

```
In [1]: 1 import numpy as np
        2 import pandas as pd
        3 import seaborn as sns
        4 import matplotlib.pyplot as plt
        5
        6 from sklearn.decomposition import PCA
        7
```

```
In [2]: 1 np.set_printoptions(suppress=True, linewidth=150, precision=2)
        2 np.random.seed(123456)
```

Clusters Specification

```
In [67]: 1 N = 1000 # Number of Entities
        2 K = 7 # Number of Clusters
        3 V = 5 # Number of Features
        4 minimum = 50 # Minimum Number of Entities in each cluster
        5
        6 remaining = N - K * minimum
        7 remaining
```

Out[67]: 650

```
In [68]: 1 ur = [np.random.uniform() for k in range(K-1)]
        2 ur.sort()
```

```
In [69]: 1 structures = []
        2 tmp = []
        3 for k in range(K-1):
        4     if k == 0:
        5         tmp.append(ur[k]-0)
        6     elif k == K-2:
        7         tmp.append(1-ur[k])
        8     else:
        9         tmp.append(ur[k+1]-ur[k])
        10
```

```
In [70]: 1 structures = [int(remaining*i+ minimum) for i in tmp]
        2 structures += [N - sum(structures)]
        3 print(structures, sum(structures))
```

[195, 173, 145, 79, 268, 79, 61] 1000

```
In [71]: 1 d1, d2 = 1, -1
        2 a = 1.0
```

```
In [72]: 1 mu_1 = np.multiply(np.random.uniform(low=d2, high=d1, size=K), 1.0)
2 mu_2 = np.multiply(np.random.uniform(low=d2, high=d1, size=K), 0.75)
3 mu_3 = np.multiply(np.random.uniform(low=d2, high=d1, size=K), 0.5)
4 mu_4 = np.multiply(np.random.uniform(low=d2, high=d1, size=K), 0.25)
5 mu_1
```

```
Out[72]: array([-0.48,  0.36, -0.27, -0.8 ,  0.64,  0.59, -0.11])
```

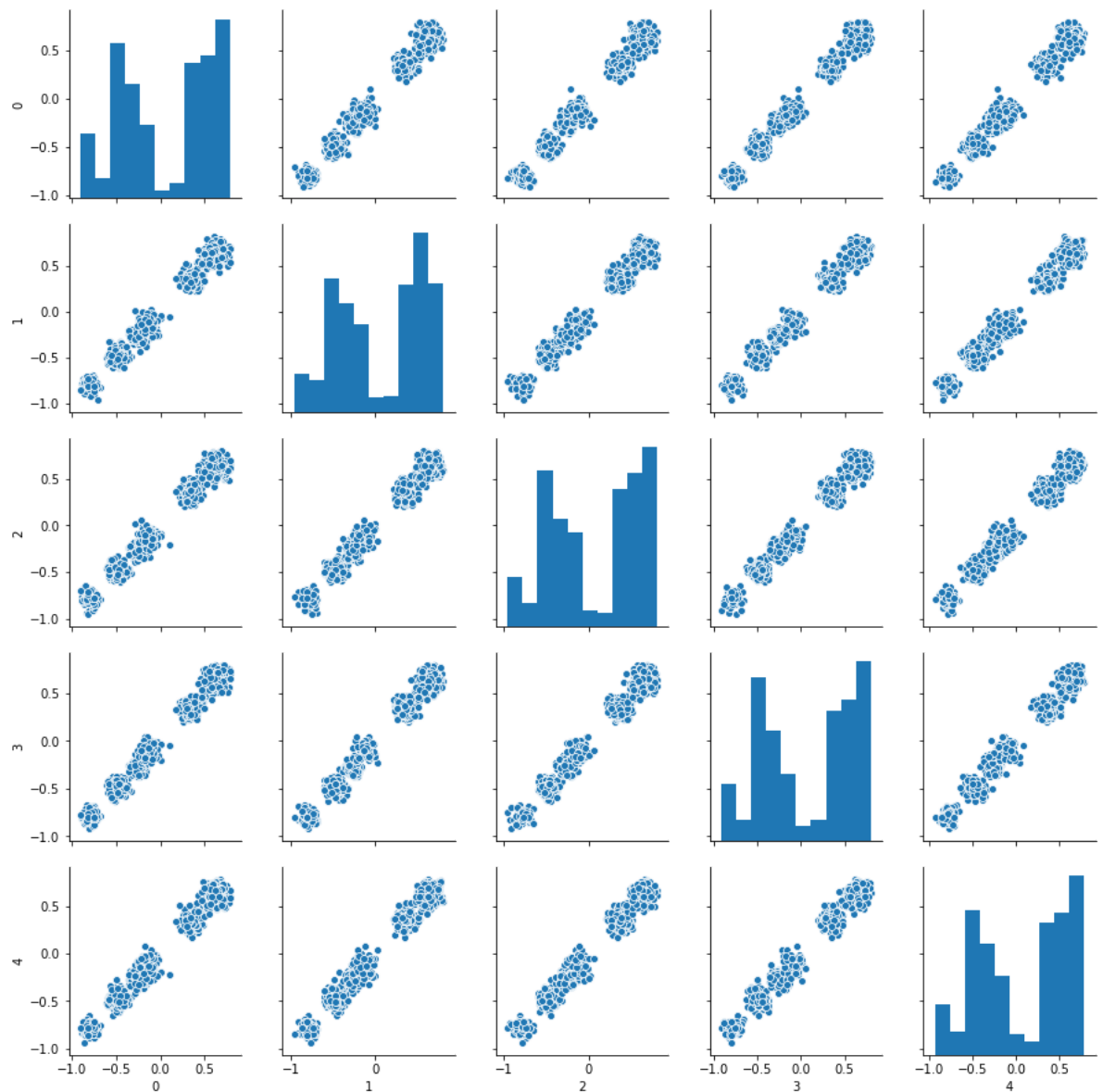
```
In [73]: 1 d1_ = 0.05*2
2 d2_ = 0.025*2
3
4 cov = np.diag(np.random.uniform(low=d2_, high=d1_, size=(K, K)))
5 cov
```

```
Out[73]: array([0.06, 0.06, 0.05, 0.05, 0.06, 0.06, 0.07])
```

