn as sns</pre>
In [18]:	<pre>fig, axs = plt.subplots(figsize = (16, 16), ncols = 3, nrows = 5) x_features = ['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'B', 'LSTAT'] for i, feature in enumerate(x_features): row = int(i/3) col = i%3 sns.regplot(x = feature, y = 'PRICE', data = boston_df, ax = axs[row][col])</pre>

In [17]: 시각화에 필요한 모듈을 임포트

In [18]: 독립 변수인 13개 피처와 종속 변수인 주택 가격, PRICE와의 회귀 관계를 보여주는 13개 그래프를 subplots()를 사용하여 5행 3열 구조로 모아서 나타냄
seaborn의 regplot()은 산점도 그래프와 선형 회귀 그래프를 함께 그려줌

2. 회귀분석(boston 주택가격 예측하기)

❖ 회귀 분석 결과를 산점도 + 선형 회귀 그래프로 시각화하기

■ 선형 회귀를 이용해 분석 모델 구축하기

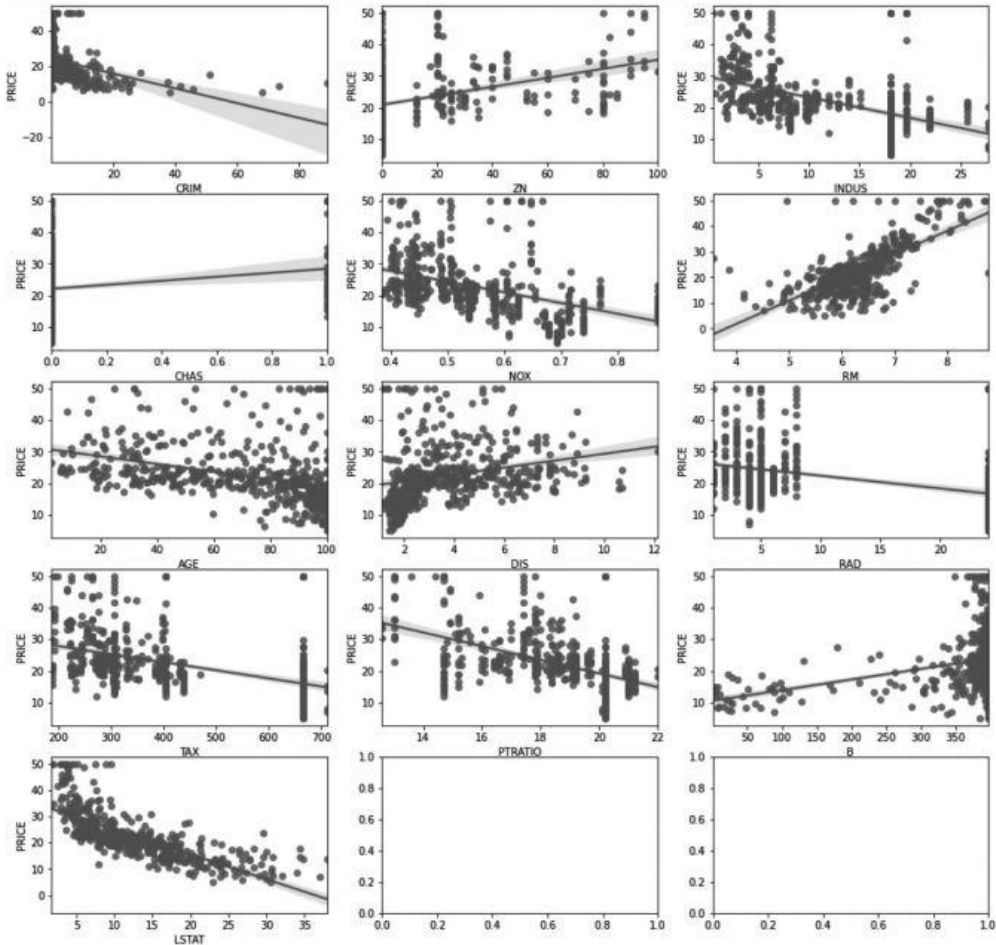


그림 10-3 13개 피처와 주택 가격의 회귀 관계를 나타낸 산점도/선형 회귀 그래프

1. 분류 분석(iris 붓꽃분류)

❖ 라이브러리 및 데이터 셋 로딩

```
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report
from sklearn import metrics
from sklearn.metrics import confusion_matrix #오차행렬은 sklearn.metrics 패키지 내 confusion_matrix로 확인 가능
```

```
iris = load_iris() #load_iris import
iris_data = iris.data #중요 데이터 iris_data 변수에 저장 후 데이터 크기 확인
iris_label = iris.target #iris_label 변수에 iris 데이터의 target 저장
```

```
print(iris)
print(iris_data)
print(iris_label)
```

❖ 훈련데이터와 테스트 데이터 분류

#(3)train, test 데이터 분리

#X:feature데이터만 / y:정답label데이터만

#X데이터셋을 머신러닝 모델에 입력 -> 모델이 내뱉는 품종 예측 결과를 정답 y와 비교하여 학습

X_train, X_test, y_train, y_test = train_test_split(iris_data, #feature(입력받는 특징 데이터)

iris_label, #label(모델이 맞춰야하는 정답값)

test_size=0.2, #test dataset 크기 조절(0.2=전체 20%를 test데이터로 사용)

random_state=7) #train데이터와 test데이터 분리시 적용되는 랜덤성)

```
print(X_train.shape)
```

```
print(X_test.shape)
```

```
print(y_train.shape)
```

```
print(y_test.shape)
```

❖ 모델 훈련 및 테스트

```
logistic_model = LogisticRegression()  
logistic_model.fit(X_train, y_train)  
y_pred = logistic_model.predict(X_test)  
scores=metrics.accuracy_score(y_test, y_pred)  
print(scores)  
  
print(classification_report(y_test, y_pred))  
print(confusion_matrix(y_test, y_pred))
```