

Unsupervised Learning

March 29, 2022

0.0.1 Principle of Maximum Variance: what is PCA supposed to do?

First of all, let us recall the principle/assumption of PCA:

1. What is the variance?
2. What is the covariance?
3. How do we compute the covariance matrix?
4. What is the meaning of the principle of maximum variance?
5. Why do we need this principle?
6. Does the principle always apply?

0.0.2 Answers:

What is the variance: variance is measures how spread our data set is, around its mean value. Taking the square root of the variance, we get standard deviation.

standard diviation : σ

variance : σ^2

What is the covariance?: covariance is a measure of the relationship between two random variables. **How do we compute the covariance matrix?:**

$$cov(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \quad (1)$$

Why do we need this principle? We need pricipel for reducing the dimensionality of a dataset.
Does the principle always apply? pricipel will not always work, if most of the correlation coefficients are smaller than 0.3, the pricipel component will not work.

0.1 Implementation: how is PCA implemented?

Here we implement the basic steps of PCA and we assemble them.

0.1.1 Importing libraries

We start importing the *numpy* library for performing matrix computations, the *pyplot* library for plotting data, and the *syntheticdata* module to import synthetic data.

```
[1]: import numpy as np
import matplotlib.pyplot as plt

import syntheticdata
```

0.1.2 Centering the Data

Implement a function with the following signature to center the data as explained in *Marsland*.

```
[2]: def center_data(A):
    m = np.mean(A, axis=0)
    X = A - m

    # INPUT:
    # A      [NxM] numpy data matrix (N samples, M features)
    #
    # OUTPUT:
    # X      [NxM] numpy centered data matrix (N samples, M features)

    return X
```

Test your function checking the following assertion on *testcase*:

```
[3]: testcase = np.array([[3., 11., 4.3], [4., 5., 4.3], [5., 17., 4.5], [4, 13., 4.4]])
answer = np.array([[-1., -0.5, -0.075], [0., -6.5, -0.075], [1., 5.5, 0.125], [0., 1.5, 0.
→025]])
np.testing.assert_array_almost_equal(center_data(testcase), answer)

#The test did work
```

0.1.3 Computing Covariance Matrix

Implement a function with the following signature to compute the covariance matrix as explained in *Marsland*.

```
[4]: def compute_covariance_matrix(A):
    C = np.cov(np.transpose(A))
    # INPUT:
    # A      [NxM] centered numpy data matrix (N samples, M features)
    #
    # OUTPUT:
    # C      [MxM] numpy covariance matrix (M features, M features)
    #
    # Do not apply centering here. We assume that A is centered before this_
    →function is called.
```

```
return C
```

Test your function checking the following assertion on *testcase*:

```
[5]: testcase = center_data(np.array([[22., 11., 5.5], [10., 5., 2.5], [34., 17., 8.5], [28.
    ↪, 14., 7]]))
answer = np.array([[580., 290., 145.], [290., 145., 72.5], [145., 72.5, 36.25]])

# Depending on implementation the scale can be different:
to_test = compute_covariance_matrix(testcase)

answer = answer/answer[0, 0]
to_test = to_test/to_test[0, 0]

np.testing.assert_array_almost_equal(to_test, answer)

#The test did work
```

0.1.4 Computing eigenvalues and eigenvectors

Use the linear algebra package of `numpy` and its function `np.linalg.eig()` to compute eigenvalues and eigenvectors. Notice that we take the real part of the eigenvectors and eigenvalues. The covariance matrix *should* be a symmetric matrix, but the actual implementation in `compute_covariance_matrix()` can lead to small round off errors that lead to tiny imaginary additions to the eigenvalues and eigenvectors. These are purely numerical artifacts that we can safely remove.

Note: If you decide to NOT use `np.linalg.eig()` you must make sure that the eigenvalues you compute are of unit length!

```
[6]: def compute_eigenvalue_eigenvectors(A):
    # INPUT:
    # A      [DxD] numpy matrix
    #
    # OUTPUT:
    # eigval  [D] numpy vector of eigenvalues
    # eigvec  [DxD] numpy array of eigenvectors
    eigval, eigvec = np.linalg.eig(A)

    # Numerical roundoff can lead to (tiny) imaginary parts. We correct that
    ↪here.
    eigval = eigval.real
    eigvec = eigvec.real

    return eigval, eigvec
```

Test your function checking the following assertion on *testcase*:

```
[7]: testcase = np.array([[2,0,0],[0,5,0],[0,0,3]])
      answer1 = np.array([2.,5.,3.])
      answer2 = np.array([[1.,0.,0.],[0.,1.,0.],[0.,0.,1.]])
      x,y = compute_eigenvalue_eigenvectors(testcase)
      np.testing.assert_array_almost_equal(x, answer1)
      np.testing.assert_array_almost_equal(y, answer2)

      #The test did work
```

0.1.5 Sorting eigenvalues and eigenvectors

Implement a function with the following signature to sort eigenvalues and eigenvectors as explained in *Marsland*.

Remember that eigenvalue $eigval[i]$ corresponds to eigenvector $eigvec[:,i]$.

```
[8]: def sort_eigenvalue_eigenvectors(eigval, eigvec):
      # INPUT:
      # eigval      [D] numpy vector of eigenvalues
      # eigvec      [DxD] numpy array of eigenvectors
      #
      # OUTPUT:
      # sorted_eigval      [D] numpy vector of eigenvalues
      # sorted_eigvec      [DxD] numpy array of eigenvectors

      indices = np.argsort(eigval)
      indices = indices[::-1]
      sorted_eigvec = eigvec[:,indices]
      sorted_eigval = eigval[indices]

      return sorted_eigval, sorted_eigvec
```

Test your function checking the following assertion on *testcase*:

```
[9]: testcase = np.array([[2,0,0],[0,5,0],[0,0,3]])
      answer1 = np.array([5.,3.,2.])
      answer2 = np.array([[0.,0.,1.],[1.,0.,0.],[0.,1.,0.]])
      x,y = compute_eigenvalue_eigenvectors(testcase)
      x,y = sort_eigenvalue_eigenvectors(x,y)
      np.testing.assert_array_almost_equal(x, answer1)
      np.testing.assert_array_almost_equal(y, answer2)

      #The test did work
```

0.1.6 PCA Algorithm

Implement a function with the following signature to compute PCA as explained in *Marsland* using the functions implemented above.

```
[10]: def pca(A,m):
    # INPUT:
    # A      [NxM] numpy data matrix (N samples, M features)
    # m      integer number denoting the number of learned features (m <= M)
    #
    # OUTPUT:
    # pca_eigvec [Mxm] numpy matrix containing the eigenvectors (M
    → dimensions, m eigenvectors)
    # P         [Nxm] numpy PCA data matrix (N samples, m features)

    K = center_data(A)
    L = compute_covariance_matrix(K)
    eigval, eigvec = compute_eigenvalue_eigenvectors(L)
    eigval, eigvec = sort_eigenvalue_eigenvectors(eigval, eigvec)

    if m > 0:
        eigvec = eigvec[:, :m]

    P = np.dot(np.transpose(eigvec), np.transpose(K))
    pca_eigvec = eigvec

    return pca_eigvec, P.T
```

Test your function checking the following assertion on *testcase*:

```
[11]: testcase = np.array([[22., 11., 5.5], [10., 5., 2.5], [34., 17., 8.5]])
x,y = pca(testcase,2)

import pickle
answer1_file = open('PCAanswer1.pkl', 'rb'); answer2_file = open('PCAanswer2.
→ pkl', 'rb')
answer1 = pickle.load(answer1_file); answer2 = pickle.load(answer2_file)

test_arr_x = np.sum(np.abs(np.abs(x) - np.abs(answer1)), axis=0)
np.testing.assert_array_almost_equal(test_arr_x, np.zeros(2))

test_arr_y = np.sum(np.abs(np.abs(y) - np.abs(answer2)))
np.testing.assert_almost_equal(test_arr_y, 0)

#The test did work
```

0.2 Understanding: how does PCA work?

We now use the PCA algorithm you implemented on a toy data set in order to understand its inner workings.

0.2.1 Loading the data

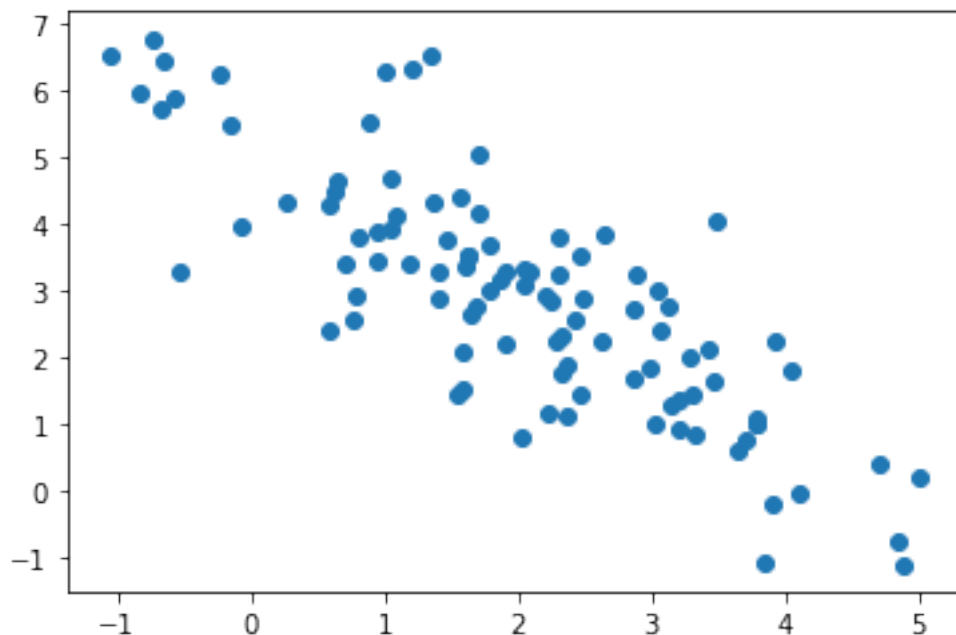
The module *syntheticdata* provides a small synthetic dataset of dimension [100x2] (100 samples, 2 features).

```
[12]: X = syntheticdata.get_synthetic_data1()
```

0.2.2 Visualizing the data

Visualize the synthetic data using the function *scatter()* from the *matplotlib* library.

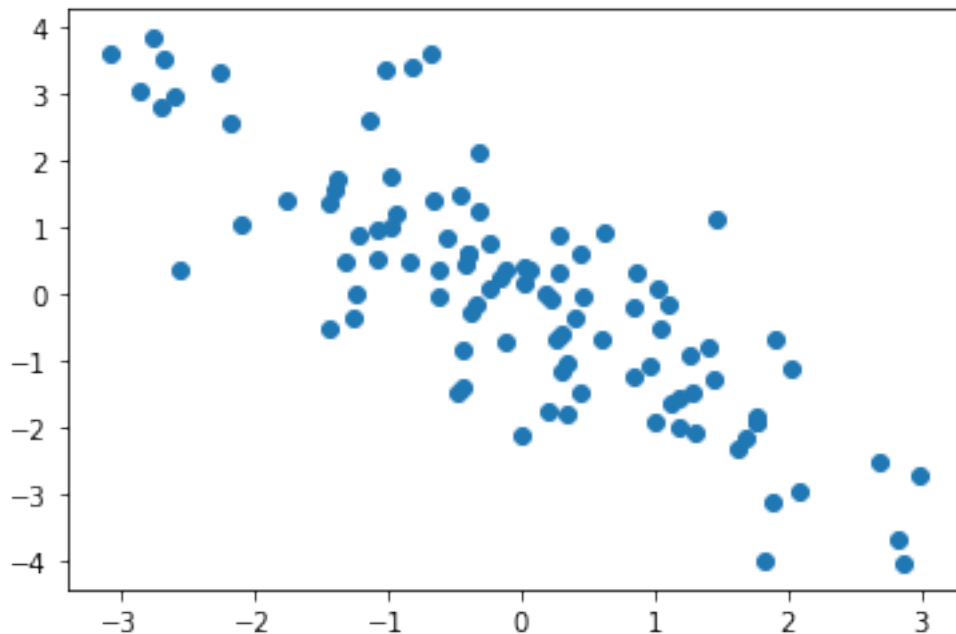
```
[13]: #Visualizing the data  
  
plt.scatter(X[:,0],X[:,1])  
plt.show()
```



0.2.3 Visualize the centered data

Notice that the data visualized above is not centered on the origin (0,0). Use the function defined above to center the data, and the replot it.

```
[14]: #Centering the data  
  
X = center_data(X)  
plt.scatter(X[:,0],X[:,1])  
plt.show()
```

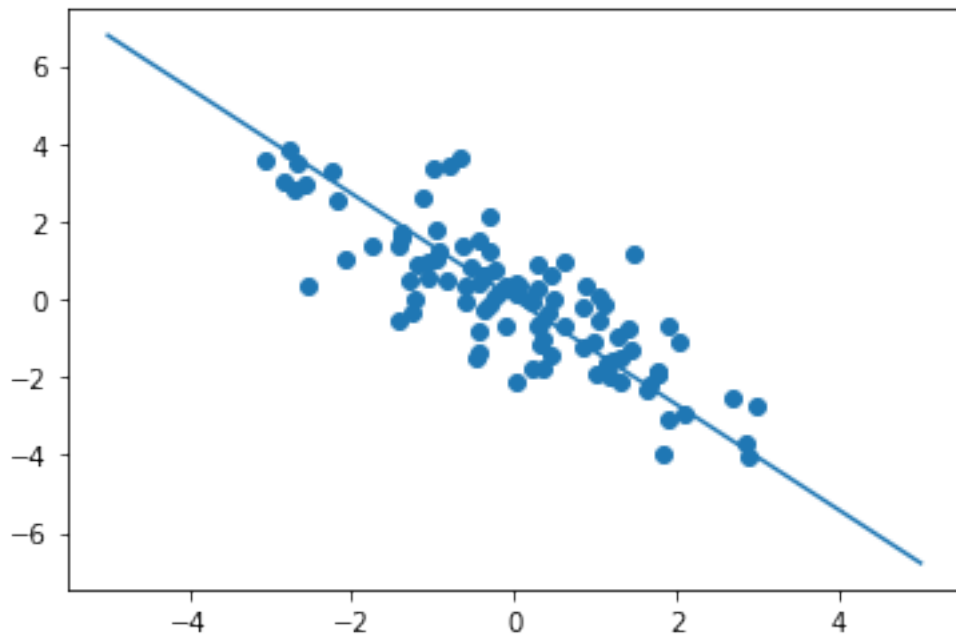


0.2.4 Visualize the first eigenvector

Visualize the vector defined by the first eigenvector. To do this you need: - Use the *PCA()* function to recover the eigenvectors - Plot the centered data as done above - The first eigenvector is a 2D vector (x0,y0). This defines a vector with origin in (0,0) and head in (x0,y0). Use the function *plot()* from matplotlib to plot a line over the first eigenvector.

```
[15]: pca_eigvec, _ = pca(X, 2)  
first_eigvec = pca_eigvec[0] #defining the vector with origin in origo  
  
plt.scatter(X[:,0],X[:,1]) #plotting the data  
  
x = np.linspace(-5, 5, 1000)  
y = first_eigvec[1]/first_eigvec[0] * x
```

```
plt.plot(x,y)
plt.show()
```

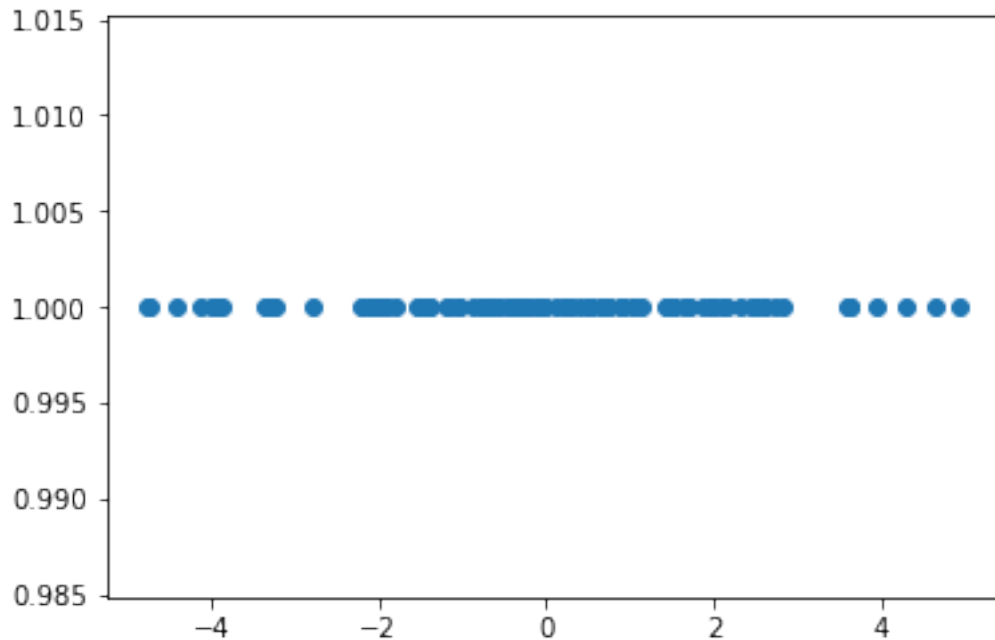


0.2.5 Visualize the PCA projection

Finally, use the `PCA()` algorithm to project on a single dimension and visualize the result using again the `scatter()` function.

```
[16]: #Visualizing the PCA projection (on a singel dimention)

_,P = pca(X, 1)
plt.scatter(P[:,0], np.ones(P.shape[0]))
plt.show()
```

0.3 Evaluation: when are the results of PCA sensible?

So far we have used PCA on synthetic data. Let us now imagine we are using PCA as a pre-processing step before a classification task. This is a common setup with high-dimensional data. We explore when the use of PCA is sensible.

