

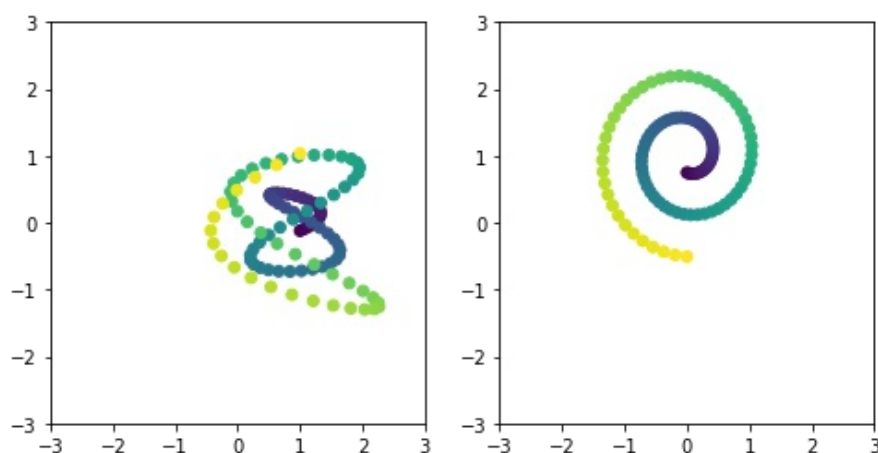
# Canonical Correlation Analysis

In this exercise, we consider canonical correlation analysis (CCA) on two simple problems, one in low dimensions and one in high dimensions. The goal is to implement the primal and dual versions of CCA to handle these two different cases. The first dataset consists of two trajectories in two dimensions. The dataset is extracted and plotted below. The first data points are shown in dark blue, and the last ones are shown in yellow.

In [1]:

```
import numpy
import matplotlib
%matplotlib inline
from matplotlib import pyplot as plt
import utils

X,Y = utils.getdata()
p1,p2 = utils.plotdata(X,Y)
```



For these two trajectories, that can be understood as two different modalities of the same data, we would like to determine under which projections they appear maximally correlated.

## Implementing Primal CCA

As stated in the lecture, the CCA problem in its primal form consists of maximizing the cross-correlation objective:

$$J(w_x, w_y) = w_x^\top C_{xy} w_y$$

subject to autocorrelation constraints  $w_x^\top C_{xx} w_x = 1$  and  $w_y^\top C_{yy} w_y = 1$ . Using the method of Lagrange multipliers, this optimization problem can be reduced to finding the first eigenvector of the generalized eigenvalue problem:

$$\begin{bmatrix} 0 & C_{xy} \\ C_{yx} & 0 \end{bmatrix} \begin{bmatrix} w_x \\ w_y \end{bmatrix} = \lambda \begin{bmatrix} C_{xx} & 0 \\ 0 & C_{yy} \end{bmatrix} \begin{bmatrix} w_x \\ w_y \end{bmatrix}$$

Your first task is to write a function that solves the CCA problem in the primal (i.e. that solves the generalized eigenvalue problem above). The function you need to implement receives two matrices  $X$  and  $Y$  of size  $N \times$

`d1` and `N × d2` respectively. It returns two vectors of size `d1` and `d2` corresponding to the projections associated to the modalities `X` and `Y`. (Hint: Note that the data matrices `X` and `Y` have not been centered yet.)

