

Problem set 1: PCA, LLE, outlier detection

On the ISIS page you can download `testsNstubs.zip` containing the files `ps1_implementation.py`, `ps1_application.py` and `ps1_tests.py`. Do not modify the names of the files or the names of the functions within these files. You are free to define additional functions within the given files. Make sure that your functions have the correct signatures.

The file `ps1_tests.py` is designed to help you debug your code. It contains test functions for each of the implementation assignments in Part 1, including runtimes of *my* code on *my* machine for comparison. Make sure that your code passes all tests. Be aware that

- a passed test does not guarantee correctness for all possible inputs
- if the test module produces a plot, you have to check if the plot looks correct.
- if your code is very slow compared to the given runtimes, you should optimize your code

You have to submit

1. a report (pdf) with your analysis
2. the implementation file `ps1_implementation.py`
3. the application file `ps1_application.py` or an Ipython notebook `ps1_application.ipynb`

Part 1: Implementation

Function stubs for these assignments have been provided in `ps1_implementation.py`.

Assignment 1 (15 pts)

Write the function `pca` with signature

$$Z, U, D = \text{pca}(X, m)$$

which receives a $d \times n$ matrix X and the number of components m to be used, and which returns the principal components, as well as the projected data points in a $m \times n$ matrix Z .

U and D should contain the principal components: U is a $d \times d$ matrix, which contains the principal directions, and D is a $1 \times d$ vector, which contains the principal values, sorted in descending order (i.e. $D_1 \geq D_2 \dots$).

Assignment 2 (15 pts)

Implement the γ -index (see paper on the ISIS page),

$$y = \text{gammaidx}(X, k)$$

which receives a $d \times n$ matrix X containing the data points and the number of neighbours k , and returns the γ -index for each datapoint in the $1 \times n$ vector y .

Assignment 3 (20 pts)

Write a function `LLE` with signature

$$Y = \text{lle}(X, m, \text{n.rule}, \text{param}, \text{tol})$$

which receives a $d \times n$ matrix X containing the data points, and which calculates an m -dimensional embedding $Y \in \mathbb{R}^{m \times n}$ using the LLE algorithm. The parameter `n.rule` determines the method (`'knn'` or `'eps-ball'`) which is used to build the neighbourhood graph where `param` is the corresponding parameter (k or ϵ , respectively). `tol` is the size of the regularization parameter.

Your function should be robust against malicious parameters. Use `raise Exception` for error reporting. In particular, you should check whether the resulting graph is connected.

Part 2: Application

Function stubs for these assignments have been provided in `ps1_application.py` (you are welcome to generate an Ipython notebook for this part).

Assignment 4 (10 pts)

Write the function `usps` which applies PCA to the `usps` data set (available on ISIS) in the following manner:

1. Load the `usps` data set.
2. Analysis of PCA:
 - (a) Calculate the PCA of this data.
 - (b) Visualize (a) all principal values, (b) the largest 25 principal values (both as a bar plot, see `bar`) and (c) the first 5 principal directions (as images, see `imshow`).
3. Consider three noise scenarios (use `numpy.random.randn`):
 - **low gaussian noise**: add Gaussian noise to the images. Select an appropriate variance such that the resulting images are *moderately* noisy.
 - **high gaussian noise**: add Gaussian noise to the images. Select an appropriate variance such that the resulting images are *very* noisy.
 - **outliers**: add Gaussian noise to *only five* images. Select an appropriate variance such that those five images are *extremely* noisy and that the noise is reflected in the spectrum.

For each of the scenarios

- (a) Calculate the PCA of this data and redo the plots of principal values. Explain the differences to the original spectrum.
- (b) Denoise the images by reconstruction from projections on the `m` largest principal components: the reconstruction y of a data point x by the m largest eigenvectors v_1, \dots, v_m of the covariance matrix is given by

$$y = \mu + \sum_{i=1}^m v_i (x^\top v_i),$$

where μ is the mean of the original data set. The reconstruction error of x is $\|x - y\|$.

The number of components `m` should be tuned by hand such that the reconstructed (denoised) images are fairly good.

- (c) For 10 examples of your choice (including the noisy ones in the **outlier** setting), plot the original image, the noisy image and its denoised reconstruction (using `imshow`).

Remark: use `subplot`, `gridspec` or `axes` to arrange multiple plots/images in a single figure.

Assignment 5 (10 pts)

In this exercise, we use the positive class of the `banana` dataset (available on ISIS) as “inliers” to which we add outliers from the negative class. We will investigate the performance of outlier detection algorithms for outlier contamination rates (i.e. percentage of outliers in the data set) of 1%, 5%, 10% and 25% relative to the positive class. Write a function `outliers_calc` which –for each of these rates– repeats the following procedure 100 times:

1. Choose a random set of outliers from the negative class of the respective size (depending on the outlier rate).
2. Add the outliers to the positive class, and compute (a) the γ -index with $k = 3$, (b) the γ -index with $k = 10$ and (c) the distance to the mean for each data point.
3. Compute the AUC (area under the ROC) for each method.

For each contamination rate and method we thus obtain 100 AUC values. Store these (use `savez`)!

Write a function `outliers_disp` which creates a plot that allows to compare the performance of the methods, using a `boxplot` to visualize the distribution of the AUC values.

Assignment 6 (20 pts)

Write `lle_visualize` which applies LLE to the data sets `fishbowl`, `swissroll` and `flatroll` (see table 1), using appropriate values for the parameters. It plots the data and the resulting embedding, visualizing the “true” embedding using a colour coding (e.g. using `scatter`).

Remarks:

- LLE is sensitive with respect to the parameters. You have to fine-tune `param` and `tol`.
- It is difficult to find good parameters for the `fishbowl` data. Try the “dense fishbowl” data set first.

data set	file name	data	true embedding
fishbowl	<code>fishbowl_data.npz</code>	<code>x_noisefree</code> (3D)	<code>z</code> (2D)
dense fishbowl	<code>fishbowl_dense.npz</code>	<code>x_noisefree</code> (3D)	<code>z</code> (2D)
swissroll	<code>swissrole_data.npz</code>	<code>x_noisefree</code> (3D)	<code>z</code> (2D)
flatroll	<code>flatrole_data.npz</code>	<code>Xflat</code> (2D)	<code>true_embedding</code> (1D)

Table 1: the data sets (available on ISIS)

Assignment 7 (10 pts)

In this exercise, you study the influence of noise on LLE, using the example of the `flatroll` data set. Write the function `lle_noise` which does the following:

1. Loads the data set.
2. Adds Gaussian noise with variance 0.2 and 1.8 to the data set (this results in 2 noisy data sets).
3. Applies LLE on both data sets, where the neighborhood graph should be constructed using k -nn. For both noise levels, try to find (a) a good value for k which unrolls the flat roll and (b) a value which is obviously too large.
4. For each of the four combinations of low/high noise level good/too large k , plots the neighbourhood graph and the resulting embedding.