

EDA + Logistic Regression + PCA

LOGISTIC REGRESSION PROJECT

Import Libraries

```
In [1]: ▶ import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)

# import libraries for plotting
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline

# ignore warnings
import warnings
warnings.filterwarnings('ignore')
```

Import Dataset

```
In [2]: ▶ df=pd.read_csv(r"C:\Users\yogay\OneDrive\Desktop\Yogita_Yadav\Data Science\Machine Learning\Logistic Regression(classifier)\Project\adult.csv\adult.csv")
```

Exploratory Data Analysis

Check Shape Of Dataset

```
In [3]: ▶ df.shape
```

```
Out[3]: (32561, 15)
```

Preview Dataset

In [4]: `df.head()`

Out[4]:

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.loss	hours.per.week	native.country	income
0	90	?	77053	HS-grad	9	Widowed	?	Not-in-family	White	Female	0	4356	40	United-States	<=50K
1	82	Private	132870	HS-grad	9	Widowed	Exec-managerial	Not-in-family	White	Female	0	4356	18	United-States	<=50K
2	66	?	186061	Some-college	10	Widowed	?	Unmarried	Black	Female	0	4356	40	United-States	<=50K
3	54	Private	140359	7th-8th	4	Divorced	Machine-op-inspct	Unmarried	White	Female	0	3900	40	United-States	<=50K
4	41	Private	264663	Some-college	10	Separated	Prof-specialty	Own-child	White	Female	0	3900	40	United-States	<=50K

View Summary Of Dataset

In [5]: `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
```

Now, the summary shows that the variables - `workclass` , `occupation` and `native.country` contain missing values. All of these variables are categorical data type. So, I will impute the missing values with the most frequent value- the mode.

Impute missing value with mode

```
In [6]: ► for col in ['workclass', 'occupation', 'native.country']:
         df[col].fillna(df[col].mode()[0], inplace=True)
```

Check Again For Missing Values

```
In [7]: ► df.isnull().sum()
```

```
Out[7]: age                0
workclass                 0
fnlwgt                   0
education                0
education.num            0
marital.status           0
occupation               0
relationship             0
race                    0
sex                     0
capital.gain             0
capital.loss             0
hours.per.week           0
native.country           0
income                  0
dtype: int64
```

setting feature vector & target variable

```
In [8]: ► X = df.drop(['income'], axis=1)
         y = df['income']
```

In [9]: `df.head()`

Out[9]:

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.loss	hours.per.week	native.country	income
0	90	?	77053	HS-grad	9	Widowed	?	Not-in-family	White	Female	0	4356	40	United-States	<=50K
1	82	Private	132870	HS-grad	9	Widowed	Exec-managerial	Not-in-family	White	Female	0	4356	18	United-States	<=50K
2	66	?	186061	Some-college	10	Widowed	?	Unmarried	Black	Female	0	4356	40	United-States	<=50K
3	54	Private	140359	7th-8th	4	Divorced	Machine-op-inspct	Unmarried	White	Female	0	3900	40	United-States	<=50K
4	41	Private	264663	Some-college	10	Separated	Prof-specialty	Own-child	White	Female	0	3900	40	United-States	<=50K

split data into separate training and test set

```
In [10]: from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 0)
```

Feature Engineering

Encode Categorical Variables

```
In [11]: from sklearn import preprocessing

categorical = ['workclass', 'education', 'marital.status', 'occupation', 'relationship', 'race', 'sex', 'native.country']
for feature in categorical:
    le = preprocessing.LabelEncoder()
    X_train[feature] = le.fit_transform(X_train[feature])
    X_test[feature] = le.transform(X_test[feature])
```

Feature Scaling

```
In [12]: from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()

X_train = pd.DataFrame(scaler.fit_transform(X_train), columns = X.columns)

X_test = pd.DataFrame(scaler.transform(X_test), columns = X.columns)
```

```
In [13]: X_train.head()
```

Out[13]:

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.loss	hours.per.week	native.country
0	0.101484	2.134215	-1.494279	-0.332263	1.133894	-0.402341	-0.600270	2.214196	0.39298	-1.430470	-0.145189	-0.217407	-1.662414	0.292864
1	0.028248	-1.279379	0.438778	0.184396	-0.423425	-0.402341	0.109933	-0.899410	0.39298	0.699071	-0.145189	-0.217407	-0.200753	0.292864
2	0.247956	0.086059	0.045292	1.217715	-0.034095	0.926666	-0.600270	-0.276689	0.39298	-1.430470	-0.145189	-0.217407	-0.038346	0.292864
3	-0.850587	-1.279379	0.793152	0.184396	-0.423425	0.926666	-0.363535	0.968753	0.39298	0.699071	-0.145189	-0.217407	-0.038346	0.292864
4	-0.044989	-1.962098	-0.853275	0.442726	1.523223	-0.402341	-0.600270	-0.899410	0.39298	0.699071	-0.145189	-0.217407	-0.038346	0.292864

Logistic Regression Model With All Features

