

Load MNIST data

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4
5 df = pd.read_csv('MNIST_train.csv')
6
7 print(df.head())
```

```
   label  pixel0  pixel1  pixel2  pixel3  pixel4  pixel5  pixel6  pixel7  \
0       1       0       0       0       0       0       0       0       0
1       0       0       0       0       0       0       0       0       0
2       1       0       0       0       0       0       0       0       0
3       4       0       0       0       0       0       0       0       0
4       0       0       0       0       0       0       0       0       0

   pixel8  ...  pixel774  pixel775  pixel776  pixel777  pixel778  pixel779  \
0       0  ...         0         0         0         0         0         0
1       0  ...         0         0         0         0         0         0
2       0  ...         0         0         0         0         0         0
3       0  ...         0         0         0         0         0         0
4       0  ...         0         0         0         0         0         0

   pixel780  pixel781  pixel782  pixel783
0         0         0         0         0
1         0         0         0         0
2         0         0         0         0
3         0         0         0         0
4         0         0         0         0

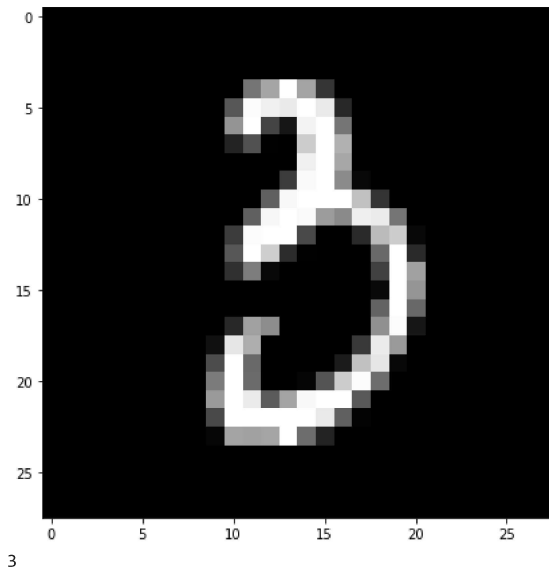
[5 rows x 785 columns]
```

```
1 label = df['label']
2 data = df.drop('label', axis=1)
```

```
1 data.shape,label.shape
```

```
((42000, 784), (42000,))
```

```
1 ### Sanity Check ###
2 plt.figure(figsize=(7,7))
3 idx= 150
4
5 image_150 = data.iloc[idx].values.reshape(28,28)
6 plt.imshow(image_150, interpolation = 'none', cmap='gray')
7 plt.show()
8 label[150]
```



Visualization using PCA

```
1 labels = label.head(15000)
2 data = data.head(15000)
3 print("the shape of sample data = ", data.shape)
```

```
the shape of sample data =  (15000, 784)
```

```
1 from sklearn.preprocessing import StandardScaler
2 standardized_data = StandardScaler().fit_transform(data)
3 print(standardized_data.shape)
```

(15000, 784)

▼ Find out why this method is not working

```
mean = data.mean() std = data.std() df_standardized = data.apply(lambda x: (x-mean)/ std) df_standardized.head()
```

```
1 #finding covariance mar=triz
2
3 sample_data = standardized_data
4 cov_matrix = np.matmul(sample_data.T, sample_data)
5 cov_matrix.shape
```

(784, 784)

```
1 from scipy.linalg import eigh
2 #TAKING TOP TWO EIGH ONLY
3 values, vectors = eigh(cov_matrix, eigvals = (782, 783)) ### top two eigvalues only calculated
4 vectors.shape, vectors.T.shape
```

((784, 2), (2, 784))

```
1 vectors = vectors.T
```

```
1 new_coordinates = np.matmul(vectors, sample_data.T)
2 vectors.shape, sample_data.shape
3 new_coordinates.shape ## Represents projected 15,000 input data points on top two eigen vectors
```

(2, 15000)

```
1 new_coordinates.shape, labels.shape
```

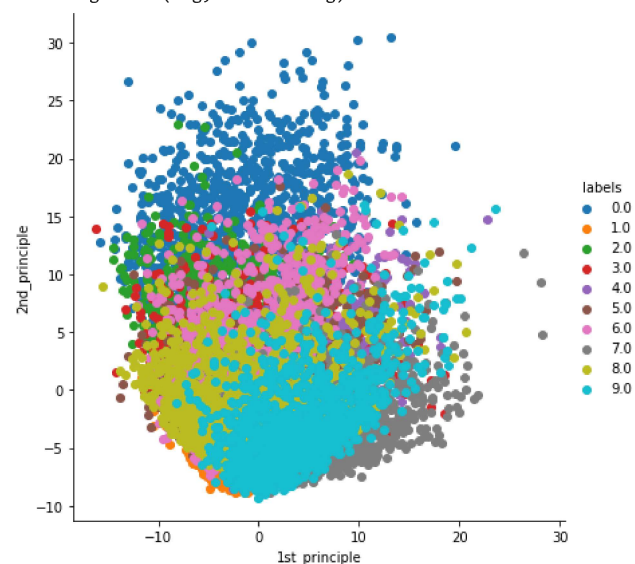
👤 ((2, 15000), (15000,))

