

eda-2

November 26, 2024

```
[ ]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
plt.style.use('dark_background')
```

```
[ ]: df = pd.read_csv('adult_with_headers.csv')
df.head()
```

```
[ ]:
age      workclass  fnlwtg  education  education_num  \
0    39      State-gov   77516   Bachelors             13
1    50  Self-emp-not-inc  83311   Bachelors             13
2    38      Private  215646   HS-grad              9
3    53      Private  234721    11th               7
4    28      Private  338409   Bachelors             13

      marital_status      occupation  relationship    race    sex  \
0      Never-married      Adm-clerical  Not-in-family  White  Male
1  Married-civ-spouse  Exec-managerial      Husband  White  Male
2      Divorced      Handlers-cleaners  Not-in-family  White  Male
3  Married-civ-spouse  Handlers-cleaners      Husband  Black  Male
4  Married-civ-spouse  Prof-specialty      Wife    Black  Female

      capital_gain  capital_loss  hours_per_week  native_country  income
0          2174           0           40   United-States  <=50K
1           0           0           13   United-States  <=50K
2           0           0           40   United-States  <=50K
3           0           0           40   United-States  <=50K
4           0           0           40         Cuba    <=50K
```

```
[ ]: df.shape
```

```
[ ]: (32561, 15)
```

```
[ ]: df.describe()
```

```
[ ]:
      count      age      fnlwgt  education_num  capital_gain  capital_loss \
mean      38.581647  1.897784e+05      10.080679      1077.648844      87.303830
std       13.640433  1.055500e+05       2.572720      7385.292085     402.960219
min       17.000000  1.228500e+04       1.000000       0.000000       0.000000
25%       28.000000  1.178270e+05       9.000000       0.000000       0.000000
50%       37.000000  1.783560e+05      10.000000       0.000000       0.000000
75%       48.000000  2.370510e+05      12.000000       0.000000       0.000000
max       90.000000  1.484705e+06      16.000000     99999.000000     4356.000000

      hours_per_week
count      32561.000000
mean       40.437456
std        12.347429
min         1.000000
25%        40.000000
50%        40.000000
75%        45.000000
max        99.000000
```

```
[ ]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32561 entries, 0 to 32560
Data columns (total 15 columns):
#   Column                Non-Null Count  Dtype
---  -
0   age                   32561 non-null  int64
1   workclass             32561 non-null  object
2   fnlwgt                32561 non-null  int64
3   education             32561 non-null  object
4   education_num         32561 non-null  int64
5   marital_status        32561 non-null  object
6   occupation            32561 non-null  object
7   relationship          32561 non-null  object
8   race                  32561 non-null  object
9   sex                   32561 non-null  object
10  capital_gain          32561 non-null  int64
11  capital_loss          32561 non-null  int64
12  hours_per_week        32561 non-null  int64
13  native_country        32561 non-null  object
14  income                32561 non-null  object
dtypes: int64(6), object(9)
memory usage: 3.7+ MB
```

```
[ ]: df['income'].value_counts()
```

```
[ ]:  <=50K    24720
      >50K     7841
      Name: income, dtype: int64
```

```
[ ]: df['sex'].value_counts()
```

```
[ ]:  Male      21790
      Female   10771
      Name: sex, dtype: int64
```

```
[ ]: df['native_country'].value_counts()
```

```
[ ]:  United-States    29170
      Mexico          643
      ?              583
      Philippines     198
      Germany         137
      Canada          121
      Puerto-Rico     114
      El-Salvador     106
      India           100
      Cuba            95
      England         90
      Jamaica         81
      South           80
      China           75
      Italy           73
      Dominican-Republic 70
      Vietnam         67
      Guatemala       64
      Japan           62
      Poland          60
      Columbia        59
      Taiwan          51
      Haiti           44
      Iran            43
      Portugal        37
      Nicaragua       34
      Peru            31
      France          29
      Greece          29
      Ecuador         28
      Ireland         24
      Hong            20
      Cambodia        19
      Trinidad&Tobago 19
      Laos            18
```

Thailand	18
Yugoslavia	16
Outlying-US(Guam-USVI-etc)	14
Honduras	13
Hungary	13
Scotland	12
Holand-Netherlands	1

