



|  |
|--|
| Experiment No. 7   |
| Apply Dimensionality Reduction on Adult Census Income Dataset and analyze the performance of the model |
| Date of Performance: 11-9-23   |
| Date of Submission: 9-10-23  |



**Aim:** Apply Dimensionality Reduction on Adult Census Income Dataset and analyze the performance of the model.

**Objective:** Able to perform various feature engineering tasks, perform dimensionality reduction on the given dataset and maximize the accuracy, Precision, Recall, F1 score.

**Theory:**

In machine learning classification problems, there are often too many factors on the basis of which the final classification is done. These factors are basically variables called features. The higher the number of features, the harder it gets to visualize the training set and then work on it. Sometimes, most of these features are correlated, and hence redundant. This is where dimensionality reduction algorithms come into play. Dimensionality reduction is the process of reducing the number of random variables under consideration, by obtaining a set of principal variables. It can be divided into feature selection and feature extraction.

**Dataset:**

Predict whether income exceeds \$50K/yr based on census data. Also known as "Adult" dataset.

Attribute Information:

Listing of attributes:

>50K, <=50K.

age: continuous.

workclass: Private, Self-emp-not-inc, Self-emp-inc, Federal-gov, Local-gov, State-gov, Without-pay, Never-worked.

fnlwgt: continuous.

education: Bachelors, Some-college, 11th, HS-grad, Prof-school, Assoc-acdm, Assoc-voc, 9th, 7th-8th, 12th, Masters, 1st-4th, 10th, Doctorate, 5th-6th, Preschool.



## Vidyavardhini's College of Engineering & Technology

### Department of Computer Engineering

---

education-num: continuous.

marital-status: Married-civ-spouse, Divorced, Never-married, Separated, Widowed, Married-spouse-absent, Married-AF-spouse.

occupation: Tech-support, Craft-repair, Other-service, Sales, Exec-managerial, Prof-specialty, Handlers-cleaners, Machine-op-inspct, Adm-clerical, Farming-fishing, Transport-moving, Priv-house-serv, Protective-serv, Armed-Forces.

relationship: Wife, Own-child, Husband, Not-in-family, Other-relative, Unmarried.

race: White, Asian-Pac-Islander, Amer-Indian-Eskimo, Other, Black.

sex: Female, Male.

capital-gain: continuous.

capital-loss: continuous.

hours-per-week: continuous.

native-country: United-States, Cambodia, England, Puerto-Rico, Canada, Germany, Outlying-US(Guam-USVI-etc), India, Japan, Greece, South, China, Cuba, Iran, Honduras, Philippines, Italy, Poland, Jamaica, Vietnam, Mexico, Portugal, Ireland, France, Dominican-Republic, Laos, Ecuador, Taiwan, Haiti, Columbia, Hungary, Guatemala, Nicaragua, Scotland, Thailand, Yugoslavia, El-Salvador, Trinidad & Tobago, Peru, Hong, Holand-Netherlands.

#### **Code:**

#### **Conclusion:**

Accuracy: Following dimensionality reduction, it exhibits approximately 0.821 accuracy.

Precision: The model shows a precision of 0.84 for the  $\leq 50K$  class and a precision of 0.72 for the  $> 50K$  class.

Model recall: The recall for the  $\leq 50K$  class is 0.95, and the recall for the  $> 50K$  class is 0.43.



## **Vidyavardhini's College of Engineering & Technology**

### **Department of Computer Engineering**

---

FI-value. The model shows an F1-score of 0.54 for the >50K class and a FI-score of 0.89 for the  $\leq 50K$  class.

# 11\_ml\_exp7.py

```
# -*- coding: utf-8 -*-  
"""11_ML_Exp7
```

Automatically generated by Colaboratory.

Original file is located at  
[https://colab.research.google.com/drive/1j-PvGfiSYWNrULq1F7girH14kama\\_ZID](https://colab.research.google.com/drive/1j-PvGfiSYWNrULq1F7girH14kama_ZID)  
"""

```
import numpy as np # linear algebra  
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)  
import os  
  
for dirname, _, filenames in os.walk('/content/adult (1).csv'):  
    for filename in filenames:  
        print(os.path.join(dirname, filename))  
  
df=pd.read_csv("/content/adult (1).csv")  
  
df.head  
  
df.columns  
  
df.shape  
  
df.info()  
  
df[df == '?'] = np.nan  
  
df.isnull().sum()  
  
for col in ['workclass', 'occupation', 'native.country']:  
    df[col].fillna(df[col].mode()[0], inplace=True)  
    df.isnull().sum()  
  
# converting categorical Columns  
df.replace({'Sex':{'male':0,'female':1}, 'Embarked':{'S':0,'C':1,'Q':2}},  
inplace=True)  
X = df.drop(['income'], axis=1)  
y = df['income']  
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)  
from sklearn import preprocessing  
categorical = ['workclass', 'education', 'marital.status', 'occupation',  
'relationship', 'race', 'sex', 'native.country']  
for feature in categorical:  
    label = preprocessing.LabelEncoder()  
    X_train[feature] = label.fit_transform(X_train[feature])  
    X_test[feature] = label.transform(X_test[feature])
```

```
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)
X_train.head()

from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score

LR = LogisticRegression()
LR.fit(X_train, y_train)
y_pred = LR.predict(X_test)
accuracy_score(y_test, y_pred)

from sklearn.decomposition import PCA
pca = PCA()
X_train = pca.fit_transform(X_train)
pca.explained_variance_ratio_

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:
label = preprocessing.LabelEncoder()
X_train[feature] = label.fit_transform(X_train[feature])
X_test[feature] = label.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)

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:
label = preprocessing.LabelEncoder()
X_train[feature] = label.fit_transform(X_train[feature])
X_test[feature] = label.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)

LR2 = LogisticRegression()
LR2.fit(X_train, y_train)

y_pred = LR2.predict(X_test)
accuracy_score(y_test, y_pred)
```

```
from sklearn.metrics import confusion_matrix
import pandas as pd
confusion = confusion_matrix(y_test, y_pred)
df_confusion = pd.DataFrame(confusion, columns=['Predicted No', 'Predicted
Yes'], index=['Actual No', 'Actual Yes'])
from sklearn.metrics import classification_report
print(classification_report(y_test, y_pred))
```