

캐글 코리아 4차 대회

학습 내용

- 라벨 인코딩 적용
- 다양한 모델 성능 비교

목차

[01. 라이브러리 임포트 및 데이터 준비](#)

[02. 데이터 전처리](#)

[03. 모델 구축하기](#)

01. 라이브러리 임포트 및 데이터 준비

[목차로 이동하기](#)

In [37]:

```
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import warnings

warnings.filterwarnings('ignore')
```

In [38]:

```
train = pd.read_csv('data/4th_kaggle/train.csv')
test = pd.read_csv('data/4th_kaggle/test.csv')
sub = pd.read_csv('data/4th_kaggle/sample_submission.csv')
```

데이터 탐색

- 컬럼명 : `df.columns`
- 행열 : `df.shape`
- 정보 : `df.info()`
- 수치 데이터 요약정보 : `df.describe()`
- 결측치 : `df.isnull().sum()`

데이터 정보

age : 나이
workclass : 고용 형태
fnlwgt : 사람 대표성을 나타내는 가중치 (final weight의 약자)
education : 교육 수준 (최종 학력)
education_num : 교육 수준 수치
marital_status: 결혼 상태
occupation : 업종
relationship : 가족 관계
race : 인종
sex : 성별
capital_gain : 양도 소득
capital_loss : 양도 손실
hours_per_week : 주당 근무 시간
native_country : 국적
income : 수익 (예측해야 하는 값, target variable)

In [39]:

```
print("학습용 데이터 : ", train.shape)  
print("테스트용 데이터 : ", test.shape)
```

학습용 데이터 : (26049, 16)
테스트용 데이터 : (6512, 15)

In [40]:

```
y = train['income']  
test['income'] = "blank"
```

In [41]:

```
all_dat = pd.concat([train, test], axis=0)  
print(all_dat.shape)
```

(32561, 16)

In [42]:

```
all_dat.income.value_counts()
```

Out[42]:

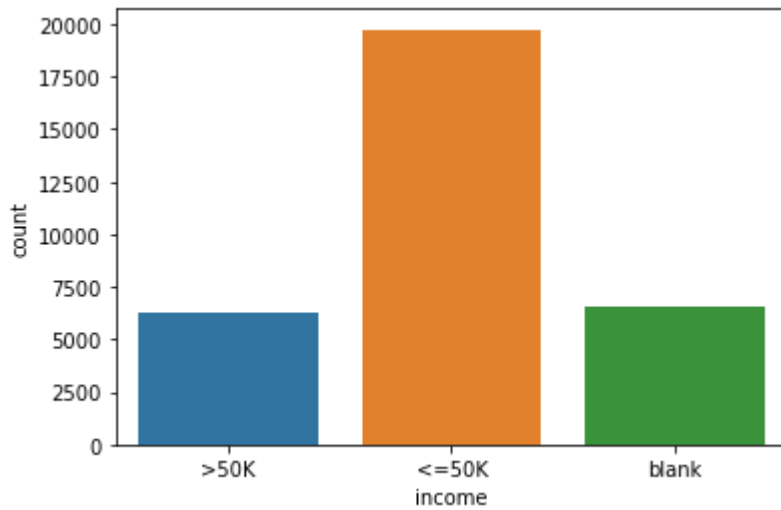
```
<=50K      19744  
blank       6512  
>50K       6305  
Name: income, dtype: int64
```

In [43]:

```
sns.countplot(x="income", data=all_dat)
```

Out[43]:

<AxesSubplot:xlabel='income', ylabel='count'>



02. 데이터 전처리

[목차로 이동하기](#)

In [44]:

```
all_dat.loc[ all_dat['income']=='>50K' , 'target'] = 1
all_dat.loc[ all_dat['income']=='<=50K' , 'target'] = 0
all_dat.loc[ all_dat['income']=='blank' , 'target'] = 999
all_dat['target'] = all_dat.target.astype("int")
```