Name: native_country, dtype: int64

```
[ ]: df['workclass'].value_counts()
```

```
[ ]: Private          22696
Self-emp-not-inc    2541
Local-gov           2093
?                   1836
State-gov           1298
Self-emp-inc        1116
Federal-gov         960
Without-pay         14
Never-worked         7
Name: workclass, dtype: int64
```

```
[ ]: df['occupation'].value_counts()
```

```
[ ]: Prof-specialty    4140
Craft-repair          4099
Exec-managerial       4066
Adm-clerical          3770
Sales                 3650
Other-service         3295
Machine-op-inspct     2002
?                     1843
Transport-moving      1597
Handlers-cleaners     1370
Farming-fishing       994
Tech-support          928
Protective-serv       649
Priv-house-serv       149
Armed-Forces          9
Name: occupation, dtype: int64
```

```
[ ]: #Dropping Education - Education No. is enough, Dropping Final Weight - Highly
↳Discrete Data so not useful.
```

```
[ ]: df = df.drop(['education', 'fnlwgt'], axis = 1)
df.head()
```

```
[ ]:  age          workclass  education_num      marital_status \
0    39          State-gov          13      Never-married
1    50  Self-emp-not-inc          13  Married-civ-spouse
2    38          Private          9      Divorced
3    53          Private          7  Married-civ-spouse
4    28          Private          13  Married-civ-spouse

      occupation  relationship  race  sex  capital_gain \
0      Adm-clerical  Not-in-family  White  Male      2174
1      Exec-managerial      Husband  White  Male          0
2  Handlers-cleaners  Not-in-family  White  Male          0
3  Handlers-cleaners      Husband  Black  Male          0
4      Prof-specialty      Wife  Black  Female          0

      capital_loss  hours_per_week  native_country  income
0              0              40  United-States  <=50K
1              0              13  United-States  <=50K
2              0              40  United-States  <=50K
3              0              40  United-States  <=50K
4              0              40      Cuba  <=50K
```

```
[ ]: #Replacing "?" with NaN
```

```
[ ]: df.replace('?', np.NaN, inplace = True)
df.head()
```

```
[ ]:  age          workclass  education_num      marital_status \
0    39          State-gov          13      Never-married
1    50  Self-emp-not-inc          13  Married-civ-spouse
2    38          Private          9      Divorced
3    53          Private          7  Married-civ-spouse
4    28          Private          13  Married-civ-spouse

      occupation  relationship  race  sex  capital_gain \
0      Adm-clerical  Not-in-family  White  Male      2174
1      Exec-managerial      Husband  White  Male          0
2  Handlers-cleaners  Not-in-family  White  Male          0
3  Handlers-cleaners      Husband  Black  Male          0
4      Prof-specialty      Wife  Black  Female          0

      capital_loss  hours_per_week  native_country  income
0              0              40  United-States  <=50K
1              0              13  United-States  <=50K
2              0              40  United-States  <=50K
3              0              40  United-States  <=50K
4              0              40      Cuba  <=50K
```