0.3.1 Loading the first set of labels

The function `get_synthetic_data_with_labels1()` from the module `syntethicdata` provides a first labeled dataset.

```
[17]: X,y = syntheticdata.get_synthetic_data_with_labels1()
```

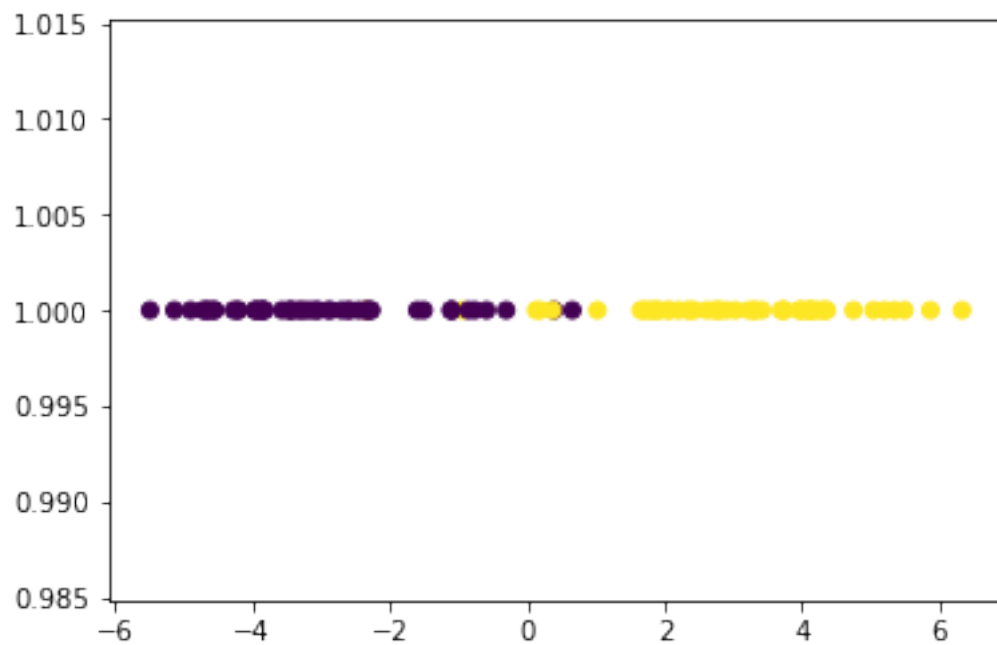
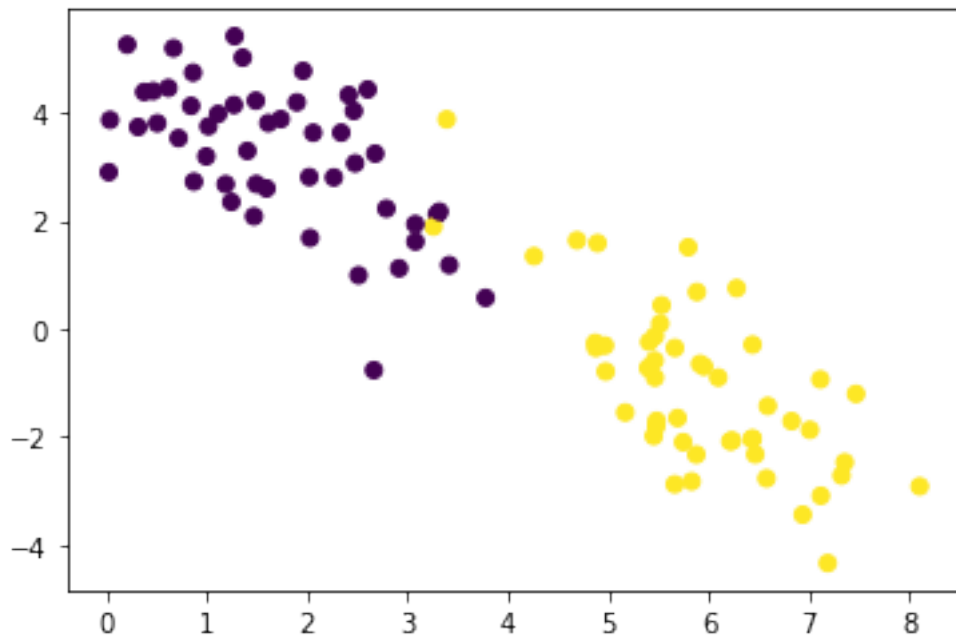
0.3.2 Running PCA

Process the data using the PCA algorithm and project it in one dimension. Plot the labeled data using `scatter()` before and after running PCA. Comment on the results.

```
[18]: plt.scatter(X[:,0],X[:,1],c=y[:,0]) #plotting the data

plt.figure()
_,P = pca(X, 1)
```

```
plt.scatter(P[:,0],np.ones(P.shape[0]),c=y[:,0]) #plotting the data in one
→dimention after running PCA
plt.show()
```



Comment: We see here that the data is “compressed” into one dimension. The two different data sets are separated.

0.3.3 Loading the second set of labels

The function `get_synthetic_data_with_labels2()` from the module `syntheticdata` provides a second labeled dataset.

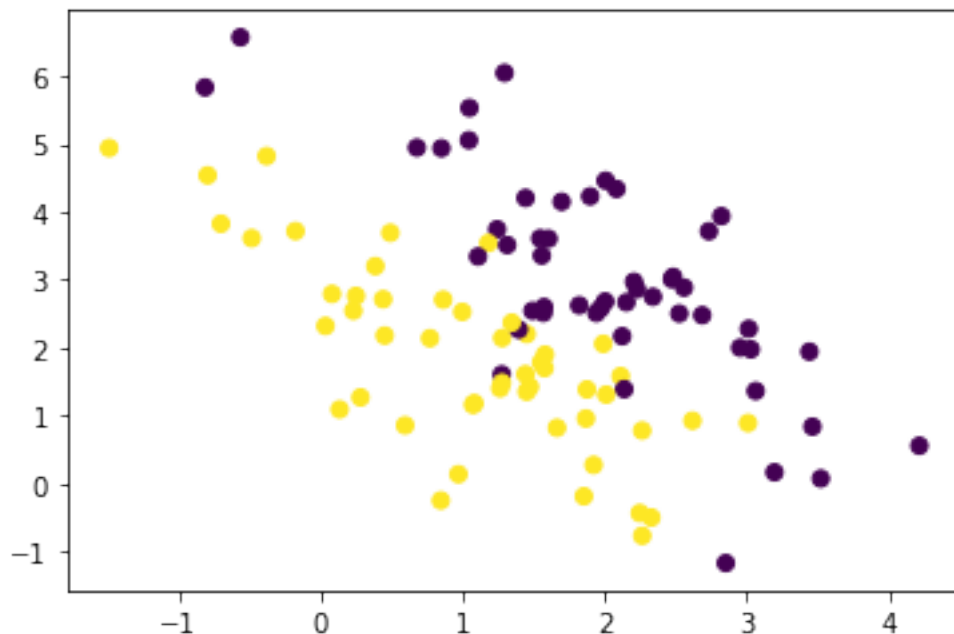
```
[19]: X,y = syntheticdata.get_synthetic_data_with_labels2()
```