In [45]:

```
all_dat.head()
```

Out[45]:

	id	age	workclass	fnlwgt	education	education_num	marital_status	occupation	relationsh
0	0	40	Private	168538	HS-grad	9	Married-civ-spouse	Sales	Husbai
1	1	17	Private	101626	9th	5	Never-married	Machine-op-inspct	Own-ch
2	2	18	Private	353358	Some-college	10	Never-married	Other-service	Own-ch
3	3	21	Private	151158	Some-college	10	Never-married	Prof-specialty	Own-ch
4	4	24	Private	122234	Some-college	10	Never-married	Adm-clerical	Not-i fam

In [46]:

```
all_dat.columns
```

Out[46]:

```
Index(['id', 'age', 'workclass', 'fnlwgt', 'education', 'education_num',  
      'marital_status', 'occupation', 'relationship', 'race', 'sex',  
      'capital_gain', 'capital_loss', 'hours_per_week', 'native_count  
ry',  
      'income', 'target'],  
      dtype='object')
```

라벨 인코딩

In [47]:

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

In [48]:

```
en_x = LabelEncoder()
all_dat['workclass_lbl'] = en_x.fit_transform(all_dat['workclass'])
all_dat.head(3)
```

Out[48]:

	id	age	workclass	fnlwgt	education	education_num	marital_status	occupation	relationsh
0	0	40	Private	168538	HS-grad	9	Married-civ-spouse	Sales	Husbai
1	1	17	Private	101626	9th	5	Never-married	Machine-op-inspct	Own-ch
2	2	18	Private	353358	Some-college	10	Never-married	Other-service	Own-ch

In [49]:

```
all_dat['education_lbl'] = en_x.fit_transform(all_dat['education'])
all_dat['marital_status_lbl'] = en_x.fit_transform(all_dat['marital_status'])
all_dat['occupation_lbl'] = en_x.fit_transform(all_dat['occupation'])
all_dat['relationship_lbl'] = en_x.fit_transform(all_dat['relationship'])
all_dat['race_lbl'] = en_x.fit_transform(all_dat['race'])
all_dat['native_country_lbl'] = en_x.fit_transform(all_dat['native_country'])
all_dat.head(3)
```

Out[49]:

	id	age	workclass	fnlwgt	education	education_num	marital_status	occupation	relationsh
0	0	40	Private	168538	HS-grad	9	Married-civ-spouse	Sales	Husbai
1	1	17	Private	101626	9th	5	Never-married	Machine-op-inspct	Own-ch
2	2	18	Private	353358	Some-college	10	Never-married	Other-service	Own-ch

3 rows × 24 columns

In [50]:

```
all_dat['sex'].unique()
```

Out[50]:

```
array(['Male', 'Female'], dtype=object)
```

In [51]:

```
mf_mapping = {"Male": 1, "Female": 2}
all_dat['sex'] = all_dat['sex'].map(mf_mapping)
all_dat.head(3)
```

Out[51]:

	id	age	workclass	fnlwgt	education	education_num	marital_status	occupation	relationship
0	0	40	Private	168538	HS-grad	9	Married-civ-spouse	Sales	Husband
1	1	17	Private	101626	9th	5	Never-married	Machine-op-inspct	Own-child
2	2	18	Private	353358	Some-college	10	Never-married	Other-service	Own-child

3 rows × 24 columns

In [52]:

```
sel_cat = ['workclass', 'education', 'marital_status',
           'occupation', 'relationship', 'race', 'native_country' ]
all_dat_n = all_dat.drop(sel_cat, axis=1)
all_dat_n
```

Out[52]:

	id	age	fnlwgt	education_num	sex	capital_gain	capital_loss	hours_per_week	income
0	0	40	168538	9	1	0	0	60	>50K
1	1	17	101626	5	1	0	0	20	<=50K
2	2	18	353358	10	1	0	0	16	<=50K
3	3	21	151158	10	2	0	0	25	<=50K
4	4	24	122234	10	2	0	0	20	<=50K
...
6507	6507	35	61343	13	1	0	0	40	blatant
6508	6508	41	32185	13	1	0	0	40	blatant
6509	6509	39	409189	3	1	0	0	40	blatant
6510	6510	35	180342	9	1	0	0	40	blatant
6511	6511	28	156819	9	2	0	0	36	blatant

