STOCKHOLM UNIVERSITY     MT 7050
DEPT. OF MATHEMATICS     EXAMINATION
Div. of Mathematical statistics     20 Feb 2024

# Re-Exam in Unsupervised Learning
## 20 Feb 2024, time 08:00-13:00

*Examinator:* Chun-Biu Li, cbli@math.su.se.
*Permitted aids:* When writing the exam, you may use any literature. Electronic devices are NOT allowed

---

NOTE: The exam consists of 4 problems with 100 points in total. Logical explanation and steps leading to the final solution must be clearly shown in order to receive full marks.

NOTE: Your answers and explanations must be to the point, **redundant writing irrelevant to the solution will result in point deduction**.

---

## Problem 1 (Basics of unsupervised learning, total 29p)

a) Consider the Gaussian mixture model (GMM) in the book "Pattern recognition and machine leanring", and suppose that the covariance matrices of all mixture components are given by $\epsilon \mathbf{I}$ such that the probability distribution function of the $k$-th Gaussian component is given by Eq. 9.41 in the book (Note: Eq. 9.41 has a small typo!). Under this setting, show that, in the limit $\epsilon \to 0$, maximizing the GMM log-likelihood Eq. 9.14 equals to minimizing the $K$-means objective function Eq. 9.1. **(18p)**

b) Show that the principal coordinates $\hat{X}_{\mathrm{MDS}} = I_{p \times N} \Lambda_{\mathrm{MDS}}^{1/2} U^\top$ is centered. **(5p)**

c) PCA and classical metric MDS are equivalent when the Euclidean distances are used. State explicitly where in PCA **(3p)** and in classical metric MDS **(3p)** the assumption of Euclidean distance is imposed. Note: Please state ONLY the relevant parts in PCA and classical metric MDS.

## Problem 2 (Graph based methods, total 30p)

For graphs with a single connected componet, the commute time distances (CTD), $c_{ij}$, expressed in terms of the eigen-values $\lambda_\alpha$ and -vectors $v_{\alpha i}$ of the normalized graph Laplacian $L_{sym}$, $c_{ij} = \mathrm{vol}(G) \sum_{\alpha=2}^{N} \frac{1}{\lambda_\alpha} \left( \frac{v_{\alpha i}}{\sqrt{d_i}} - \frac{v_{\alpha j}}{\sqrt{d_j}} \right)^2$, has the form of squared Euclidean distance, where $\mathrm{vol}(G)$ is volume of the graph, $d_i$ is the degree of the $i$-th node, with $i = 1, \cdots, N$ and $\alpha = 2, \cdots, N$. This suggests that one can embed the data points in a Euclidean space with the Cartesian coordinates $x_{\alpha i} = v_{\alpha i} \sqrt{\frac{\mathrm{vol}(G)}{\lambda_\alpha d_i}}$, called the CTD embedding. Here $\alpha$ labels the directions and $i$ labels the data point.

a) Show that $E(x_\alpha) = 0$ for $\alpha > 1$ with the weight of each data given by $P(i) = d_i/\text{vol}(G)$. **(10p)**

b) With the same weights in part a, find the covariance matrix $E(x_\alpha x_{\alpha'})$ for $\alpha, \alpha' > 1$. **(10p)**

c) Since CTD depends on the volume of the graph $\text{vol}(G)$, bigger graphs with more nodes and connections have larger values of CTD. This may be a problem if one wants to compare two graphs with the same statistical properties (e.g. both of them are random graphs) but with different number of nodes and connections. Propose a way to modify the graph distance in terms of CTD that is not sensitive to the graph and at the same time keeping the graph distance invariant under rescaling of graph weights as in Past c. Justify your answer. **(5p)**.

d) Draw one example where the mutual $k$NN graph construction with single connected component may end up with very large $k$ **(2p)**, then propose a solution for it **(3p)**.

## Problem 3 (Local linear embedding, total 27p)

This problem follows the notation in the paper "Nonlinear dimensionality reduction by locally linear embedding".

a) Show that the weights $W_{ij}^{min}$ that minimize the cost function $\epsilon(W) = \sum_i \left| \overrightarrow{X}_i - \sum_j W_{ij} \overrightarrow{X}_j \right|^2$ (i.e., Eq. 1 in the paper) subject to the constraints $\sum_j W_{ij} = 1$ are invariant under orthogonal transformation **(3p)** and rescaling **(3p)** of the data coordinates $\overrightarrow{X}_i$, $i = 1, \cdots, N$.

b) Consider the constrained least squares problem in solving the weights for a given data point $\overrightarrow{X}$, one minimizes $\epsilon(W) = \left| \overrightarrow{X} - \sum_j W_j \overrightarrow{\eta}_j \right|^2$ subject to $\sum_j W_j = 1$ where $\overrightarrow{\eta}_j$ are neighbors of $\overrightarrow{X}$. Show that the cost function $\epsilon(W)$ can be written as the quadratic form $\epsilon(W) = \sum_{j,k} W_j C_{jk} W_k$, where the scalar product matrix is defined by $C_{jk} = (\overrightarrow{X} - \overrightarrow{\eta}_j) \cdot (\overrightarrow{X} - \overrightarrow{\eta}_k)$. **(5p)**

c) Discuss in what situation that the matrix $C^{-1}$ in part b does not exist (i.e., when $C$ is a singular matrix) and what are the implications in terms of the intrinsic dimension of the data structure locally around $\overrightarrow{X}$. **(6p)**

d) Consider the eigenvector problem where the $N \times N$ weight matrix $W$ is given, one minimizes $\phi(Y) = \sum_i \left| \overrightarrow{Y}_i - \sum_j W_{ij} \overrightarrow{Y}_j \right|^2$ subject to the constraints $\sum_i \overrightarrow{Y}_i = 0$ and $\sum_i \overrightarrow{Y}_i \overrightarrow{Y}_i^T = NI$. Show that the cost function $\phi(Y)$ can be written as $\phi(Y) = \text{Tr}(Y^T M Y)$ where $M = (I - W)^T (I - W)$ and $Y$ is the $N \times d$ data matrix in the lower dimensional space. **(10p)**

## Problem 4 (Validation Methods, total 14p)

This problem refers to the lecture note on validation methods

a) Name TWO limitations of the Silhouette plot and coefficent to validate clustering results and propose solutions for each of them. **(6p)**

b) Consider the following scenario: There are $N$ data points. One assumes that the number of clusters equal 2 and exactly $N/2$ data points are randomly picked and assigned to cluster 1, the rest of the $N/2$ data points are then assigned to cluster 2. What is the value of the Silhouette coefficient for this assignment? Please show your argument clearly. **(5p)**

c) Draw a figure to show an example when the Silhouette coefficient of a data point $s(i)$ can approach $-1$. **(3p)**

*Good Luck!*