```
1 # Data massaging
2
3 new_coordinates = np.vstack((new_coordinates, labels)).T #Note do not execute this cell again and again will cause error unnecessarily.
4 dataframe = pd.DataFrame(data = new_coordinates, columns=('1st_principle', '2nd_principle', 'labels'))
5 dataframe.head() ##Transformed data i.e projected data points
```

	1st_principle	2nd_principle	labels	🔗
0	-5.558661	-5.043558	1.0	
1	6.193635	19.305278	0.0	
2	-1.909878	-7.678775	1.0	
3	5.525748	-0.464845	4.0	
4	6.366527	26.644289	0.0	

```
1 import seaborn as sns
2 sns.FacetGrid(dataframe, hue= 'labels', size=6).map(plt.scatter,
3                                                    '1st_principle',
4                                                    '2nd_principle').add_legend()
5 plt.show()
```

/usr/local/lib/python3.8/dist-packages/seaborn/axisgrid.py:337: UserWarning: The `size` parameter is deprecated. Use `figsize` or `height` instead.
warnings.warn(msg, UserWarning)



PCA using Scikit-learn

```
1 from sklearn import decomposition
2 pca = decomposition.PCA()
```

```
1 pca.n_components = 2
2 pca_data = pca.fit_transform(sample_data)
3
4 pca_data.shape
```

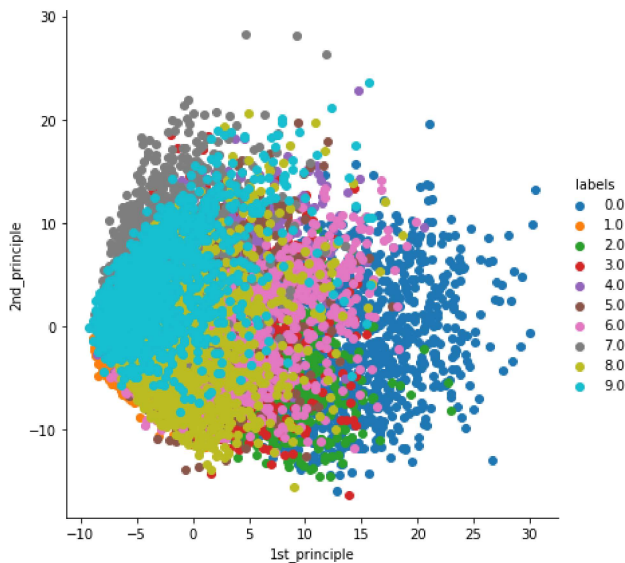
```
(15000, 2)
```

```
1 pca_data = np.vstack((pca_data.T, labels)).T
2 pca_df = pd.DataFrame(data=pca_data, columns=('1st_principle', '2nd_principle', 'labels'))
3 pca_df.head()
```

	1st_principle	2nd_principle	labels
0	-5.043562	-5.558364	1.0
1	19.305334	6.192572	0.0
2	-7.678767	-1.910249	1.0
3	-0.464817	5.525317	4.0
4	26.644314	6.366089	0.0

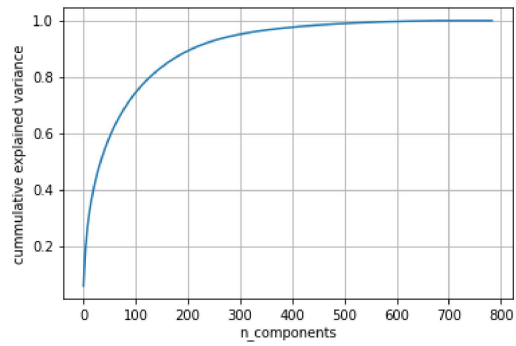
```
1 sns.FacetGrid(pca_df, hue='labels', size=6).map(plt.scatter,
2                                                  '1st_principle',
3                                                  '2nd_principle').add_legend()
4 plt.show()
```

```
/usr/local/lib/python3.8/dist-packages/seaborn/axisgrid.py:337: UserWarning: The `size` parameter is deprecated.
warnings.warn(msg, UserWarning)
```



PCA FOR DIMENTIONALITY REDUCTION ----> Checking what numbers of eigen vectors(Principle components) are required w.r.t data variance(Percentage of information to keep while performing dimentionality reduction)

```
1 pca.n_components = 784
2 pca_data = pca.fit_transform(sample_data)
3
4 percentage_var_explained = pca.explained_variance_/np.sum(pca.explained_variance_)
5
6 cum_var_explained = np.cumsum(percentage_var_explained)    #cummulative sum of variance ratios
7
8 plt.plot(cum_var_explained)
9 plt.grid()
10 plt.ylabel('cummulative explained variance')
11 plt.xlabel('n_components')
12 plt.show()
13 percentage_var_explained.shape, cum_var_explained.shape
```



((784,), (784,))