32561 rows × 17 columns

In [53]:

```
X_cat = all_dat_n.drop(['target'],axis=1)
y = all_dat_n['target']
```

In [54]:

```
train_n = all_dat_n.loc[ (all_dat_n['target']==0) | (all_dat_n['target']==1) , : ]
test_n = all_dat_n.loc[ all_dat_n['target']==999 , : ]
```

In [55]:

```
print(train_n.shape, test_n.shape)
```

```
(26049, 17) (6512, 17)
```

In [56]:

```
train_n.head(3)
```

Out[56]:

	id	age	fnlwgt	education_num	sex	capital_gain	capital_loss	hours_per_week	income	target
0	0	40	168538	9	1	0	0	60	>50K	0
1	1	17	101626	5	1	0	0	20	<=50K	0
2	2	18	353358	10	1	0	0	16	<=50K	0

In [57]:

```
test_n.head(3)
```

Out[57]:

	id	age	fnlwgt	education_num	sex	capital_gain	capital_loss	hours_per_week	income	target
0	0	28	67661	10	2	0	0	40	blank	999
1	1	40	37869	9	1	0	0	50	blank	999
2	2	20	109952	10	1	0	0	25	blank	999

In [58]:

```
sel_cat = ['income']

train_n = train_n.drop(['income'], axis=1)
test_n = test_n.drop(['income', 'target'], axis=1)

print(train_n.shape, test_n.shape)
```

```
(26049, 16) (6512, 15)
```

In [59]:

```
train_n.columns
```

Out[59]:

```
Index(['id', 'age', 'fnlwgt', 'education_num', 'sex', 'capital_gain',  
      'capital_loss', 'hours_per_week', 'target', 'workclass_lbl',  
      'education_lbl', 'marital_status_lbl', 'occupation_lbl',  
      'relationship_lbl', 'race_lbl', 'native_country_lbl'],  
      dtype='object')
```

In [60]:

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

In [61]:

```
sel = ['age', 'education_num', 'sex']  
  
X_tr_all = train_n[sel]  
y_tr_all = train_n['target']  
X_test_all = test_n[sel]  
  
X_train, X_test, y_train, y_test = train_test_split(X_tr_all,  
                                                    y_tr_all,  
                                                    test_size=0.3,  
                                                    random_state=77)
```

03. 모델 구축하기

[목차로 이동하기](#)

로지스틱 모델 만들기

In [62]:

```
from sklearn.linear_model import LogisticRegression
```

In [63]:

```
model = LogisticRegression()  
model.fit(X_train, y_train)  
  
model.score(X_train, y_train), model.score(X_test, y_test),
```

Out[63]:

```
(0.7960403641548756, 0.7901471529110684)
```

다른 모델 확인해 보기

In [64]:

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.ensemble import GradientBoostingClassifier
```

In [65]:

```
model_list = [RandomForestClassifier(), AdaBoostClassifier(), GradientBoostingClassifier()]

for model in model_list:
    m = model
    m.fit(X_train, y_train)

    ac_tr = model.score(X_train, y_train)
    ac_test = model.score(X_test, y_test)

    print(ac_tr, ac_test)
```

```
0.8174838214324888 0.7856685860524633
0.8044861248217615 0.7982085732565579
0.8076669957222771 0.7980806142034549
```

최종 모델

In [66]:

```
model = GradientBoostingClassifier()
model.fit(X_train, y_train)

pred = model.predict(X_test_all)
```

In [67]:

```
sub['prediction'] = pred
sub.to_csv("secondSub4th_gb.csv", index=False)
```

In [68]:

```
### score : 0.80939
```