

AI 보안 기술개발 교육

머신러닝(Machine Learning)의 이해

AI 보안 기술개발 교육

머신러닝의 이해

1. Scikit-Learn 라이브러리
2. 데이터 분할
3. 교차 검증
4. 탐색적 데이터 분석

Scikit-Learn 라이브러리

- 머신러닝 알고리즘을 구현한 오픈소스 라이브러리 중 가장 유명한 라이브러리 중 하나
- 일관되고 간결한 API가 강점이며, 문서화가 잘되어 있음
- 알고리즘은 파이썬 클래스로 구현되고, 데이터 셋은 Numpy 배열, Pandas DataFrame 및 SciPy 희소행렬 등 사용가능



scikit-learn
Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Scikit-Learn 데이터 표현 – Feature Matrix

- 표본(sample)은 데이터셋이 설명하는 개별 객체를 나타냄
- 특징(feature)은 각 표본을 연속적인 수치값, 부울값, 이산값으로 표현하는 개별 관측치를 의미
- 표본 : 행렬의 행
- 행의 개수 : `n_samples`
- 특징(feature) : 행렬의 열
- 열의 개수 : `n_features`
- 관례적으로 특징행렬은 변수 `X`에 저장
- `[n_samples, n_features]` 형태의 2차원 배열 구조를 사용
(주로 Numpy 배열, Pandas DataFrame, SciPy 희소행렬을 사용)

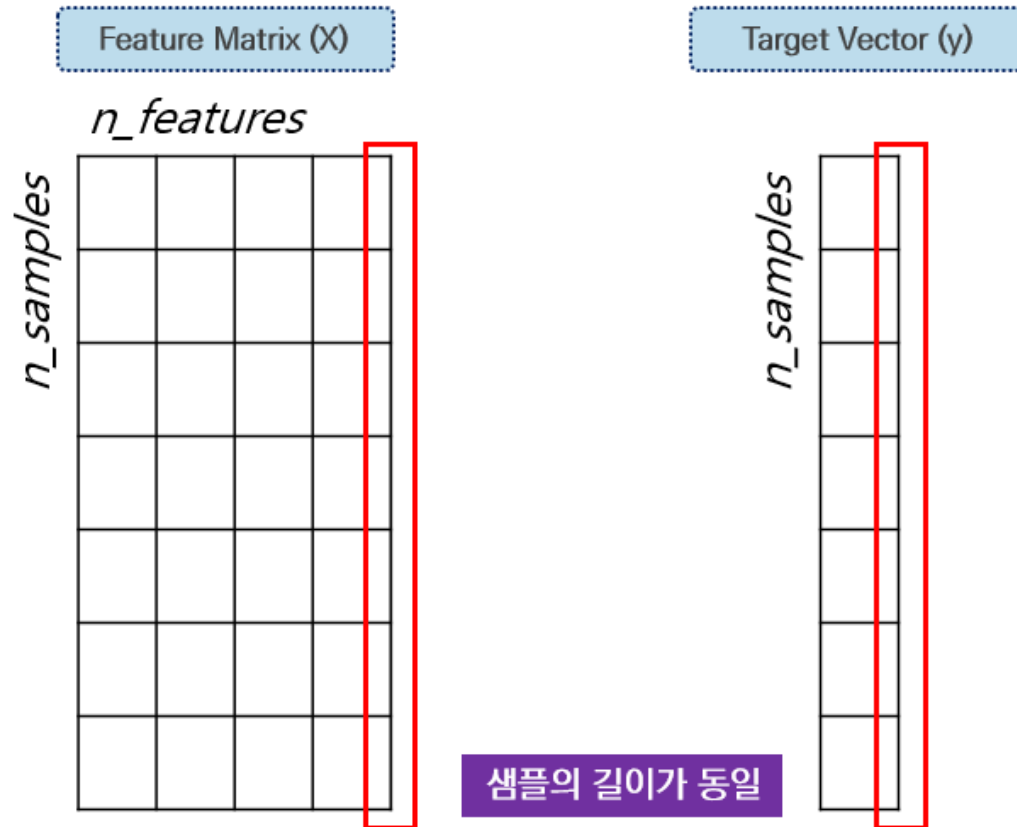
Scikit-Learn 데이터 표현 – Target Vector

- 연속적인 수치값, 이산 클래스/레이블을 가짐
- 길이 : `n_samples`
- 관례적으로 대상벡터는 변수 `y`에 저장
- 1차원 배열 구조를 사용 (주로 Numpy 배열, Pandas Series를 사용)
- 특징 행렬로부터 예측하고자 하는 값의 벡터
- 종속 변수, 출력 변수, 결과 변수, 반응 변수 라고도 함

Scikit-Learn 데이터 표현 – Target Vector

- 연속적인 수치값, 이산 클래스/레이블을 가짐
- 길이 : `n_samples`
- 관례적으로 대상벡터는 변수 `y`에 저장
- 1차원 배열 구조를 사용 (주로 Numpy 배열, Pandas Series를 사용)
- 특징 행렬로부터 예측하고자 하는 값의 벡터
- 종속 변수, 출력 변수, 결과 변수, 반응 변수 라고도 함

Scikit-Learn 데이터 표현 - Layout



Scikit-Learn 데이터 표현 - Layout

```
import seaborn as sd
```

```
iris = sd.load_dataset('iris')
```

```
iris.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 150 entries, 0 to 149  
Data columns (total 5 columns):  
sepal_length    150 non-null float64  
sepal_width     150 non-null float64  
petal_length    150 non-null float64  
petal_width     150 non-null float64  
species         150 non-null object  
dtypes: float64(4), object(1)  
memory usage: 5.9+ KB
```

```
iris.head()
```

