

```
In [5]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline

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

```
In [6]: df = pd.read_csv(r"C:\Users\santo\OneDrive\Desktop\Data science\Jan 2026\3rd - KNN\projects\LOGISTIC REGRESSION , PCA, EDA\adu
```

```
In [7]: df.shape
```

```
Out[7]: (32561, 15)
```

```
In [8]: df.head()
```

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



```
In [9]: 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
```

```
In [10]: df[df=='?'] = np.nan
```

```
In [11]: 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        30725 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       30718 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   31978 non-null   object  
 14  income            32561 non-null   object  
dtypes: int64(6), object(9)
memory usage: 3.7+ MB
```

```
In [12]: for col in ['workclass', 'occupation', 'native.country']:
    df[col].fillna(df[col].mode()[0], inplace=True)
```

```
In [13]: df.isnull().sum()
```

```
Out[13]: 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
```

```
In [14]: X = df.drop(['income'], axis=1)  
  
y = df['income']
```

```
In [15]: X.head()
```

```
Out[15]:   age  workclass  fnlwgt  education  education.num  marital.status  occupation  relationship  race  sex  capital.gain  capital.loss  hours.per.week  
0    90     Private    77053    HS-grad           9        Widowed    Prof-specialty  Not-in-family  White  Female           0        4356  
1    82     Private   132870    HS-grad           9        Widowed    Exec-managerial  Not-in-family  White  Female           0        4356  
2    66     Private   186061  Some-college          10        Widowed    Prof-specialty  Unmarried    Black  Female           0        4356  
3    54     Private   140359    7th-8th            4        Divorced  Machine-op-inspct  Unmarried    White  Female           0        3900  
4    41     Private   264663  Some-college          10       Separated    Prof-specialty  Own-child    White  Female           0        3900
```



```
In [16]: 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)
```

```
In [18]: 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])
```

```
In [19]: 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 [20]: X_train.head()
```

```
Out[20]:
```

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relationship	race	sex	capital.gain	ca
0	0.101484	2.600478	-1.494279	-0.332263	1.133894	-0.402341	-0.782234	2.214196	0.39298	-1.430470	-0.145189	-
1	0.028248	-1.884720	0.438778	0.184396	-0.423425	-0.402341	-0.026696	-0.899410	0.39298	0.699071	-0.145189	-
2	0.247956	-0.090641	0.045292	1.217715	-0.034095	0.926666	-0.782234	-0.276689	0.39298	-1.430470	-0.145189	-
3	-0.850587	-1.884720	0.793152	0.184396	-0.423425	0.926666	-0.530388	0.968753	0.39298	0.699071	-0.145189	-
4	-0.044989	-2.781760	-0.853275	0.442726	1.523223	-0.402341	-0.782234	-0.899410	0.39298	0.699071	-0.145189	-

```
In [21]: 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: {:.4f}'.format(accuracy_score(y_test, y_pred)))
```

Logistic Regression accuracy score with all the features: 0.8218

```
In [22]: from sklearn.decomposition import PCA
pca = PCA()
X_train = pca.fit_transform(X_train)
pca.explained_variance_ratio_
```

```
Out[22]: array([0.14757168, 0.10182915, 0.08147199, 0.07880174, 0.07463545,
   0.07274281, 0.07009602, 0.06750902, 0.0647268 , 0.06131155,
   0.06084207, 0.04839584, 0.04265038, 0.02741548])
```

```
In [23]: 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: {:.4f}'.format(accuracy_score(y_test, y_pred)))
```

Logistic Regression accuracy score with the first 13 features: 0.8213

```
In [24]: 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: {:.4f}'.format(accuracy_score(y_test, y_pred)))
```

Logistic Regression accuracy score with the first 12 features: 0.8227

```
In [25]: 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: {:.4f}'.format(accuracy_score(y_test, y_pred)))

```

Logistic Regression accuracy score with the first 11 features: 0.8186

```

In [26]: 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

```

In [28]: 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()
```

AttributeError

Traceback (most recent call last)

Cell In[28], line 3

```
1 plt.figure(figsize=(8,6))
2 plt.plot(np.cumsum(pca.explained_variance_ratio_))
----> 3 plt.Xlim(0,14,1)
4 plt.Xlabel('Number of components')
5 plt.ylabel('Cumulative explained variance')
```

AttributeError: module 'matplotlib.pyplot' has no attribute 'Xlim'

In []: