

**Experiment No: 03**

● **Aim:** To implement PCA / SVD / LDA.

● **Theory:**

**Principal Component Analysis (PCA):**

Principal Component Analysis is an unsupervised learning algorithm that is used for the dimensionality reduction in machine learning. It is a statistical process that converts the observations of correlated features into a set of linearly uncorrelated features with the help of orthogonal transformation. These new transformed features are called the Principal Components. It is one of the popular tools that is used for exploratory data analysis and predictive modeling. It is a technique to draw strong patterns from the given dataset by reducing the variances.

PCA works by considering the variance of each attribute because the high attribute shows the good split between the classes, and hence it reduces the dimensionality. Some real-world applications of PCA are image processing, movie recommendation system, optimizing the power allocation in various communication channels. It is a feature extraction technique, so it contains the important variables and drops the least important variable.

**Steps for PCA algorithm:**

Getting the dataset

Firstly, we need to take the input dataset and divide it into two subparts X and Y, where X is the training set, and Y is the validation set.

**Representing data into a structure**

Now we will represent our dataset into a structure. Such as we will represent the two-dimensional matrix of independent variable X. Here each row corresponds to the data items, and the column corresponds to the Features. The number of columns is the dimensions of the dataset.

**Standardizing the data**

In this step, we will standardize our dataset. Such as in a particular column, the features with high variance are more important compared to the features with lower variance.

If the importance of features is independent of the variance of the feature, then we will divide each data item in a column with the standard deviation of the column. Here we will name the matrix as Z.

**Calculating the Covariance of Z**

To calculate the covariance of Z, we will take the matrix Z, and will transpose it. After transpose, we will multiply it by Z. The output matrix will be the Covariance matrix of Z.

## Calculating the Eigen Values and Eigen Vectors

Now we need to calculate the eigenvalues and eigenvectors for the resultant covariance matrix Z. Eigenvectors or the covariance matrix are the directions of the axes with high information. And the coefficients of these eigenvectors are defined as the eigenvalues.

## Sorting the Eigen Vectors

In this step, we will take all the eigenvalues and will sort them in decreasing order, which means from largest to smallest. And simultaneously sort the eigenvectors accordingly in matrix P of eigenvalues. The resultant matrix will be named as P\*.

## Calculating the new features Or Principal Components

Here we will calculate the new features. To do this, we will multiply the P\* matrix to the Z. In the resultant matrix Z\*, each observation is the linear combination of original features. Each column of the Z\* matrix is independent of each other.

## Remove less or unimportant features from the new dataset.

The new feature set has occurred, so we will decide here what to keep and what to remove. It means, we will only keep the relevant or important features in the new dataset, and unimportant features will be removed out.

## Applications of Principal Component Analysis

- PCA is mainly used as the dimensionality reduction technique in various AI applications such as computer vision, image compression, etc.
- It can also be used for finding hidden patterns if data has high dimensions. Some fields where PCA is used are Finance, data mining, Psychology, etc

