

Generating and Comparing Different Datasets

```
In [44]: import pandas as pd
import numpy as np

import matplotlib.pyplot as plt
import matplotlib
import seaborn as sns

from sklearn.datasets import make_classification
from sklearn.linear_model import LassoCV, LogisticRegression
from sklearn.decomposition import PCA
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

from sklearn.preprocessing import StandardScaler
```

Make Two Feature Datasets For Visualization:

Changing class separation changes the difficulty of the classification task. The data points no longer remain easily separable in case of lower class separation.

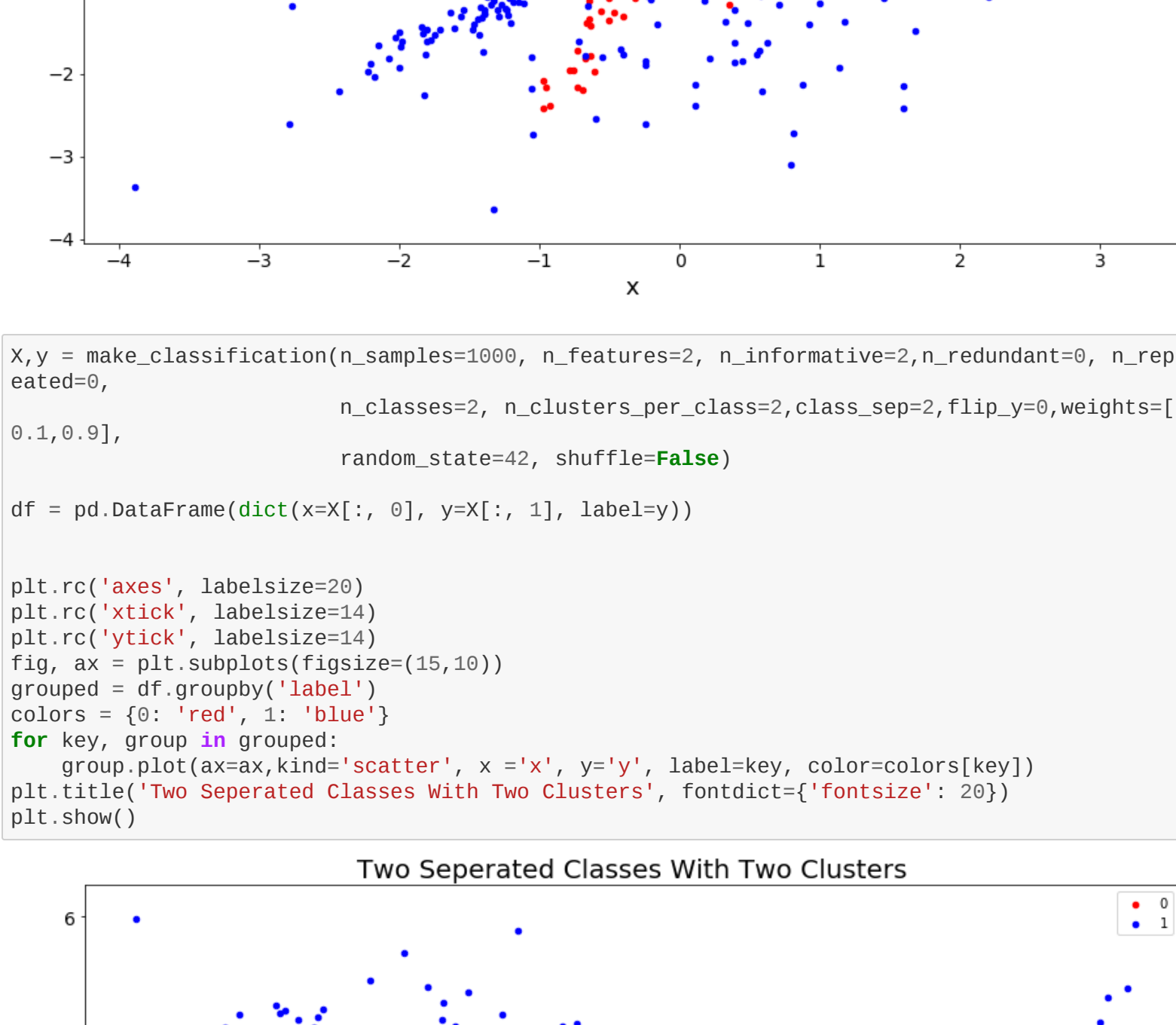
The sklearn make_classification API helps us create datasets with different distributions. We are focusing on class_sep hyperparameter to evaluate the performances of Lasso models.

Below first have a look at how make_classification API works.

Let's have two datasets with class separations of 0.2 and 2. Both datasets have two clusters per class. They both have 2 features for us to be able to visualize the data.