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa

Feature Matrix와 Target Vector를 나누는 작업 필요

Scikit-Learn 데이터 표현 - Layout

```
X = iris.drop('species', axis=1)
X.head()
```

“species” 열이 삭제된 X 데이터프레임(Feature Matrix) 생성

	sepal_length	sepal_width	petal_length	petal_width
0	5.1	3.5	1.4	0.2
1	4.9	3.0	1.4	0.2
2	4.7	3.2	1.3	0.2
3	4.6	3.1	1.5	0.2
4	5.0	3.6	1.4	0.2

```
y = iris['species']
y.head()
```

“species” 열을 추출하여 y Series(Target Vector) 생성

```
0    setosa
1    setosa
2    setosa
3    setosa
4    setosa
Name: species, dtype: object
```

Scikit-Learn 데이터 표현 - Layout

```
X = iris.drop('species', axis=1)
X.head()
```

“species” 열이 삭제된 X 데이터프레임(Feature Matrix) 생성

	sepal_length	sepal_width	petal_length	petal_width
0	5.1	3.5	1.4	0.2
1	4.9	3.0	1.4	0.2
2	4.7	3.2	1.3	0.2
3	4.6	3.1	1.5	0.2
4	5.0	3.6	1.4	0.2

```
y = iris['species']
y.head()
```

“species” 열을 추출하여 y Series(Target Vector) 생성

```
0    setosa
1    setosa
2    setosa
3    setosa
4    setosa
Name: species, dtype: object
```

Scikit-Learn을 활용한 머신러닝 수행 절차

- 1 데이터 준비
- 2 모델 클래스 선택
- 3 모델 인스턴스 생성과 하이퍼파라미터 선택
- 4 특징 행렬과 대상 벡터 준비
- 5 모델을 데이터에 적합
- 6 새로운 데이터를 이용해 예측
- 7 모델 평가

Scikit-Learn을 활용한 머신러닝 수행 절차

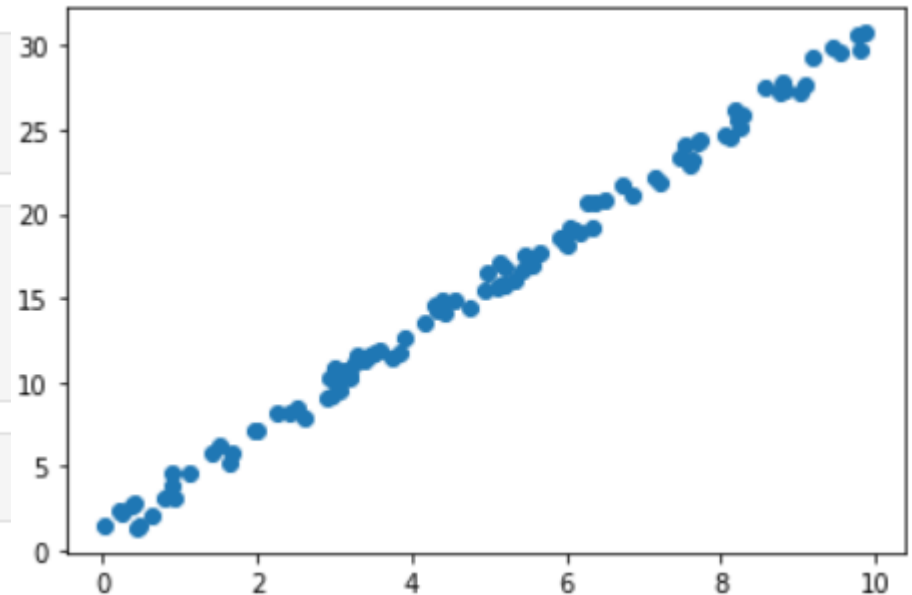
1

데이터 준비

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

```
rs = np.random.RandomState(10)
X = 10 * rs.rand(100)
y = 3 * X + 2 * rs.rand(100)
```

```
plt.scatter(X, y)
```



Scikit-Learn을 활용한 머신러닝 수행 절차

2

모델 클래스 선택

3

모델 인스턴스 생성과 하이퍼파라미터 선택

입력데이터(x), 출력데이터(y)가 모두 연속형 수치 데이터이므로 그에 맞는 분석 모델을 선택

```
from sklearn.linear_model import LinearRegression  
regr = LinearRegression()
```

선형회귀 객체(인스턴스) 생성 - 디폴트

```
from sklearn.linear_model import LinearRegression  
regr = LinearRegression(fit_intercept = True)
```

선형회귀 객체(인스턴스) 생성 -
(fit_intercept=True라는 하이퍼파라미터를 제공)

Scikit-Learn을 활용한 머신러닝 수행 절차

4

특징 행렬과 대상 벡터 준비

```
X = x.reshape(-1, 1)
print(X.shape, y.shape)
```

(100, 1) (100,)

5

모델을 데이터에 적합

```
regr.fit(X, y)
```

X, y에 맞는 선형회귀 모델을 적합(모델 생성)

```
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

```
regr.coef_
```

모델의 기울기

```
array([2.9855087])
```

```
regr.intercept_
```

모델의 y 절편

```
0.9878534341975644
```

Scikit-Learn을 활용한 머신러닝 수행 절차

6

새로운 데이터를 이용해 예측

```
x_new = np.linspace(-1, 11, num=100)
```

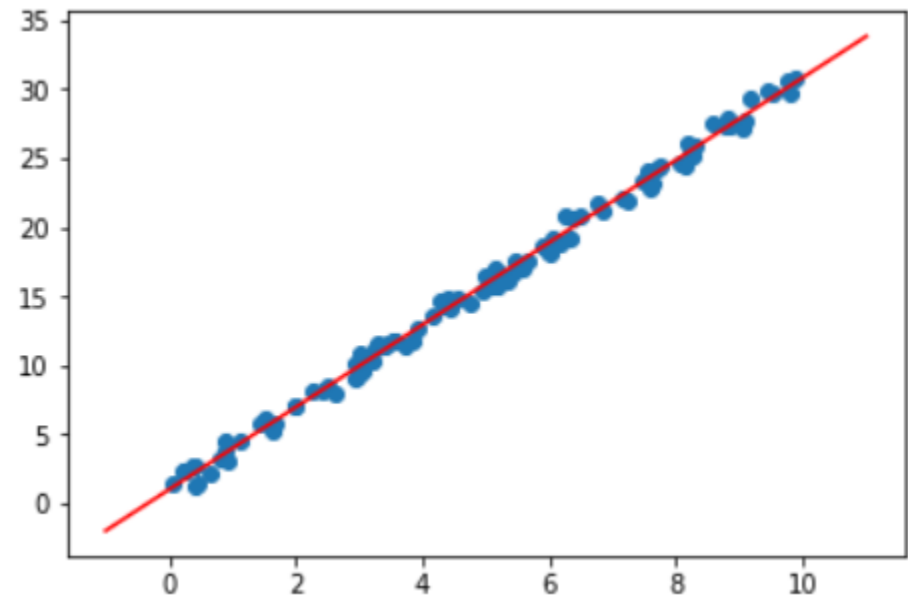
```
X_new = x_new.reshape(-1, 1)  
X_new.shape
```

```
(100, 1)
```

```
y_pred = regr.predict(X_new)
```

새로 입력된 X_new에 대한
모델 예측값(y_pred) 생성

```
plt.plot(X_new, y_pred, c='red')  
plt.scatter(X, y)
```



Scikit-Learn을 활용한 머신러닝 수행 절차

7

모델 평가

```
from sklearn.metrics import mean_squared_error
```

```
rmse = np.sqrt(mean_squared_error(y, y_pred))
```

```
print(rmse)
```

```
13.708237122486333
```


데이터 레이블링

```
import seaborn as sd
import pandas as pd
```

```
iris = sd.load_dataset('iris')
```

```
X = iris.drop('species', axis=1)
y = iris['species']
```

```
iris['species'].value_counts()
```

```
virginica    50
setosa       50
versicolor  50
```

```
Name: species, dtype: int64
```

```
from sklearn.preprocessing import LabelEncoder
```

```
encoder = LabelEncoder()
y = encoder.fit_transform(y)
```

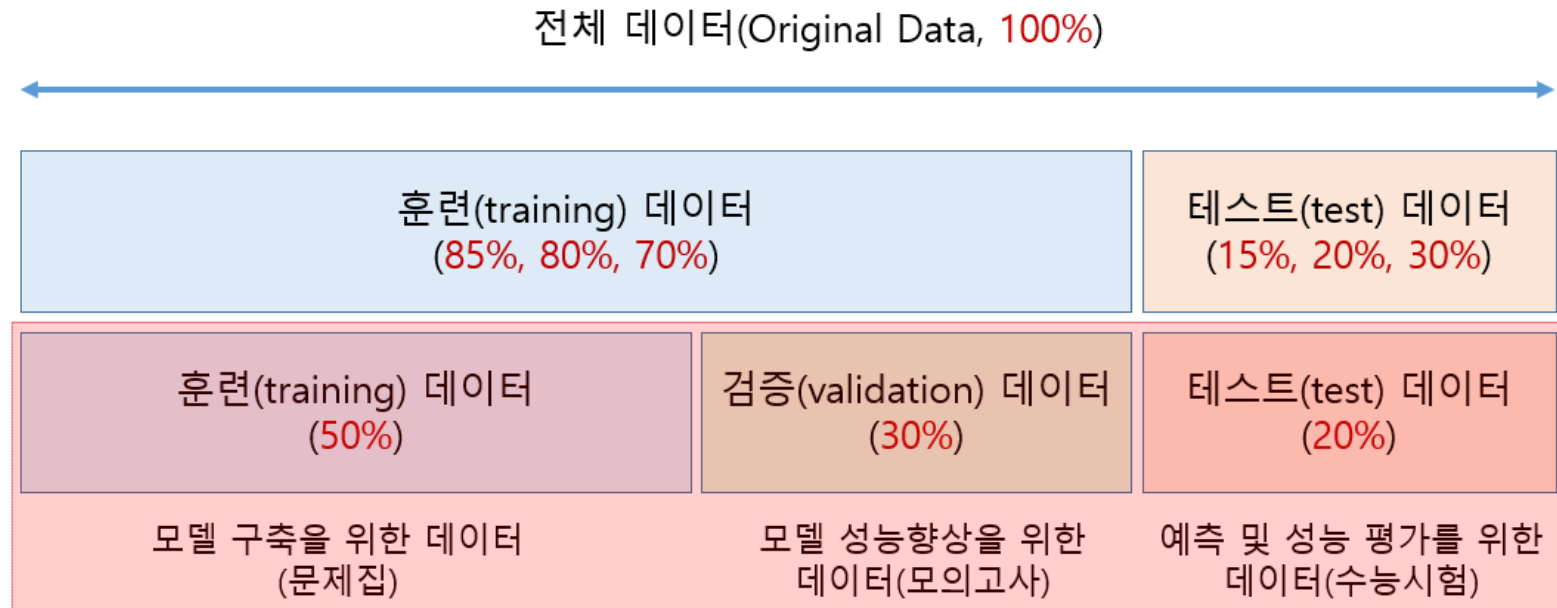
```
y
```

