

12주차 과제

Classifier



담당교수님: 윤장혁 교수님

과목명: Data Analytics

이름: 박민성

학번: 201611145

전공: 산업공학과

제출일: 2020.06.05

Object

- 다양한 binary classifier를 학습하여 결과를 확인
 - ✓ 여러 classifier 모델을 테스트해보고 결과를 함께 첨부 (2개 이상)
 - ✓ 최종적으로 선택한 모델은 test accuracy 기준 55% 이상을 만족

Classifier Model 1 - Logistic Regression Model

```
C: > Users > minisong > Desktop > 12주차 2.py > ...
1  import pandas as pd
2  import seaborn as sns
3  import matplotlib.pyplot as plt
4  import numpy as np
5  df=pd.read_csv('./data_week12.csv')
6  #sns.countplot(x='Class variable (0 or 1).', data=df)
7  #plt.show()
8  from sklearn.linear_model import LogisticRegression
9  from sklearn.model_selection import train_test_split
10 Y=df['Class variable (0 or 1).']
11 X=df[['Number of times pregnant.','Plasma glucose concentration a 2 hours
12 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.3)
13 log_clf = LogisticRegression()
14 log_clf.fit(X_train,Y_train)
15 print(log_clf.score(X_test, Y_test)*100, "%")
```

76.62337662337663 %

전체 자료 중 70%를 학습시키고 30%로 테스트를 진행한 결과 정확도는 약 76.6%가 나왔다.

Classifier Model 2 – Naïve Bayes Model

```
> Users > minisong > Desktop > 12주차 1.py > ...
1  import pandas as pd
2  from sklearn.model_selection import train_test_split
3  from sklearn.naive_bayes import GaussianNB
4  from sklearn.metrics import accuracy_score
5  df=pd.read_csv('./data_week12.csv')
6  Y=df['Class variable (0 or 1).']
7  X=df[['Number of times pregnant.','Plasma glucose concentration a 2 hours
8  X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.3)
9  model = GaussianNB()
10 model.fit(X_train,Y_train)
11 accuracy = model.score(X_test,Y_test)*100
12 print(accuracy,"%")
```

77.05627705627705 %

마찬가지로 전체 자료 중 70%를 학습시키고 30%로 테스트를 진행한 결과 정확도는 약 77.1%가 나왔다.

Final Selection Model – Naïve Bayes Model

두 모델 다 test accuracy 가 55% 이상으로 조건에 부합하지만 최종적으로는 정확도가 조금 더 높은 Naïve Bayes Model을 선정하기로 하였다.