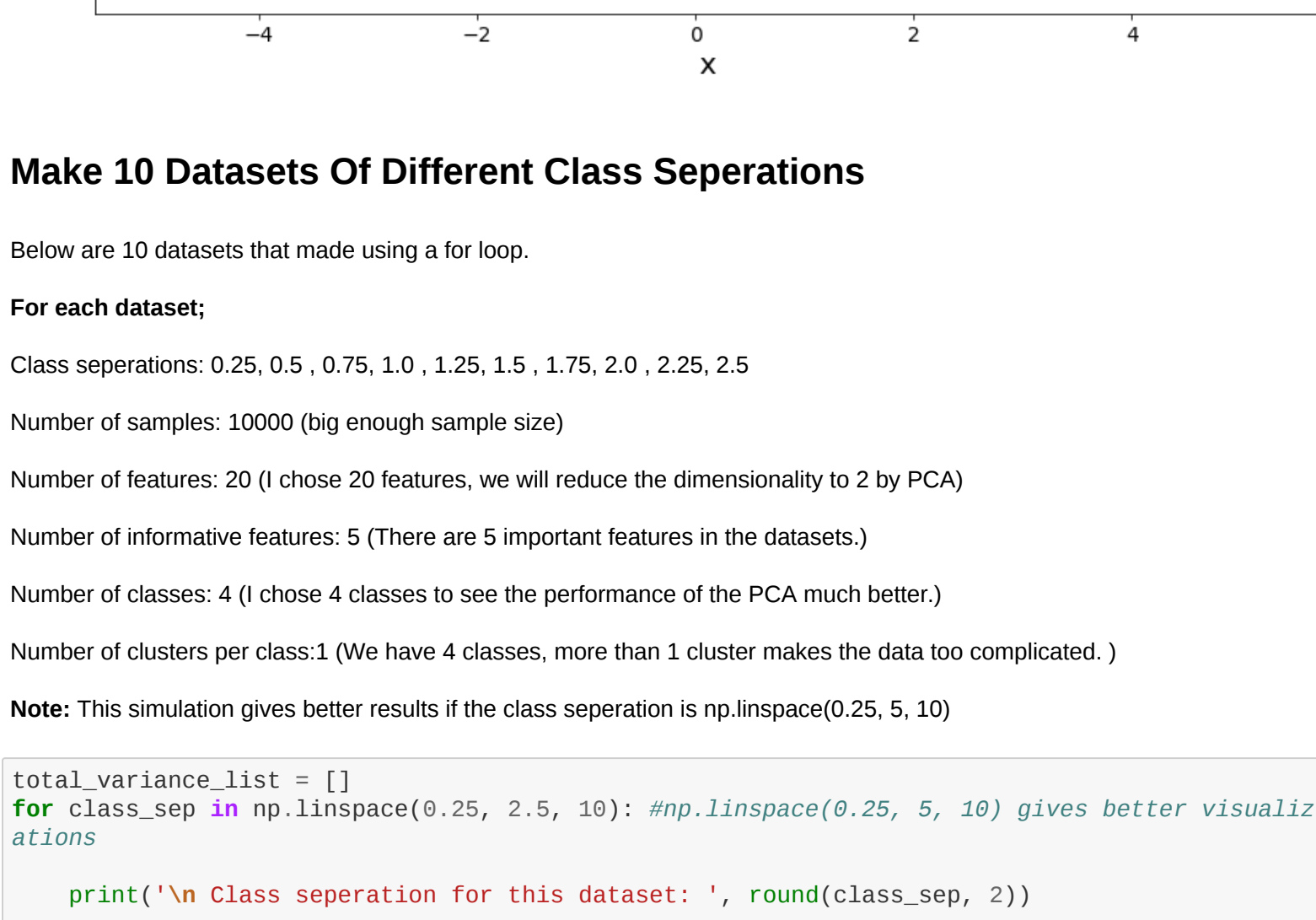
```
In [17]: X, y = make_classification(n_samples=1000, n_features=2, n_informative=2, n_redundant=0, n_repeated=0,
                                n_classes=2, n_clusters_per_class=2, class_sep=0.2, flip_y=0, weights=[0.1, 0.9],
                                random_state=42, shuffle=False)
df = pd.DataFrame(dict(x=X[:, 0], y=X[:, 1], label=y))

plt.rc('axes', labelsiz=10)
plt.rc('tick', labelsiz=4)
plt.rc('y', labelsiz=14)
fig, ax = plt.subplots(figsize=(15, 10))
grouper = df.groupby('label')
colors = {'red', 1: 'blue'}
for key, group in grouper:
    group.plot(ax=ax, kind='scatter', x='x', y='y', label=key, color=colors[key])
plt.title('Two Non-separated Classes With Two Clusters', fontdict={'fontsize': 20})
plt.show()
```



```
In [16]: X, y = make_classification(n_samples=1000, n_features=2, n_informative=2, n_redundant=0, n_repeated=0,
                                n_classes=2, n_clusters_per_class=2, class_sep=2, flip_y=0, weights=[0.1, 0.9],
                                random_state=42, shuffle=False)
df = pd.DataFrame(dict(x=X[:, 0], y=X[:, 1], label=y))

plt.rc('axes', labelsiz=20)
plt.rc('tick', labelsiz=4)
plt.rc('y', labelsiz=14)
fig, ax = plt.subplots(figsize=(15, 10))
grouper = df.groupby('label')
colors = {'red', 1: 'blue'}
for key, group in grouper:
    group.plot(ax=ax, kind='scatter', x='x', y='y', label=key, color=colors[key])
plt.title('Two Separated Classes With Two Clusters', fontdict={'fontsize': 20})
plt.show()
```



Make 10 Datasets Of Different Class Separations

Below are 10 datasets that made using a loop.

For each dataset:

Class separations: 0.25, 0.5, 0.75, 1.0, 1.25, 1.5, 1.75, 2.0, 2.25, 2.5

Number of samples: 10000 (big enough sample size)

Number of features: 20 (I chose 20 features, we will reduce the dimensionality to 2 by PCA)

Number of informative features: 5 (There are 5 important features in the datasets)

Number of classes: 4 (I chose 4 classes to see the performance of the PCA much better)

Number of clusters per class: 1 (We have 4 classes, more than 1 cluster makes the data too complicated.)

