Exercises for the course
## Machine Learning 2
Summer semester 2021

Abteilung Maschinelles Lernen
Institut für Softwaretechnik und theoretische Informatik
Fakultät IV, Technische Universität Berlin
Prof. Dr. Klaus-Robert Müller
Email: klaus-robert.mueller@tu-berlin.de

# Exercise Sheet 1

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

The Locally Linear Embedding (LLE) method takes as input a collection of data points $\vec{x}_1, \ldots, \vec{x}_N \in \mathbb{R}^d$ 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:

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

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 $\eta = (\vec{\eta}_1, \ldots, \vec{\eta}_K)$ 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:

$$\min_{w} \quad w^\top C w \qquad \text{subject to} \quad w^\top \mathbb{1} = 1.$$

where $C = (\mathbb{1}\vec{x}^\top - \eta)(\mathbb{1}\vec{x}^\top - \eta)^\top$ is the covariance matrix associated to the data point $\vec{x}$ and $\mathbb{1}$ 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:

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

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

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

SNE is an embedding algorithm that operates by minimizing the Kullback-Leibler divergence between two discrete probability distributions $p$ and $q$ representing the input space and the embedding space respectively. In 'symmetric SNE', 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:

$$C = D_{\mathrm{KL}}(p \,\|\, q)$$
$$= \sum_{i=1}^{N} \sum_{j=1}^{N} p_{ij} \log \left( \frac{p_{ij}}{q_{ij}} \right)$$

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

$$\frac{\partial C}{\partial q_{ij}} = -\frac{p_{ij}}{q_{ij}}. \tag{1}$$

(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

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

(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

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

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

Download the programming files on ISIS and follow the instructions.

# 1 Symmetries in LLE

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

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

$\operatorname{argmin} \mathcal{E}'(w) = \operatorname{argmin} \alpha^2 \mathcal{E}(w)$

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

$= \sum_{i=1}^{N} \| \vec{x}_i + \vec{v} - \sum_j w_{ij} \vec{x}_j - \sum_j w_{ij} \vec{v} \|^2$

$= \sum_{i=1}^{N} \| \vec{x}_i + \vec{v} - \sum_j w_{ij} \vec{x}_j - \vec{v} \underbrace{\sum_j w_{ij}}_{= 1} \|^2$

$= \sum_{i=1}^{N} \| \vec{x}_i + \vec{v} - \vec{v} - \sum_j w_{ij} \vec{x}_j \|^2$

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

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

$= \sum_{i=1}^{N} \left( U(\vec{x}_i - \sum_j w_{ij} \vec{x}_j) \right)^T \left( U(\vec{x}_i - \sum_j w_{ij} \vec{x}_j) \right)$

$= \sum_{i=1}^{N} (\vec{x}_i - \sum_j w_{ij} \vec{x}_j)^T \underbrace{U^T \cdot U}_{= \mathbb{1}} (\vec{x}_i - \sum_j w_{ij} \vec{x}_j)$

$= \sum_{i=1}^{N} (\vec{x}_i - \sum_j w_{ij} \vec{x}_j)^T (\vec{x}_i - \sum_j w_{ij} \vec{x}_j)$

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

## 2 Closed Form for LLE

a)  $\varepsilon_i(w) = \| \vec{x} - \sum\limits_{j=1}^{k} w_j \vec{\eta}_j \|^2 = \| \vec{x}\, w^T \mathbb{1} - \eta^T w \|^2$

$= \| (\vec{x}\, \mathbb{1}^T\, w - \eta^T w) \|^2 = \| (\vec{x}\, \mathbb{1}^T - \eta^T)\, w \|^2$

$= \big( (\vec{x}\, \mathbb{1}^T - \eta^T)\, w \big)^T \big( (\vec{x}\, \mathbb{1}^T - \eta^T)\, w \big)$

$= w^T (\vec{x}\, \mathbb{1}^T - \eta^T)^T (\vec{x}\, \mathbb{1}^T - \eta^T)\, w$

$= w^T (\mathbb{1}\, \vec{x}^T - \eta)(\vec{x}\, \mathbb{1}^T - \eta^T)\, w$

$= w^T (\mathbb{1}\, \vec{x}^T - \eta)(\mathbb{1}\, \vec{x}^T - \eta)^T\, w \quad = w^T C\, w$

b)  $\mathcal{L}(w, \lambda) = w^T C\, w + \lambda(w^T \mathbb{1} - 1) = w^T C\, w + \lambda w^T \mathbb{1} - \lambda$

$\dfrac{\partial \mathcal{L}(w, \lambda)}{\partial \lambda} = w^T \mathbb{1} - 1 \overset{!}{=} 0 \iff w^T \mathbb{1} = 1$

$\dfrac{\partial \mathcal{L}(w, \lambda)}{\partial w} = \dfrac{\partial w^T C\, w}{\partial w} + \dfrac{\partial \lambda \mathbb{1}^T w}{\partial w}$

$= \dfrac{\partial w^T C\, w}{\partial w^T} + \dfrac{\partial w^T C\, w}{\partial w} + \lambda \mathbb{1}^T$

$= \dfrac{\partial w C^T w^T}{\partial w} + \dfrac{\partial w^T C\, w}{\partial w} + \lambda \mathbb{1}^T$

$= C^T w^T + C\, w + \lambda \mathbb{1}^T = w^T C^T + w^T C + \lambda \mathbb{1}^T = w^T(C + C^T) + \lambda \mathbb{1}^T = 0$

$C \text{ symmetric} \implies C^T = C$

$\implies w^T \cdot 2C + \lambda \mathbb{1}^T = 0 \iff 2(w^T C) + \lambda \mathbb{1}^T = 2(C^T w)^T + \lambda \mathbb{1}^T$