1 Replacing NaN with Forward Fill

```
[ ]: df.fillna(method = 'ffill', inplace = True)
```

2 Label Encoding

```
[ ]: from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df['workclass'] = le.fit_transform(df['workclass'])
df['marital.status'] = le.fit_transform(df['marital.status'])
df['occupation'] = le.fit_transform(df['occupation'])
df['relationship'] = le.fit_transform(df['relationship'])
df['race'] = le.fit_transform(df['race'])
df['sex'] = le.fit_transform(df['sex'])
df['native_country'] = le.fit_transform(df['native_country'])
df['income'] = le.fit_transform(df['income'])
df.head()
```

```
[ ]:
age  workclass  education_num  marital_status  occupation \
0   39         7           13   Never-married         1
1   50         6           13  Married-civ-spouse         4
2   38         4           9    Divorced              6
3   53         4           7  Married-civ-spouse         6
4   28         4           13  Married-civ-spouse        10

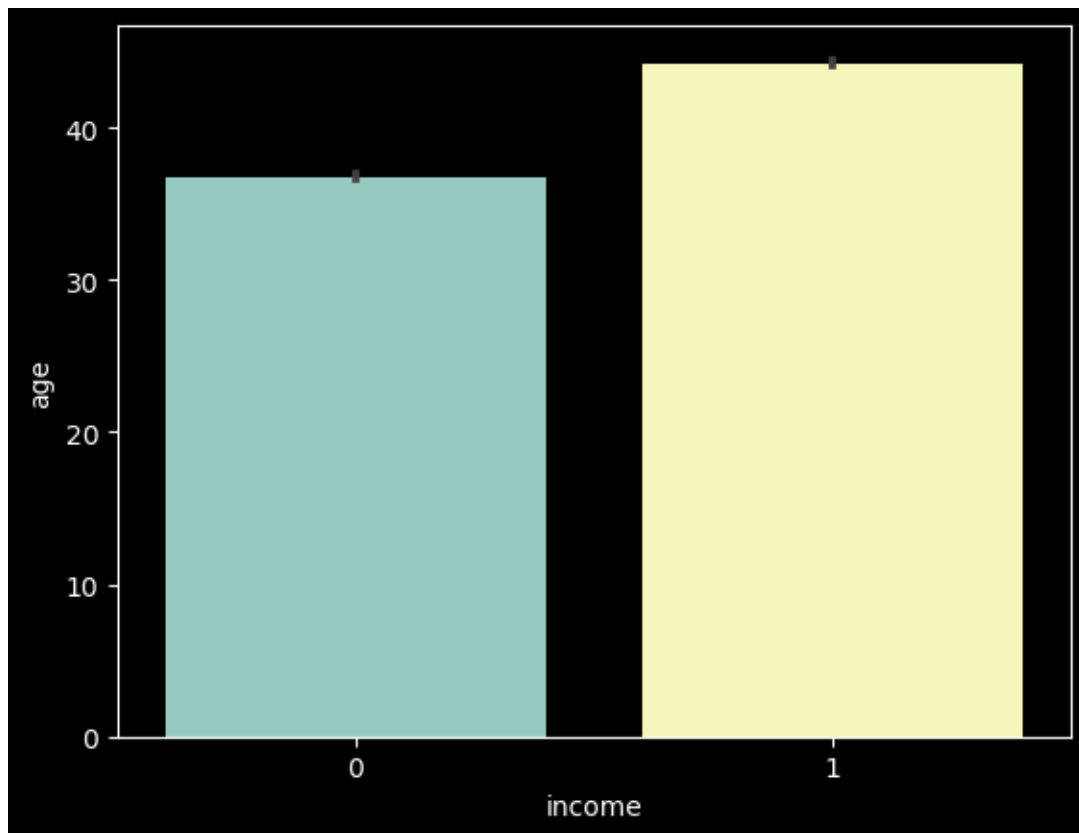
relationship  race  sex  capital_gain  capital_loss  hours_per_week \
0             1    4    1         2174             0             40
1             0    4    1             0             0             13
2             1    4    1             0             0             40
3             0    2    1             0             0             40
4             5    2    0             0             0             40

native_country  income  marital.status
0             39         0             4
1             39         0             2
2             39         0             0
3             39         0             2
4             5         0             2
```

3 People with more age earn more

```
[ ]: sns.barplot(x = 'income', y = 'age', data = df)
```

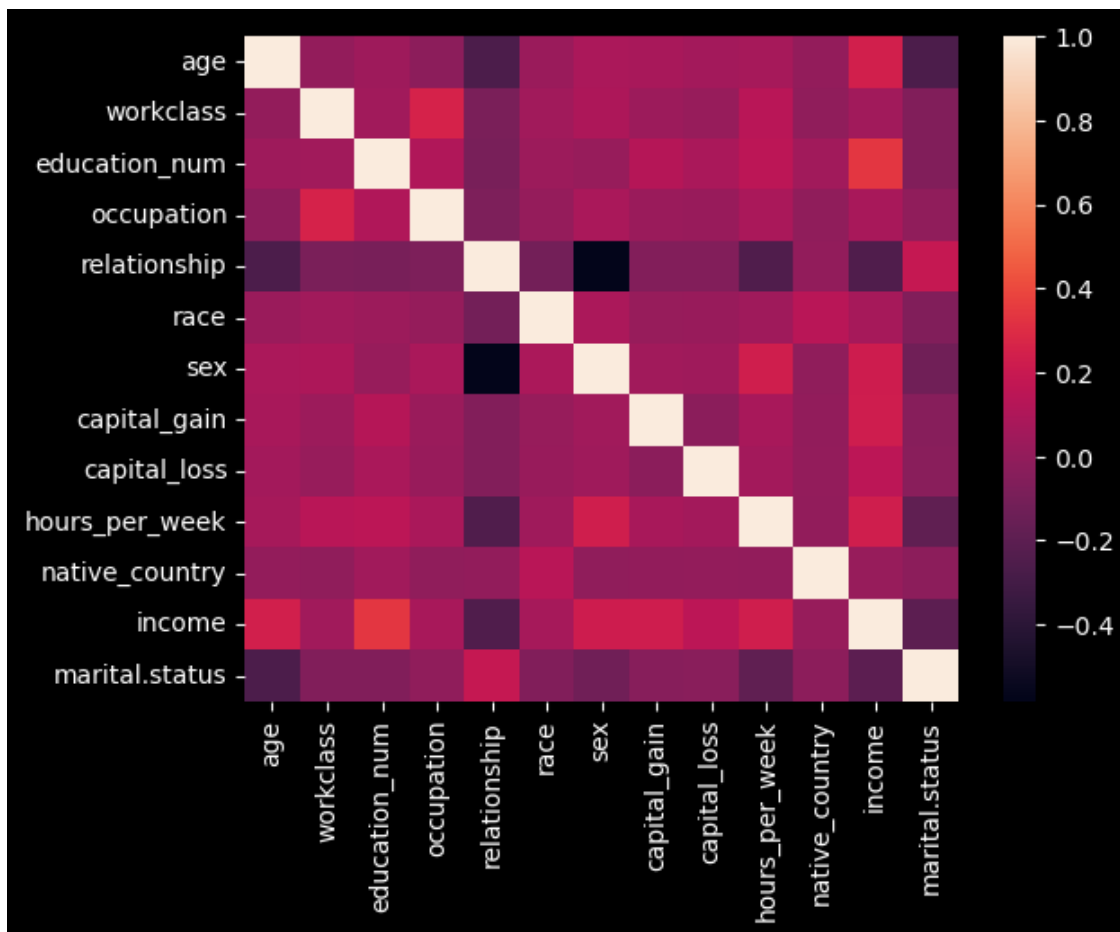
```
[ ]: <Axes: xlabel='income', ylabel='age'>
```



```
[ ]: sns.heatmap(df.corr())
```

C:\Users\excel\AppData\Local\Temp\ipykernel_25540\58359773.py:1: FutureWarning:
The default value of numeric_only in DataFrame.corr is deprecated. In a future
version, it will default to False. Select only valid columns or specify the
value of numeric_only to silence this warning.
sns.heatmap(df.corr())

```
[ ]: <Axes: >
```



```
[ ]: X = df.drop(['income'],axis = 1)
     y = df['income']
```

```
[ ]: print('Shape of X =', X.shape)
     print('Shape of y =', y.shape)
```

Shape of X = (32561, 13)

Shape of y = (32561,)

4 Train Test Split

```
[ ]: 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=42)
```

```
[ ]: # Instantiate Gaussian Naive Bayes model
     gnb = GaussianNB()
```



```
# Fit the model
gnb.fit(X_train, y_train)
```

```
[ ]: GaussianNB()
```

```
[ ]: from sklearn.metrics import classification_report
      from sklearn.metrics import confusion_matrix
      from sklearn.metrics import accuracy_score
```

```
[ ]: y_pred = gnb.predict(X_test)
      print(classification_report(y_test, y_pred))
      print(confusion_matrix(y_test, y_pred))
      print(accuracy_score(y_test, y_pred)*100)
```

	precision	recall	f1-score	support
0	0.82	0.95	0.88	4942
1	0.67	0.32	0.44	1571
accuracy			0.80	6513
macro avg	0.74	0.64	0.66	6513
weighted avg	0.78	0.80	0.77	6513

```
[[4698 244]
 [1065 506]]
79.90173499155534
```

```
[ ]:
```