## Implementation of Principal Component Analysis (PCA):

```
Principal Component Analysis

In [1]: from sklearn.datasets import load_digits
import pandas as pd

In [2]: dataset = load_digits()
dataset.keys()

Out[2]: dict_keys(['data', 'target', 'frame', 'feature_names', 'target_names', 'images', 'DESCR'])

In [3]: dataset.data.shape

Out[3]: (1797, 64)

In [4]: dataset.data[0]

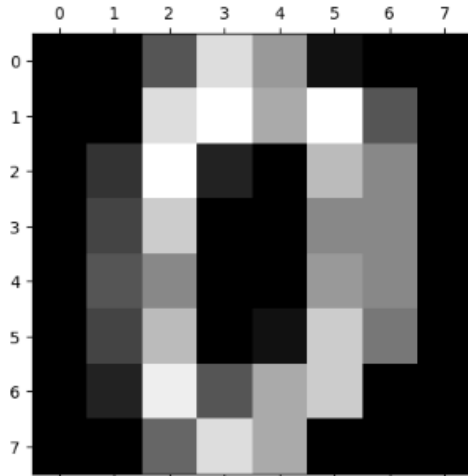
Out[4]: array([ 0.,  0.,  5., 13.,  9.,  1.,  0.,  0.,  0.,  0., 13., 15., 10.,
 15.,  5.,  0.,  0.,  3., 15.,  2.,  0., 11.,  8.,  0.,  0.,  4.,
 12.,  0.,  0.,  8.,  8.,  0.,  0.,  5.,  8.,  0.,  0.,  9.,  8.,
  0.,  0.,  4., 11.,  0.,  1., 12.,  7.,  0.,  0.,  2., 14.,  5.,
 10., 12.,  0.,  0.,  0.,  0.,  6., 13., 10.,  0.,  0.,  0.])

In [5]: dataset.data[0].reshape(8,8)

Out[5]: array([[ 0.,  0.,  5., 13.,  9.,  1.,  0.,  0.],
 [ 0.,  0., 13., 15., 10., 15.,  5.,  0.],
 [ 0.,  3., 15.,  2.,  0., 11.,  8.,  0.],
 [ 0.,  4., 12.,  0.,  0.,  8.,  8.,  0.],
 [ 0.,  5.,  8.,  0.,  0.,  9.,  8.,  0.],
 [ 0.,  4., 11.,  0.,  1., 12.,  7.,  0.],
 [ 0.,  2., 14.,  5., 10., 12.,  0.,  0.],
 [ 0.,  0.,  6., 13., 10.,  0.,  0.,  0.]])
```

```
In [6]: from matplotlib import pyplot as plt
%matplotlib inline
plt.gray()
plt.imshow(dataset.data[0].reshape(8,8))
```

```
Out[6]: <matplotlib.image.AxesImage at 0x23cfdfe67d0>
<Figure size 640x480 with 0 Axes>
```



```
In [7]: dataset.target[:5]
```

```
Out[7]: array([0, 1, 2, 3, 4])
```

```
In [8]: df = pd.DataFrame(dataset.data, columns=dataset.feature_names)
df.head()
```

```
Out[8]:
```

|   | pixel_0_0 | pixel_0_1 | pixel_0_2 | pixel_0_3 | pixel_0_4 | pixel_0_5 | pixel_0_6 | pixel_0_7 | pixel_1_0 | pixel_1_1 | ... | pixel_6_6 | pixel_6_7 | pixel_7_0 | pixel_7_1 | pix |
|---|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----|-----------|-----------|-----------|-----------|-----|
| 0 | 0.0       | 0.0       | 5.0       | 13.0      | 9.0       | 1.0       | 0.0       | 0.0       | 0.0       | 0.0       | ... | 0.0       | 0.0       | 0.0       | 0.0       | 0.0 |
| 1 | 0.0       | 0.0       | 0.0       | 12.0      | 13.0      | 5.0       | 0.0       | 0.0       | 0.0       | 0.0       | ... | 0.0       | 0.0       | 0.0       | 0.0       | 0.0 |
| 2 | 0.0       | 0.0       | 0.0       | 4.0       | 15.0      | 12.0      | 0.0       | 0.0       | 0.0       | 0.0       | ... | 5.0       | 0.0       | 0.0       | 0.0       | 0.0 |
| 3 | 0.0       | 0.0       | 7.0       | 15.0      | 13.0      | 1.0       | 0.0       | 0.0       | 0.0       | 8.0       | ... | 9.0       | 0.0       | 0.0       | 0.0       | 0.0 |
| 4 | 0.0       | 0.0       | 0.0       | 1.0       | 11.0      | 0.0       | 0.0       | 0.0       | 0.0       | 0.0       | ... | 0.0       | 0.0       | 0.0       | 0.0       | 0.0 |

