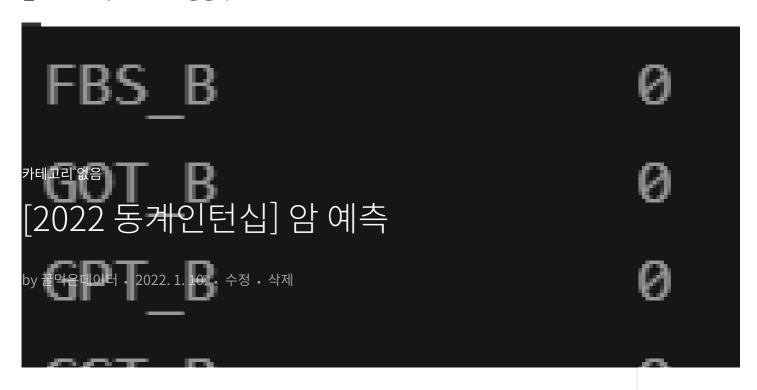
## 데이터스트





홈 태그 방명록



import pandas as pd
data=pd.read\_csv('C:/Users/dudtj/0

data=pd.read\_csv('C:/Users/dudtj/OneDrive - 숭실대학교 - Soongsi I University/Desktop/CL/python/동계인턴십\_Data\_1000/phenotype\_10 00.txt',engine="python",sep=" ")

print(data.info())