Note: This simulation gives better results if the class separation is np.linspace(0.25, 5, 10)

```
In [23]: total_variance_list = []
for class_sep in np.linspace(0.25, 2.5, 10): #np.linspace(0.25, 5, 10) gives better visualizations
    print('\n Class separation for this dataset: ', round(class_sep, 2))

    # Generate 10 datasets with different class separations
    X, y = make_classification(class_sep=class_sep, n_samples=10000, n_features=20,
                              n_informative=5, n_classes=4, n_clusters_per_class=1)

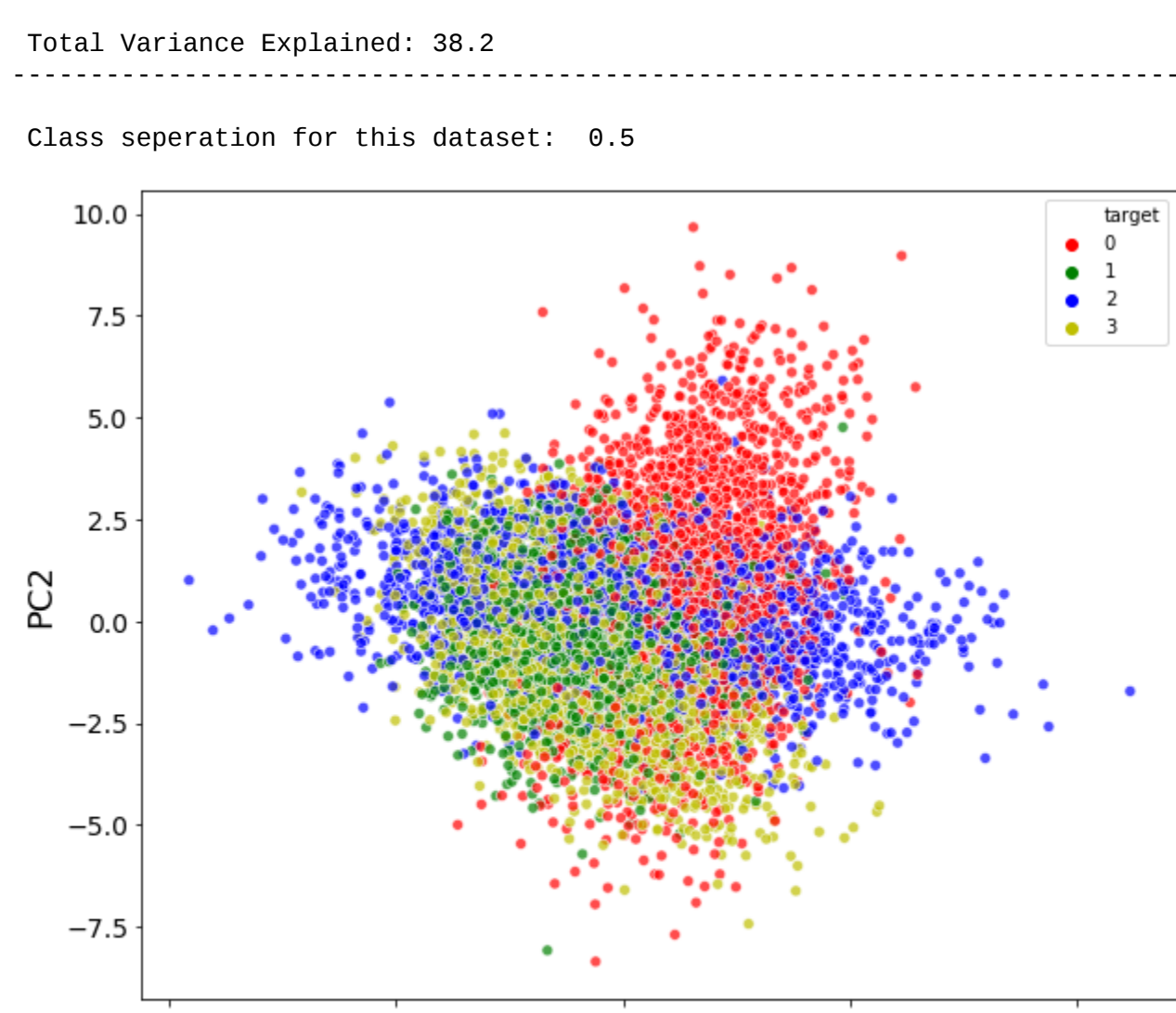
    # reduce dimensions by PCA
    pca = PCA(n_components=2)
    principalComponents = pca.fit_transform(X)

    # Make a dataframe using the new components
    df = pd.DataFrame(data=principalComponents, columns=['PC1', 'PC2'])
    # Add the target column to the dataframe
    df['target'] = y

    # plot the new components of the dataset
    fig = plt.figure(figsize=(10, 8))
    sns.scatterplot(x=df['PC1'], y=df['PC2'], hue=df['target'], palette=['r', 'g', 'b', 'y'], alpha=0.7)
    plt.show()

    # print the explained variance
    total_variance_explained = round(np.sum(pca.explained_variance_ratio_) * 100, 2)
    total_variance_list.append(total_variance_explained)
    print('variance of each component', pca.explained_variance_ratio_)
    print('\n Total Variance Explained: ', total_variance_explained)
    print('=' * 100)
```

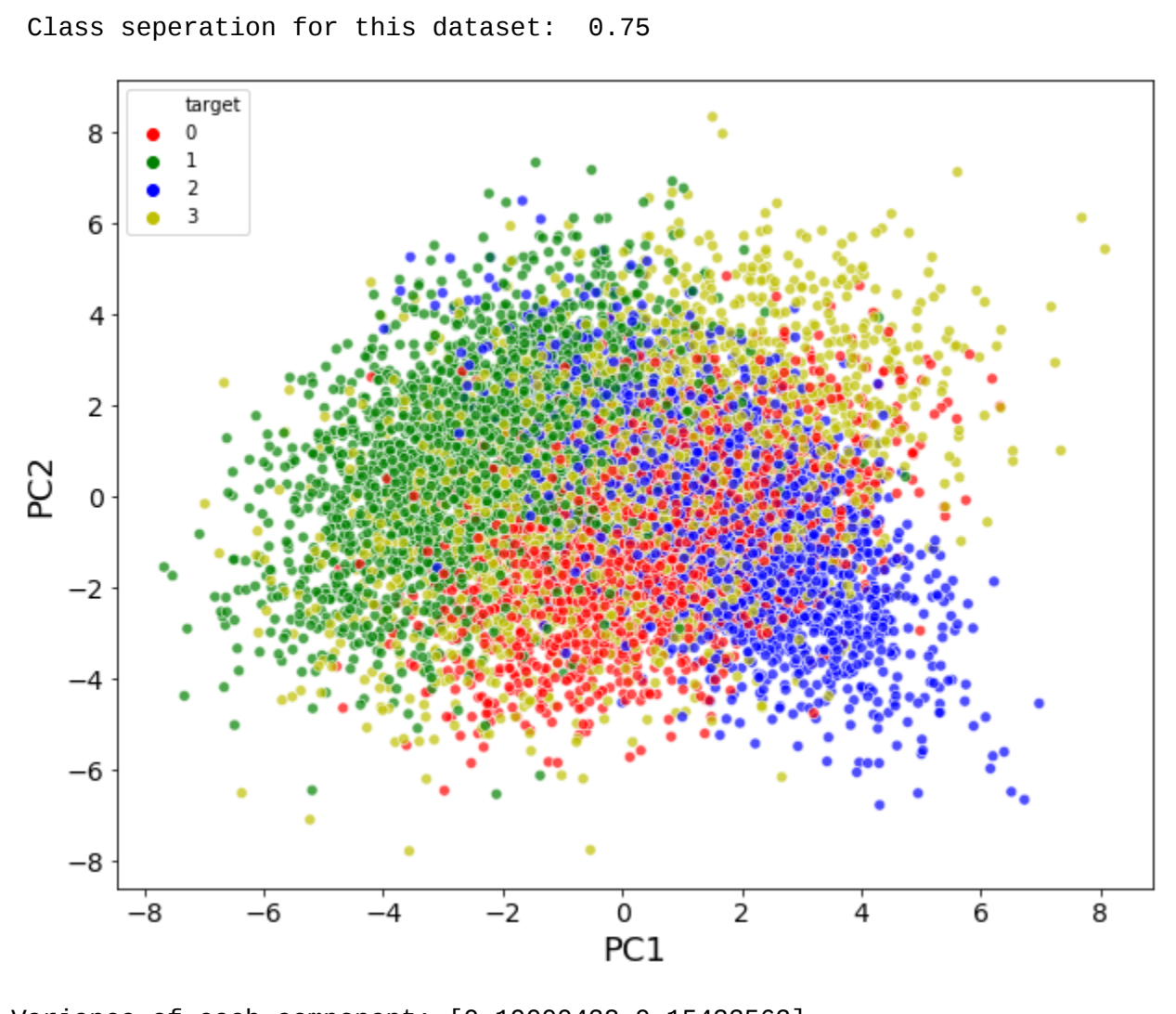
Class separation for this dataset: 0.25



Variance of each component: [0.24315515 0.13888963]