Create Y matrix with a=1 (few or no intermix)

```
In [74]: 1 Y1 = np.zeros([N, V])
2 interval = 0
3 for k in range(K):
4     for i in range(interval, structures[k]+interval):
5         Y1[i, :] = np.random.normal(loc=mu_1[k], scale=cov[k], size=V)
6     interval += structures[k]
```

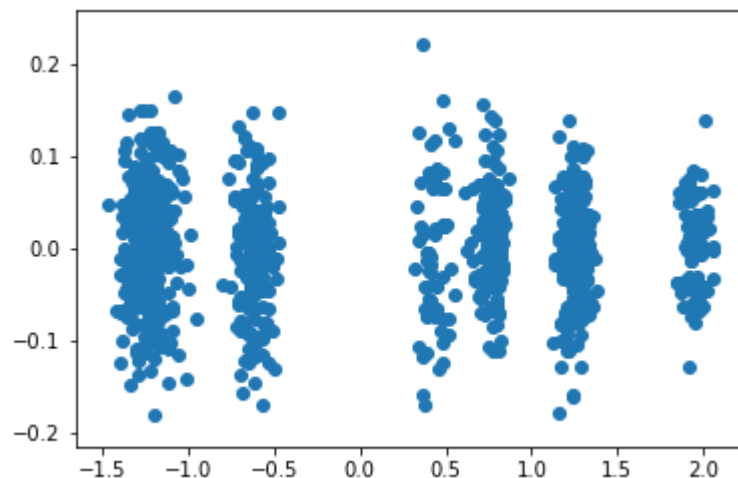
```
In [75]: 1 t1 = pd.DataFrame(Y1) # Create a Pandas DataFrame to plot scatter
          2 sns.pairplot(t1)
          3 plt.show()
```



Compute PCA as it is mentioned in the paper

```
In [76]: 1 pca = PCA(n_components=2)
          2 Y_r1 = pca.fit(Y1).transform(Y1)
```

```
In [77]: 1 plt.figure()
2 colors = ['navy', 'turquoise', 'darkorange']
3 lw = 2
4 target_names = list(range(V))
5
6 plt.scatter(Y_r1[:, 0], Y_r1[:,1])
7 plt.show()
```

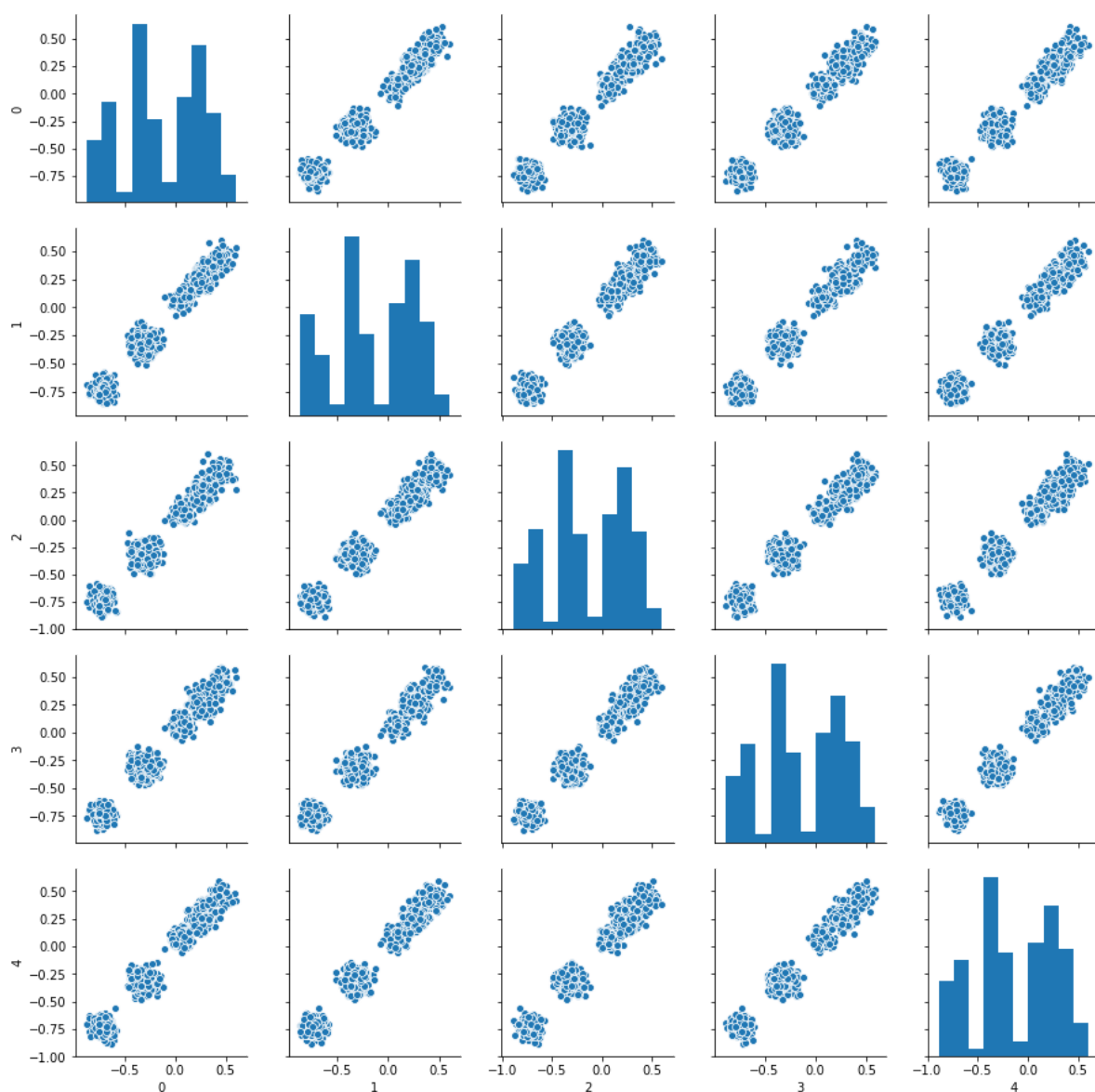


The Problem is' in this scatter plot shows just 5 clusters instead of 7 clusters.

