

Machine Learning 2 – Group ESHG

Assignment 01

Willi Gierke, Arik Elimelech, Mehmed Halilovic, Leon Sixt

May 2, 2017

Exercise 1: Symmetries (40 P)

In the third paragraph of page 3, it is claimed that the optimal weights W_{ij} which minimize the cost function \mathcal{E} are independent with respect to scaling, translation and rotation of the original data \vec{X}_i . Prove that. That is, *prove* that the minimum (or minima) of \mathcal{E} is (or are) invariant under the following symmetries:

- (i) Replace all \vec{X}_i with $\alpha\vec{X}_i$, for an $\alpha \in \mathbb{R}^+ \setminus \{0\}$,
- (ii) Replace all \vec{X}_i with $\vec{X}_i + \vec{v}$, for a vector $\vec{v} \in \mathbb{R}^D$,
- (iii) Replace all \vec{X}_i with $U \cdot \vec{X}_i$, where U is an orthogonal $D \times D$ matrix (this additionally includes mirror symmetries).

Briefly explain why you obtain a statement for all possible rotations (e.g. rotating around arbitrary fixed points) by proving statements (i)-(iii), despite the fact that multiplying with U as in (iii) always leaves the origin (i.e., the zero vector) fixed (and only the origin in general).

Exercise 2: Lagrange Multipliers (30 P)

In the second paragraph on page 3, it is stated that finding the optimal W_{ij} is a least-squares problem, which is shown to have an explicit analytic solution in Appendix A. In the following, assume the notation of Appendix A. For abbreviation (and clarity of notation), additionally write $w = (w_1, \dots, w_K)^\top$ for the weight vector which is optimized, $\eta = (\eta_1, \dots, \eta_K)^\top$ for the $(K \times D)$ -matrix of nearest neighbors of \vec{x} , $\mathbf{1} = (1, \dots, 1)^\top$ for the K -dimensional vector of ones, and $C = (\mathbf{1}\vec{x}^\top - \eta)(\mathbf{1}\vec{x}^\top - \eta)^\top$ for the local covariance matrix at \vec{x} . We would like to work out the following claims from Appendix A:

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

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

In particular, prove equation (3) on page 9.

- (ii) *Show* by using the Lagrangian method for constrained optimization that the minimum of the optimization problem is explicitly given by

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

- (iii) *Show* that the minimum 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: Kullback-Leibler Divergence (30 P)

The objective of t-SNE is based on minimization of the Kullback-Leibler divergence between two probability distributions p and q .

$$C = D_{\text{KL}}(P||Q) = \sum_j p_j \log \left(\frac{p_j}{q_j} \right)$$

where $\sum_j p_j = 1$ and $\sum_j q_j = 1$. Minimization of such quantity also intervenes in various probabilistic machine learning models. In this exercise, we derive the gradient of the Kullback-Leibler divergence, both with respect to the probability distribution itself, and a reparameterization of it.

- *Show* that

$$\frac{\partial C}{\partial q_i} = -\frac{p_i}{q_i}.$$

- The probability q_i that has to be optimized for all i is now reparameterized as $q_i = \frac{e^{x_i}}{\sum_k e^{x_k}}$. Here, x_i can be interpreted as the unnormalized log-probability associated to q_i . *Show* that

$$\frac{\partial C}{\partial x_i} = -p_i + q_i$$

- *Explain* which of these two gradients is the most appropriate for practical use in a learning algorithm. Motivate your choice (1) in terms of stability or boundedness of the gradient, and (2) in terms of ability to produce a valid probability distribution.

Task 1**i**

$$\begin{aligned}
\epsilon(W) &= \sum_i |\vec{X}_i - \sum_j W_{ij} \vec{X}_j|^2 \\
\tilde{\epsilon}(W) &= \sum_i |\alpha \vec{X}_i - \sum_j W_{ij} \alpha \vec{X}_j|^2 \\
\tilde{\epsilon}(W) &= \sum_i |\alpha (\vec{X}_i - \sum_j W_{ij} \vec{X}_j)|^2 \\
\tilde{\epsilon}(W) &= \sum_i \alpha^2 |\vec{X}_i - \sum_j W_{ij} \vec{X}_j|^2 \\
\tilde{\epsilon}(W) &= \alpha^2 \sum_i |\vec{X}_i - \sum_j W_{ij} \vec{X}_j|^2
\end{aligned}$$

Since α^2 in $\tilde{\epsilon}(W)$ does not influence the weights of the maximum we get the same result as in $\epsilon(W)$

ii

$$\begin{aligned}
\epsilon(W) &= \sum_i |\vec{X}_i - \sum_j W_{ij} \vec{X}_j|^2 \\
\tilde{\epsilon}(W) &= \sum_i |(\vec{X}_i + \vec{v}) - \sum_j W_{ij} (\vec{X}_j + \vec{v})|^2 \\
\tilde{\epsilon}(W) &= \sum_i |(\vec{X}_i + \vec{v}) - (\sum_j W_{ij} \vec{X}_j) - (\sum_j W_{ij} \vec{v})|^2 \\
\tilde{\epsilon}(W) &= \sum_i |(\vec{X}_i + \vec{v}) - (\sum_j W_{ij} \vec{X}_j) - \vec{v} (\sum_j W_{ij})|^2 \\
\tilde{\epsilon}(W) &= \sum_i |\vec{X}_i + \vec{v} - (\sum_j W_{ij} \vec{X}_j) - \vec{v}|^2 \\
\tilde{\epsilon}(W) &= \sum_i |\vec{X}_i - (\sum_j W_{ij} \vec{X}_j)|^2 = \epsilon(W)
\end{aligned}$$

iii

$$\begin{aligned}
\epsilon(W) &= \sum_i |\vec{X}_i - \sum_j w_{ij} \vec{X}_j|^2 \\
\tilde{\epsilon}(W) &= \sum_i |U\vec{X}_i - \sum_j w_{ij} U\vec{X}_j|^2 \\
\tilde{\epsilon}(W) &= \sum_i |U\vec{X}_i - U(\sum_j w_{ij} \vec{X}_j)|^2 \\
\tilde{\epsilon}(W) &= \sum_i |U(\vec{X}_i - \sum_j w_{ij} \vec{X}_j)|^2 \\
\tilde{\epsilon}(W) &= \sum_i |\vec{X}_i - \sum_j w_{ij} \vec{X}_j|^2 = \epsilon(W)
\end{aligned}$$

