



Experiment No. 7
Apply Dimensionality Reduction on Adult Census Income Dataset and analyze the performance of the model
Date of Performance: 04/09/2023
Date of Submission: 13/09/2023



**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:

- 1)Accuracy: The accuracy of your model, which stands at 0.8227, can be influenced by how well reduced features capture data patterns. This means that your model is correct in approximately 82.27% of its predictions when using the initial 12 features.
- 2)Precision: Precision quantifies the ratio of true positive predictions to all positive predictions. The process of dimensionality reduction may result in a less precise distinction between true and false positives. Your report indicates a precision of 0.72 for the ">50K" class.
- 3)Recall: Recall evaluates the model's ability to correctly identify actual positive instances. Reducing dimensionality may lead to the omission of some positive cases. In the case of the ">50K" class, the recall value is 0.43.
- 4)F1 Score: The F1 score represents the harmonic mean of precision and recall. It decreases if either precision or recall is affected. Your ">50K" class has an F1 score of 0.54.

```
confusion matrix
[[7012  398]
 [1334 1025]]

      precision    recall  f1-score   support

    <=50K      0.84      0.95      0.89      7410
    >50K      0.72      0.43      0.54      2359

 accuracy              0.82      9769
 macro avg      0.78      0.69      0.72      9769
weighted avg      0.81      0.82      0.81      9769
```

Logistic Regression accuracy score with the first 12 features: 0.8227

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import io
from sklearn.metrics import accuracy_score, precision_score, f1_score, confusion_matrix, classification_report
from sklearn.model_selection import cross_val_score
from sklearn.metrics import mean_squared_error

import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
```

```
file = ('/content/adult.csv')
df = pd.read_csv(file)
```

```
df.shape
```

(32561, 15)

```
df.head()
```

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relatio
0	90	?	77053	HS-grad	9	Widowed	?	Not-in
1	82	Private	132870	HS-grad	9	Widowed	Exec-managerial	Not-in
2	66	?	186061	Some-college	10	Widowed	?	Unm
Machine-								

```
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[df == '?'] = np.nan
```

```
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
```

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

```
df.isnull().sum()
```

```
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
```

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

y = df['income']
```

```
X.head()
```

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relatio
0	90	Private	77053	HS-grad	9	Widowed	Prof-specialty	Not-in
1	82	Private	132870	HS-grad	9	Widowed	Exec-managerial	Not-in
2	66	Private	186061	Some-college	10	Widowed	Prof-specialty	Unm

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

	age	workclass	fnlwgt	education	education.num	marital.status	occupation
0	0.101484	2.600478	-1.494279	-0.332263	1.133894	-0.402341	-0.782234
1	0.028248	-1.884720	0.438778	0.184396	-0.423425	-0.402341	-0.026696
2	0.247956	-0.090641	0.045292	1.217715	-0.034095	0.926666	-0.782234
3	-0.850587	-1.884720	0.793152	0.184396	-0.423425	0.926666	-0.530388
4	-0.044989	-2.781760	-0.853275	0.442726	1.523223	-0.402341	-0.782234

```
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("confusion matrix\n",confusion_matrix(y_test,y_pred))
print(classification_report(y_test, y_pred))

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

confusion matrix					
[[6987 423]					
[1319 1040]]					
	precision	recall	f1-score	support	
<=50K	0.84	0.94	0.89	7410	
>50K	0.71	0.44	0.54	2359	
accuracy			0.82	9769	
macro avg	0.78	0.69	0.72	9769	
weighted avg	0.81	0.82	0.81	9769	

Logistic Regression accuracy score with all the features: 0.8217

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

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])
```

```
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("confusion matrix\n",confusion_matrix(y_test,y_pred))
print(classification_report(y_test, y_pred))

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

confusion matrix					
[[6984 426]					
[1320 1039]]					
	precision	recall	f1-score	support	
<=50K	0.84	0.94	0.89	7410	
>50K	0.71	0.44	0.54	2359	
accuracy			0.82	9769	
macro avg	0.78	0.69	0.72	9769	
weighted avg	0.81	0.82	0.81	9769	
Logistic Regression accuracy score with the first 13 features: 0.8213					

```
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("confusion matrix\n",confusion_matrix(y_test,y_pred))
print(classification_report(y_test, y_pred))

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

confusion matrix				
[[7012 398]				
[1334 1025]]				
	precision	recall	f1-score	support
<=50K	0.84	0.95	0.89	7410
>50K	0.72	0.43	0.54	2359
accuracy			0.82	9769
macro avg	0.78	0.69	0.72	9769
weighted avg	0.81	0.82	0.81	9769
Logistic Regression accuracy score with the first 12 features: 0.8227				