$= (2(C^T w)^T + \lambda \mathbb{1}^T)^T = 2C\, w + \lambda \mathbb{1} \iff 2C\, w = -\lambda \mathbb{1}$

$\iff C\, w = -\tfrac{\lambda}{2} \mathbb{1} \iff w = -\tfrac{\lambda}{2} C^{-1} \mathbb{1} \iff = w^T \mathbb{1} \iff \mathbb{1}^T w = -\tfrac{\lambda}{2} \mathbb{1}^T C^{-1} \mathbb{1}$

$\iff 1 = -\tfrac{\lambda}{2} \mathbb{1}^T C^{-1} \mathbb{1} \iff \lambda = -\dfrac{2}{\mathbb{1}^T C^{-1} \mathbb{1}}$

$\iff w = \dfrac{C^{-1} \mathbb{1}}{\mathbb{1}^T C^{-1} \mathbb{1}}$

c) $Cw = \underline{1}, \quad w^T \underline{1} = 1$

$\Rightarrow w = c^{-1}\underline{1}$

$\Longleftrightarrow \dfrac{w}{w^T \underline{1}} = \dfrac{c^{-1}\underline{1}}{\underline{1}^T C^{-1}\underline{1}}$

## 3 SNE and Kullback-Leibler Divergence

a) 
$$\frac{\partial \mathcal{L}}{\partial q_{ij}} = \frac{\partial}{\partial q_{ij}} \sum_{i=1}^{N} \sum_{j=1}^{N} P_{ij} \log\left(\frac{P_{ij}}{q_{ij}}\right) = \frac{\partial}{\partial q_{ij}} P_{ij} \log\left(\frac{P_{ij}}{q_{ij}}\right)$$

$$= \frac{\partial}{\partial q_{ij}} P_{ij}\left(\log(P_{ij}) - \log(q_{ij})\right) = -P_{ij}\frac{1}{q_{ij}} = -\frac{P_{ij}}{q_{ij}}$$

b) 
$$\frac{\partial}{\partial z_{ij}} \sum_{i=1}^{N} \sum_{j=1}^{N} P_{ij} \log\left(\frac{P_{ij} \sum_{h=1}^{N} \sum_{l=1}^{N} e^{z_{hl}}}{e^{z_{ij}}}\right)$$

$$= \frac{\partial}{\partial z_{ij}} \sum_{i=1}^{N} \sum_{j=1}^{N} P_{ij}\left(\log\left(P_{ij} \sum_{h=1}^{N} \sum_{l=1}^{N} e^{z_{hl}}\right) - \log e^{z_{ij}}\right)$$

$$= \frac{\partial}{\partial z_{ij}} \sum_{i=1}^{N} \sum_{j=1}^{N} P_{ij}\left(\log(P_{ij}) + \log\left(\sum_{h=1}^{N} \sum_{l=1}^{N} e^{z_{hl}}\right) - z_{ij}\right)$$

$$= \frac{\partial}{\partial z_{ij}} \sum_{i=1}^{N} \sum_{j=1}^{N} P_{ij} \log\left(\sum_{h=1}^{N} \sum_{l=1}^{N} e^{z_{hl}}\right) + \sum_{i=1}^{N} \sum_{j=1}^{N} (-P_{ij} z_{ij})$$

$$= \frac{\partial}{\partial z} \sum_{i=1}^{N} \sum_{j=1}^{N}\left(\log\left(\sum_{h=1}^{N} \sum_{l=1}^{N} e^{z_{hl}}\right) - P_{ij} z_{ij}\right)$$

$$= \frac{\partial}{\partial z_{ij}} \log\left(\sum_{h=1}^{N} \sum_{l=1}^{N} e^{z_{hl}}\right) - P_{ij} z_{ij} = q_{ij} - P_{ij} = -P_{ij} + q_{ij}$$

c) $q_{ij} \leftarrow q_{ij} - \rho \dfrac{\partial \mathcal{L}}{\partial q_{ij}} \rightarrow -\dfrac{P_{ij}}{q_{ij}}$   If the probability of the embedded space gets very close to 0, the gradient becomes very large $\rightarrow$ not stable, not bounded

$z_{ij} \leftarrow z_{ij} - \rho \dfrac{\partial \mathcal{L}}{\partial z_{ij}} \rightarrow -P_{ij} + q_{ij} \rightarrow$ The gradient will never become infinite $\rightarrow$ stable and bounded

d)

$$\frac{\partial C}{\partial \vec{u}_i} = \sum_{j=1}^{N} \frac{\partial C}{\partial z_{\cdot j}} \cdot \frac{\partial z_j}{\partial \vec{u}_i} + \frac{\partial C}{\partial z_{j\cdot}} \cdot \frac{\partial z_j}{\partial \vec{u}_i}$$

$$= \sum_{j=1}^{N} (-p_{\cdot j} + q_{\cdot j}) \frac{\partial - \|\vec{u}_i - \vec{u}_j\|^2}{\partial \vec{u}_i} + (-p_{j\cdot} + q_{j\cdot}) \frac{\partial - \|\vec{u}_i - \vec{u}_j\|^2}{\partial \vec{u}_i}$$

$$= \sum_{j=1}^{N} (-p_{\cdot j} + q_{\cdot j})(-2(\vec{u}_i - \vec{u}_j)) + (-p_{j\cdot} + q_{j\cdot})(-2(\vec{u}_i - \vec{u}_j))$$

$$= \sum_{j=1}^{N} 2(-(p_{\cdot j} - q_{\cdot j})(-2(\vec{u}_i - \vec{u}_j)) = \sum_{j=1}^{N} 4(p_{\cdot j} - q_{\cdot j})(\vec{u}_i - \vec{u}_j)$$