Total Variance Explained: 38.2

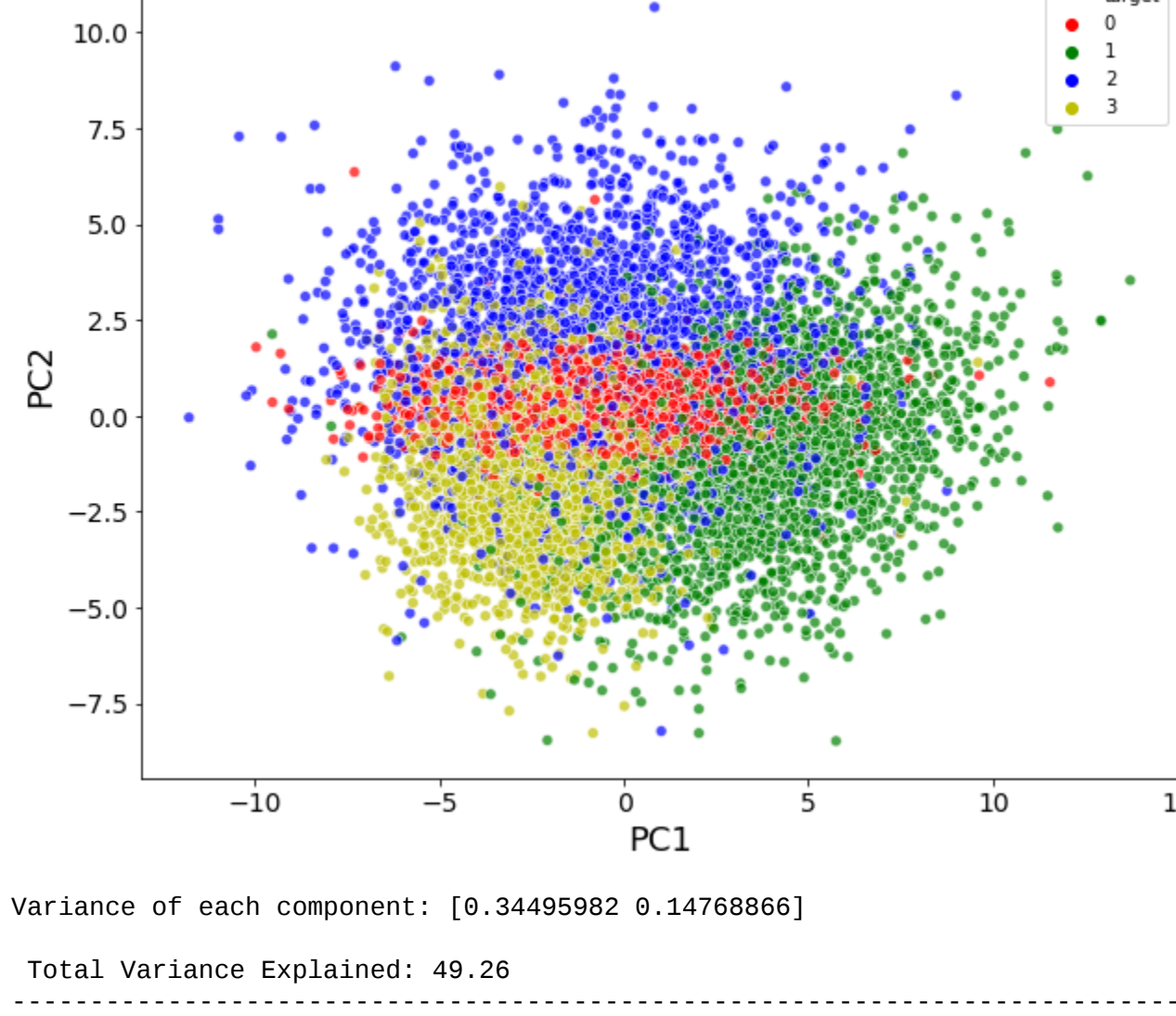
Class separation for this dataset: 0.5



Variance of each component: [0.18557843 0.15182282]

