

---

**Task 3: Principal Component Analysis**

---

---

3a)

---

Normalizing can convert the original data into pure, dimensionless values and remove the unit restriction on data, then it is easy to compare and weight indicator among data with different units or scales. After Normalizing the data is easy to handle. One Advantage is that the convergence speed of the model improves. For example, if there is only two features, and the range of component 1 is 1000 and the range of component 2 is 5. The speed of iteration is very low during the optimization. But after normalizing the speed will be faster. Another advantage is that the accuracy of the model improves. In the last example, when we need to calculate about the distance, the second feature make much less contribute to the result than the first one.

Code:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 # get data
4 iris = np.loadtxt("./dataSets/iris.txt", delimiter=',')
5 iris_set = iris[:, 0:4]
6 iris_label = iris[:, 4]
7 # Normalizing
8 iris_norm = (iris_set - np.mean(iris_set, axis=0)) / np.std(iris_set, axis=0)
9 print('mean:', np.mean(iris_norm, axis=0))
10 print('var:', np.var(iris_norm, axis=0))
```

---

3b)

---

We need only two components to in order to explain at least 95% of the dataset variance. From the picture below, we can know that the cumulative variance of first two components is 95.8010%.

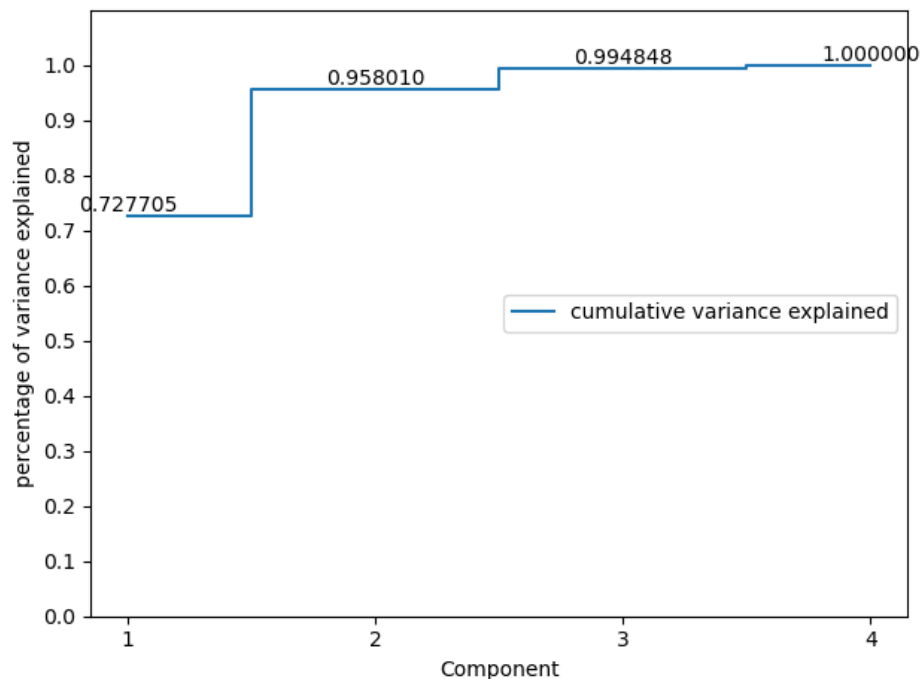


Figure 1: cumulative variance explained

Code:

```

1 #PCA
2 cov_matrix = np.cov(iris_norm, rowvar=0)
3 eigenvalues, eigenvectors = np.linalg.eig(cov_matrix)
4 eigenvalues_sorted = np.sort(eigenvalues)[-1::-1]
5 total = np.sum(eigenvalues_sorted)
6 var_explained = np.zeros(4)
7 for i, item in zip(range(4), eigenvalues_sorted):
8     var_explained[i] = item/total
9 cum_var_explained = np.cumsum(var_explained)
10 x = ['1', '2', '3', '4']
11 plt.step(x, cum_var_explained, where='mid', label='cumulative variance explained')
12 for a, b in zip(x, cum_var_explained):
13     plt.text(a, b, '%f' % b, ha='center', va='bottom')
14 plt.ylim((0, 1.1))
15 plt.yticks(np.arange(0, 1.1, 0.1))
16 plt.xlabel('Component')
17 plt.ylabel('percentage of variance explained')
18 plt.legend(loc='center right')
19 plt.show()

```

3c)

From the picture below, we can know that the data can be clearly distinguished by using 2 components, specially, through component1 the data is distinguished. Two components can explain 95% of the dataset variance. After PCA we spared two feature, i.e., 50% data, but the data can be still good distinguished and only 5% dataset variance is lost.

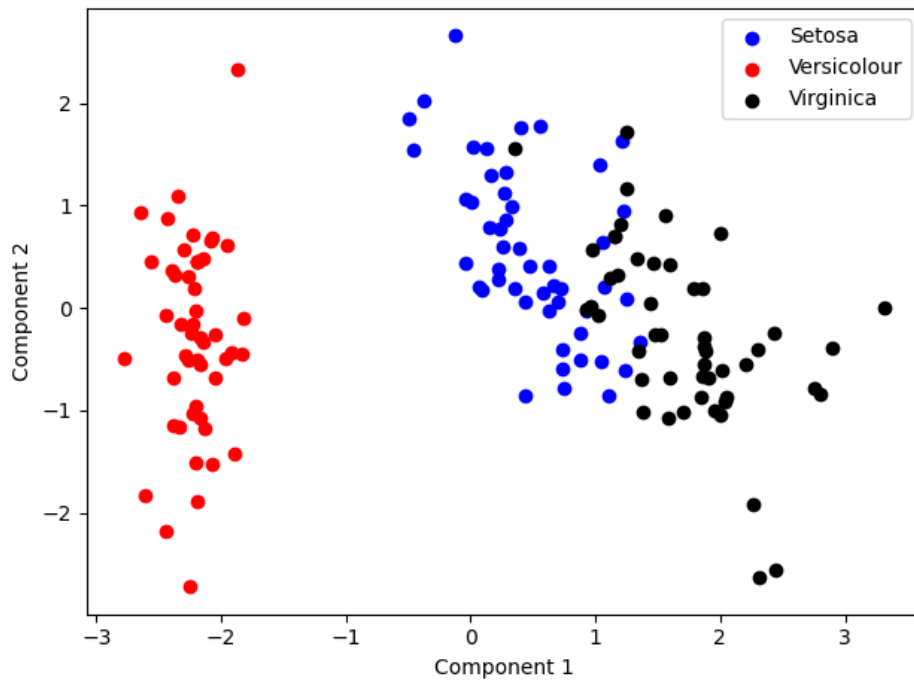


Figure 2: lower dimensional projection

Code:

```

1 vector_n = eigenvectors[:,0:2]
2 iris_lowDim = iris_norm@vector_n
3 name=['Setosa', 'Versicolour', 'Virginica']
4 print(iris_lowDim.shape)
5 for lab, col in zip(range(3), ('blue', 'red', 'black')):
6     plt.scatter(iris_lowDim[iris_label==lab, 0], iris_lowDim[iris_label==lab, 1],
7                 label=name[lab], color=col)
8 plt.xlabel('Component 1')
9 plt.ylabel('Component 2')
10 plt.legend()
11 plt.show()

```

3d)

N. of components	x1	x2	x3	x4
1	0.22925036	0.18006108	0.29805579	0.31692192
2	0.11261513	0.19281695	0.13314032	0.12717403
3	0.07919279	0.23615959	0.13302896	0.1257637
4	0.04925365	0.23658978	0.1322316	0.10227491

Code:

```

1 # reconstruct
2 def nrmse(eigvec, n):
3     vec_n = eigvec[:, 0:n]
4     lowDim = iris_norm @ vec_n
5     std = np.std(iris_set, axis=0)

```

```

6   iris_reconstruct = lowDim@vec_n.T+np.mean(iris_set, axis=0)
7   k = np.max(iris_set, axis=0)-np.min(iris_set, axis=0)
8   return np.sqrt(np.mean(np.square(iris_reconstruct-iris_set), axis=0))/k
9   for n in range(4):
10  print(nrmse(eigenvectors, n))

```

3e)

1. Explain the difference between PCA and ZCA whitening.

PCA whitening: It ensures the variance of each dimension of the data is 1.

ZCA whitening: It ensures that the variance of each dimension of the data is the same. It does a rotation operation on the basis of PCA whitening to make the data after whitening closer to the original data

2. State the equation(s) to compute the ZCA whitening parameters, given the data.

There are  $m$  samples, and the feature dimension of each sample is  $n$ . We have a sample set feature matrix  $X$  of  $n$  rows and  $m$  columns. And then we calculate Zero-average each row of  $X$ . The covariance matrix of  $X$ :

$$\Sigma = \frac{1}{m} X X^T \quad (1)$$

Calculate the eigenvalues and corresponding eigenvectors of the covariance matrix:

$$\frac{1}{m} X X^T = U \Lambda U^T \quad (2)$$

Rotate data

$$X_{rot} = U^T X \quad (3)$$

Then we get the PCA whitening

$$X_{PCAwhite,i} = \frac{X_{rot,i}}{\sqrt{\lambda_i + e}}, e = 1.0e^{-5} \quad (4)$$

ZCA whitening:

$$X_{ZCAwhite,i} = U X_{PCAwhite,i} \quad (5)$$

3. State the equation(s) to whiten a (new) data example  $x$ , given the ZCA parameters.

4. Compute and report the ZCA whitening parameters for the unnormalized IRIS data (including numerical values!).

```

1   import numpy as np
2   # get data
3   iris = np.loadtxt("./dataSets/iris.txt", delimiter=',')
4   iris_set = iris[:, 0:4]
5   iris_label = iris[:, 4]
6
7
8   def ZCA_whitening(x):
9       avg = np.mean(x, axis=0)
10      x = x - avg
11      sigma = np.cov(x, rowvar=0)
12      u,s,v = np.linalg.svd(sigma)
13      xRot = x@u
14      xPCAwhite = xRot / np.sqrt(s + 1e-5)
15      xZCAwhite = (u@xPCAwhite.T).T

```

# Übung TUDaExercise

LastName, FirstName: \_\_\_\_\_

EnrollmentID: \_\_\_\_\_

```
16     return xPCAwhite, xZCAwhite
17
18
19 xPCAwhite, xZCAwhite = ZCA_whitening(iris_set)
20 print(xZCAwhite)
```

result:

```
[[-1.01364052e-02  5.43753800e-01 -1.24393454e+00 -5.95425420e-01]
 [-7.66248148e-02 -7.85557354e-01 -1.43413894e+00  3.41837272e-01]
 [-7.03817057e-01 -7.68193655e-02 -1.28900583e+00  3.41850724e-01]
 [-1.13281235e+00 -1.11237957e-01 -7.86956953e-01 -7.29890731e-01]
 [-3.64283465e-01  9.41999744e-01 -1.83465916e+00 -6.4652558e-01]
 [-1.79831899e-01  1.62886291e+00 -1.88743144e+00  3.82451224e-01]
 [-1.25687616e+00  6.59057714e-01 -9.21713357e-01 -2.86753435e-01]
 [-2.97371385e-01  4.21581333e-01 -1.81524529e+00 -7.41285204e-01]
 [-1.38895543e+00 -6.16444559e-01 -9.89647358e-01 -4.9386826e-01]
 [-3.38880150e-01 -3.41188433e-01 -9.51559822e-01 -1.89861812e+00]
 [-3.38869574e-01  9.54480622e-01 -1.24353363e+00 -7.53319477e-01]
 [-9.79219846e-01  6.97654794e-01 -5.77312663e-01 -1.81812173e+00]
 [-9.93732945e-01 -6.37188677e-01 -1.10958384e+00 -8.79868860e-01]
 [-1.42577806e+00 -4.25251386e-01 -1.87821845e+00 -4.59727831e-01]
 [-1.74837383e+00  1.22635826e+00 -2.85189487e+00 -1.58652878e-01]
 [-7.98147356e-01  2.68898651e+00 -1.48623978e+00 -5.49986765e-02]
 [-6.68181090e-01  1.28824737e+00 -1.86861974e+00  5.86598471e-01]
 [-4.74644539e-02  4.89944919e-01 -1.44619696e+00 -7.39220812e-02]
 [-1.07568388e+00  1.89458628e+00 -1.33898740e+00 -6.18161728e-01]
 [-3.58592410e-01  1.48795725e+00 -9.9197777e-01 -4.37689173e-01]
 [-5.77173795e-01  2.17249947e-01 -1.11768802e+00  9.96413857e-01]
 [-2.26414572e-01  1.13046221e+00 -1.28119198e+00  9.77832557e-02]
 [-9.94173799e-01  9.72423342e-01 -1.31897570e+00  1.31649426e-02]
 [-5.56223648e-02  3.58144970e-02 -1.44655268e+00  3.89921584e-01]
 [-1.34681544e+00  9.58516444e-01 -2.76023444e-03 -1.62490380e+00]
 [-4.14197289e-02 -7.04289375e-01 -1.16978320e+00 -7.28230871e-01]
 [-1.34681544e+00  9.58516444e-01 -2.76023444e-03 -1.62490380e+00]
 [-3.45808776e-01  4.80917439e-01 -1.22640851e+00  1.95415846e-02]
 [-1.67628989e-01  5.36147984e-01 -1.17286406e+00 -7.26357996e-01]
 [-1.84676275e-01  1.45587877e-01 -1.45317792e+00 -4.64288482e-01]
 [-1.06987955e-01  1.84842287e-01 -6.28932934e-01 -9.47831995e-01]
 [-6.95631340e-01 -2.13483650e-01 -8.38174316e-01 -8.56695857e-01]
 [-8.96485931e-01 -6.42755998e-02 -1.98891413e+00  3.71114869e-01]
 [-3.47857776e-01  2.41287368e+00 -4.48882674e-01 -1.52471487e+00]
 [-4.67338211e-01  2.29131769e+00 -1.12448934e+00 -7.82775658e-01]
 [-3.38880150e-01 -3.41188433e-01 -9.51559822e-01 -1.89861812e+00]
 [-2.58618685e-01 -4.46652615e-01 -1.76922596e+00 -2.68615376e-02]
 [-1.29154962e+00  7.85687696e-02 -1.92644416e+00 -2.83853802e-01]
 [-3.38880150e-01 -3.41188433e-01 -9.51559822e-01 -1.89861812e+00]
 [-1.35323772e+00 -3.99712296e-01 -1.81685111e+00  3.45417891e-01]
 [-1.75913840e-02  3.27821544e-01 -1.13753279e+00 -7.83877155e-01]
```

(a) data1.