5 rows x 64 columns

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

```
Out[9]:
```

|       | pixel_0_0 | pixel_0_1   | pixel_0_2   | pixel_0_3   | pixel_0_4   | pixel_0_5   | pixel_0_6   | pixel_0_7   | pixel_1_0   | pixel_1_1   | ... | pixel_6_6   | pix    |
|-------|-----------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-----|-------------|--------|
| count | 1797.0    | 1797.000000 | 1797.000000 | 1797.000000 | 1797.000000 | 1797.000000 | 1797.000000 | 1797.000000 | 1797.000000 | 1797.000000 | ... | 1797.000000 | 1797.0 |
| mean  | 0.0       | 0.303840    | 5.204786    | 11.835838   | 11.848080   | 5.781859    | 1.382270    | 0.129861    | 0.005565    | 1.993879    | ... | 3.725097    | 0.     |
| std   | 0.0       | 0.907192    | 4.754826    | 4.248842    | 4.287388    | 5.686418    | 3.325775    | 1.037383    | 0.094222    | 3.198180    | ... | 4.919408    | 0.     |
| min   | 0.0       | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    | ... | 0.000000    | 0.     |
| 25%   | 0.0       | 0.000000    | 1.000000    | 10.000000   | 10.000000   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    | ... | 0.000000    | 0.     |
| 50%   | 0.0       | 0.000000    | 4.000000    | 13.000000   | 13.000000   | 4.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    | ... | 1.000000    | 0.     |
| 75%   | 0.0       | 0.000000    | 9.000000    | 15.000000   | 15.000000   | 11.000000   | 0.000000    | 0.000000    | 0.000000    | 3.000000    | ... | 7.000000    | 0.     |
| max   | 0.0       | 8.000000    | 16.000000   | 16.000000   | 16.000000   | 16.000000   | 16.000000   | 15.000000   | 2.000000    | 16.000000   | ... | 16.000000   | 13.    |

8 rows x 64 columns

```
In [10]: X = df
y = dataset.target
```

```
In [11]: y
```

```
Out[11]: array([0, 1, 2, ..., 8, 9, 8])
```

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

```
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
X_scaled
```

```
Out[12]: array([[ 0.          , -0.33501649, -0.04308102, ..., -1.14664746,
        -0.5056698 , -0.19600752],
       [ 0.          , -0.33501649, -1.09493684, ...,  0.54856067,
        -0.5056698 , -0.19600752],
       [ 0.          , -0.33501649, -1.09493684, ...,  1.56568555,
        1.6951369 , -0.19600752],
       ...,
       [ 0.          , -0.33501649, -0.88456568, ..., -0.12952258,
        -0.5056698 , -0.19600752],
       [ 0.          , -0.33501649, -0.67419451, ...,  0.8876023 ,
        -0.5056698 , -0.19600752],
       [ 0.          , -0.33501649,  1.00877481, ...,  0.8876023 ,
        -0.26113572, -0.19600752]])
```

```
In [13]: from sklearn.model_selection import train_test_split
```

```
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=30)
```

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

```
model = LogisticRegression()
model.fit(X_train, y_train)
model.score(X_test, y_test)
```

```
Out[14]: 0.9722222222222222
```

```
In [20]: X.shape
```

```
Out[20]: (1797, 64)
```

#### Use PCA to reduce dimensions

```
In [15]: from sklearn.decomposition import PCA
```

```
pca = PCA(0.95)
X_pca = pca.fit_transform(X)
X_pca.shape
```

```
Out[15]: (1797, 29)
```

```
In [21]: pca.explained_variance_ratio_
```

```
Out[21]: array([0.14890594, 0.13618771, 0.11794594, 0.08409979, 0.05782415,
        0.0491691 , 0.04315987, 0.03661373, 0.03353248, 0.03078806,
        0.02372341, 0.02272697, 0.01821863, 0.01773855, 0.01467101,
        0.01409716, 0.01318589, 0.01248138, 0.01017718, 0.00905617,
        0.00889538, 0.00797123, 0.00767493, 0.00722904, 0.00695889,
        0.00596081, 0.00575615, 0.00515158, 0.0048954 ])
```

```
In [22]: pca.n_components_
```

```
Out[22]: 29
```

You can see that both combined retains  $0.14+0.13=0.27$  or 27% of important feature information

```
In [24]: X_train_pca, X_test_pca, y_train, y_test = train_test_split(X_pca, y, test_size=0.2, random_state=30)
```

```
In [25]: model = LogisticRegression(max_iter=1000)
model.fit(X_train_pca, y_train)
model.score(X_test_pca, y_test)
```

```
Out[25]: 0.9694444444444444
```

Let's now select only two components

```
In [26]: pca = PCA(n_components=2)
X_pca = pca.fit_transform(X)
X_pca.shape

Out[26]: (1797, 2)

In [27]: X_pca

Out[27]: array([[ -1.25946509,  21.27488304],
 [  7.95759075, -20.76870055],
 [  6.99193132, -9.95598473],
 ...,
 [ 10.80130987, -6.96024689],
 [-4.87210787, 12.42395562],
 [-0.34434859,  6.36555467]])
```

```
In [29]: pca.explained_variance_ratio_

Out[29]: array([0.14890594, 0.13618771])
```