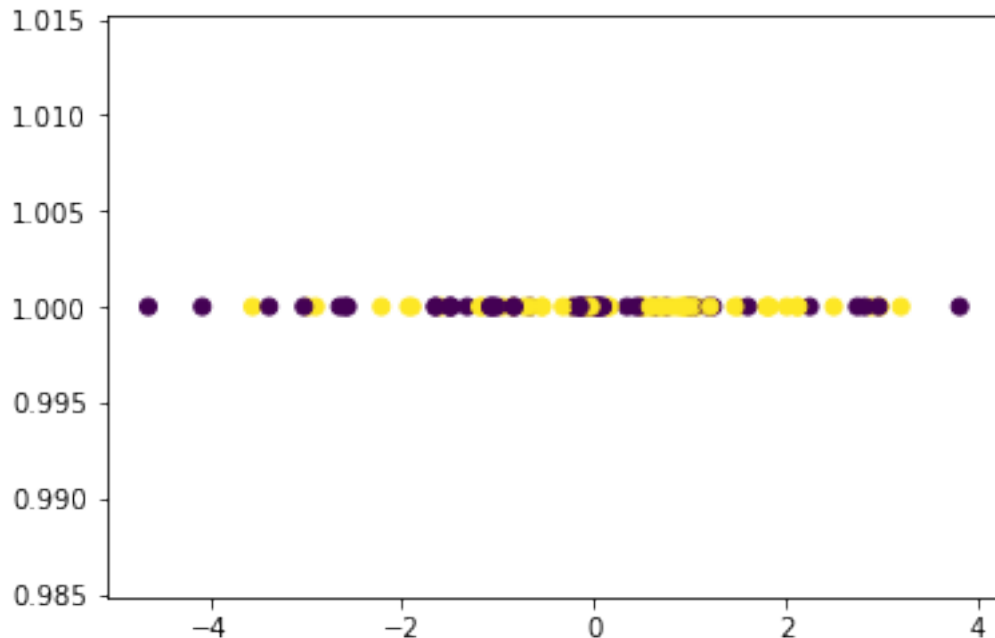
0.3.4 Running PCA

As before, process the data using the PCA algorithm and project it in one dimension. Plot the labeled data using `scatter()` before and after running PCA. Comment on the results.

```
[20]: plt.scatter(X[:,0],X[:,1],c=y[:,0]) #plotting the data

plt.figure()
_,P = pca(X, 1)
plt.scatter(P[:,0],np.ones(P.shape[0]),c=y[:,0]) #plotting the data in one_
→dimension after running PCA
plt.show()
```





Comment: Again, the data is “compressed” in the second plot, but this time the datasets are not that well separated and looks a bit disorganized.

How would the result change if you were to consider the second eigenvector? Or if you were to consider both eigenvectors?

Answer: The dataset would be more organized if we used the second eigenvector.

Using both eigenvectors, we get higher variance

0.4 Case study 1: PCA for visualization

We now consider the *iris* dataset, a simple collection of data ($N=150$) describing iris flowers with four ($M=4$) features. The features are: Sepal Length, Sepal Width, Petal Length and Petal Width. Each sample has a label, identifying each flower as one of 3 possible types of iris: Setosa, Versicolour, and Virginica.

Visualizing a 4-dimensional dataset is impossible; therefore we will use PCA to project our data in 2 dimensions and visualize it.

0.4.1 Loading the data

The function `get_iris_data()` from the module `syntheticdata` returns the *iris* dataset. It returns a data matrix of dimension $[150 \times 4]$ and a label vector of dimension $[150]$.

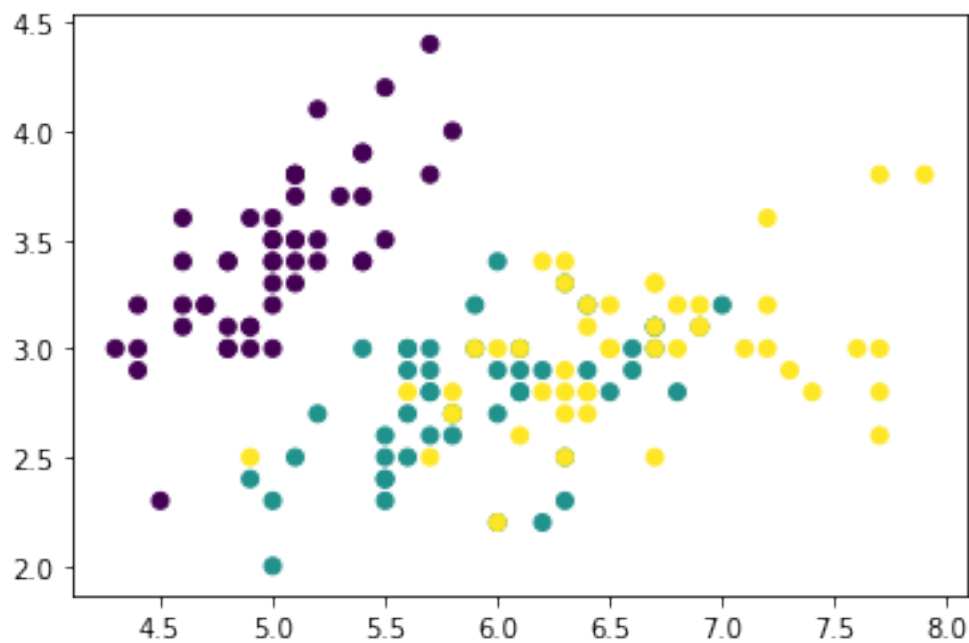
```
[21]: X,y = syntheticdata.get_iris_data()
```

0.4.2 Visualizing the data by selecting features

Try to visualize the data (using label information) by randomly selecting two out of the four features of the data. You may try different pairs of features.

```
[22]: import random
r = random.sample(range(0,3), 2) #making a variable which makes two random
    ↪ indices

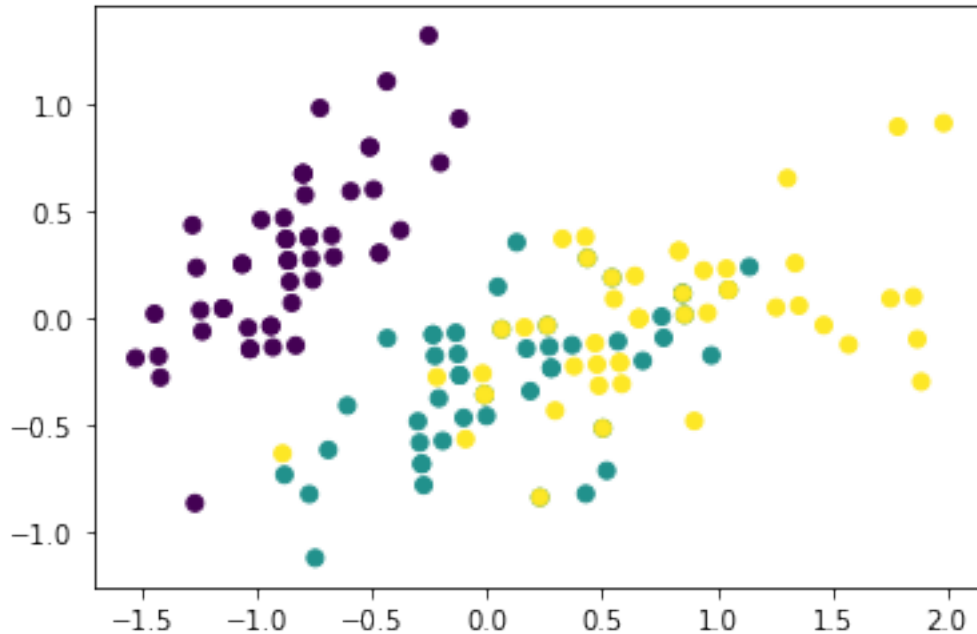
X = X[:,[r[0], r[1]]]
plt.scatter(X[:,0],X[:,1],c=y) #plotting
plt.show()
```



0.4.3 Visualizing the data by PCA

Process the data using PCA and visualize it (using label information). Compare with the previous visualization and comment on the results.

```
[23]: _,P = pca(X, 2)
plt.scatter(P[:,0],P[:,1],c=y) #plotting after running PCA
plt.show()
```



Comment: Enter your comment here.

0.5 Case study 2: PCA for compression

We now consider the *faces in the wild (lfw)* dataset, a collection of pictures ($N=1280$) of people. Each pixel in the image is a feature ($M=2914$).

0.5.1 Loading the data

The function `get_lfw_data()` from the module `syntheticdata` returns the *lfw* dataset. It returns a data matrix of dimension $[1280 \times 2914]$ and a label vector of dimension $[1280]$. It also returns two parameters, h and w , reporting the height and the width of the images (these parameters are necessary to plot the data samples as images). Beware, it might take some time to download the data. Be patient :)

```
[24]: X,y,h,w = syntheticdata.get_lfw_data()
```

0.5.2 Inspecting the data

Choose one datapoint to visualize (first coordinate of the matrix X) and use the function `imshow()` to plot and inspect some of the pictures.