```
In [14]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score

logreg = LogisticRegression()
logreg.fit(X_train, y_train)
y_pred = logreg.predict(X_test)

print('Logistic Regression accuracy score with all the features: {0:0.4f}'.format(accuracy_score(y_test, y_pred)))
```

Logistic Regression accuracy score with all the features: 0.8204

Logistic Regression With PCA

```
In [15]: ▶ from sklearn.decomposition import PCA  
pca = PCA()  
X_train = pca.fit_transform(X_train)  
pca.explained_variance_ratio_
```

```
Out[15]: array([0.15112277, 0.10122703, 0.09056424, 0.0802928 , 0.07708238,  
               0.07350038, 0.06774638, 0.06602885, 0.06115879, 0.06007244,  
               0.05358847, 0.04835632, 0.04181168, 0.02744748])
```

comment

We can see that approximately 97.25% of variance is explained by the first 13 variables.

Only 2.75% of variance is explained by the last variable. So, we can assume that it carries little information.

So, I will drop it, train the model again and calculate the accuracy.

Logistic Regression With First 13 Features

```
In [18]: X = df.drop(['income', 'native.country'], axis=1)
y = df['income']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 0)

categorical = ['workclass', 'education', 'marital.status', 'occupation', 'relationship', 'race', 'sex']
for feature in categorical:
    le = preprocessing.LabelEncoder()
    X_train[feature] = le.fit_transform(X_train[feature])
    X_test[feature] = le.transform(X_test[feature])

X_train = pd.DataFrame(scaler.fit_transform(X_train), columns = X.columns)

X_test = pd.DataFrame(scaler.transform(X_test), columns = X.columns)

logreg = LogisticRegression()
logreg.fit(X_train, y_train)
y_pred = logreg.predict(X_test)

print('Logistic Regression accuracy score with the first 13 features: {0:0.4f}'.format(accuracy_score(y_test, y_pred)))
```

Logistic Regression accuracy score with the first 13 features: 0.8209

Logistic Regression With First 12 Features

```
In [19]: X = df.drop(['income', 'native.country', 'hours.per.week'], axis=1)
y = df['income']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 0)

categorical = ['workclass', 'education', 'marital.status', 'occupation', 'relationship', 'race', 'sex']
for feature in categorical:
    le = preprocessing.LabelEncoder()
    X_train[feature] = le.fit_transform(X_train[feature])
    X_test[feature] = le.transform(X_test[feature])

X_train = pd.DataFrame(scaler.fit_transform(X_train), columns = X.columns)

X_test = pd.DataFrame(scaler.transform(X_test), columns = X.columns)

logreg = LogisticRegression()
logreg.fit(X_train, y_train)
y_pred = logreg.predict(X_test)

print('Logistic Regression accuracy score with the first 12 features: {0:0.4f}'.format(accuracy_score(y_test, y_pred)))
```

Logistic Regression accuracy score with the first 12 features: 0.8209

Logistic Regression With First 11 Features

```
In [20]: X = df.drop(['income', 'native.country', 'hours.per.week', 'capital.loss'], axis=1)
y = df['income']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 0)

categorical = ['workclass', 'education', 'marital.status', 'occupation', 'relationship', 'race', 'sex']
for feature in categorical:
    le = preprocessing.LabelEncoder()
    X_train[feature] = le.fit_transform(X_train[feature])
    X_test[feature] = le.transform(X_test[feature])

X_train = pd.DataFrame(scaler.fit_transform(X_train), columns = X.columns)
X_test = pd.DataFrame(scaler.transform(X_test), columns = X.columns)

logreg = LogisticRegression()
logreg.fit(X_train, y_train)
y_pred = logreg.predict(X_test)

print('Logistic Regression accuracy score with the first 11 features: {0:0.4f}'.format(accuracy_score(y_test, y_pred)))
```

Logistic Regression accuracy score with the first 11 features: 0.8184

Comment

We can see that accuracy has significantly decreased to 0.8184 if I drop the last three features.

Select right number of dimensions

