SLT coding exercise #1

# Locally Linear Embedding
https://gitlab.vis.ethz.ch/vwegmayr/slt-coding-exercises

Due on Monday, March 6th, 2017

Carl Rynegardh
16-909-327

# Contents

# The Model

The model section is intended to allow you to recapitulate the essential ingredients used in Locally Linear Embedding. Write down the *necessary* equations to specify Locally Linear Embedding and and shortly explain the variables that are involved. This section should only introduce the equations, their solution should be outlined in the implementation section.

Hard limit: One page

---

LLE is a type of unsupervised learning used for the problem of non linear dimensionality reduction. We want to minimize the reconstruction error with regard to the weights W.

Reconstruction error:

$$\mathcal{E}(W) = \sum_i \left| X_i - \sum_j W_{ij} X_j \right|^2$$

In the above formula $W_{ij}$ stands for the contribution from the j-th data point to the i-th data point. $X_i$ stands for the i-ith data point with D dimensions. Minimizing with regards to the weights, we also have two constraints: $\sum_j W_{ij} = 1$ and that i-th data point only should be reconstructed with the help of it's $K$ neighbors.

To receive the embedded coordinate in d, $d << D$, dimensions the following cost function is minimized with regards to Y, which is the low dimensional vector X is mapped to:

$$\Phi(Y) = \sum_i \left| Y_i - \sum_j W_{ij} Y_j \right|^2 = \sum_{ij} M_{ij}(Y_i \cdot Y_j)$$

---

# The Questions

This is the core section of your report, which contains the tasks for this exercise and your respective solutions. Make sure you present your results in an illustrative way by making use of graphics, plots, tables, etc. so that a reader can understand the results with a single glance. Check that your graphics have enough resolution or are vector graphics. Consider the use of GIFs when appropriate.
Hard limit: Two pages

## (a) Get the data

For this exercise we will work with the MNIST data set. In order to learn more about it and download it, go to http://yann.lecun.com/exdb/mnist/.

## (b) Locally linear embedding

Implement the LLE algorithm and apply it to the MNIST data set. Provide descriptive visualizations for 2D & 3D embedding spaces. Is it possible to see clusters?

## (c) Cluster structure

Investigate the cluster structure of the data. Can you observe block structures in the $M$ matrix (use matrix plots)? Also plot the singular values of $M$. Do you notice something? Can you think of ways to determine the optimal embedding dimension?
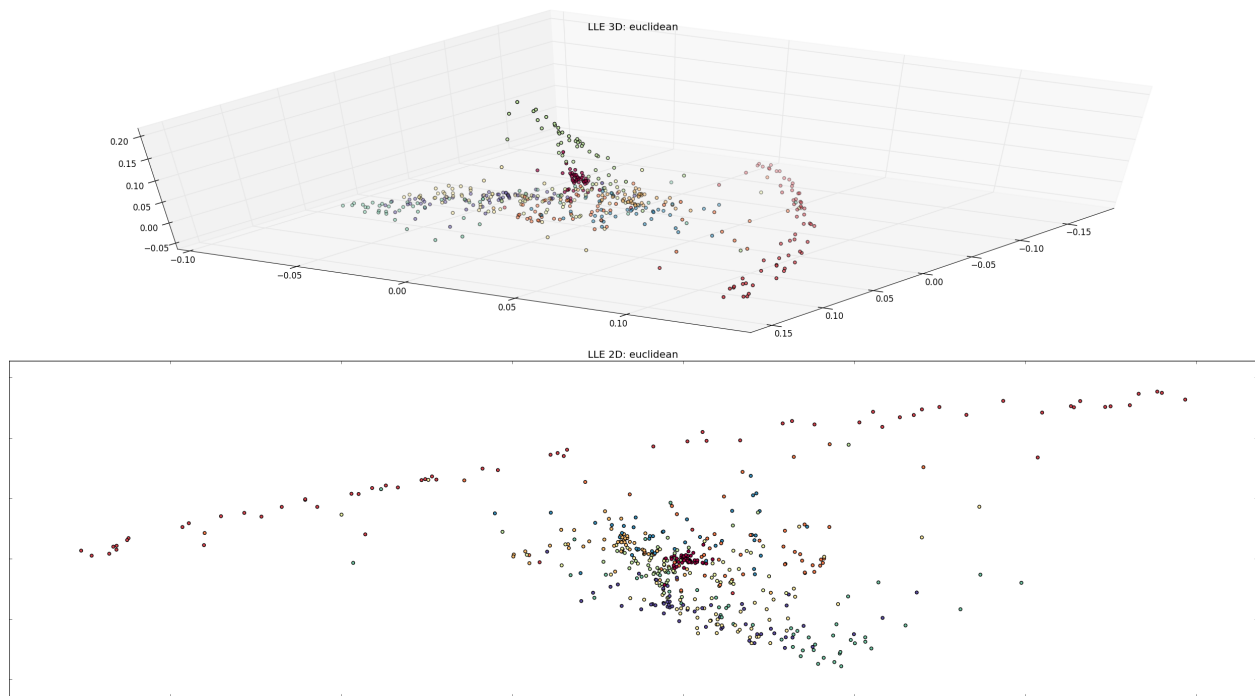
## (d) Nearest Neighbors

Investigate the influence of the choice of how many nearest neighbors you take into account. Additionally, try different metrics to find the nearest neighbors (we are dealing with images!).

## (e) Linear manifold interpolation

Assume you pick some point in the embedding space. How can you map it back to the original (high dimensional) space? Investigate how well this works for points within and outside the manifold (does it depend on the dimensionality of the embedding space?) Try things like linearly interpolating between two embedding vectors and plot the sequence of images along that line. What happens if you do that in the original space?

---

**(b):** In figure 1, we see both 3D and 2D of the embedded space. In 2D we can see some clustering and in 3D we see even more. While it might not be perfect, we can at least see some clustering.
**(c):** There seem to be no structure in the M matrix. It might seem like there are structure in the diagonal but that is natural since that is the same number multiplied with itself(Fig 2).
**(d):** Looking at figure 4-9 one can easily see that different metrics gives different clustering and with more neighbors the structure get more complex. I can imagine that it is important to especially try with different number of neighbors to find a good clustering.
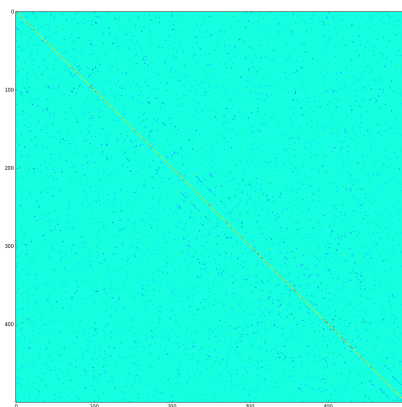
---

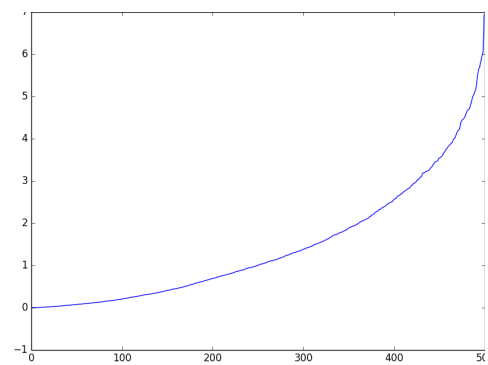(a) 2D LLE, 3D LLE plot

Figure 1



Figure 2: The M matrix
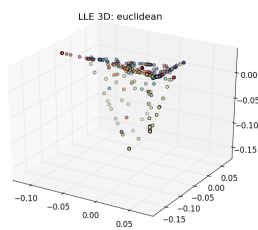


Figure 3: Singular values of the M matrix

Figure 4: k = 3, metric = euclidean
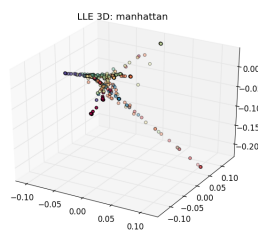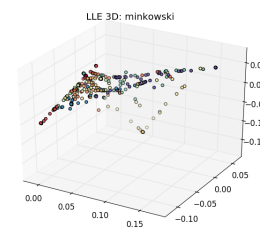


Figure 5: k = 3, metric = manhattan
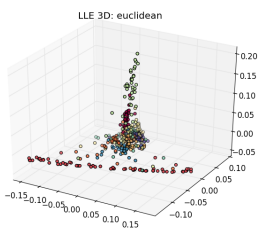


Figure 6: k = 3, metric = minkowski
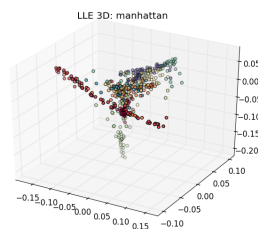


Figure 7: k = 15, metric = euclidean



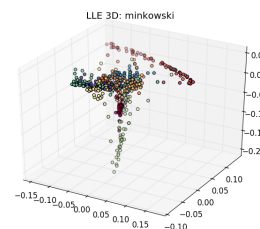Figure 8: k = 15, metric = manhattan



Figure 9: k = 15, metric = minkowski

# The Implementation

In the implementation section you give a concise insight to the practical aspects of this coding exercise. It mainly mentions the optimization methods used to solve the model equations. Did you encounter numerical or efficiency problems? If yes, how did you solve them? Provide the link to your git branch of this coding exercise.

Hard limit: One page

Your Answer

# Your Page

Your page gives you space to include ideas, observations and results which do not fall into the categories provided by us. You can also use it as an appendix to include things which did not have space in the other sections.

No page limit.

Your Answer
YOUR GIT BRANCH