```
[[-1.10828130e-01  4.97588223e-01 -1.51726710e+00  9.18182940e-02]
 [-3.74211874e-01 -2.67388481e+00 -1.94927622e+00  5.58087621e-01]
 [-1.54235738e+00  2.87668808e-01 -8.42143338e-01 -4.53834870e-01]
 [-3.64986477e-01  5.96985888e-01 -1.54397547e+00  9.28739640e-01]
 [-8.10324352e-01  1.78196389e+00 -4.288897570e-01 -7.65147329e-02]
 [-3.19876847e-01 -7.44886455e-01 -1.51418389e+00  8.31382187e-02]
 [-5.18117956e-01  1.62872802e+00 -5.96359620e-01 -1.12137314e+00]
 [-1.10588464e+00  1.85494307e-01 -8.93368701e-01 -5.88639196e-01]
 [-2.58289493e-01  1.84896841e+00 -1.12124613e+00 -7.98647725e-01]
 [-8.05248991e-02  3.89413888e-02 -1.29555681e+00 -4.85135891e-01]
 [-2.82288648e+00  5.98388871e-02  1.25416142e-01 -5.81317547e-01]
 [-6.25389632e-01  3.99472961e-01  2.78165439e-01  8.87385483e-02]
 [-1.62961924e+00 -2.91965932e-02  3.45284596e-01 -4.87854015e-01]
 [-5.04892212e-01 -1.80981592e+00  2.59886465e-02  2.77361154e-01]
 [-1.16184077e+00 -8.22877535e-01 -6.57996265e-03  1.31841739e-01]
 [-1.02956848e+00 -4.57353322e-02  1.18288323e+00 -9.2832914e-01]
 [-4.37222151e-02  9.17817777e-01  6.63861662e-01  6.65843018e-02]
 [-1.53418415e+00 -1.38682156e+00  9.98573686e-02 -2.91836775e-02]
 [-1.27168497e+00 -5.06133397e-01  3.62687272e-01 -8.48464189e-01]
 [-1.56285606e+00 -4.52154731e-01  3.44968620e-01  6.33905410e-01]
 [-1.12865991e+00 -2.52141817e+00  1.46695916e-02 -1.81869280e-01]
 [-2.17689385e-01  4.83795575e-03  1.27622262e-01  6.88497355e-01]
 [-8.74583837e-01 -2.42487434e+00 -6.57894546e-02 -9.26700331e-01]
 [-2.12254891e-01 -1.89459515e-04  9.65141969e-01 -7.55843317e-01]
 [-3.03328871e-01 -4.38074353e-01 -3.48894877e-01  8.80877563e-01]
 [-1.64416851e+00 -2.21037547e-01 -1.74748147e-01 -3.27115317e-02]
 [-1.42381204e+00  5.48578974e-01  1.87455765e+00 -1.10268636e-01]
 [-2.78862766e-01 -6.38378138e-01  8.86992594e-01 -1.47266138e+00]
 [-1.01134675e+00 -2.44826894e+00 -3.54798470e-01  5.44717434e-01]
 [-3.66688380e-01 -1.27633951e+00  2.88684132e-01 -5.53675389e-01]
 [-1.02846972e+00  9.71786889e-01  8.54894554e-01 -7.31827551e-01]
 [-7.08989333e-01 -8.58743987e-01 -2.73854919e-01  2.32285080e-01]
 [-5.18297479e-01 -1.28395475e+00  5.57286278e-01 -3.88968726e-01]
 [-1.92351200e-01 -1.96257829e-01  1.28278893e+00 -1.66584151e+00]
 [-1.07889738e+00 -5.77875472e-01  2.71893692e-02 -3.16321051e-01]
 [-1.45894822e+00 -4.38163987e-01 -1.39414533e-01 -1.62386122e-02]
 [-1.71847797e+00 -8.78848243e-01  2.15532336e-01 -6.42998911e-01]
 [-1.11698746e+00 -1.64427858e-01  2.91662493e-01  2.52845432e-01]
 [-2.10131928e-01 -1.33546328e-01  4.98453778e-01  9.28530121e-02]
 [-2.78441928e-01 -1.36121980e+00 -3.17819588e-01 -2.42617914e-01]
```

(b) data2.

```
[[-4.29933103e-01 -1.57261976e+00  1.38648113e-01 -3.34933966e-01]
 [-3.44573654e-01 -1.48556475e+00  1.39562995e-01 -6.14177081e-01]
 [-4.11683268e-02 -9.11895684e-01  1.57564818e-02 -1.05133372e-01]
 [-7.17489289e-01 -2.72884289e-01  1.28243138e+00 -5.31588285e-01]
 [-1.98337220e+00  7.37698551e-01  1.31913265e+00 -1.84924168e-01]
 [-6.45682822e-01  1.33127552e+00  7.38962774e-01  3.05312114e-01]
 [-1.31463407e+00 -1.39847843e-02  2.03844327e-01 -1.57989264e-01]
 [-1.24419845e+00 -2.21847869e+00 -1.78948805e-01 -2.33856152e-01]
 [-1.08931815e+00  3.88381214e-01  7.85647978e-01 -2.64233446e-01]
 [-6.94811789e-01 -1.20246362e+00  1.99888415e-01  1.69743376e-01]
 [-3.31584962e+00 -4.97133846e-01  1.26245325e+00 -1.74610755e+00]
 [-1.84527182e-01  2.16542884e-01  8.58738228e-01 -6.07391782e-01]
 [-1.34486174e-02 -1.12862795e+00  1.22168229e-01 -2.53584987e-01]
 [-1.15976428e+00  1.78426749e+00 -1.89384821e-01  6.28332684e-02]
 [-8.47926291e-01 -5.15723344e-01  6.38143958e-01 -3.05867282e-01]
 [-8.89153614e-01  3.54584280e-01  9.78978529e-01 -9.18669368e-01]
 [-7.57265774e-01 -2.91883663e-03  6.89744221e-01 -3.753654932e-01]
 [-5.19363133e-01 -3.88758959e-01  2.71684365e-01 -3.98977112e-01]
 [-6.64913251e-01 -1.58612553e+00 -8.48897072e-01  1.08882834e+00]
 [-5.48418498e-01 -3.93558869e-01  4.89452784e-01 -1.19287618e-01]
 [-1.21886282e+00  1.56393827e+00  1.35716797e+00  1.77873864e+00]
 [-1.16498531e+00 -2.45111299e-01  9.28225185e-01  8.38264315e-01]
 [-1.28483259e+00  2.46831878e-02  7.33689499e-01  5.87827968e-01]
 [-6.08378013e-01  3.78848344e-01  1.65174397e+00 -5.75516878e-01]
 [-2.34222439e+00 -1.48561434e-01  1.47575644e+00 -7.22154756e-01]
 [-2.83481375e+00 -4.15551824e-01  1.09127987e+00  8.68487121e-01]
 [-1.33781832e+00 -4.11196688e-02  1.78237242e+00 -1.61885936e+00]
 [-6.48814471e-01 -1.84183561e+00  1.28149371e+00 -6.15489975e-01]
 [-9.01918521e-01  1.71891258e+00  7.18798876e-01  1.74295959e+00]
 [-3.58084370e-01  5.57592831e-01  2.96721617e-01  1.31282175e+00]
 [-2.69128181e-01 -6.38562262e-01  5.73215383e-01  6.57713760e-01]
 [-9.34642247e-01 -3.94529838e-02  3.27121462e-01  1.28488552e+00]
 [-1.07683826e+00 -8.98485889e-01  4.72986778e-01  1.59232801e+00]
 [-1.87294845e+00 -2.18465977e-01 -4.12313863e-03  3.19179512e+00]
 [-5.43665695e-02  6.64632928e-01  1.98943180e-01  2.31683474e+00]
 [-1.68214108e-02  4.85635858e-01  1.38876523e+00 -3.52489246e-01]
 [-1.2686275e+00  2.52863582e+00  2.84819622e+00 -8.36858786e-01]
```

(c) data3.

```
[ 2.19353327e+00 -9.95499872e-01  1.18128445e+00 -1.13365344e-01]
 [-1.59658897e-01 -1.82437994e+00  8.56564752e-01 -5.41786874e-01]
 [ 8.55383855e-01  5.39649514e-01  3.68936150e-01  1.69228202e+00]
 [-1.53728221e+00  1.99777685e-02  6.62778295e-01  1.59582572e+00]
 [ 2.13788454e+00 -4.88687877e-01  1.57517823e+00 -1.26897289e+00]
 [ 4.41162848e-01 -3.38889121e-01  1.24332775e-01  9.47925867e-01]
 [ 1.26687887e-01  1.14887382e+00  1.89698588e+00  6.88889953e-01]
 [ 1.14122136e+00  7.85876247e-01  1.58544867e+00 -1.21883288e+00]
 [ 1.89189676e-01 -5.26717869e-01  1.48216524e-01  1.05984735e+00]
 [-4.82877488e-01  2.62168980e-01  6.29764923e-01  7.11841102e-01]
 [-1.17645997e-01 -1.81632240e-01  8.35721318e-01  9.68138677e-01]
 [ 1.58025984e+00  1.12213961e-01  1.42934649e+00 -1.66898126e+00]
 [ 1.99486123e+00 -5.84842927e-01  9.84155351e-01 -6.40898034e-01]
 [ 2.11782931e+00  2.18627237e+00  1.62886918e+00 -1.11762441e+00]
 [-8.03179482e-02 -2.35441130e-01  6.33468894e-01  1.44163422e+00]
 [-9.95688170e-03 -1.98988537e-01  1.28478320e+00 -9.54916487e-01]
 [-1.82916380e+00 -1.28662946e-01  2.44449901e+00 -2.41476046e+00]
 [ 2.79359489e+00 -4.76385554e-01 -1.78418697e-02  1.28948249e+00]
 [-8.52888644e-01  1.57361777e+00  8.72958048e-01  2.84445991e+00]
 [-3.91022811e-01  8.83898987e-01  1.51808661e+00 -4.43546184e-01]
 [-5.59578864e-01  2.69774884e-01  5.58699288e-01  8.76773478e-01]
 [ 1.24215804e+00  8.27194922e-02  9.84382168e-02  1.38986518e+00]
 [ 5.49998028e-01  2.84328178e-01  1.22939285e-01  2.35519877e+00]
 [-1.68366836e+00 -2.85759939e-01 -8.86163530e-01  2.95965345e+00]
 [-1.16498531e+00 -2.45111299e-01  9.28225185e-01  8.38264315e-01]
 [ 3.31828779e-01  8.08117071e-01  8.68938913e-01  1.25843312e+00]
 [ 2.75928881e-01  9.24837461e-01  2.87944162e-01  2.52682141e+00]
 [ 1.89638876e+00 -3.13372625e-01 -5.35184786e-01  2.73654582e+00]
 [ 5.45322176e-01 -1.41223642e+00 -5.84777848e-02  1.33478476e+00]
 [ 4.24836458e-01  3.71736195e-02  3.16171462e-01  1.21737918e+00]
 [-9.25333798e-01  1.54887868e+00  8.18783585e-01  1.93614917e+00]
 [-1.28621264e+00  6.25196245e-01  1.26185968e+00  2.52664157e-01]
```

(d) data4.

Figure 3: ZCA

3f)

In general, PCA is suitable for linear dimensionality reduction of the data. Kernel Principal Component Analysis enables non-linear dimensionality reduction of data and is used to deal with linear distinguishable data set. The basic idea of KPCA is that for a input matrix  $X$ , by using a nonlinear mapping we map all the samples in  $X$  to a high or even infinite dimensional space to making it linearly distinguishable, and then process the data in this high-dimensional space using PCA.

These are the steps for the implementation of the kernel PCA algorithm:

1. Choose a kernel mapping  $k(\mathbf{x}_m, \mathbf{x}_n)$ .
  2. Obtain  $\mathbf{K}$  based on training data  $\{\mathbf{x}_n, (n = 1, \dots, N)\}$ .
  3. Solve eigenvalue problem of  $\mathbf{K}$  to get  $\lambda_i$  and  $\mathbf{a}_i$ .
  4. For each given data point  $\mathbf{x}$ , obtain its principal components in the feature space:  $(f(\mathbf{x}) \cdot \phi_i) = \sum_{n=1}^N a_n^{(i)} k(\mathbf{x}, \mathbf{x}_n)$
  5. Do whatever processing (e.g., feature selection, classification) in the feature space.
- (Reference: <http://fourier.eng.hmc.edu/e161/lectures/kernelPCA/node4.html>)

Limits: KPCA costs more time than PCA