Notice that `imshow` receives as a first argument an image to be plot; the image must be provided as a rectangular matrix, therefore we reshape a sample from the matrix X to have height h and

width w . The parameter *cmap* specifies the color coding; in our case we will visualize the image in black-and-white with different gradations of grey.

```
[25]: plt.imshow(X[0,:].reshape((h, w)), cmap=plt.cm.gray)
      plt.show()
```



0.5.3 Implementing a compression-decompression function

Implement a function that first uses PCA to project samples in low-dimensions, and then reconstruct the original image.

Hint: Most of the code is the same as the previous `PCA()` function you implemented. You may want to refer to *Marsland* to check out how reconstruction is performed.

```
[26]: def encode_decode_pca(A,m):
      # INPUT:
      # A      [NxM] numpy data matrix (N samples, M features)
      # m      integer number denoting the number of learned features (m <= M)
      #
      # OUTPUT:
      # Ahat   [NxM] numpy PCA reconstructed data matrix (N samples, M features)

      K2 = center_data(A)
      L2 = compute_covariance_matrix(K2)
      eigval, eigvec = compute_eigenvalue_eigenvectors(L2)
      eigval, eigvec = sort_eigenvalue_eigenvectors(eigval, eigvec)
```

```

if m > 0:
    eigvec = eigvec[:, :m]

x = np.dot(np.transpose(eigvec), np.transpose(K2))
Ahat = np.transpose(np.dot(eigvec, x)) + np.mean(K2, axis=0)
return Ahat

```

0.5.4 Compressing and decompressing the data

Use the implemented function to encode and decode the data by projecting on a lower dimensional space of dimension 200 ($m=200$).

```

[27]: Xhat = encode_decode_pca(X, 200)
      print(Xhat)

```

```

[[ 1.44817193e+02  1.50202728e+02  1.54983446e+02 ... -5.04860528e+01
  -3.79366459e+01 -2.67548960e+01]
 [-3.05458848e+01 -3.71682028e+01 -4.37690377e+01 ...  5.65153447e+00
  -5.11652878e-02  2.01507418e+01]
 [-2.50659506e+00  2.72109516e+00  1.81170277e+01 ...  5.78473422e+01
   6.98044830e+01  8.77943082e+01]
 ...
 [-9.69113245e-01 -1.21564438e+01 -3.22364705e+01 ... -7.45073156e+01
  -6.71533978e+01 -4.98508130e+01]
 [-5.15908904e+01 -3.86735736e+01 -2.02154690e+01 ...  8.01005322e+01
   7.75474246e+01  5.80646470e+01]
 [-6.65815802e+01 -6.87320832e+01 -7.38834750e+01 ... -8.56664355e+01
  -6.48000869e+01 -5.54935392e+01]]

```

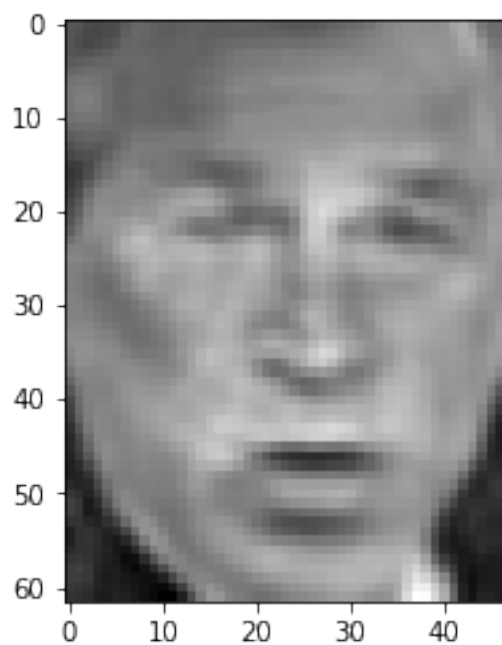
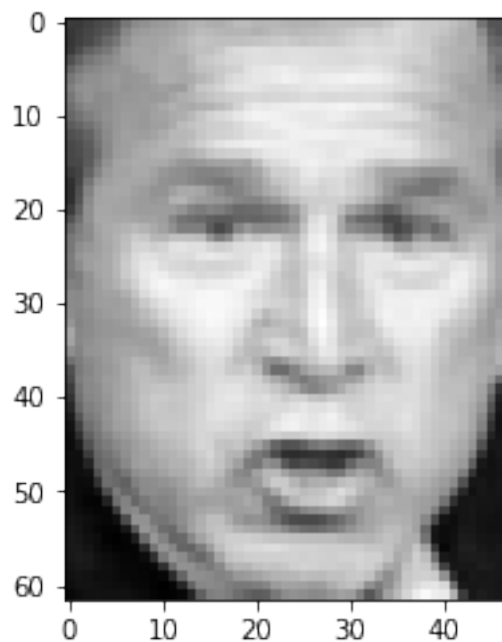
0.5.5 Inspecting the reconstructed data

Use the function *imshow* to plot and compare original and reconstructed pictures. Comment on the results.

```

[28]: plt.imshow(X[20,:].reshape((h, w)), cmap=plt.cm.gray)
      plt.figure()
      plt.imshow(Xhat[20,:].reshape((h, w)), cmap=plt.cm.gray)
      plt.show()

```

Comment: Enter your comment here.

0.5.6 Evaluating different compressions

Use the previous setup to generate compressed images using different values of low dimensions in the PCA algorithm (e.g.: 100, 200, 500, 1000). Plot and comment on the results.

```
[29]: plt.imshow(X[0,:].reshape((h, w)), cmap=plt.cm.gray)

plt.figure()
Xhat = encode_decode_pca(X,100)
plt.imshow(Xhat[0,:].reshape((h, w)), cmap=plt.cm.gray)

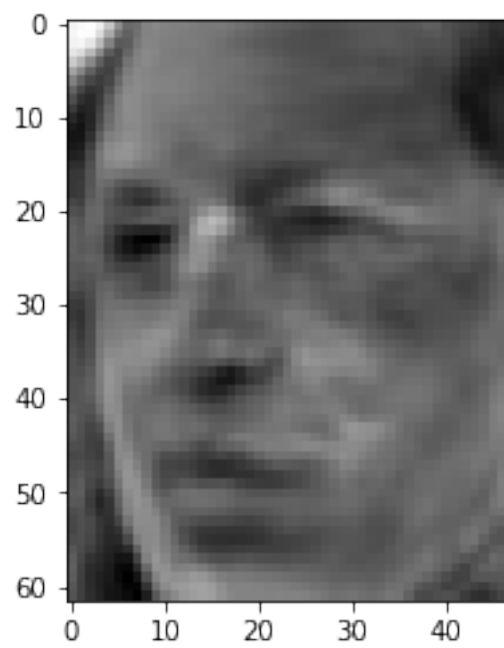
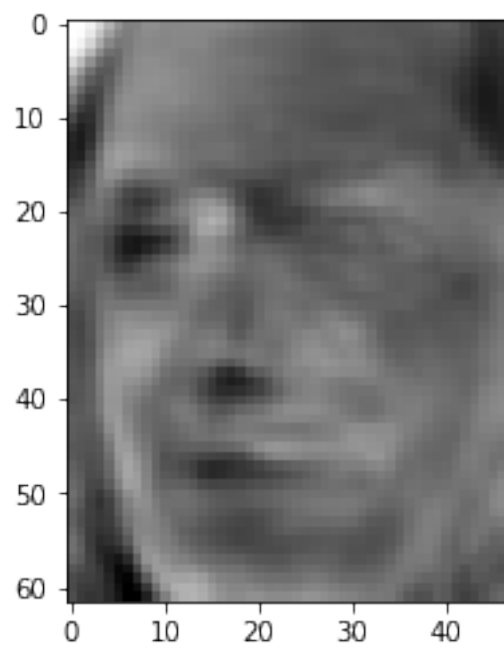
plt.figure()
Xhat = encode_decode_pca(X,200)
plt.imshow(Xhat[0,:].reshape((h, w)), cmap=plt.cm.gray)

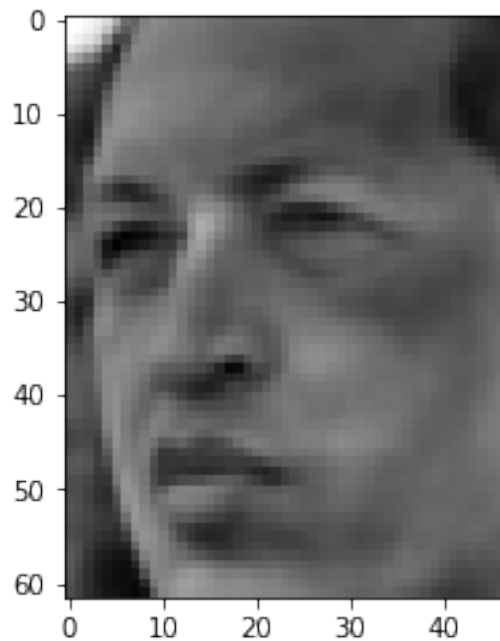
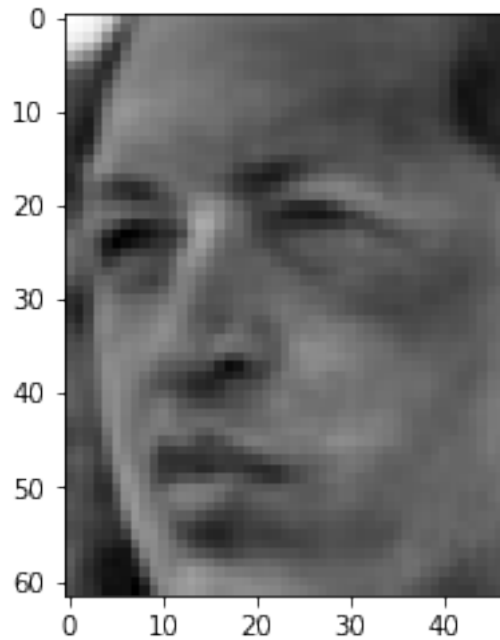
plt.figure()
Xhat = encode_decode_pca(X,500)
plt.imshow(Xhat[0,:].reshape((h, w)), cmap=plt.cm.gray)

plt.figure()
Xhat = encode_decode_pca(X,1000)
plt.imshow(Xhat[0,:].reshape((h, w)), cmap=plt.cm.gray)
plt.show()

#plotting with different m-value
```







