

## Exercise Sheet 1

### Exercise 1: Symmetries in LLE (25 P)

The Locally Linear Embedding (LLE) method takes as input a collection of data points  $\vec{x}_1, \dots, \vec{x}_N \in \mathbb{R}^d$  and embeds them in some low-dimensional space. LLE operates in two steps, with the first step consisting of minimizing the objective

$$\mathcal{E}(w) = \sum_{i=1}^N \left\| \vec{x}_i - \sum_j w_{ij} \vec{x}_j \right\|^2$$

where  $w$  is a collection of reconstruction weights subject to the constraint  $\forall i : \sum_j w_{ij} = 1$ , and where  $\sum_j$  sums over the  $K$  nearest neighbors of the data point  $\vec{x}_i$ . The solution that minimizes the LLE objective can be shown to be invariant to various transformations of the data.

Show that invariance holds in particular for the following transformations:

- (a) Replacement of all  $\vec{x}_i$  with  $\alpha \vec{x}_i$ , for an  $\alpha \in \mathbb{R}^+ \setminus \{0\}$ ,
- (b) Replacement of all  $\vec{x}_i$  with  $\vec{x}_i + \vec{v}$ , for a vector  $\vec{v} \in \mathbb{R}^d$ ,
- (c) Replacement of all  $\vec{x}_i$  with  $U \vec{x}_i$ , where  $U$  is an orthogonal  $d \times d$  matrix.

### Exercise 2: Closed form for LLE (25 P)

In the following, we would like to show that the optimal weights  $w$  have an explicit analytic solution. For this, we first observe that the objective function can be decomposed as a sum of as many subobjectives as there are data points:

$$\mathcal{E}(w) = \sum_{i=1}^N \mathcal{E}_i(w) \quad \text{with} \quad \mathcal{E}_i(w) = \left\| \vec{x}_i - \sum_j w_{ij} \vec{x}_j \right\|^2$$

Furthermore, because each subobjective depends on different parameters, they can be optimized independently. We consider one such subobjective and for simplicity of notation, we rewrite it as:

$$\mathcal{E}_i(w) = \left\| \vec{x} - \sum_{j=1}^K w_j \vec{\eta}_j \right\|^2$$

where  $\vec{x}$  is the current data point (we have dropped the index  $i$ ), where  $\eta = (\vec{\eta}_1, \dots, \vec{\eta}_K)$  is a matrix of size  $K \times d$  containing the  $K$  nearest neighbors of  $\vec{x}$ , and  $w$  is the vector of size  $K$  containing the weights to optimize and subject to the constraint  $\sum_{j=1}^K w_j = 1$ .

- (a) Prove that the optimal weights for  $\vec{x}$  are found by solving the following optimization problem:

$$\min_w w^\top C w \quad \text{subject to} \quad w^\top \mathbf{1} = 1.$$

where  $C = (\mathbf{1} \vec{x}^\top - \eta)(\mathbf{1} \vec{x}^\top - \eta)^\top$  is the covariance matrix associated to the data point  $\vec{x}$  and  $\mathbf{1}$  is a vector of ones of size  $K$ .

- (b) Show using the method of Lagrange multipliers that the minimum of the optimization problem found in (a) is given analytically as:

$$w = \frac{C^{-1} \mathbf{1}}{\mathbf{1}^\top C^{-1} \mathbf{1}}.$$

- (c) Show that the optimal  $w$  can be equivalently found by solving the equation  $Cw = \mathbf{1}$  and then rescaling  $w$  such that  $w^\top \mathbf{1} = 1$ .

**Exercise 3: t-SNE and Kullback-Leibler Divergence (25 P)**

The t-SNE embedding algorithm operates by minimizing the Kullback-Leibler divergence between two discrete probability distributions  $p$  and  $q$  representing the input space and the embedding space respectively. These discrete distributions assign to each pair of data points  $(i, j)$  in the dataset the probability scores  $p_{ij}$  and  $q_{ij}$  respectively, corresponding to how close the two data points are in the input and embedding spaces. Once the exact probability functions are defined, the embedding algorithm proceeds by optimizing the function:

$$\begin{aligned} C &= D_{\text{KL}}(p \parallel q) \\ &= \sum_{i=1}^N \sum_{j=1}^N p_{ij} \log \left( \frac{p_{ij}}{q_{ij}} \right) \end{aligned}$$

where  $p$  and  $q$  are subject to the constraints  $\sum_{i=1}^N \sum_{j=1}^N p_{ij} = 1$  and  $\sum_{i=1}^N \sum_{j=1}^N q_{ij} = 1$ . Specifically, the algorithm minimizes  $q$  which itself is a function of the coordinates in the embedded space. Optimization is typically performed using gradient descent.

In this exercise, we derive the gradient of the Kullback-Leibler divergence, first with respect to the probability scores  $q_{ij}$ , and then with respect to the embedding coordinates of which  $q_{ij}$  is a function.

(a) *Show that*

$$\frac{\partial C}{\partial q_{ij}} = -\frac{p_{ij}}{q_{ij}}. \quad (1)$$

(b) The probability matrix  $q$  is now reparameterized using a ‘softargmax’ function:

$$q_{ij} = \frac{\exp(z_{ij})}{\sum_{k=1}^N \sum_{l=1}^N \exp(z_{kl})}$$

The new variables  $z_{ij}$  can be interpreted as unnormalized log-probabilities. *Show that*

$$\frac{\partial C}{\partial z_{ij}} = -p_{ij} + q_{ij}. \quad (2)$$

(c) *Explain* which of the two gradients, (1) or (2), is the most appropriate for practical use in a gradient descent algorithm. Motivate your choice, first in terms of the stability or boundedness of the gradient, and second in terms of the ability to maintain a valid probability distribution during training.

(d) The scores  $z_{ij}$  are now reparameterized as

$$z_{ij} = -\|\vec{y}_i - \vec{y}_j\|^2$$

where the coordinates  $\vec{y}_i, \vec{y}_j \in \mathbb{R}^h$  of data points in embedded space now appear explicitly. *Show* using the chain rule for derivatives that

$$\frac{\partial C}{\partial \vec{y}_i} = \sum_{j=1}^N 4(p_{ij} - q_{ij}) \cdot (\vec{y}_i - \vec{y}_j).$$

**Exercise 4: Programming (25 P)**

Download the programming files on ISIS and follow the instructions.