In [2]:

```
import numpy as np
from scipy.linalg import eig

X,Y = utils.getdata()
p1,p2 = utils.plotdata(X,Y)

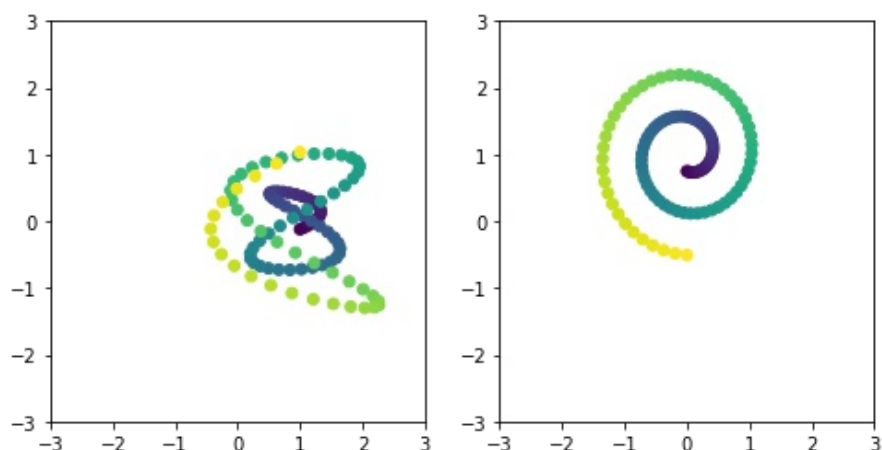
def CCAprimal(X,Y):
    # centering data
    X=X-np.mean(X,axis=0)
    Y=Y-np.mean(Y,axis=0)
    # calculating covariance matrices
    N = X.shape[0]
    x_dim = X.shape[1]
    y_dim = Y.shape[1]

    C_xx= (X.T.dot(X))/N
    C_yy= (Y.T.dot(Y))/N
    C_xy= (X.T.dot(Y))/N
    C_yx= (Y.T.dot(X))/N

    C_xx_yy = np.zeros((x_dim+y_dim,x_dim+y_dim))
    C_xx_yy[0:x_dim,0:x_dim]= C_xx
    C_xx_yy[x_dim:y_dim+x_dim,x_dim:y_dim+x_dim]= C_yy # C = [[C_xx,0],[0,C_yy]]

    C_xy_yx = np.zeros((x_dim+y_dim,x_dim+y_dim))
    C_xy_yx[0:x_dim,x_dim:y_dim+x_dim]= C_xy
    C_xy_yx[x_dim:y_dim+x_dim,0:x_dim]= C_yx # C = [[C_xx,0],[0,C_yy]]

    # solving the eigenvector problem
    eigvals,eigenvecs=eig(C_xy_yx, b=C_xx_yy)
    sol_ind = np.argmax(eigvals)
    w = eigenvecs[:,sol_ind].T
    wx = w[0:x_dim]
    wx = wx/np.linalg.norm(wx)
    wy= w[x_dim:x_dim+y_dim]
    wy= wy/np.linalg.norm(wy)
    return wx.reshape((2,)),wy.reshape((2,))
```



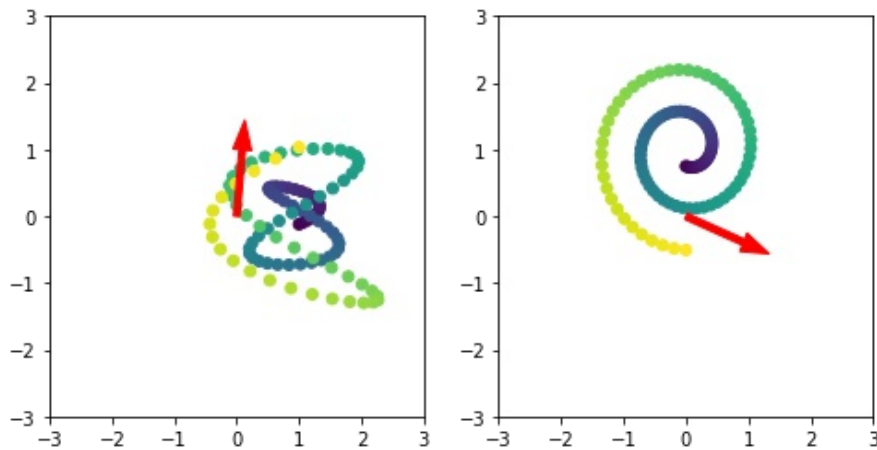
The function can now be called with our dataset. The learned projection vectors  $w_x$  and  $w_y$  are plotted as red arrows.

In [3]:

```
wx,wy = CCAprimal(X,Y)

p1,p2 = utils.plotdata(X,Y)
p1.arrow(0,0,1*wx[0],1*wx[1],color='red',width=0.1)
```

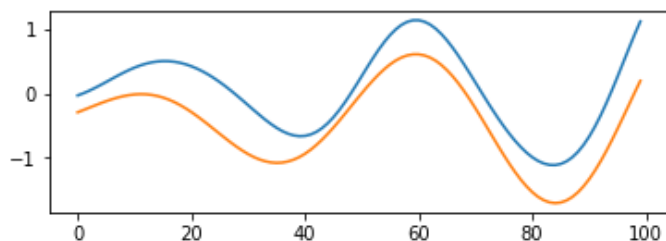
```
p2.arrow(0,0,1*wy[0],1*wy[1],color='red',width=0.1)
plt.show()
```



In each modality, the arrow points in a specific direction (note that the optimal CCA directions are defined up to a sign flip of both  $w_x$  and  $w_y$ ). Furthermore, we can verify CCA has learned a meaningful solution by projecting the data on it.

In [4]:

```
plt.figure(figsize=(6,2))
plt.plot(numpy.dot(X,wx))
plt.plot(numpy.dot(Y,wy))
plt.show()
```



Clearly, the data is correlated in the projected space.