You can see that both combined retains  $0.14+0.13=0.27$  or 27% of important feature information

```
In [30]: X_train_pca, X_test_pca, y_train, y_test = train_test_split(X_pca, y, test_size=0.2, random_state=30)

model = LogisticRegression(max_iter=1000)
model.fit(X_train_pca, y_train)
model.score(X_test_pca, y_test)

Out[30]: 0.6083333333333333
```

## **Singular Value Decomposition (SVD):**

The Singular Value Decomposition (SVD) of a matrix is a factorization of that matrix into three matrices. It has some interesting algebraic properties and conveys important geometrical and theoretical insights about linear transformations. It also has some important applications in data science. In this article, I will try to explain the mathematical intuition behind SVD and its geometrical meaning.

### **Mathematics behind SVD:**

The SVD of  $m \times n$  matrix  $A$  is given by the formula  $A = U \Sigma V^T$

#### **where:**

- $U$ :  $m \times m$  matrix of the orthonormal eigenvectors of  $AA^T$ .
- $V^T$ : transpose of a  $n \times n$  matrix containing the orthonormal eigenvectors of  $A^T A$ .
- $\Sigma$ : diagonal matrix with  $r$  elements equal to the root of the positive eigenvalues of  $AA^T$  or  $A^T A$  (both matrices have the same positive eigenvalues anyway).

## Implementation of Singular Value Decomposition (SVD):

jupyter Singular Value Decomposition (SVD) Last Checkpoint: 09/19/2023 (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted | Python 3 (ipykernel)

```
In [1]: from skimage.color import rgb2gray
from skimage import data
import matplotlib.pyplot as plt
import numpy as np
from scipy.linalg import svd
```

C:\ProgramData\anaconda3\Lib\site-packages\paramiko\transport.py:219: CryptographyDeprecationWarning: Blowfish has been deprecated  
"class": algorithms.Blowfish,

Singular Value Decomposition

```
In [ ]: X = np.array([[3, 3, 2], [2, 3, -2]])
print(X)
# perform SVD
U, singular, V_transpose = svd(X)
# print different components
print("U: ", U)
print("Singular array", singular)
print("V^T", V_transpose)
```

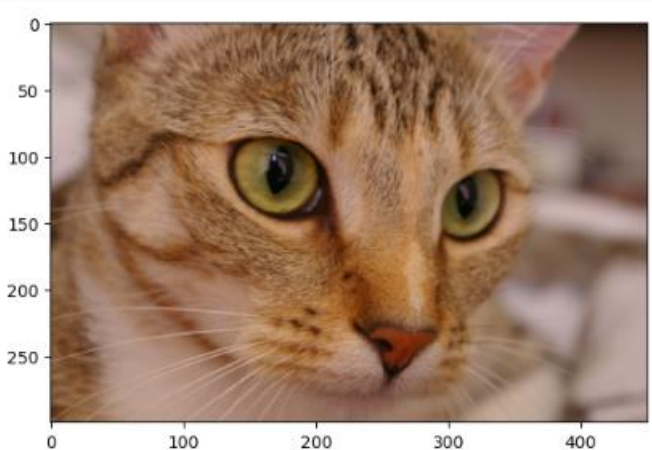
Calculate Pseudo inverse

```
In [3]: singular_inv = 1.0 / singular
# create m x n matrix of zeroes and put singular values in it
s_inv = np.zeros(X.shape)
s_inv[0][0] = singular_inv[0]
s_inv[1][1] = singular_inv[1]
# calculate pseudoinverse
M = np.dot(np.dot(V_transpose.T, s_inv.T), U.T)
print(M)
```

```
[[ 0.11462451  0.04347826]
 [ 0.07114625  0.13043478]
 [ 0.22134387 -0.26086957]]
```

SVD on image compression

```
In [4]: cat = data.chelsea()
plt.imshow(cat)
# convert to grayscale
gray_cat = rgb2gray(cat)
```