Now let us repeat the procedure with $a=0.75$

```
In [78]: 1 Y2 = np.zeros([N, V])
2 interval = 0
3 for k in range(K):
4     for i in range(interval, structures[k]+interval):
5         Y2[i, :] = np.random.normal(loc=mu_2[k], scale=cov[k], size=
6         interval += structures[k])
```

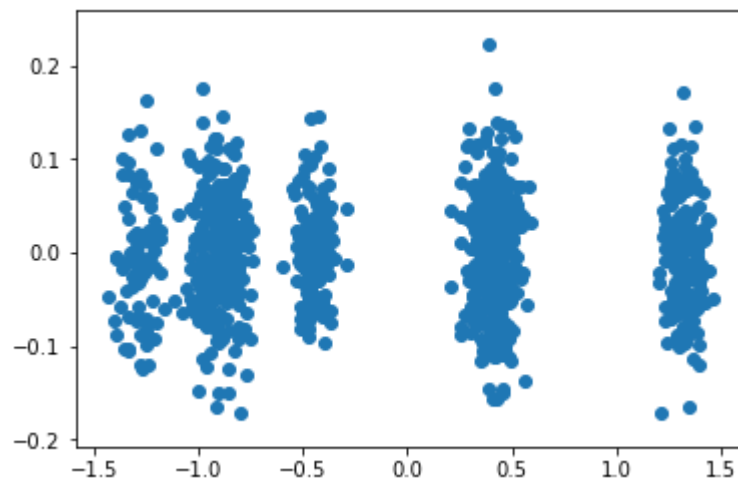
```
In [79]: 1 t2 = pd.DataFrame(Y2) # Create a Pandas DataFrame to plot scatter
          2 sns.pairplot(t2)
          3 plt.show()
```



Compute PCA as it is used in the paper

```
In [80]: 1 pca = PCA(n_components=2)
          2 Y_r2 = pca.fit(Y2).transform(Y2)
```

```
In [81]: 1 plt.figure()
2 colors = ['navy', 'turquoise', 'darkorange']
3 lw = 2
4 target_names = list(range(V))
5
6 plt.scatter(Y_r2[:, 0], Y_r2[:,1])
7 plt.show()
```