Total Variance Explained: 33.66

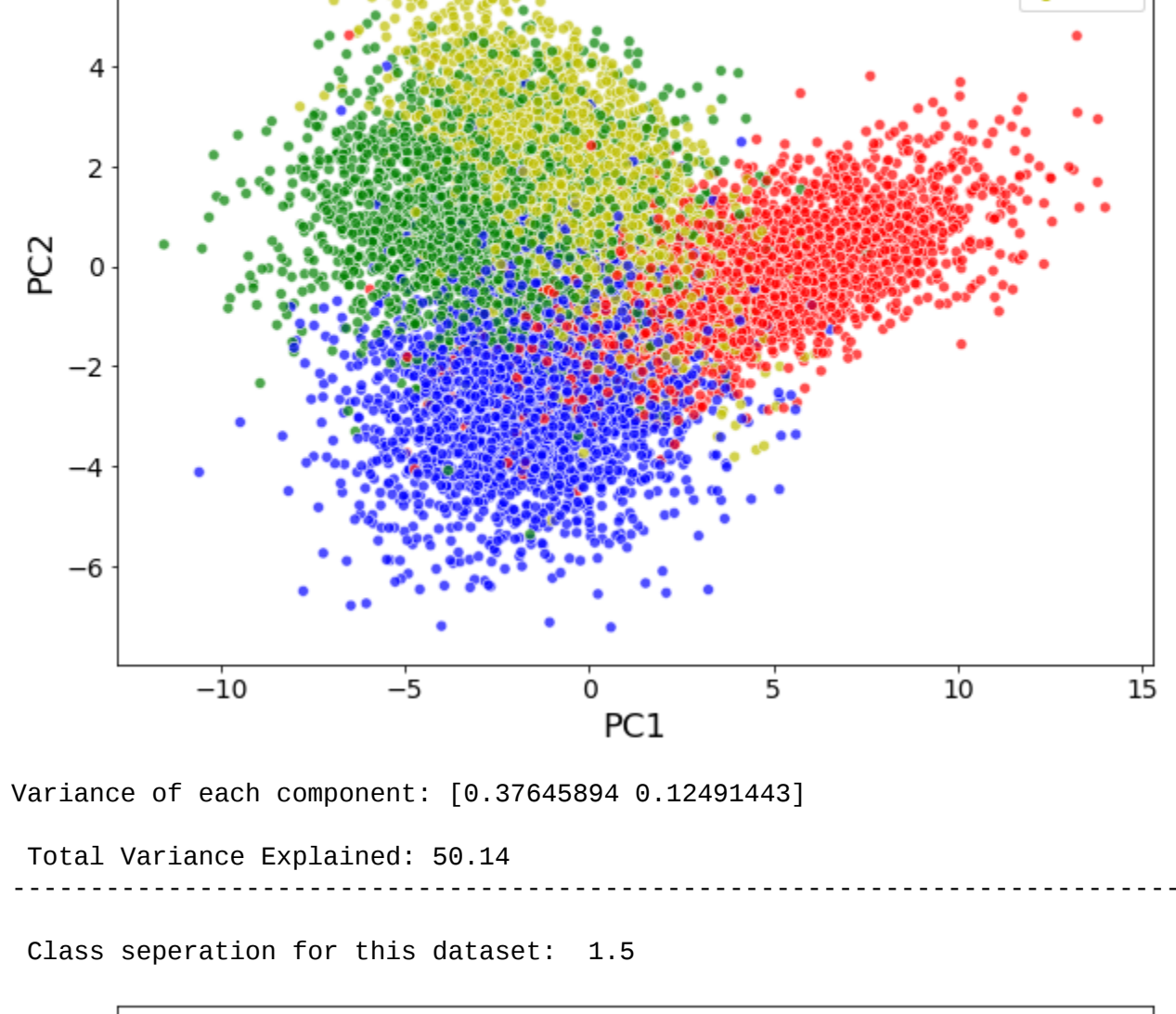
Class separation for this dataset: 0.75



Variance of each component: [0.19098428 0.15433563]

Total Variance Explained: 34.43

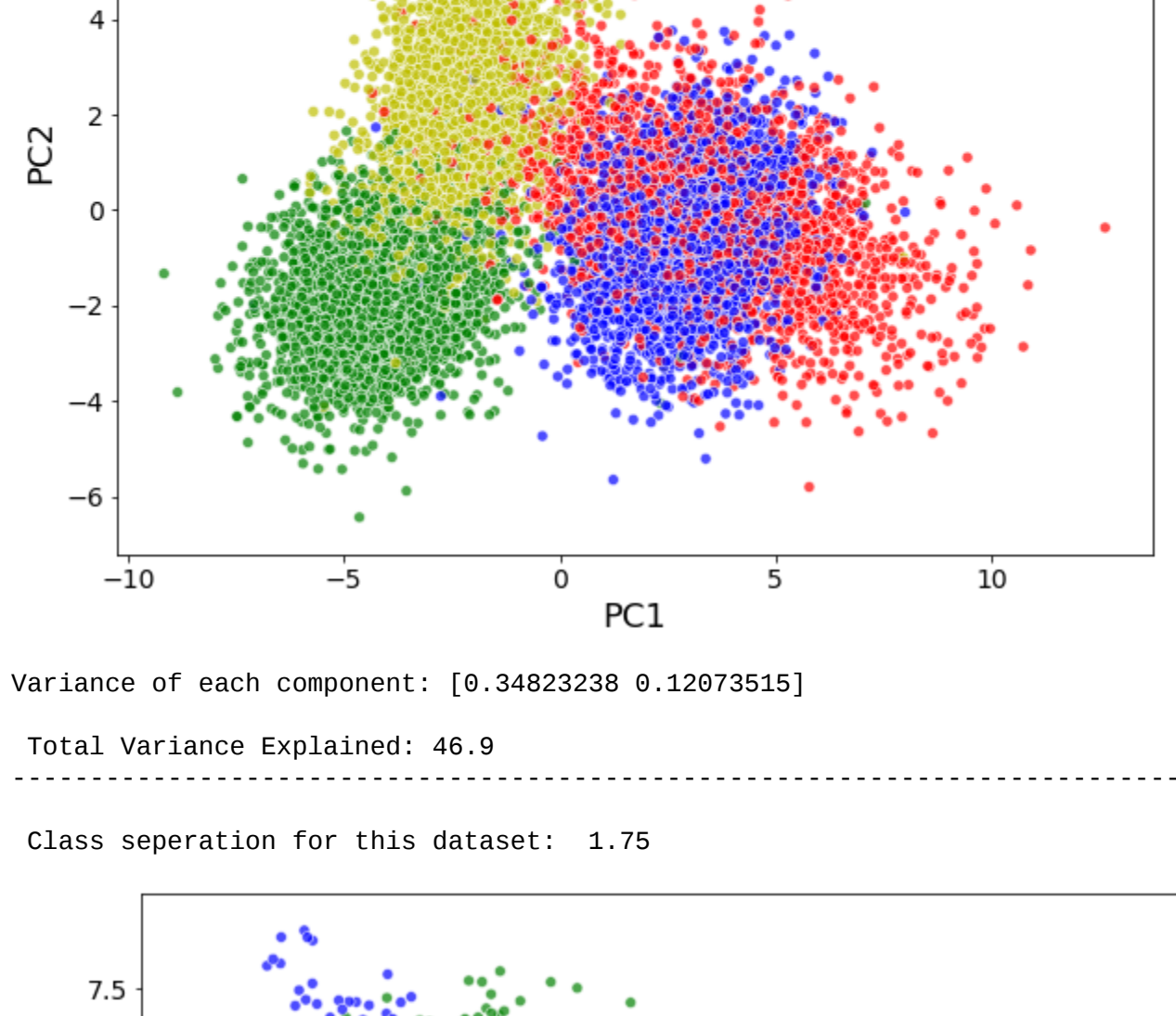
Class separation for this dataset: 1.0



Variance of each component: [0.34495982 0.14768866]