Comment: We see that the plot before running the PCA is less blurry. After running the PCA the first image is very blurry, but it gets less and less blurry with higher m-value. The last picture with $m=1000$ is less blurry than the three pictures above, and looks like the first picture (before the PCA was run).

0.6 Master Students: PCA Tuning

If we use PCA for compression or decompression, it may be not trivial to decide how many dimensions to keep. In this section we review a principled way to decide how many dimensions to keep.

The number of dimensions to keep is the only *hyper-parameter* of PCA. A method designed to decide how many dimensions/eigenvectors is the *proportion of variance*:

$$\text{POV} = \frac{\sum_{i=1}^m \lambda_i}{\sum_{j=1}^M \lambda_j},$$

where λ are eigenvalues, M is the dimensionality of the original data, and m is the chosen lower dimensionality.

Using the *POV* formula we may select a number M of dimensions/eigenvalues so that the proportion of variance is, for instance, equal to 95%.

Implement a new PCA for encoding and decoding that receives in input not the number of dimensions for projection, but the amount of proportion of variance to be preserved.

```
[30]: #def encode_decode_pca_with_pov(A,p):  
      # INPUT:  
      # A      [NxM] numpy data matrix (N samples, M features)  
      # p      float number between 0 and 1 denoting the POV to be preserved  
      #  
      # OUTPUT:  
      # Ahat [NxM] numpy PCA reconstructed data matrix (N samples, M features)  
      # m     integer reporting the number of dimensions selected  
  
      #return Ahat,m
```

Import the *lfw* dataset using the *get_lfw_data()* in *syntheticdata*. Use the implemented function to encode and decode the data by projecting on a lower dimensional space such that $\text{POV}=0.9$. Use the function *imshow* to plot and compare original and reconstructed pictures. Comment on the results.

```
[31]: #X,y,h,w = syntheticdata.get_lfw_data()
```

```
[32]: #Xhat,m = encode_decode_pca_with_pov(X,None)
```

```
[33]: #plt.imshow(X[0,:].reshape((h, w)), cmap=plt.cm.gray)  
      #plt.figure()  
      #plt.imshow(Xhat[0,:].reshape((h, w)), cmap=plt.cm.gray)
```

Comment: Enter your comment here.

1 K-Means Clustering (Bachelor and master students)

In this section you will use the *k-means clustering* algorithm to perform unsupervised clustering. Then you will perform a qualitative assesment of the results.

1.0.1 Importing scikit-learn library

We start importing the module *cluster.KMeans* from the standard machine learning library *scikit-learn*.

```
[34]: from sklearn.cluster import KMeans
```

1.0.2 Loading the data

We will use once again the *iris* data set. The function *get_iris_data()* from the module *syntheticdata* returns the *iris* dataset. It returns a data matrix of dimension [150x4] and a label vector of dimension [150].

```
[35]: X,y = syntheticdata.get_iris_data()
```

1.0.3 Projecting the data using PCA

To allow for visualization, we project our data in two dimensions as we did previously. This step is not necessary, and we may want to try to use *k-means* later without the PCA pre-processing. However, we use PCA, as this will allow for an easy visualization.

```
[36]: _,P = pca(X, 2)
```

1.0.4 Running k-means

We will now consider the *iris* data set as an unlabeled set, and perform clustering to this unlabeled set. We can compare the results of the clustering to the labeled calsses.

Use the class *KMeans* to fit and predict the output of the *k-means* algorithm on the projected data. Run the algorithm using the following values of $k = \{2, 3, 4, 5\}$.

```
[37]: KM = KMeans(2)
      yhat2 = KM.fit_predict(P)

      KM = KMeans(3)
      yhat3 = KM.fit_predict(P)

      KM = KMeans(4)
      yhat4 = KM.fit_predict(P)
```

```
KM = KMeans(5)
yhat5 = KM.fit_predict(P)
```

1.0.5 Qualitative assessment

Plot the results of running the k-means algorithm, compare with the true labels, and comment.