## 분류 전체보기 🔯

파이썬 머신러닝 완벽 가 이드

AI 데이터 연구단

공지사항

최근글 : 인기글

[2022 동계 인턴십]암… 2022.01.12



[동계인턴 십] 암 예··· 2022.01.11



[2022 동계인턴십] 암 예측 2022.01.10

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 66 columns):
    Column
                    Non-Null Count Dtype
                                                34
                                                   PCAN89
                                                                   18 non-null
                                                                                    float64
                                                   FCAN80
                                                                   30 non-null
                                                                                    float64
                    1000 non-null
                                    object
                                                  FCAN81
                                                                   69 non-null
                                                                                    float64
                                               36
    IID
                    1000 non-null
                                    object
                                                   FCAN82
                                                                   1 non-null
                                                                                    float64
     AGE_B
                    1000 non-null
                                    int64
                                                   FCAN83
                                                                                    float64
                                               38
                                                                   23 non-null
    SMOK B
                    990 non-null
                                     float64
                                                   FCAN84
                                                                   20 non-null
                                                                                    float64
    SMOKA MOD B
                    679 non-null
                                    float64
                                               40 FCAN86
                                                                   35 non-null
                                                                                    float64
    ALCO_B
                    965 non-null
                                     float64
                                                                   11 non-null
                                                   FCAN89
                                                                                    float64
    ALCO_AMOUNT_B 915 non-null
                                               42
                                                   FFV1
                                                                   569 non-null
                                                                                    float64
    EXER B
                    927 non-null
                                    float64
                                                                   569 non-null
                                                                                    float64
                                                  FVC
    MDM B
                    46 non-null
                                     float64
                                                                                    float64
                                               44
                                                   BIL
                                                                   948 non-null
    MHTN B
                                     float64
                    120 non-null
                                                   WBC
                                                                   928 non-null
                                                                                    float64
10
    MLPD_B
                    28 non-null
                                    float64
                                               46
                                                   CREAT
                                                                   962 non-null
                                                                                    float64
    PHTN B
                    153 non-null
                                     float64
                                               47
                                                   STOMA
                                                                   1000 non-null
                                                                                    int64
    PDM R
                    67 non-null
                                     float64
                                               48
                                                   COLON
                                                                   1000 non-null
                                                                                    int64
    PLPD B
                    36 non-null
                                     float64
                                                   LIVER
                                                                   1000 non-null
                                                                                    int64
                    998 non-null
                                     float64
    нт в
                                               50
                                                                   1000 non-null
                                                   LUNG
                                                                                    int64
    WT B
                    998 non-null
                                     float64
                                                   PROST
                                                                   1000 non-null
                                                                                    int64
    WAIST B
                    971 non-null
                                     float64
                                                   THROI
                                                                   1000 non-null
                                                                                    int64
    SBP B
                    988 non-null
                                     float64
                                                   BREAC
                                                                   1000 non-null
                                                                                    int64
    DBP B
                    989 non-null
                                     float64
                                                   RECTM
                                                                   1000 non-null
                                                                                    int64
19 CHO B
                    995 non-null
                                     float64
                                                                   1000 non-null
                                                                                    float64
                                                   SCOLON
                    995 non-null
                                     float64
 20 LDL B
                                                                   1000 non-null
                                                                                    float64
                                                   SRECTM
     TG B
                    995 non-null
                                     float64
                                                                                    float64
                                               57
                                                   SPROST
                                                                   1000 non-null
    HDL B
                    987 non-null
                                               58
                                                   STHROI
                                                                   1000 non-null
                                                                                    float64
   FBS_B
                    998 non-null
                                     float64
                                               59
                                                   SBREAC
                                                                   1000 non-null
                                                                                    float64
 24
    GOT B
                    996 non-null
                                     float64
                                                                                    float64
                                               60
                                                   SLUNG
                                                                   1000 non-null
25 GPT B
                    996 non-null
                                     float64
                                                   SSTOMA
                                                                   1000 non-null
                                                                                    float64
 26 GGT_B
                    992 non-null
                                     float64
                                                   SLIVER
                                                                   1000 non-null
                                                                                    float64
    URIC_B
                    929 non-null
                                     float64
                                               63 SEX1
                                                                   1000 non-null
                                                                                   int64
 28
    PCAN80
                    4 non-null
                                     float64
                                                   CRC
                                                                   1000 non-null
                                                                                    int64
    PCAN81
                    15 non-null
                                     float64
                                               65 SCRC
                                                                   1000 non-null
                                                                                    float64
    PCAN82
                    0 non-null
                                     float64
 30
                                              dtypes: float64(53), int64(11), object(2)
    PCAN83
                    1 non-null
                                     float64
                                              memory usage: 515.8+ KB
    PCAN84
 32
                    8 non-null
                                     float64
                                              None
```

위 그림과 같이 null값이 데이터의 절반 이상을 차지하는 변수들이 있다. 데이터 보존을 위해 non-null값이 900 이상인 데이터들만 가져와 다시 df로 저장하였다. (1000개의 데이터이므로 90%만 뽑았다.)

이렇게 할 수 있는 이유는 drop시킨 변수들은 암 예측에 중요하지 않다고 판단한 변수들이었기 때문이다.

data.columns

df=data.loc[:,['AGE\_B','SMOK\_B','ALCO\_B','ALCO\_AMOUNT\_B','EXER\_
B','HT\_B','WT\_B','WAIST\_B','SBP\_B','DBP\_B','CHO\_B','LDL\_B','TG\_
B','HDL\_B','FBS\_B','GOT\_B','GPT\_B','GGT\_B','URIC\_B','BIL','WBC',
'CREAT','STOMA','COLON','LIVER','LUNG','PROST','THROI','BREAC',
'RECTM','SCOLON','SRECTM','SPROST','STHROI','SBREAC','SLUNG','SS
TOMA','SLIVER','SEX1','CRC','SCRC']]



[FIND-A] 금 융경제학… 2022.01.09



Sleep AI Challen… 2022.01.07



최근댓글

태그

사이킷런, 머신러닝, 파이썬머신러닝완벽가이 드, 수면다원검사, 수면, AI데이터연구단

\_\_\_

전체 방문자

3

Today: 0 Yesterday: 0

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 41 columns):
                  Non-Null Count Dtype
    Column
#
                 1000 non-null
                                int64
0
    AGE B
            990 non-null
                                float64
    SMOK B
1
                 965 non-null
                                float64
2
    ALCO B
    ALCO AMOUNT B 915 non-null
                               float64
 3
            927 non-null
4
    EXER B
                                float64
 5
                 998 non-null
                               float64
    HT B
                 998 non-null
 6
    WT B
                                float64
                971 non-null
988 non-null
7
                                float64
    WAIST B
8
    SBP B
                               float64
                989 non-null
 9
    DBP B
                                float64
                995 non-null
                                float64
 10 CHO B
                 995 non-null
 11
    LDL B
                                float64
                 995 non-null
 12 TG B
                                float64
                987 non-null
 13 HDL B
                                float64
                998 non-null
 14 FBS B
                                float64
 15 GOT B
                 996 non-null
                               float64
 16 GPT B
                996 non-null
                                float64
 17 GGT B
                 992 non-null
                                float64
                929 non-null
 18 URIC B
                                float64
                 948 non-null
 19 BIL
                                float64
 20 WBC
                 928 non-null
                                float64
 21 CREAT
                  962 non-null
                                float64
                  1000 non-null
 22 STOMA
                                int64
                  1000 non-null
 23 COLON
                               int64
 24 LIVER
                  1000 non-null
                                int64
 25 LUNG
                  1000 non-null
                               int64
                  1000 non-null
                                int64
 26 PROST
                  1000 non-null
 27 THROI
                                int64
 28
                  1000 non-null
                                int64
    BREAC
                  1000 non-null
                                int64
 29 RECTM
 30 SCOLON
                  1000 non-null
                                float64
                  1000 non-null
                                float64
 31
    SRECTM
                  1000 non-null
                                float64
 32 SPROST
 33 STHROI
                  1000 non-null
                               float64
                  1000 non-null
 34 SBREAC
                                float64
 35 SLUNG
                  1000 non-null
                                float64
 36 SSTOMA
                  1000 non-null
                                float64
                  1000 non-null
                                float64
 37
    SLIVER
```

```
39 CRC 1000 non-null int64
40 SCRC 1000 non-null float64
dtypes: float64(30), int64(11)
memory usage: 320.4 KB
None
```

이와 같이 null값이 많지 않은 변수들을 df에 저장한 뒤 dropna()함수를 통해 그 데이터를 날리고 검정을 하였다. 그 이유는 제거를 한 뒤 데이터의 양은 747개로 충분하다고 생각하여 이를 test와 train set으로 분리하여 진행하였다.

```
df.dropna(inplace=True)
print(df.isnull().sum())
```