```
array([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, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
       2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
       2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2], dtype=int64)
```

TargetVector 에 대한 Labeling 필요

데이터 분할

머신러닝/딥러닝 학습데이터는 훈련(training) 데이터, 검증(validation) 데이터 및 테스트(test)로 분할(split)하여 사용



훈련 데이터 : 모델의 훈련 및 가중치 업데이트 등의 목적으로 사용

검증 데이터 : 훈련된 모델의 평가 및 최종 모델을 선정하기 위해 사용

테스트 데이터 : 모델의 예측 및 평가를 위해 사용

데이터 분할

```
from sklearn.model_selection import train_test_split
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state=25)
```

```
print(X_train.shape, X_test.shape, y_train.shape, y_test.shape)
```

```
(120, 4) (30, 4) (120,) (30,)
```

```
pd.Series(y_train).value_counts()
```

```
2    42  
0    41  
1    37
```

```
dtype: int64
```

데이터 분할

- 랜덤 층화 샘플링

- 각 층별 배분된 표본을 단순임의추출의 방법으로 표본을 추출

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, shuffle=True, stratify=y, random_state=25)
```

```
pd.Series(y_train).value_counts()
```

```
2    40  
1    40  
0    40
```

```
dtype: int64
```

교차 검증(Cross Validation)

- 교차검증의 개념 및 절차

- ① 교차 검증 1단계에서는 데이터를 학습용과 테스트용으로 나눔
- ② 모델의 테스트 성능을 기록
- ③ 교차 검증의 매 단계마다 다른 파티션으로 위의 작업을 수행
- ④ 모델의 최종 성능은 매 단계의 테스트 성능을 평균 계산

교차 검증은 모델의 변동성을 줄여주며 오버피팅 방지 효과

교차 검증을 통해 모든 데이터를 학습용 데이터로 사용할 수 있음

교차 검증(Cross Validation)

- k 폴드 교차 검증(K-fold Cross Validation)

- 데이터를 무작위로 k개의 동일한 크기인 폴드로 분할 (보통 k값으로 3, 5, 10을 많이 사용)
- 각 시행 단계에서 특정 폴드를 테스트용으로, 나머지는 학습용으로 사용
- 각 폴드를 테스트 세트로 한 번씩 사용하고 이 과정을 k번 반복 시행함
- 최종적으로 모델 성능의 평균을 계산

- $k = 5$, repeat = 5

반복시행	폴드1	폴드2	폴드3	폴드4	폴드5
1	테스트	학습	학습	학습	학습
2	학습	테스트	학습	학습	학습
3	학습	학습	테스트	학습	학습
4	학습	학습	학습	테스트	학습
5	학습	학습	학습	학습	테스트

교차 검증(Cross Validation)

```
from sklearn.datasets import load_iris
iris = load_iris()
X = iris.data
y = iris.target
```

```
from sklearn.neighbors import KNeighborsClassifier
model = KNeighborsClassifier(n_neighbors=3)
```

```
from sklearn.model_selection import cross_val_score
cross_val_score(model, X, y, cv = 5)
```

```
array([0.96666667, 0.96666667, 0.93333333, 0.96666667, 1.        ])
```

교차 검증(Cross Validation)

- 단일 관측치 제거 방식(LOOCV)

- Leave-one-out cross validation

검증을 시행할 때 마다 한 지점을 제외한 모든 지점에서 훈련

- 매 시행 단계에서 테스트 샘플을 고정하는 방식

- 데이터를 n 개의 서브세트로 분할하고, n 개 중 1개를 테스트용으로 두고 $n-1$ 개로 학습을 수행

- 데이터 크기가 n 이면 n 번의 교차 검증을 수행

교차 검증(Cross Validation)

```
from sklearn.model_selection import LeaveOneOut
scores = cross_val_score(model, X, y, cv = LeaveOneOut())
scores
```

```
array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,  
       1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,  
       1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,  
       1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,  
       1., 1., 0., 1., 0., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1.,  
       1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,  
       1., 1., 1., 1., 0., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,  
       0., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1.,  
       1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.] )
```

Grid Search

```
from sklearn.datasets import load_iris
import numpy as np
iris = load_iris()
X = iris.data
y = iris.target
```

```
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import GridSearchCV
```

```
params = {"fit_intercept": [True, False],
          "normalize": [True, False]}
grid = GridSearchCV(LinearRegression(), params, iid=True, cv=7)
```