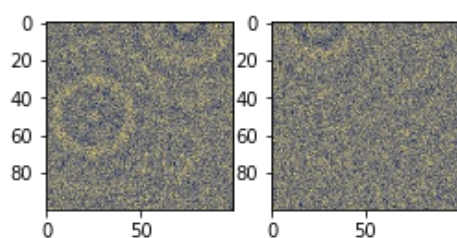
## Implementing Dual CCA

In the second part of the exercise, we consider the case where the data is high dimensional (with  $d \gg N$ ). Such high-dimensionality occurs for example, when input data are images. We consider the scenario where sources emit spatially, and two (noisy) receivers measure the spatial field at different locations. We would like to identify signal that is common to the two measured locations, e.g. a given source emitting at a given frequency. We first load the data and show one example.

In [5]:

```
X,Y = utils.getHDdata()

utils.plotHDdata(X[0],Y[0])
plt.show()
```



Several sources can be perceived, however, there is a significant level of noise. Here again, we will use CCA to find subspaces where the two modalities are maximally correlated. In this example, because there are many more dimensions than there are data points, it is more advantageous to solve CCA in the dual. Your task is to

more dimensions than there are data points, it is more advantageous to solve CCA in the dual. Your task is to implement a CCA dual solver that receives two data matrices of size  $N \times d_1$  and  $N \times d_2$  respectively as input, and returns the associate CCA directions (two vectors of respective sizes  $d_1$  and  $d_2$ ).

In [6]:

```
def CCAdual(X,Y):

    # centering data
    X=X-np.mean(X,axis=0)
    Y=Y-np.mean(Y,axis=0)
    # calculating covariance matrices
    N = X.shape[0]
    x_dim = X.shape[1]
    y_dim = Y.shape[1]

    A = (X.dot(X.T))
    B = (Y.dot(Y.T))

    C_xx_yy = np.zeros((2*N,2*N))
    # regularization , see https://www.di.ens.fr/~fbach/kernelICA-jmlr.pdf , page 11
    A_reg= A + 0.05*np.identity(N)
    B_reg = B + 0.05*np.identity(N)

    C_xx_yy[0:N,0:N]= A_reg.dot(A_reg)
    C_xx_yy[N:2*N,N:2*N]= B_reg.dot(B_reg)    # C = [[C_xx,0],[0,C_yy]]

    C_xy_yx = np.zeros((2*N,2*N))
    C_xy_yx[0:N,N:2*N]= A.dot(B)
    C_xy_yx[N:2*N,0:N]= B.dot(A)    # C = [[C_xx,0],[0,C_yy]]

    # solving the eigenvector problem
    eigvals,eigenvecs=eig(C_xy_yx, b=C_xx_yy)

    sol_ind = np.argmax(eigvals)

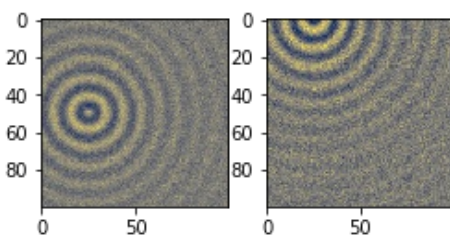
    alpha = eigenvecs[:,sol_ind].T
    alpha_x = alpha[0:N]
    alpha_y= alpha[N:2*N]
    wx = X.T.dot(alpha_x)
    wy = Y.T.dot(alpha_y)
    return wx,wy
```

We now call the function we have implemented with a training sequence of 100 pairs of images. Because the returned solution is of same dimensions as the inputs, it can be rendered in a similar fashion.

In [7]:

```
wx,wy = CCAdual(X[:100],Y[:100])

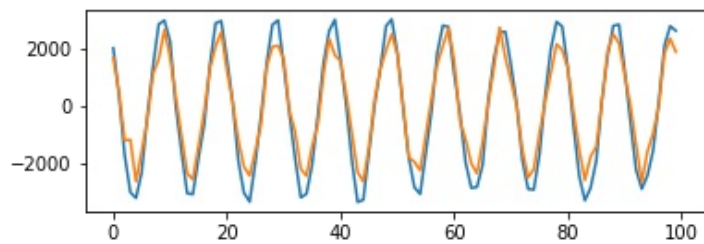
utils.plotHDdata(wx,wy)
plt.show()
```



Here, we can clearly see a common factor that has been extracted between the two fields, specifically a point source emitting at a particular frequency. A test sequence of 100 pairs of images can now be projected on these two filters:

In [8]:

```
plt.figure(figsize=(6,2))  
plt.plot(numpy.dot(X[100:],wx))  
plt.plot(numpy.dot(Y[100:],wy))  
plt.show()
```



**Clearly the two projected signals are correlated and the input noise has been strongly reduced.**