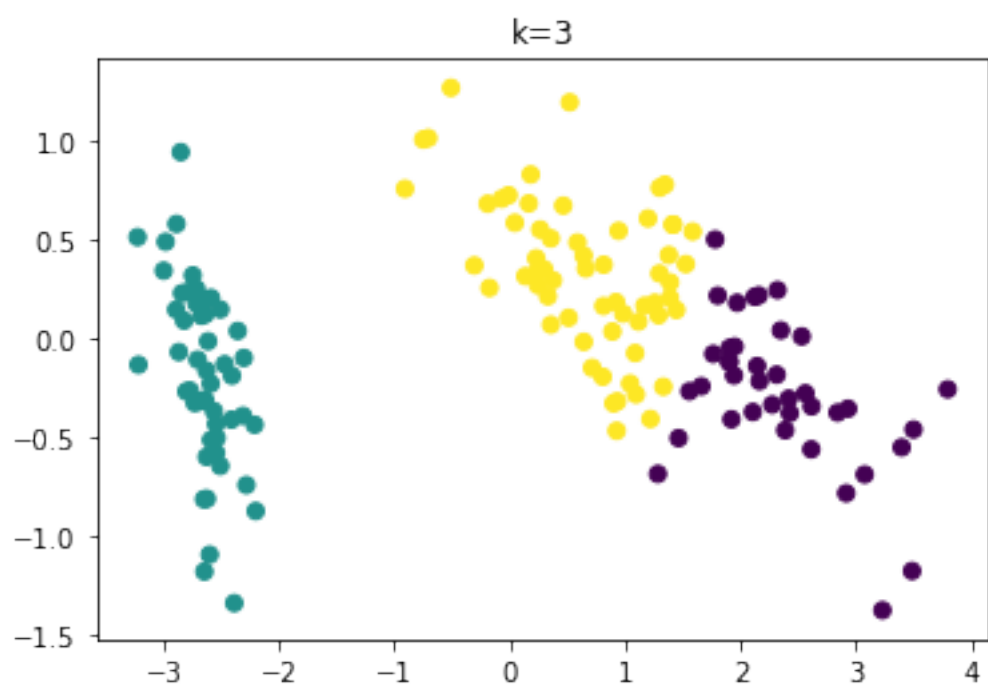
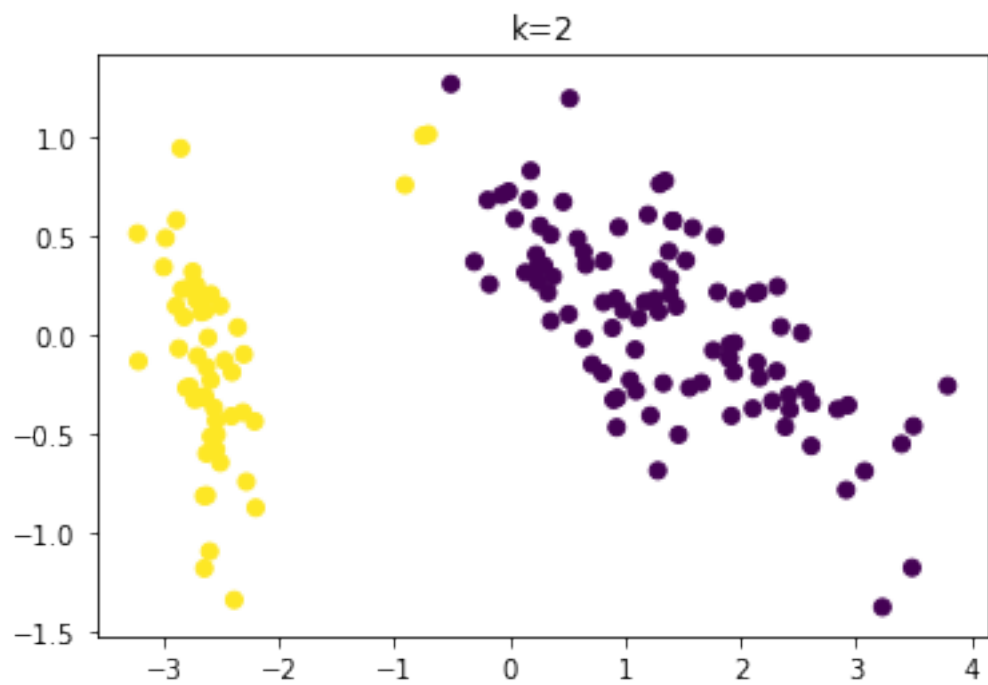
```
[38]: plt.scatter(P[:,0],P[:,1],c=yhat2)
plt.title('k=2')
plt.figure()

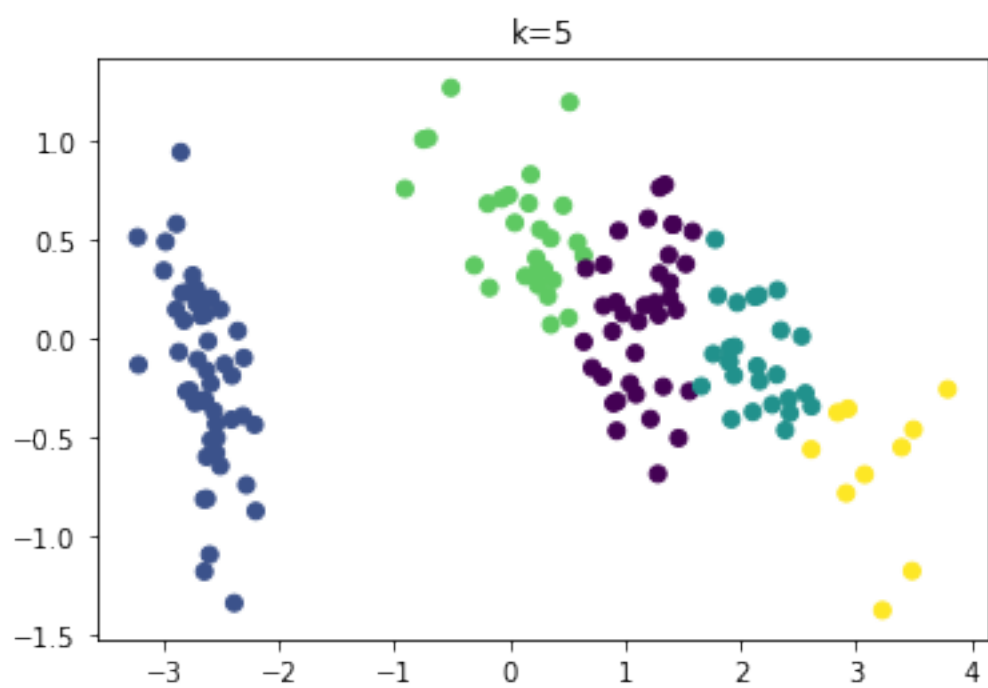
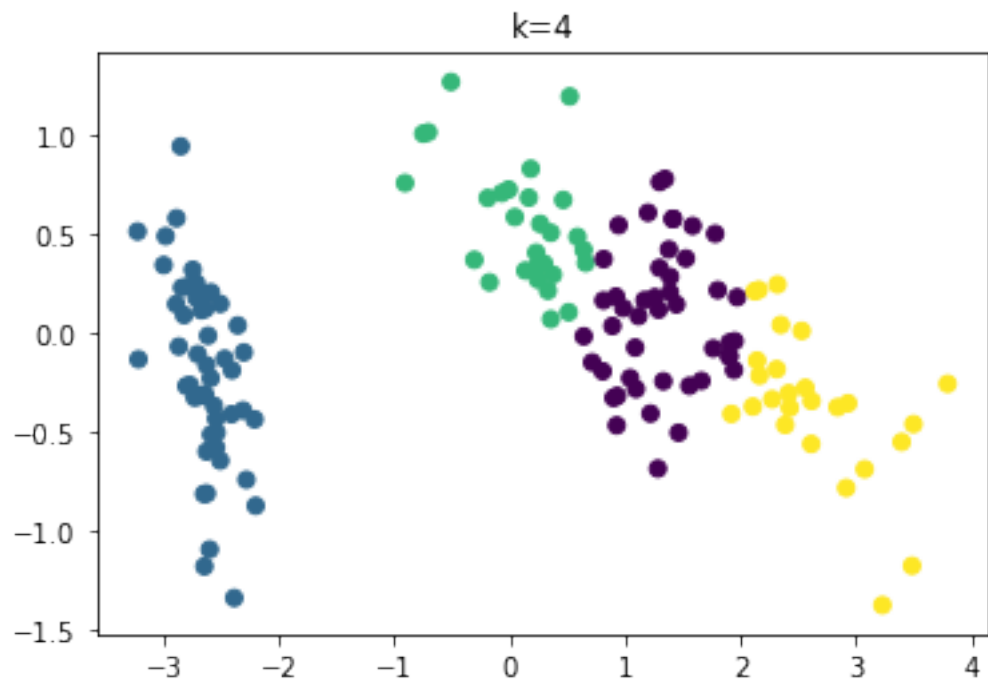
# Repeat for k=3, 4 and 5:
plt.scatter(P[:,0],P[:,1],c=yhat3)
plt.title('k=3')
plt.figure()

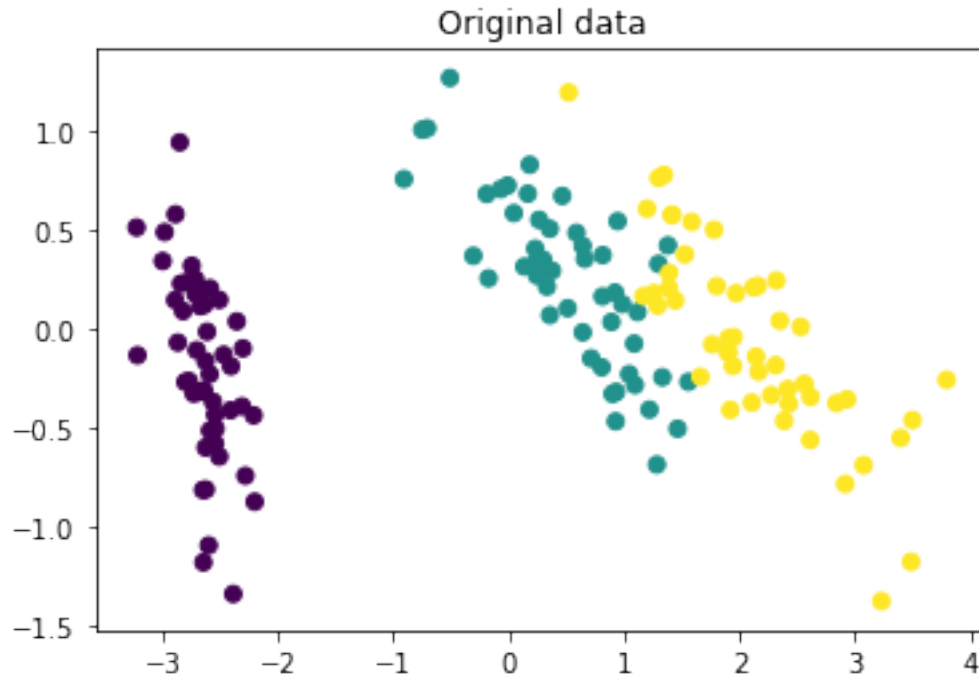
plt.scatter(P[:,0],P[:,1],c=yhat4)
plt.title('k=4')
plt.figure()

plt.scatter(P[:,0],P[:,1],c=yhat5)
plt.title('k=5')
plt.figure()

plt.scatter(P[:,0],P[:,1],c=y)
plt.title('Original data')
plt.show()
```







Comment: After we ran KMeans algorithm, the with $k=2$, $k=3$, $k=4$ and $k=5$, we see the data has the same coordinates but are more separated into 5 sets, unlike the original data with three sets.