```
In [21]: X = df.drop(['income'], axis=1)
y = df['income']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 0)

categorical = ['workclass', 'education', 'marital.status', 'occupation', 'relationship', 'race', 'sex', 'native.country']
for feature in categorical:
    le = preprocessing.LabelEncoder()
    X_train[feature] = le.fit_transform(X_train[feature])
    X_test[feature] = le.transform(X_test[feature])

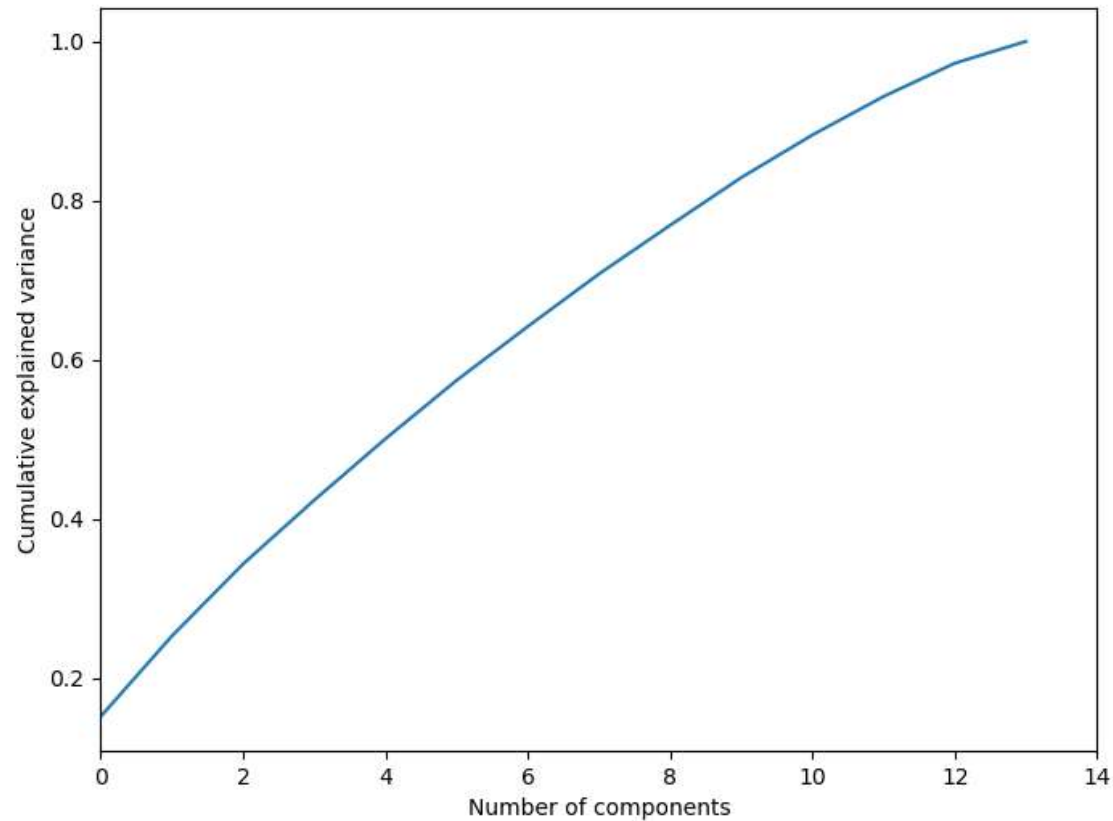
X_train = pd.DataFrame(scaler.fit_transform(X_train), columns = X.columns)

pca= PCA()
pca.fit(X_train)
cumsum = np.cumsum(pca.explained_variance_ratio_)
dim = np.argmax(cumsum >= 0.90) + 1
print('The number of dimensions required to preserve 90% of variance is',dim)
```

The number of dimensions required to preserve 90% of variance is 12

Plot explained variance ratio with number of dimensions

```
In [22]: ▶ plt.figure(figsize=(8,6))  
plt.plot(np.cumsum(pca.explained_variance_ratio_))  
plt.xlim(0,14,1)  
plt.xlabel('Number of components')  
plt.ylabel('Cumulative explained variance')  
plt.show()
```



Comment

The above plot shows that almost 90% of variance is explained by the first 12 components.