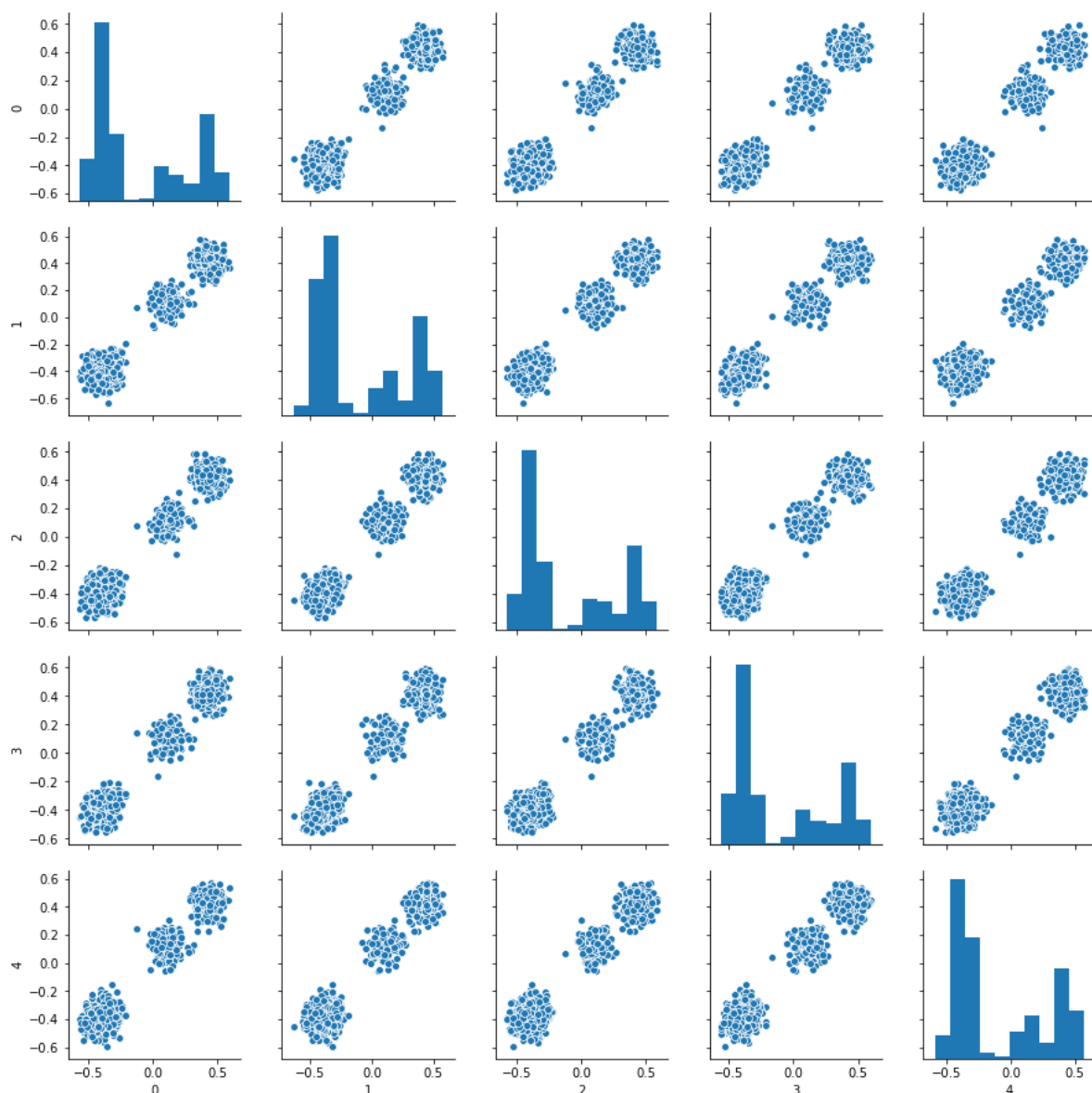


Now the Five clusters are intermixed as we expected!

Now let us repeat the procedure with $a=0.5$

```
In [82]: 1 Y3 = np.zeros([N, V])
2 interval = 0
3 for k in range(K):
4     for i in range(interval, structures[k]+interval):
5         Y3[i, :] = np.random.normal(loc=mu_3[k], scale=cov[k], size=V)
6         interval += structures[k]
```

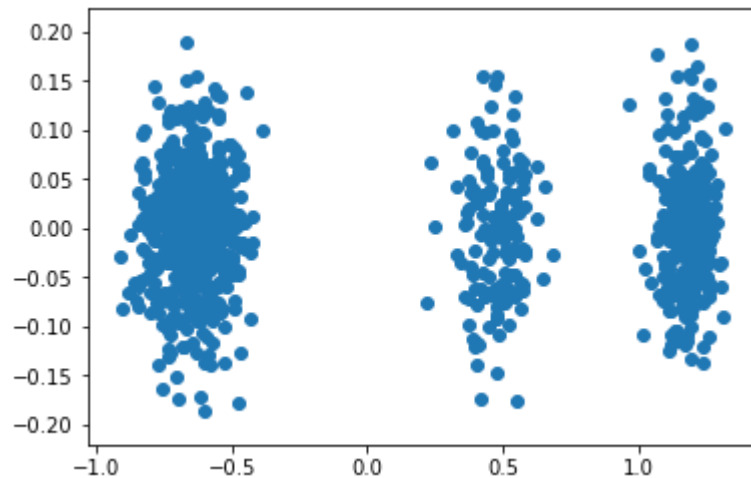
```
In [83]: 1 t3 = pd.DataFrame(Y3) # Create a Pandas DataFrame to plot scatter
          2 sns.pairplot(t3)
          3 plt.show()
```



Compute PCA as it is used in the paper

```
In [84]: 1 pca = PCA(n_components=2)
          2 Y_r3 = pca.fit(Y3).transform(Y3)
```

```
In [85]: 1 plt.figure()
2 colors = ['navy', 'turquoise', 'darkorange']
3 lw = 2
4 target_names = list(range(V))
5
6 plt.scatter(Y_r3[:, 0], Y_r3[:,1])
7 plt.show()
```



Now the clusters are more intermixed as it is expected!

```
In [ ]: 1
```