### Calculate the SVD and plot the image

```
In [6]: U, S, V_T = svd(gray_cat, full_matrices=False)
S = np.diag(S)
fig, ax = plt.subplots(5, 2, figsize=(8, 20))

curr_fig = 0
for r in [5, 10, 70, 100, 200]:
    cat_approx = U[:, :r] @ S[0:r, :r] @ V_T[:, :r, :]
    ax[curr_fig][0].imshow(cat_approx, cmap='gray')
    ax[curr_fig][0].set_title("k = "+str(r))
    ax[curr_fig, 0].axis('off')
    ax[curr_fig][1].set_title("Original Image")
    ax[curr_fig][1].imshow(gray_cat, cmap='gray')
    ax[curr_fig, 1].axis('off')
    curr_fig += 1
plt.show()
```

k = 5



Original Image



k = 10



Original Image



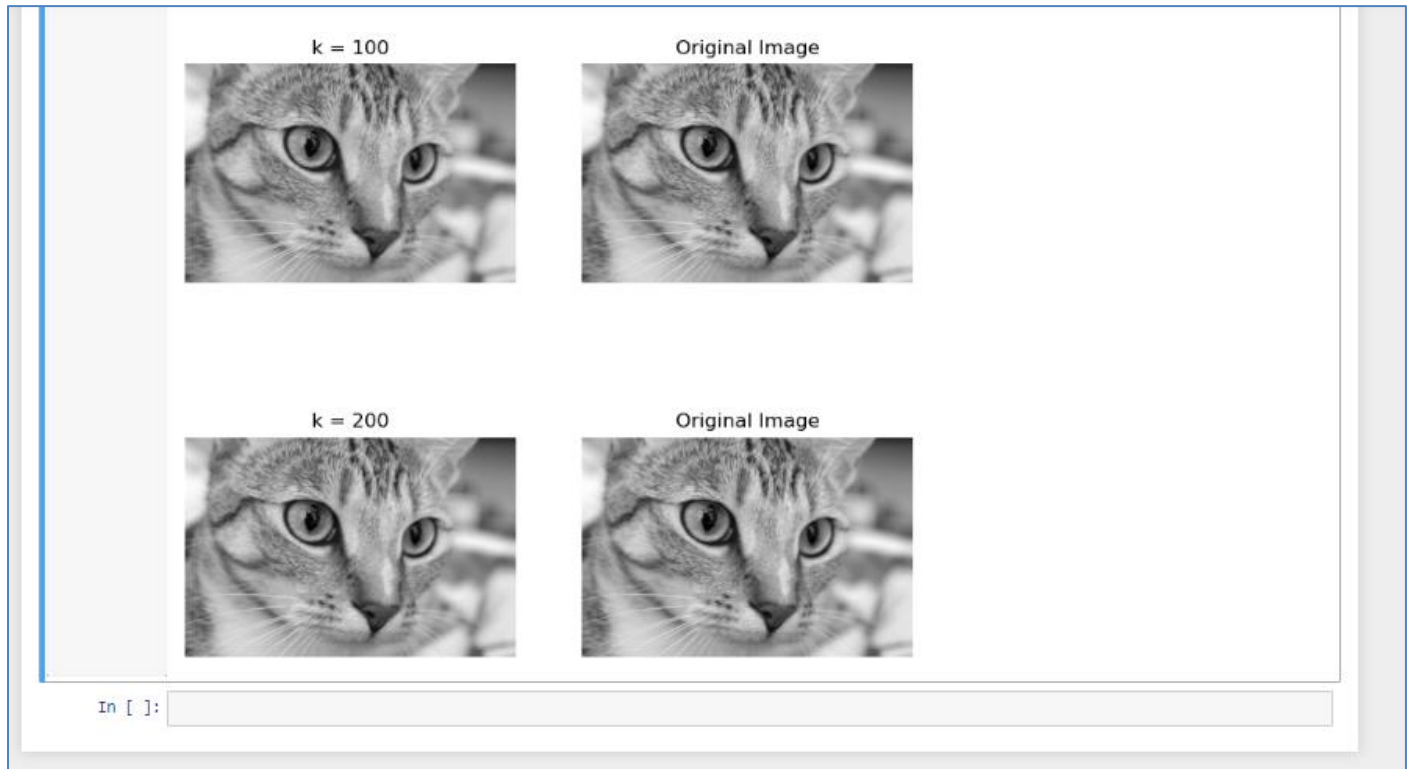
k = 70



Original Image







### ● **Conclusion:**

Both PCA and SVD are valuable tools for simplifying and analysing high-dimensional data, which can lead to improved data understanding and more efficient machine learning models. The choice between PCA and SVD depends on the specific problem and the mathematical and computational tools available.