Total Variance Explained: 49.28

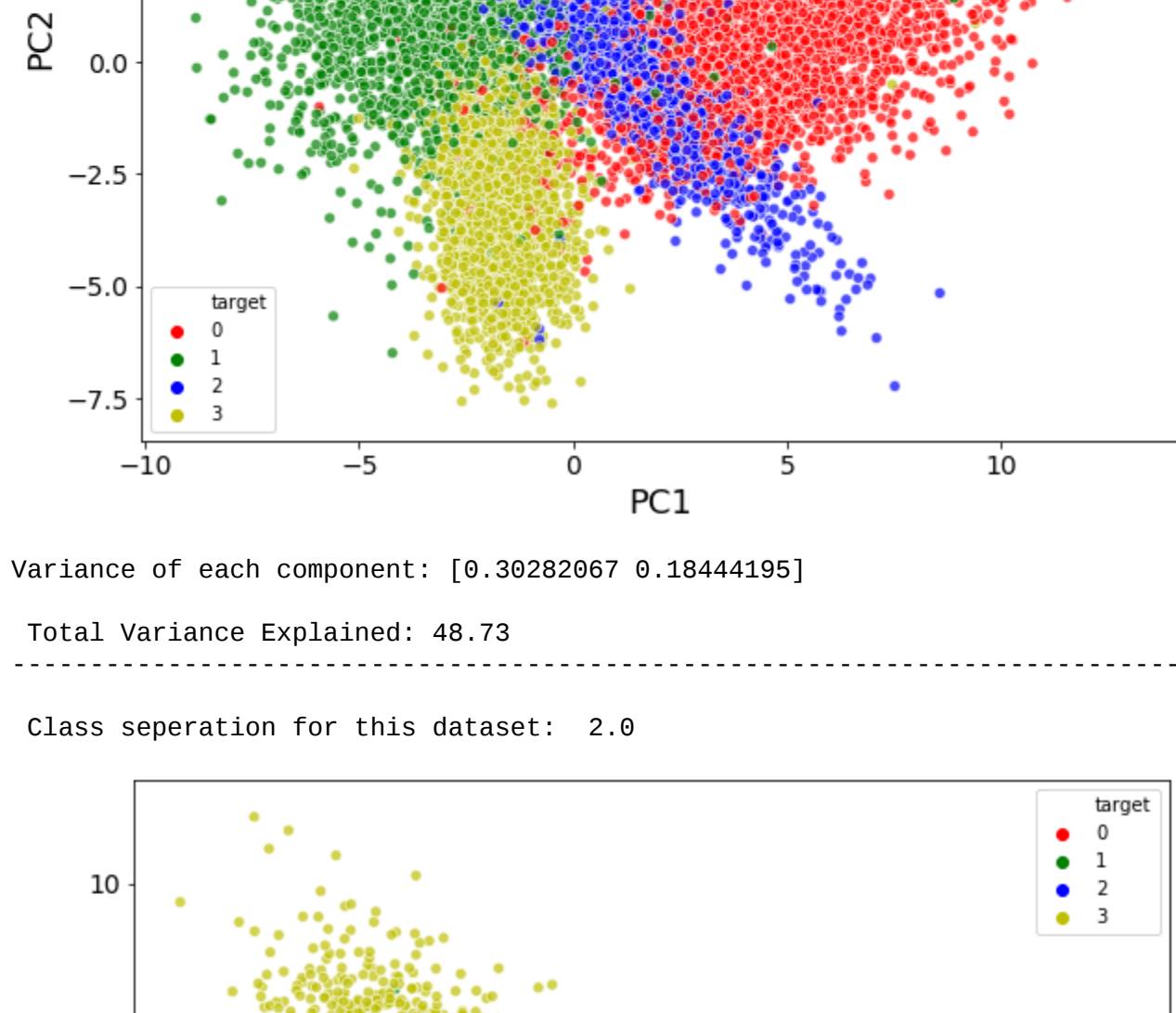
Class separation for this dataset: 1.25



Variance of each component: [0.37845894 0.12491443]

Total Variance Explained: 58.14

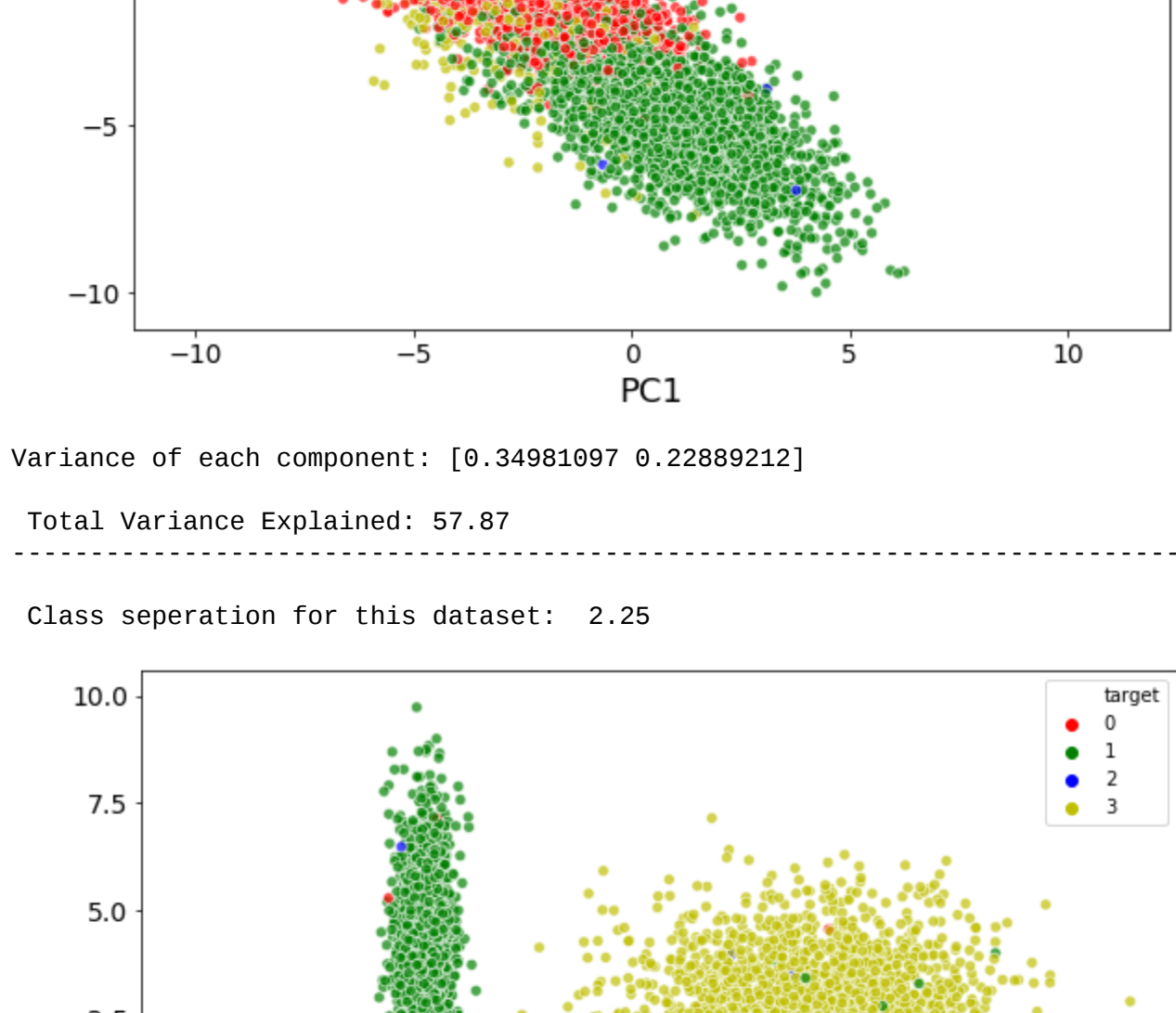
Class separation for this dataset: 1.5



Variance of each component: [0.34823238 0.12072515]