AGE_B	0
SMOK_B	0
ALCO_B	0
ALCO_AMOUNT_B	0
EXER_B	0
нт_в	0
WT_B	0
WAIST_B	0
SBP_B	0
DBP_B	0
СНО_В	0
LDL_B	0
TG_B	0
HDL_B	0
FBS_B	0
GOT_B	0
GPT_B	0
GGT_B	0
URIC_B	0
BIL	0
WBC	0
CREAT	0
STOMA	0
COLON	0
LIVER	0
SLIVER	0
SEX1	0
CRC	0
SCRC	0
dtype: int64	

```
from sklearn.preprocessing import LabelEncoder
def format_features(df):
    features=['STOMA','COLON','LIVER','LUNG','PROST','THROI','BR
EAC','RECTM']
    for feature in features:
        le=LabelEncoder()
        le=le.fit(df[feature])
        df[feature]=le.transform(df[feature])
    return df
```

features를 암종류로 놓고 y를 'LUNG'암으로 두고 X를 이 features를 제거 한 변수로 둔 뒤 진행하였다.

```
features=['STOMA','COLON','LIVER','LUNG','PROST','THROI','BREAC'
,'RECTM']
y_df = df['LUNG']
X_df =df.drop(features, axis=1)
```

train set과 test set의 비율은 8:2 비율로 진행하였다.

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X_df, y_df, test_
size=0.2)
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
dt clf = DecisionTreeClassifier()
rf_clf = RandomForestClassifier()
Ir_clf = LogisticRegression()
#DecisionTreeClassifier 학습/예측/평가
dt_clf.fit(X_train,y_train)
dt_pred = dt_clf.predict(X_test)
print("DecisionTreeClassifier 정확도 :",accuracy_score(y_test,dt
_pred))
#RandomForestClassifier 학습/예측/평가
rf_clf.fit(X_train,y_train)
rf_pred = rf_clf.predict(X_test)
print("RandomForestClassifier 정확도 :",accuracy_score(y_test,rf
_pred))
#LogisticRegression 학습/예측/평가
Ir_clf.fit(X_train,y_train)
Ir_pred = Ir_clf.predict(X_test)
print("LogisticRegression 정확도 :",accuracy_score(y_test, |r_pre
d))
DecisionTreeClassifier 정확도: 0.986666666666667
RandomForestClassifier 정확도: 0.95333333333333333
LogisticRegression 정확도 : 0.9666666666666667
C:\Users\dudtj\anaconda3\lib\site-packages\sklearn\linear model\ logistic.py:763:
ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
Increase the number of iterations (max_iter) or scale the data as shown in:
   https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
   https://scikit-learn.org/stable/modules/linear model.html#logistic-regression
 n_iter_i = _check_optimize_result(
```

위와 같이 DecisionTree나 RandomForest, LogisticRegression이 모두 95%보다 높은 정확도를 볼 수 있었다.

```
import numpy as np
from sklearn.model_selection import KFold
def exec_kfold(clf,folds=5):
   kfold=KFold(n_splits=folds)
   scores=[]
```

## 교차 검증을 통한 정확도 또한 95%를 넘는 정확도로 상당히 높은 수치다.

```
from sklearn.model_selection import cross_val_score
scores=cross_val_score(dt_clf,X_df,y_df,cv=5)
for iter_count,accuracy in enumerate(scores):
    print("교차 검증 정확도:",iter_count,accuracy)
print("평균 정확도:",np.mean(scores))
```

```
from sklearn.model_selection import GridSearchCV parameters= {'max_depth':[2,3,5,10], 'min_samples_split':[2,3,5], 'min_samples_leaf':[1,5,8]} grid_dclf=GridSearchCV(dt_clf,param_grid=parameters,scoring='acc uracy',cv=5) grid_dclf.fit(X_train,y_train) print("GridSearchCV 최적 하이퍼 파라미터:",grid_dclf.best_params__)
```

```
print("GridSearchCV 최고 정확도:",grid_dclf.best_score_)
best_dclf = grid_dclf.best_estimator_

dpredictions = best_dclf.predict(X_test)
accuracy = accuracy_score(y_test,dpredictions)
print("테스트 세트에서의 DecisionTreeClassifier 정확도:",accuracy)
```

```
GridSearchCV 최적 하이퍼 파라미터: {'max_depth': 2, 'min_samples_leaf': 1, 'min_samples_split': 2}
GridSearchCV 최고 정확도: 0.9849439775910364
테스트 세트에서의 DecisionTreeClassifier 정확도: 0.973333333333333333
```

최적 하이퍼 파라미터는 max\_depth가 2 ,min\_samples\_leaf가 1, min\_samples\_split이 2가 나왔고, 이에 대한 최고 정확도는 98.49%가 나왔다.

♡ 공감 🖒 🛍 ····

## 댓글 0

여러분의 소중한 댓글을 입력해주세요.

등록

( 1 2 **3** 4 5 6 7 ··· 15 )

TEL. 02.1234.5678 / 경기 성남시 분당구 판교역로 © Kakao Corp.