In the last step we used the fact that an orthogonal matrix does not change the length of a vector because

$$|U\vec{v}|^2 = \langle U\vec{v}, U\vec{v} \rangle = \langle \vec{v}, U^T U \vec{v} \rangle = \langle \vec{v}, \vec{v} \rangle = |\vec{v}|^2$$

General Rotations: We get a rotation around arbitrary points by combining translation and rotation. If we want to rotate around point A , we can first do a translation so that point A becomes the origin. Then we rotate around the origin and after that we reverse the translation.

Task 2**i**

$$\begin{aligned}
\epsilon &= |\vec{x} - \sum_j w_j \vec{\eta}_j|^2 \\
\epsilon &= \left| \underbrace{\sum_j w_j \vec{x}}_{=1} - \sum_j w_j \vec{\eta}_j \right|^2 \\
\epsilon &= \left| \sum_j w_j (\vec{x} - \vec{\eta}_j) \right|^2 \\
\epsilon &= \left(\sum_j w_j (\vec{x} - \vec{\eta}_j) \right) \left(\sum_k w_k (\vec{x} - \vec{\eta}_k) \right) \\
\epsilon &= \sum_j \sum_k w_j w_k (\vec{x} - \vec{\eta}_j) (\vec{x} - \vec{\eta}_k) \\
\epsilon &= \sum_{j,k} w_j w_k C_{jk} \\
\epsilon &= \sum_j w_j \sum_k C_{jk} w_k \\
\epsilon &= \sum_j w_j (Cw)_j
\end{aligned}$$

where $(Cw)_j = \sum_k C_{jk} w_k$ and $C_{jk} = (\vec{x} - \vec{\eta}_k)(\vec{x} - \vec{\eta}_j) = (\vec{x} - \vec{\eta}_j)(\vec{x} - \vec{\eta}_k) = C_{kj}$

$$\epsilon = \vec{w}^T C \vec{w}$$

ii

$$\begin{aligned}
L(w, \lambda) &= w^T C w + \lambda(w^T \mathbf{1} - 1) \\
\frac{\partial L}{\partial \lambda} &= w^T \mathbf{1} - 1 \\
\frac{\partial L}{\partial w_i} &= \frac{\partial}{\partial w_i}
\end{aligned}$$

iii

Task 3**i**

$$\begin{aligned}
C &= \sum_j p_j \log\left(\frac{p_j}{q_j}\right) \\
\frac{\partial C}{\partial q_i} &= p_i * \frac{\partial}{\partial q_i} \log\left(\frac{p_i}{q_i}\right) \\
\frac{\partial C}{\partial q_i} &= p_i * \frac{q_i}{p_i} * \left(-\frac{p_i}{q_i^2}\right) \\
\frac{\partial C}{\partial q_i} &= -\frac{p_i}{q_i}
\end{aligned}$$

ii

$$\frac{\partial C}{\partial x_i} = \frac{\partial C}{\partial \vec{q}} \frac{\partial \vec{q}}{\partial x_i}$$

As we just have shown

$$\left(\frac{\partial C}{\partial \vec{q}}\right)_j = \frac{\partial C}{\partial q_j} = -\frac{p_j}{q_j} \quad .$$

For the second factor we get

$$\begin{aligned}
\left(\frac{\partial \vec{q}}{\partial x_i}\right)_j &= \frac{\partial q_j}{\partial x_i} = \frac{\partial}{\partial x_i} \frac{e^{x_j}}{\sum_k e^{x_k}} \\
&= \frac{e^{x_j} \sum_k e^{x_k} \delta_{ij} - e_i^x e_j^x}{\left[\sum_k e^{x_k}\right]^2} \\
&= \frac{e^{x_j} \delta_{ij}}{\sum_k e^{x_k}} - \frac{e^{x_i}}{\sum_k e^{x_k}} \frac{e^{x_j}}{\sum_k e^{x_k}} \\
&= q_j \delta_{ij} - q_i q_j \\
\Rightarrow \frac{\partial C}{\partial \vec{q}} \frac{\partial \vec{q}}{\partial x_i} &= \sum_j \left(-\frac{p_j}{q_j}\right) (q_j \delta_{ij} - q_i q_j) \\
&= \sum_j \left(-\frac{p_j q_j \delta_{ij}}{q_j}\right) + \left(\frac{p_j q_i q_j}{q_j}\right) \\
&= \sum_j (-p_j \delta_{ij}) + (p_j q_i) \\
&= \sum_j p_j q_i - \sum_j p_j \delta_{ij} \\
&= q_i \sum_j p_j - p_i = q_i - p_i
\end{aligned}$$

iii

Stability/Boundedness: The first derivative has q_i in the denominator which can become close to zero. In that case the gradient would go to infinity. The second derivative only consists of a sum which performs better in these aspects.

Validity: To perform a gradient descent we can only vary the projected points x_i so it's more meaningful to calculate the derivative with respect to x_i .