Total Variance Explained: 46.9

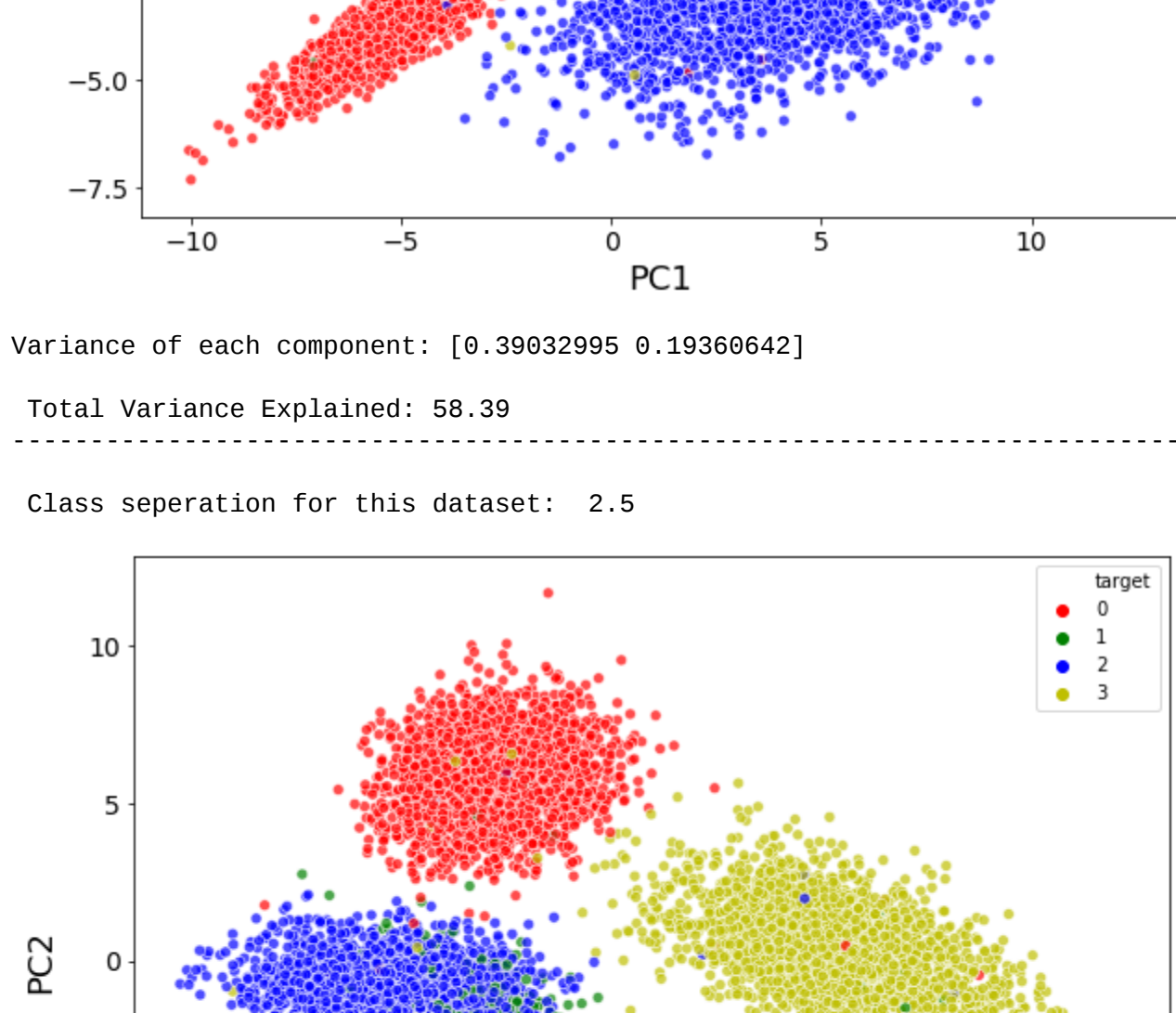
Class separation for this dataset: 1.75



Variance of each component: [0.30282867 0.18444195]

Total Variance Explained: 48.73

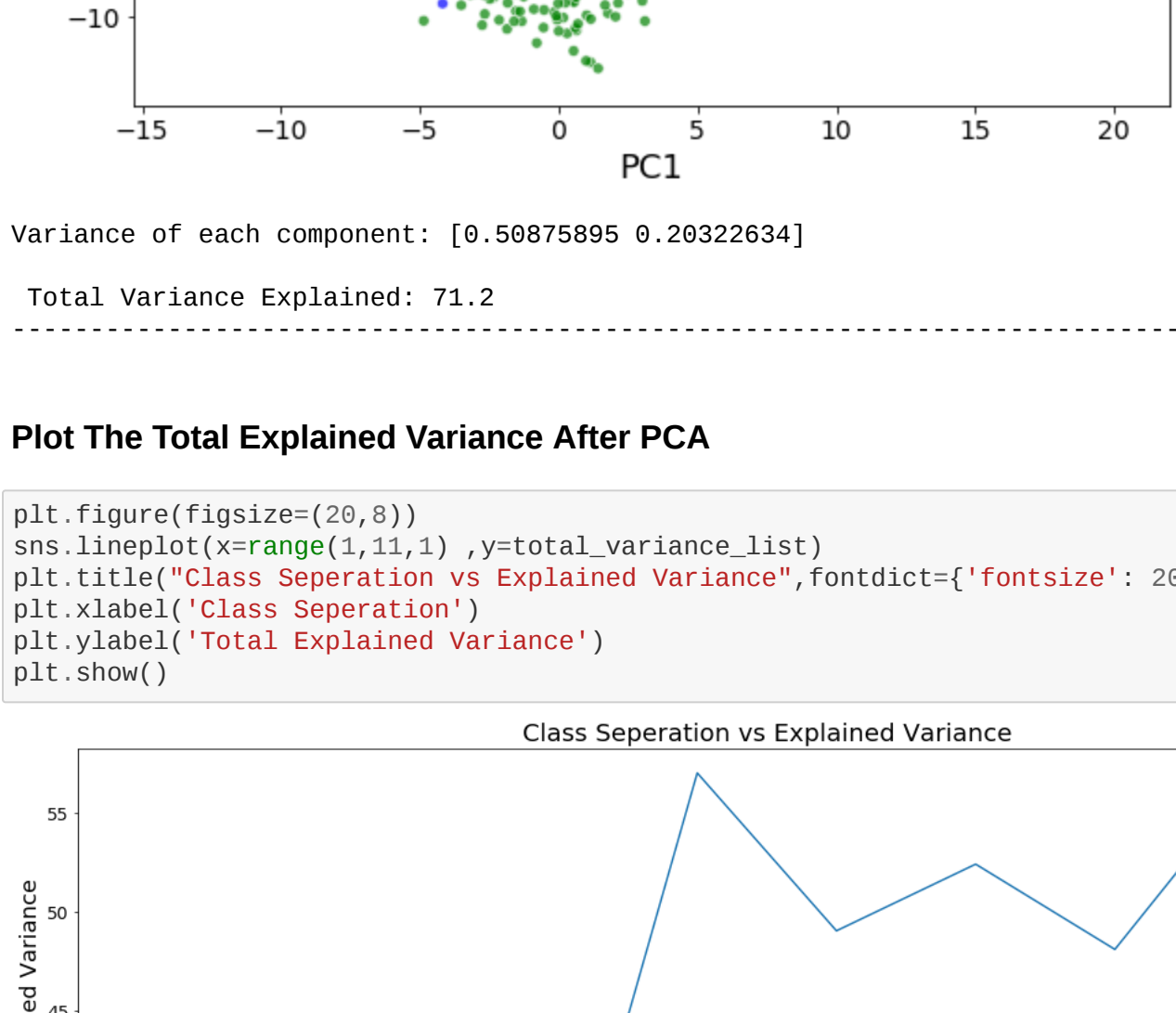
Class separation for this dataset: 2.0



Variance of each component: [0.34981097 0.22889212]

Total Variance Explained: 57.87

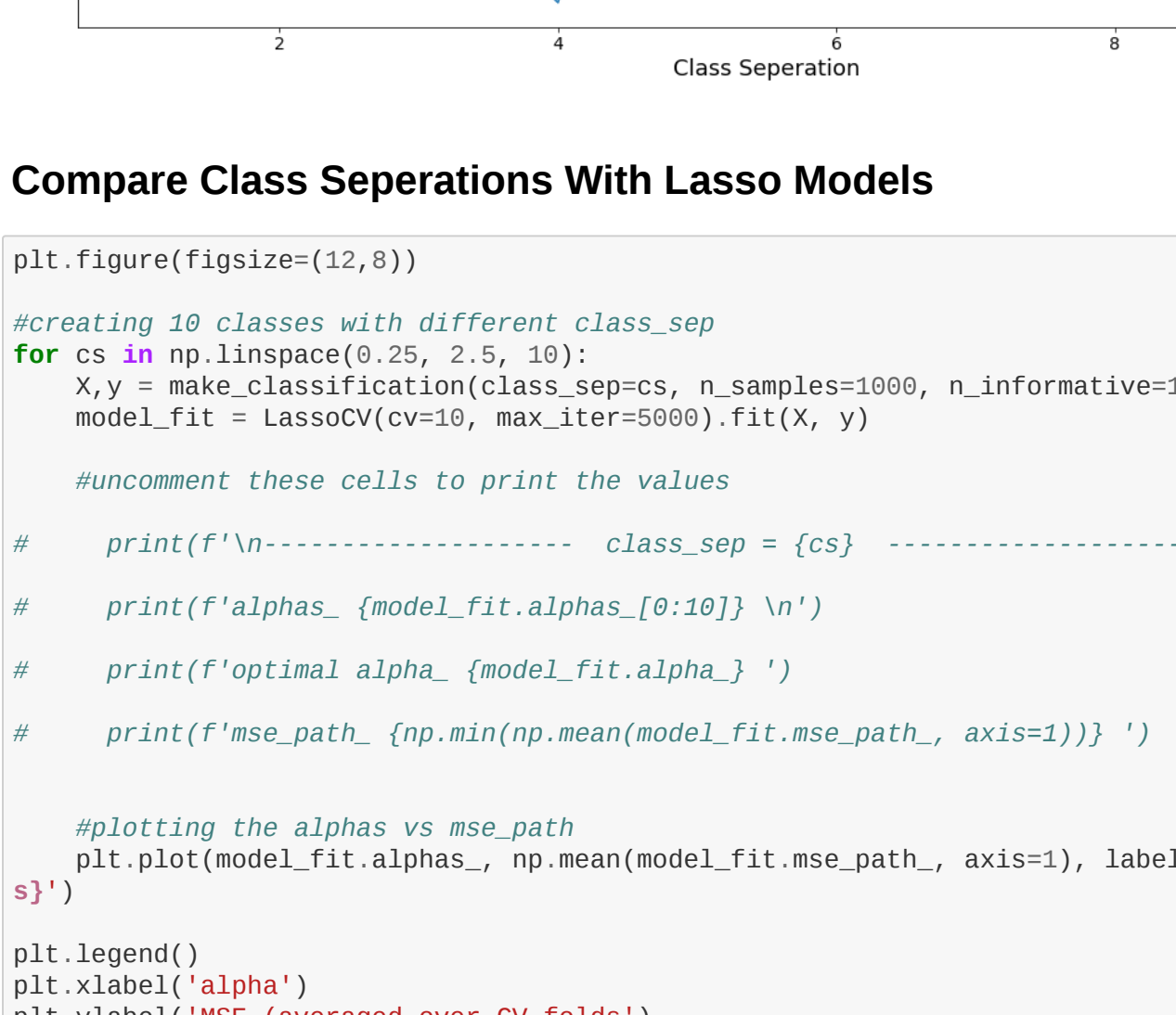
Class separation for this dataset: 2.25



Variance of each component: [0.39032995 0.19369842]

Total Variance Explained: 58.39

Class separation for this dataset: 2.5



Variance of each component: [0.58875895 0.28322634]

Total Variance Explained: 71.2

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5

Class separation for this dataset: 2.5