

# 01. 기본- 결정트리(decision tree)

- date : 22/05

## 데이터 셋 다운로드

- UCI : <https://www.kaggle.com/uciml/pima-indians-diabetes-database>
- 피마 인디언 당뇨병 발병 예측

(가) decision tree는 classification(분류)와 regression(회귀) 문제에 널리 사용하는 모델이다.

(나) 스무고개 놀이의 질문과 비슷하다.

```
In [46]: # 라이브러리 불러오기
import pandas as pd
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
```

## Data Fields

구분		설명
Pregnancies		임신
Glucose		포도당
BloodPressure		혈압
SkinThickness		피부두께
Insulin		인슐린
BMI	BMI	111
Diabetes Pedigree Function		당뇨병혈통기능
Age		나이
Outcome		결과

```
In [47]: pima = pd.read_csv("diabetes.csv")
```

```
In [48]: pima.columns
```

```
Out[48]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
               'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
              dtype='object')
```

```
In [49]: pima.head()
```

```
Out[49]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
0	6	148	72	35	0	33.6	0.627	50
1	1	85	66	29	0	26.6	0.351	31
2	8	183	64	0	0	23.3	0.672	32
3	1	89	66	23	94	28.1	0.167	21

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
4	0	137	40	35	168	43.1	2.288	33

## Feature Selection

In [50]: `pima.columns`

Out[50]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'], dtype='object')

In [51]: `pima.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Pregnancies                          768 non-null    int64
1   Glucose                             768 non-null    int64
2   BloodPressure                       768 non-null    int64
3   SkinThickness                      768 non-null    int64
4   Insulin                            768 non-null    int64
5   BMI                                 768 non-null    float64
6   DiabetesPedigreeFunction            768 non-null    float64
7   Age                                768 non-null    int64
8   Outcome                            768 non-null    int64
dtypes: float64(2), int64(7)
memory usage: 54.1 KB
```

In [52]: `pima.head(3)`

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
0	6	148	72	35	0	33.6	0.627	50
1	1	85	66	29	0	26.6	0.351	31
2	8	183	64	0	0	23.3	0.672	32

In [53]: `pima.columns`

Out[53]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'], dtype='object')

```
In [54]: # 데이터 셋 (feature와 target 변수로 나누기)
feature_cols = ['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
               'BMI', 'DiabetesPedigreeFunction', 'Age']
X = pima[feature_cols] # Features
y = pima.Outcome       # Target variable
```

## 데이터 나누기

In [55]: `# 데이터 셋 나누기`  
`X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=`

In [56]: `print(X_test.columns)`

```
print(X_train.columns)
print(y_train.shape)
```

```
Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
      'BMI', 'DiabetesPedigreeFunction', 'Age'],
      dtype='object')
Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
      'BMI', 'DiabetesPedigreeFunction', 'Age'],
      dtype='object')
(537,)
```

```
In [57]: # 의사결정 트리 모델 생성 및 학습
model = DecisionTreeClassifier(max_depth=5, random_state=0).fit(X_train,y_train)

# 예측
y_pred = model.predict(X_test)
```

## 모델 평가

```
In [58]: from sklearn import metrics
```

```
In [59]: # Model Accuracy, 얼마나 정확한가? 정확도
print("Accuracy:",metrics.accuracy_score(y_test, y_pred))
```

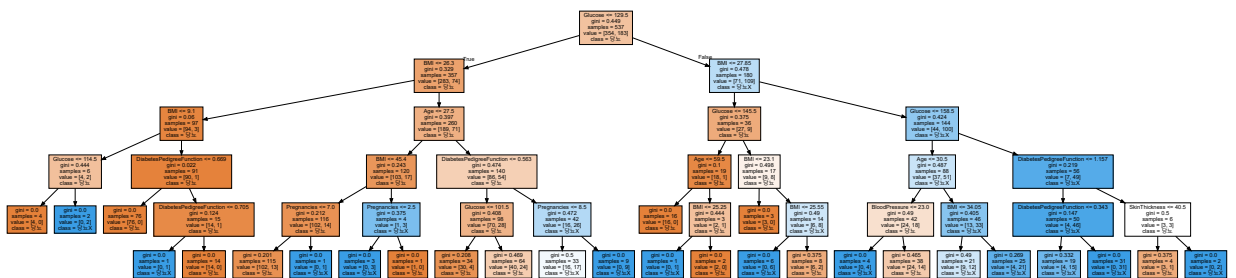
Accuracy: 0.7575757575757576

## 시각화 1

```
In [60]: from sklearn.tree import export_graphviz
export_graphviz(model, out_file="tree.dot",
                class_names=['당뇨', '당뇨X'],
                feature_names = feature_cols,
                impurity = True, # gini 계수
                filled=True)    # color

import graphviz
with open("tree.dot") as f:
    dot_graph = f.read()

display(graphviz.Source(dot_graph))
```



## 시각화 2

- png파일로 만들기

```
In [61]: import pydotplus
from IPython.display import Image
```

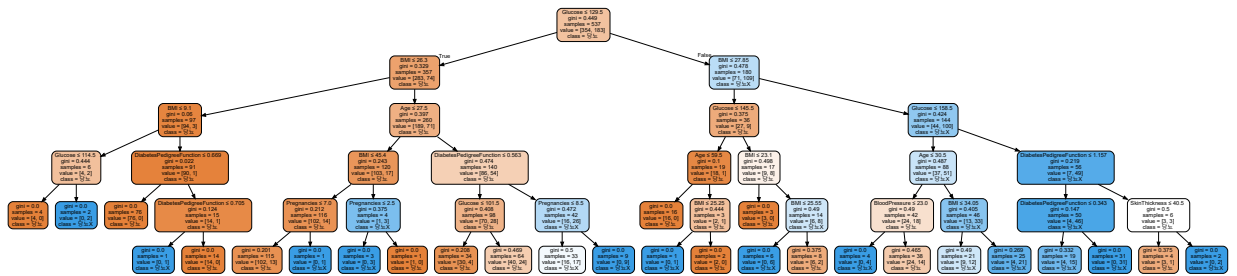
```
In [62]: import graphviz

# model : 모델명,
# class_n : 클래스명,
# feature_n : 특징 이름
```

```
def tree_plot(model, class_n, feature_n):
    export_graphviz(model, out_file="tree.dot",
                    class_names = class_n,
                    feature_names = feature_n,
                    impurity = True, # gini 계수
                    filled=True,
                    rounded=True,
                    special_characters=True) # color

    with open("tree.dot") as f:
        dot_graph = f.read()
        display(graphviz.Source(dot_graph))

tree_plot(model, ['당뇨', '당뇨X'], feature_cols)
```



## 모델 성능 개선

```
In [63]: model = DecisionTreeClassifier(criterion="entropy",
                                       max_depth=3,
                                       random_state=0) # 의사결정트리 모델

model.fit(X_train,y_train) # 학습
y_pred = model.predict(X_test) # 데이터 셋 예측

# 정확도 확인
print("Accuracy:", metrics.accuracy_score(y_test, y_pred))
```

Accuracy: 0.7705627705627706