Grid Search

```
grid.fit(X, y)
```

```
GridSearchCV(cv=7, error_score='raise-deprecating',  
             estimator=LinearRegression(copy_X=True, fit_intercept=True,  
                                         n_jobs=None, normalize=False),  
             iid=True, n_jobs=None,  
             param_grid={'fit_intercept': [True, False],  
                          'normalize': [True, False]},  
             pre_dispatch='2*n_jobs', refit=True, return_train_score=False,  
             scoring=None, verbose=0)
```

```
grid.best_params_
```

```
{'fit_intercept': True, 'normalize': False}
```

```
model = grid.best_estimator_  
model
```

```
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

탐색적 데이터 분석(Exploratory Data Analysis)

- 탐색적 데이터 분석(Exploratory Data Analysis)이란?

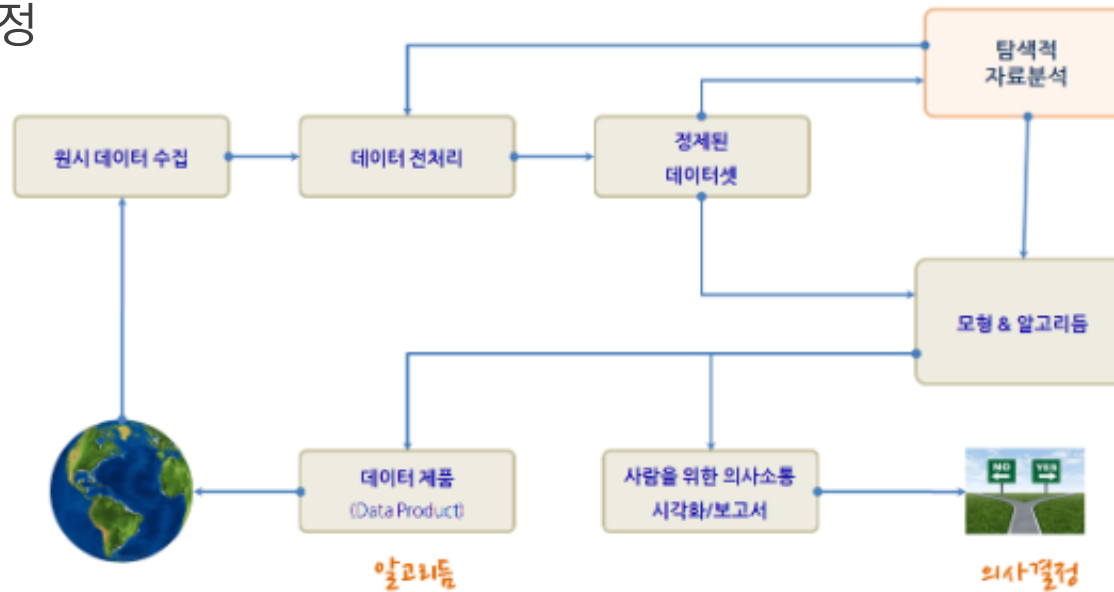
존 튜키라는 미국의 저명한 통계학자가 창안한 자료 분석 방법론

데이터가 가지고 있는 본연의 특징과 의미를 찾는 것을 목적으로 함

데이터를 다양한 각도에서 관찰하고 이해하는 과정

데이터를 분석하기 전에 통계적인 방법이나 시각화 도구를 활용하여 데이터를 직관적으로

파악해 보는 과정



출처 : <https://statklee.github.io/ml/ml-eda.html>

EDA – 와인 품질 데이터

- 관측값 : 총 6,497건 (레드 와인: 1,599건, 화이트 와인: 4,898건)
- 입력변수 : 12개 (고정산, 휘발산, 구연산, 잔여당, 염화물, 무수아황산, 총이산화황, 밀도, 산성도, 황산염, 알콜도수와 같은 와인의 물리화학적 특성들과 red, white의 와인 타입)
- 출력변수 : 1개 (와인품질평가점수, 가장 낮은 품질 1점 ~ 가장 높은 품질 10점)
- 데이터 수집
 - <https://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-red.csv>
 - <https://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-white.csv>

EDA – 와인 품질 데이터

▪ 데이터 파일 read 및 데이터프레임 구성 확인

```
redwine = pd.read_csv("data/winequality-red.csv", sep=";", header=0)
redwine["type"] = "red"
redwine.head()
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality	type
0	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5	red
1	7.8	0.88	0.00	2.6	0.098	25.0	67.0	0.9968	3.20	0.68	9.8	5	red
2	7.8	0.76	0.04	2.3	0.092	15.0	54.0	0.9970	3.26	0.65	9.8	5	red
3	11.2	0.28	0.56	1.9	0.075	17.0	60.0	0.9980	3.16	0.58	9.8	6	red
4	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5	red

```
redwine.shape
```

```
(1599, 13)
```

EDA – 와인 품질 데이터

▪ 데이터 파일 read 및 데이터프레임 구성 확인

```
whitewine = pd.read_csv("data/winequality-white.csv", sep=";", header=0)
whitewine["type"] = "white"
whitewine.head()
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality	type
0	7.0	0.27	0.36	20.7	0.045	45.0	170.0	1.0010	3.00	0.45	8.8	6	white
1	6.3	0.30	0.34	1.6	0.049	14.0	132.0	0.9940	3.30	0.49	9.5	6	white
2	8.1	0.28	0.40	6.9	0.050	30.0	97.0	0.9951	3.26	0.44	10.1	6	white
3	7.2	0.23	0.32	8.5	0.058	47.0	186.0	0.9956	3.19	0.40	9.9	6	white
4	7.2	0.23	0.32	8.5	0.058	47.0	186.0	0.9956	3.19	0.40	9.9	6	white

```
whitewine.shape
```

```
(4898, 13)
```

EDA – 와인 품질 데이터

- red/white 데이터프레임 병합 및 컬럼명 변경

```
wine = redwine.append(whitewine)
```

```
wine.shape
```

```
(6497, 13)
```

```
wine.columns = wine.columns.str.replace(' ', '_')  
wine.head()
```

	fixed_acidity	volatile_acidity	citric_acid	residual_sugar	chlorides	free_sulfur_dioxide	total_sulfur_dioxide	density	pH	sulphates	alcohol	quality	type
0	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5	red
1	7.8	0.88	0.00	2.6	0.098	25.0	67.0	0.9968	3.20	0.68	9.8	5	red
2	7.8	0.76	0.04	2.3	0.092	15.0	54.0	0.9970	3.26	0.65	9.8	5	red
3	11.2	0.28	0.56	1.9	0.075	17.0	60.0	0.9980	3.16	0.58	9.8	6	red
4	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5	red

EDA – 와인 품질 데이터

▪ 요약통계 확인

```
wine.describe()
```

	fixed_acidity	volatile_acidity	citric_acid	residual_sugar	chlorides	free_sulfur_dioxide	total_sulfur_dioxide	density	pH	sulphates
count	6497.000000	6497.000000	6497.000000	6497.000000	6497.000000	6497.000000	6497.000000	6497.000000	6497.000000	6497.000000
mean	7.215307	0.339666	0.318633	5.443235	0.056034	30.525319	115.744574	0.994697	3.218501	0.531268
std	1.296434	0.164636	0.145318	4.757804	0.035034	17.749400	56.521855	0.002999	0.160787	0.148806
min	3.800000	0.080000	0.000000	0.600000	0.009000	1.000000	6.000000	0.987110	2.720000	0.220000
25%	6.400000	0.230000	0.250000	1.800000	0.038000	17.000000	77.000000	0.992340	3.110000	0.430000
50%	7.000000	0.290000	0.310000	3.000000	0.047000	29.000000	118.000000	0.994890	3.210000	0.510000
75%	7.700000	0.400000	0.390000	8.100000	0.065000	41.000000	156.000000	0.996990	3.320000	0.600000
max	15.900000	1.580000	1.660000	65.800000	0.611000	289.000000	440.000000	1.038980	4.010000	2.000000

EDA – 와인 품질 데이터

▪ 출력변수 quality 확인

```
pd.DataFrame(wine.quality.value_counts())
```

quality	
6	2836
5	2138
7	1079
4	216
8	193
3	30
9	5

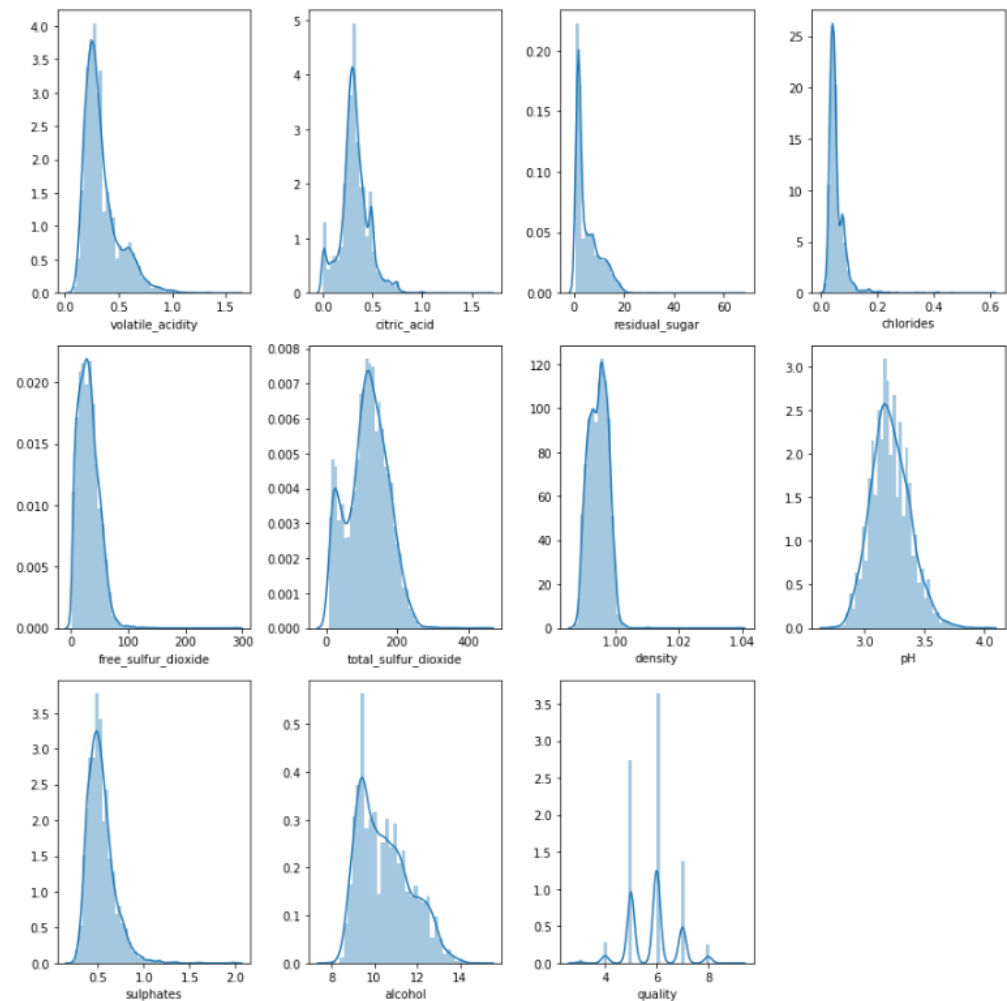
```
wine.groupby('type')['quality'].value_counts().unstack('type')
```

type	red	white
quality		
3	10.0	20.0
4	53.0	163.0
5	681.0	1457.0
6	638.0	2198.0
7	199.0	880.0
8	18.0	175.0
9	NaN	5.0

EDA – 와인 품질 데이터

■ 각 변수별 분포

```
plt.figure(figsize=(12,12))
for i in range(1,12):
    plt.subplot(3,4,i)
    sns.distplot(wine.iloc[:,i])
plt.tight_layout()
plt.show()
```



EDA – 와인 품질 데이터

■ 상관분석

```
wine_corr = wine.corr()  
wine_corr
```

상관계수 R은 $-1 \leq R \leq 1$ 범위
1에 가까우면 양의 상관관계
-1에 가까우면 음의 상관관계
0에 가까우면 상관관계 없음

	fixed_acidity	volatile_acidity	citric_acid	residual_sugar	chlorides	free_sulfur_dioxide	total_sulfur_dioxide	density
fixed_acidity	1.000000	0.219008	0.324436	-0.111981	0.298195	-0.282735	-0.329054	0.458910
volatile_acidity	0.219008	1.000000	-0.377981	-0.196011	0.377124	-0.352557	-0.414476	0.271296
citric_acid	0.324436	-0.377981	1.000000	0.142451	0.038998	0.133126	0.195242	0.096154
residual_sugar	-0.111981	-0.196011	0.142451	1.000000	-0.128940	0.402871	0.495482	0.552517
chlorides	0.298195	0.377124	0.038998	-0.128940	1.000000	-0.195045	-0.279630	0.362615
free_sulfur_dioxide	-0.282735	-0.352557	0.133126	0.402871	-0.195045	1.000000	0.720934	0.025717
total_sulfur_dioxide	-0.329054	-0.414476	0.195242	0.495482	-0.279630	0.720934	1.000000	0.032395
density	0.458910	0.271296	0.096154	0.552517	0.362615	0.025717	0.032395	1.000000
pH	-0.252700	0.261454	-0.329808	-0.267320	0.044708	-0.145854	-0.238413	0.011686
sulphates	0.299568	0.225984	0.056197	-0.185927	0.395593	-0.188457	-0.275727	0.259478
alcohol	-0.095452	-0.037640	-0.010493	-0.359415	-0.256916	-0.179838	-0.265740	-0.686745
quality	-0.076743	-0.265699	0.085532	-0.036980	-0.200666	0.055463	-0.041385	-0.305858

EDA – 와인 품질 데이터

■ 상관분석 – 양의 상관관계

```
pd.DataFrame(wine_corr.loc[wine_corr.quality > 0, 'quality'])
```

	quality
citric_acid	0.085532
free_sulfur_dioxide	0.055463
pH	0.019506
sulphates	0.038485
alcohol	0.444319
quality	1.000000

양의 상관 관계를 가지는 변수들 (구연산, 무수아황산, 산성도, 황산염, 알콜도수)

EDA – 와인 품질 데이터

▪ 상관분석 – 음의 상관관계

```
pd.DataFrame(wine_corr.loc[wine_corr.quality < 0, 'quality'])
```

	quality
fixed_acidity	-0.076743
volatile_acidity	-0.265699
residual_sugar	-0.036980
chlorides	-0.200666
total_sulfur_dioxide	-0.041385
density	-0.305858

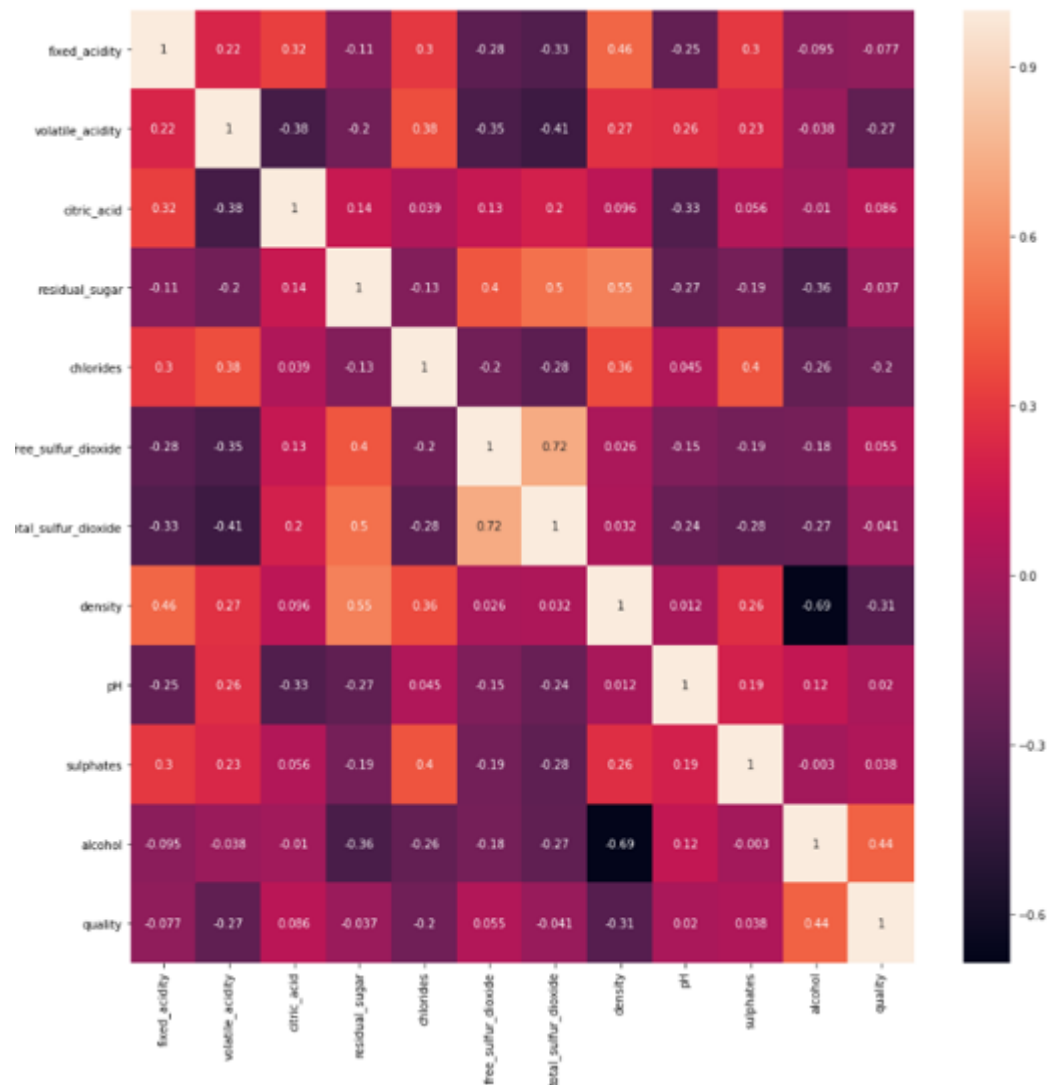
음의 상관 관계를 가지는 변수들 (고정산, 휘발산, 잔여당, 염화물, 총이산화황, 밀도)

EDA – 와인 품질 데이터

상관분석 – heat map

```
import seaborn as sns

plt.figure(figsize=(15,15))
sns.heatmap(wine_corr,color = "k", annot=True)
```

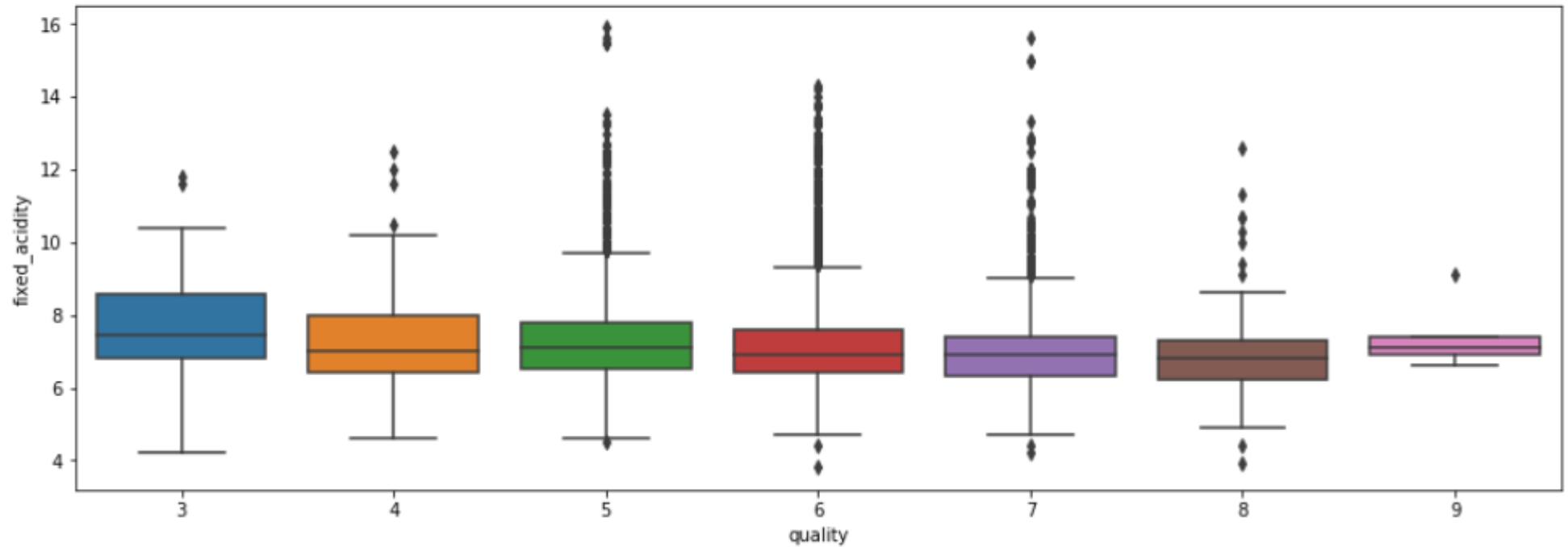


EDA – 와인 품질 데이터

■ 상관분석 – box plot

```
plt.figure(figsize=(15,5))  
sns.boxplot(x="quality", y="fixed_acidity", data=wine)
```

<matplotlib.axes._subplots.AxesSubplot at 0x21038e85240>



Q&A



Thank you