2 Quantitative Assessment of K-Means (Bachelor and master students)

We used k-means for clustering and we assessed the results qualitatively by visualizing them. However, we often want to be able to measure in a quantitative way how good the clustering was. To do this, we will use a classification task to evaluate numerically the goodness of the representation learned via k-means.

Reload the *iris* dataset. Import a standard `LogisticRegression` classifier from the module `sklearn.linear_model`. Use the k-means representations learned previously (`yhat2, ..., yhat5`) and the true label to train the classifier. Evaluate your model on the training data (we do not have a test set, so this procedure will assess the model fit instead of generalization) using the `accuracy_score()` function from the `sklearn.metrics` module. Plot a graph showing how the accuracy score varies when changing the value of k . Comment on the results.

- Train a Logistic regression model using the first two dimensions of the PCA of the iris data set as input, and the true classes as targets.
- Report the model fit/accuracy on the training set.
- For each value of K :
 - One-Hot-Encode the classes outputted by the K-means algorithm.

- Train a Logistic regression model on the K-means classes as input vs the real classes as targets.
- Calculate model fit/accuracy vs. value of K.
- Plot your results in a graph and comment on the K-means fit.

```
[39]: from sklearn.linear_model import LogisticRegression
      from sklearn import metrics

      X,y = syntheticdata.get_iris_data()
      clf = LogisticRegression().fit(P, y)
      print(f"score accuracy: {clf.score(P, y)}")
```

score accuracy: 0.9133333333333333

```
/srv/conda/envs/notebook/lib/python3.6/site-
packages/sklearn/linear_model/logistic.py:432: FutureWarning: Default solver
will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
  FutureWarning)
/srv/conda/envs/notebook/lib/python3.6/site-
packages/sklearn/linear_model/logistic.py:469: FutureWarning: Default
multi_class will be changed to 'auto' in 0.22. Specify the multi_class option to
silence this warning.
  "this warning.", FutureWarning)
```

```
[44]: k2 = [2,3,4,5]

      accSco = []

      from sklearn.preprocessing import OneHotEncoder
      from sklearn.metrics import accuracy_score

      enc = OneHotEncoder(categories="auto")

      for i in k2:
          KM = KMeans(i)
          yhat = KM.fit_predict(P).reshape(-1, 1)
          OneHotEncode_yhat = enc.fit_transform(yhat).toarray()
          clf = LogisticRegression().fit(yhat, y)
          y_pred = clf.predict(yhat)
          AS = accuracy_score(y, y_pred)
          accSco.append(AS)

      #yhat2 = [1,0,0,0]
      #yhat3 = [0,1,0,0]
      #yhat4 = [0,0,1,0]
```

```

#yhat5 = [0,0,0,1]

print(f"k=2: {accSco[0]}")
print(f"k=3: {accSco[1]}")
print(f"k=4: {accSco[2]}")
print(f"k=5: {accSco[3]}")
print(accSco)
plt.plot(k2, accSco)
plt.xlabel("k")
plt.ylabel("score")
plt.show()

```

```

/srv/conda/envs/notebook/lib/python3.6/site-
packages/sklearn/linear_model/logistic.py:432: FutureWarning: Default solver
will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
  FutureWarning)
/srv/conda/envs/notebook/lib/python3.6/site-
packages/sklearn/linear_model/logistic.py:469: FutureWarning: Default
multi_class will be changed to 'auto' in 0.22. Specify the multi_class option to
silence this warning.
  "this warning.", FutureWarning)
/srv/conda/envs/notebook/lib/python3.6/site-
packages/sklearn/linear_model/logistic.py:432: FutureWarning: Default solver
will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
  FutureWarning)
/srv/conda/envs/notebook/lib/python3.6/site-
packages/sklearn/linear_model/logistic.py:469: FutureWarning: Default
multi_class will be changed to 'auto' in 0.22. Specify the multi_class option to
silence this warning.
  "this warning.", FutureWarning)
/srv/conda/envs/notebook/lib/python3.6/site-
packages/sklearn/linear_model/logistic.py:432: FutureWarning: Default solver
will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
  FutureWarning)
/srv/conda/envs/notebook/lib/python3.6/site-
packages/sklearn/linear_model/logistic.py:469: FutureWarning: Default
multi_class will be changed to 'auto' in 0.22. Specify the multi_class option to
silence this warning.
  "this warning.", FutureWarning)
/srv/conda/envs/notebook/lib/python3.6/site-
packages/sklearn/linear_model/logistic.py:432: FutureWarning: Default solver
will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
  FutureWarning)
/srv/conda/envs/notebook/lib/python3.6/site-
packages/sklearn/linear_model/logistic.py:469: FutureWarning: Default
multi_class will be changed to 'auto' in 0.22. Specify the multi_class option to

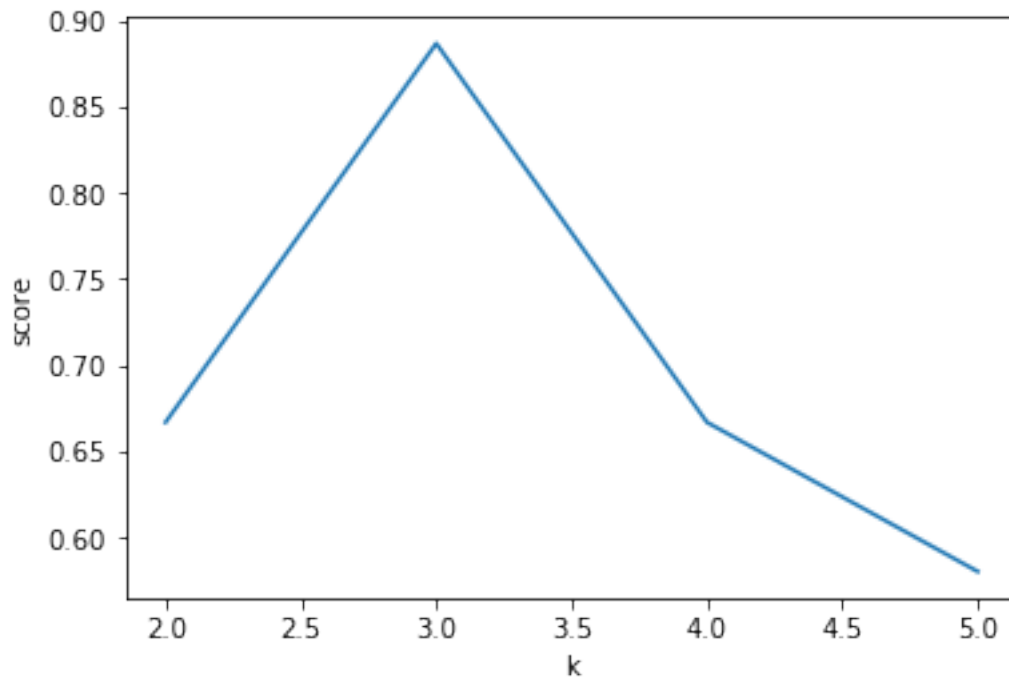
```

```

silence this warning.
    "this warning.", FutureWarning)

k=2: 0.6666666666666666
k=3: 0.8866666666666667
k=4: 0.6666666666666666
k=5: 0.58
[0.6666666666666666, 0.8866666666666667, 0.6666666666666666, 0.58]

```



Comment: I see that $k=3$ gives the best results, which is correct. Above when we plot for $k=2$, $k=3$, $k=4$ and $k=5$, the k which is closest to our original data is **3**. I think the scores for $k=2$, $k=4$ and $k=5$, seems to have pretty good accuracy, even though they are a bit low in my opinion.

3 Conclusions

In this notebook we studied **unsupervised learning** considering two important and representative algorithms: **PCA** and **k-means**.

First, we implemented the PCA algorithm step by step; we then run the algorithm on synthetic data in order to see its working and evaluate when it may make sense to use it and when not. We then considered two typical uses of PCA: for **visualization** on the *iris* dataset, and for **compression-decompression** on the *lfw* dataset.

We then moved to consider the k-means algorithm. In this case we used the implementation provided by *scikit-learn* and we applied it to another prototypical unsupervised learning problem:

clustering; we used *k-means* to process the *iris* dataset and we evaluated the results visually.

In the final part, we considered two additional questions that may arise when using the above algorithms. For PCA, we considered the problem of **selection of hyper-parameters**, that is, how we can select the hyper-parameter of our algorithm in a reasonable fashion. For k-means, we considered the problem of the **quantitative evaluation** of our results, that is, how can we measure the performance or usefulness of our algorithms.

[]: