

Some Studies on unsupervised multi-view clustering methods

Project Report

Prasun De

under

Dr. Swagatam Das

Associate Professor

Electronics and Communication Sciences Unit

Indian Statistical Institute

203 B T Road

Kolkata 700108

Indian Statistical Institute

India

August 11, 2024

Abstract

We study various approaches to unsupervised clustering of multi-dimensional data using both linear and non - linear algorithms .

Introduction

Clustering is a process which partitions a given data set into homogeneous groups based on given features such that similar objects are kept in a group whereas dissimilar objects are in different groups . It deals with finding structure in a collection of unlabeled data. As one of the classical clustering algorithms , **k-means** provides an intuitive and effective way to perform clustering. In specific, the k-means clustering is composed of (i) calculating k prototypes (i.e., centres of k clusters) given an assignment of samples to clusters and (ii) updating the assignment matrix by minimizing the sum-of-squares cost given the prototypes. These two steps are alternately performed until convergence. Due to its conceptual simplicity, easy-implementation and high efficiency, k-means clustering has been intensively studied and extended ([10] ; [9]; [11]; [12]). As an important extension, **kernel k-means** first maps data onto a high-dimensional space through a feature mapping and then conducts a standard k-means clustering in that space .In many practical applications of clustering, samples are represented by multiple groups of features extracted from different information sources. For example, three kinds of feature representations (colour, shape and texture) are extracted to distinguish one flower from another. These different sources usually provide complementary information, and it is better to let learning algorithms optimally combine them in order to obtain excellent clustering. This line of research is known as **multiple kernel clustering** in the literature.

1 k-means clustering

Let $\{x_i\}_{i=1}^n \subseteq \mathcal{X}$ be the set of n data-points .The objective of the **k-means** algorithm is to minimize the sum-of-squares loss over the *cluster assignment* matrix $\mathbf{Z} \in \{0, 1\}^{n \times k}$. We can formulate it as the following optimization problem ,

$$\begin{aligned}
& \min_{Z \in \{0,1\}^{n,k}} \sum_{i=1, c=1}^{n,k} Z_{ic} |x_i - \mu_c|^2 \text{ such that } \sum_{c=1}^k Z_{ic} = 1 \text{ where} \\
& \sum_{i=1}^n Z_{ic} = n_c, \text{ and } \mu_c = \frac{\sum_{i=1}^n Z_{ic} x_i}{n_c} \text{ for } c = 1, \dots, k
\end{aligned} \tag{1}$$

We use the following algorithm to obtain **local optima** of the cost function . We have assumed *Lloyd's conditions* , that is , we will never encounter empty clusters during any stage of the procedure . Also note that k is either known , or we use G-means algorithm to obtain the most likely value of k . The k -means algorithm is highly sensitive to the initial (random) choice of the centroid matrix $\vec{\mu} = [\mu_1 : \mu_2 : \dots : \mu_k]$.

We shall describe the k means algorithm applied on the data-set X , with known number of clusters k .

Algorithm 1: kmeans(X , k)

Input: X , k

Output: Z , $\vec{\mu}$

- 1 Randomly initialize $\vec{\mu} = [\mu_1 : \mu_2 : \dots : \mu_k]$
 - 2 Given $\vec{\mu}$ obtain Z by setting $Z_{ic} = 1$ if $c = \operatorname{argmin} \{d(\mu_j, x_i)\}_{j=1}^n$ and 0 for other values of c (given i)
 - 3 Given Z obtain $\vec{\mu}$ by setting $\mu_c = \frac{\sum_{i=1}^n Z_{ic} x_i}{n_c}$ where $n_c = \sum_{i=1}^n Z_{ic}$ for $c = 1, \dots, k$
 - 4 Repeat Steps 2 and 3 until convergence of $\vec{\mu}$
-

2 Learning the k in k -means

When clustering a dataset, the right number k of clusters to use is often not obvious, and choosing k automatically is a difficult problem. The **G-means** algorithm [2] is an effective method to tackle this issue .

G-means repeatedly makes decisions based on a statistical test for the data assigned to each center. If the data currently assigned to a k -means center appear to be Gaussian, then we represent that data with only one center. However, if the same data do not appear to be Gaussian, then we use multiple centers to model the data properly. The algorithm will run

k -means multiple times (up to k times when finding k centers), so the time complexity is at most $O(k)$ times that of k -means.

In the following algorithm : X denotes the data matrix , α is the level of significance of the statistical test . Widely used statistical tests for this purpose are the *Anderson - Darling test* and *Kolmogorov - Smirnoff test*

Algorithm 2: Gmeans(X , α)

Input: X , α

Output: k

- 1 Let C be initial set of centers .(Usually $\{C \leftarrow \bar{x}\}$)
 - 2 $C \leftarrow k\text{-means}(C, X)$
 - 3 Let $\{x_i \text{---class}(x_i) = j\}$ be the set of datapoints assigned to center j .
 - 4 Use statistical test to detect if each $\{x_i \text{---class}(x_i) = j\}$ follow a Gaussian distribution at confidence level α
 - 5 If the data look Gaussian, keep c_j . Otherwise replace c_j with two centers.
 - 6 Repeat from step 2 until no more centers are added.
-

For the **Anderson-Darling test** , the data can be then tested for uniformity with a distance test [4] . The formula for the test statistic A to assess if data $\{Y_1 < \dots < Y_n\}$ (note that the data must be put in order) comes from a CDF F is

$$A^2 = -n - S ,$$

where

$$S = \sum_{i=1}^n \frac{2i-1}{n} [\ln(F(Y_i)) + \ln(1 - F(Y_{n+1-i}))] .$$

The test statistic can then be compared against the critical values of the theoretical distribution. Note that in this case no parameters are estimated in relation to the cumulative distribution function F .

3 Kernel K-means Clustering (KKM)

Let $\{x_i\}_{i=1}^n \subseteq \mathcal{X}$ be the set of n data-points , and $\phi(\cdot) : x \in \mathcal{X} \rightarrow \mathcal{H}$ be a feature mapping which maps x into a reproducing Kernel Hilbert Space \mathcal{H} .

The objective of the **k-means** algorithm is to minimize the sum-of-squares loss over the *cluster assignment* matrix $\mathbf{Z} \in \{0, 1\}^{n \times k}$. We can formulate it as the following optimization problem ,

$$\min_{\mathbf{Z} \in \{0,1\}^{n,k}} \sum_{i=1, c=1}^{n,k} Z_{ic} \|\phi(x_i) - \mu_c\|^2 \text{ such that } \sum_{c=1}^k Z_{ic} = 1 \text{ where} \quad (2)$$

$$\sum_{i=1}^n Z_{ic} = n_c, \text{ and } \mu_c = \frac{\sum_{i=1}^n Z_{ic} \phi(x_i)}{n_c} \text{ for } c = 1, \dots, k$$

We will now deduce an equivalent optimization problem from this . Let us define the **kernel matrix** \mathbf{K} by $K_{ij} = \phi(x_i)^T \phi(x_j)$ for $i, j \in \{1, \dots, n\}$. Now , we observe that ,

$$\begin{aligned} & \sum_{i=1, c=1}^{n,k} Z_{ic} \|\phi(x_i) - \mu_c\|^2 \\ &= \sum_{i=1, c=1}^{n,k} Z_{ic} (\phi(x_i) - \mu_c)^T (\phi(x_i) - \mu_c) \\ &= \sum_{i=1, c=1}^{n,k} Z_{ic} (\phi(x_i)^T - \mu_c^T) (\phi(x_i) - \mu_c) \\ &= \sum_{i=1, c=1}^{n,k} Z_{ic} (\phi(x_i)^T \phi(x_i) - \phi(x_i)^T \mu_c - \mu_c^T \phi(x_i) + \mu_c^T \mu_c) \\ &= \sum_{i=1}^n \left(\sum_{c=1}^k Z_{ic} \right) \phi(x_i)^T \phi(x_i) - \sum_{i=1, c=1}^{n,k} Z_{ic} (\phi(x_i)^T \mu_c + \mu_c^T \phi(x_i) - \mu_c^T \mu_c) \\ &= \sum_{i=1}^n K_{ii} - \sum_{i=1, c=1}^{n,k} Z_{ic} (\phi(x_i)^T \mu_c + \mu_c^T \phi(x_i) - \mu_c^T \mu_c) \\ &= \sum_{i=1}^n K_{ii} - \sum_{c=1}^k \sum_{i=1}^n Z_{ic} (\phi(x_i)^T \mu_c + \mu_c^T \phi(x_i) - \mu_c^T \mu_c) \end{aligned}$$

Now , observe the second term in the summation on the previous line . We shall first simplify the inner-most summation that runs over i . For this , we first observe a technical result :

If \mathbf{A} is a $k \times n$ matrix, and \mathbf{B} is a $n \times n$ matrix then :

$$(\mathbf{ABA}^T)_{ij} = \sum_{s=1}^n \sum_{t=1}^n a_{is} b_{st} a_{jt}$$

This can be seen easily by just expanding the LHS of the claimed equality . Now, we use this fact to proceed in our calculations :

$$\begin{aligned}
& \sum_{i=1}^n \mathbf{Z}_{ic} (\phi(x_i)^T \mu_c + \mu_c^T \phi(x_i) - \mu_c^T \mu_c) \\
&= \sum_{i=1}^n \mathbf{Z}_{ic} (\phi(x_i)^T \mu_c + \mu_c^T \phi(x_i) - \mu_c^T \mu_c) - \frac{\left(\sum_{j=1}^n \mathbf{Z}_{jc} \phi(x_j)^t \right) \left(\sum_{j=1}^n \mathbf{Z}_{jc} \phi(x_j) \right)}{n_c^2} \\
&= \sum_{i=1}^n \mathbf{Z}_{ic} \left[\phi(x_i)^T \cdot \frac{\sum_{j=1}^n \mathbf{Z}_{jc} \phi(x_j)}{n_c} + \frac{\sum_{j=1}^n \mathbf{Z}_{jc} \phi(x_j)}{n_c} \cdot \phi(x_i) \right] \\
&\quad - \frac{\left(\sum_{j=1}^n \mathbf{Z}_{jc} \phi(x_j)^T \right) \left(\sum_{j=1}^n \mathbf{Z}_{jc} \phi(x_j) \right)}{n_c^2} \\
&= \sum_{i=1}^n \mathbf{Z}_{ic} \left[\sum_{j=1}^n \frac{\mathbf{Z}_{jc} K(x_i, x_j)}{n_c} \right] + \sum_{i=1}^n \mathbf{Z}_{ic} \left[\sum_{j=1}^n \frac{\mathbf{Z}_{jc} K(x_i, x_j)}{n_c} \right] \\
&\quad - \sum_{i=1}^n \mathbf{Z}_{ic} \sum_{j=1}^n \sum_{s=1}^n \frac{\mathbf{Z}_{sc} \mathbf{Z}_{jc} K(x_j, x_s)}{n_c^2} \\
&= 2 \cdot \sum_{i=1}^n \sum_{j=1}^n \frac{\mathbf{Z}_{ic} \mathbf{Z}_{jc} K(x_i, x_j)}{n_c} - \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^n \frac{\mathbf{Z}_{ic} \mathbf{Z}_{jc} \mathbf{Z}_{sc} K(x_j, x_s)}{n_c^2} \\
&= 2 \cdot \sum_{i=1}^n \sum_{j=1}^n \frac{\mathbf{Z}_{ic} \mathbf{Z}_{jc} K(x_i, x_j)}{n_c} - \sum_{j=1}^n \sum_{s=1}^n \sum_{i=1}^n \frac{\mathbf{Z}_{ic}}{n_c^2} \cdot \mathbf{Z}_{jc} \mathbf{Z}_{sc} K(x_j, x_s)
\end{aligned}$$

$$= \sum_{i=1}^n \sum_{j=1}^n \frac{Z_{ic} Z_{jc} K(x_i, x_j)}{n_c}$$

Now , observe that :

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^n \frac{Z_{ic} Z_{jc} K(x_i, x_j)}{n_c} &= \sum_{i=1}^n \sum_{j=1}^n \frac{Z_{ic}}{\sqrt{n_c}} \cdot K(x_i, x_j) \cdot \frac{Z_{jc}}{\sqrt{n_c}} \\ &= (AKA^T)_{cc} \end{aligned}$$

Here , \mathbf{A} is a $k \times n$ matrix such that : $A_{ij} = \frac{Z_{ji}}{\sqrt{n_i}}$. Here , we have used the technical result that we had mentioned earlier . So , we get that :

$$\sum_{i=1, c=1}^{n, k} Z_{ic} \|\phi(x_i) - \mu_c\|^2 = Tr(K) - \sum_{c=1}^k (AKA^T) \quad (*)$$

Now , let $\mathbf{L} = \begin{bmatrix} \frac{1}{n_1} & 0 & 0 & \dots & 0 \\ 0 & 0 & \frac{1}{n_2} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \frac{1}{n_k} \end{bmatrix}$ Then , we observe that :

$$\begin{aligned} ZL^{1/2} &= \begin{bmatrix} z_{11} & z_{12} & z_{13} & \dots & z_{1k} \\ z_{21} & z_{22} & z_{23} & \dots & z_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ z_{n1} & x_{n2} & x_{n3} & \dots & x_{nk} \end{bmatrix} \times \begin{bmatrix} \frac{1}{\sqrt{n_1}} & 0 & 0 & \dots & 0 \\ 0 & \frac{1}{\sqrt{n_2}} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \frac{1}{\sqrt{n_k}} \end{bmatrix} \\ &= \begin{bmatrix} \frac{z_{11}}{\sqrt{n_1}} & \frac{z_{12}}{\sqrt{n_2}} & \frac{z_{13}}{\sqrt{n_3}} & \dots & \frac{z_{1k}}{\sqrt{n_k}} \\ \frac{z_{21}}{\sqrt{n_1}} & \frac{z_{22}}{\sqrt{n_2}} & \frac{z_{23}}{\sqrt{n_3}} & \dots & \frac{z_{2k}}{\sqrt{n_k}} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{z_{n1}}{\sqrt{n_1}} & \frac{z_{n2}}{\sqrt{n_2}} & \frac{z_{n3}}{\sqrt{n_3}} & \dots & \frac{z_{nk}}{\sqrt{n_k}} \end{bmatrix} \end{aligned}$$

So , we can take $\mathbf{A} = (ZL^{1/2})^T$. Since \mathbf{L} is symmetric , $\mathbf{A} = L^{1/2} Z^T$. Plugging this value of \mathbf{A} in $*$, we find that the objective function is equivalent to :

$$Tr(K) - \sum_{c=1}^k \left(L^{1/2} Z^T K Z L^{1/2} \right)_{cc} = Tr(K) - Tr(L^{1/2} Z^T K Z L^{1/2})$$

We can say that the minimization problem :

$$\min_{Z \in \{0,1\}^{n,k}} \sum_{i=1, c=1}^{n,k} Z_{ic} \|\phi(x_i) - \mu_c\|^2 \text{ such that } \sum_{c=1}^k Z_{ic} = 1 \text{ where}$$

$$\sum_{i=1}^n Z_{ic} = n_c, \text{ and } \mu_c = \frac{\sum_{i=1}^n Z_{ic} \phi(x_i)}{n_c} \text{ for } c = 1, \dots, k$$

is equivalent to :

The optimization problem can thus be rewritten as :

$$\min_{Z \in \{0,1\}^{n,k}} \text{Tr}(\mathbf{K}) - \text{Tr}(\mathbf{L}^{1/2} \mathbf{Z}^T \mathbf{K} \mathbf{Z} \mathbf{L}^{1/2}) \text{ s.t. } \mathbf{Z} \mathbf{1}_k = \mathbf{1}_n \quad (3)$$

Here , \mathbf{K} is the kernel matrix , $\mathbf{L} = \text{diag}[n_1^{-1}, \dots, n_k^{-1}]$, and $\mathbf{1}_l \in \mathbf{R}^l$ is a column vector with all entries 1 . The variable \mathbf{Z} is discrete, which makes the optimization problem very difficult to solve. However, we can approximate this problem through relaxing \mathbf{Z} to take arbitrary real values. Specifically, by defining $\mathbf{H} = \mathbf{Z} \mathbf{L}^{1/2}$ and letting \mathbf{H} take real values , we obtain a relaxed version of the above problem .

$$\min_{H \in \mathbf{R}^{n \times k}} \text{Tr}(\mathbf{K}(\mathbf{I}_n - \mathbf{H} \mathbf{H}^T)) \text{ s.t. } \mathbf{H}^T \mathbf{H} = \mathbf{I}_k \quad (4)$$

where \mathbf{I}_k is the $k \times k$ identity matrix . Noting that $\mathbf{Z}^T \mathbf{Z} = \mathbf{L}^{-1}$, we have $\mathbf{L}^{1/2} \mathbf{Z}^T \mathbf{Z} \mathbf{L}^{1/2} = \mathbf{I}_k$. This leads us to the orthogonality constraint on \mathbf{H} . Finally, one can **obtain the optimal \mathbf{H} by taking the k eigenvectors that correspond to the k largest eigenvalues of \mathbf{K}** . We shall show it with the help of the following machinery that was developed in [2] .

Lemma from [8]

Let $\mathbf{y} \in \mathbf{R}^q$, $\mathbf{B}^T \in \mathbf{R}^{q \times p}$, such that $1 \leq q \leq p$. Consider the orthogonal linear transformation $\mathbf{x} \mapsto \mathbf{y} = \mathbf{B}^T \mathbf{x}$. Let $\Sigma_{\mathbf{y}} = \mathbf{B}^T \Sigma \mathbf{B}$ be the variance - covariance matrix for \mathbf{y} . Then, $\text{Tr}(\Sigma_{\mathbf{y}})$ is maximized by taking $\mathbf{B} = \mathbf{A}_q$, \mathbf{A}_q consists of the first q columns of \mathbf{A} . Note that \mathbf{A} is an orthogonal matrix whose k -th column is the k -th eigenvector of Σ .

Proof of Lemma :

Let β_k be the k -th column of \mathbf{B} . Now, the columns of \mathbf{A} form a basis for \mathbf{R}^p . So, there exists $\{c_{jk}\}_{j,k=1,1}^{j,k=p,q}$ such that $\beta_k = \sum_{j=1}^p c_{jk} \alpha_j$ for $k = 1, \dots, q$. So, we may say that $\mathbf{B} = \mathbf{A}\mathbf{C}$. So,

$$\begin{aligned} \mathbf{B}^T \Sigma \mathbf{B} &= \mathbf{C}^T \mathbf{A}^T \Sigma \mathbf{A} \mathbf{C} \\ &= \mathbf{C}^T \Lambda \mathbf{C} \\ &= \sum_{j=1}^p \lambda_j \mathbf{c}_j \mathbf{c}_j^T \end{aligned}$$

Here, \mathbf{c}_j is the j -th row of \mathbf{C}

$$\begin{aligned} \text{Tr}(\mathbf{B}^T \Sigma \mathbf{B}) &= \sum_{j=1}^p \lambda_j \text{Tr}(\mathbf{c}_j \mathbf{c}_j^T) \\ &= \sum_{j=1}^p \lambda_j \text{Tr}(\mathbf{c}_j^T \mathbf{c}_j) \\ &= \sum_{j=1}^p \lambda_j \mathbf{c}_j^T \mathbf{c}_j \\ &= \sum_{j=1}^p \sum_{k=1}^q \lambda_j c_{jk}^2 \dots (i) \end{aligned}$$

Now, $\mathbf{C} = \mathbf{A}^T \mathbf{B}$, and the orthogonality of \mathbf{A} with orthonormal property of the columns of \mathbf{B} implies that $\mathbf{C}^T \mathbf{C} = \mathbf{B}^T \mathbf{A} \mathbf{A}^T \mathbf{B} = \mathbf{B}^T \mathbf{B} = \mathbf{I}_q$. So,

$\sum_{j=1}^q \sum_{k=1}^q c_{jk}^2 = q \dots (ii)$ and the columns of \mathbf{C} are also orthonormal . \mathbf{C} can be thought of as the first q columns of a $p \times p$ orthogonal matrix , say , \mathbf{D} . Rows of \mathbf{D} are orthonormal , so satisfy that $\mathbf{d}_j^T \mathbf{d}_j = 1$ for $j = 1, \dots, p$. As rows of \mathbf{C} consist of first q elements of rows of \mathbf{D} , $\mathbf{c}_j^T \mathbf{c}_j \leq 1$. So , $\sum_{k=1}^q c_{jk}^2 \leq 1 \dots (iii)$. Now , $\sum_{k=1}^q c_{jk}^2$ is coefficient of λ_j in (i) . So , the sum of these coeff. is q from (ii) . None of these can exceed q from (iii) . Now , since $\lambda_1 \geq \dots \geq \lambda_p$, it is clear that $\sum_{j=1}^p (\sum_{k=1}^q c_{jk}^2) \lambda_j$ will be maximized if we can find a set of c_{jk} for which :

$$\sum_{k=1}^q c_{jk}^2 = \begin{cases} 1 & j = 1, \dots, q \\ 0 & j = q + 1, \dots, p \end{cases}$$

If $\mathbf{B}^T = \mathbf{A}_q^T$, then

$$c_{jk} = \begin{cases} 1 & 1 \leq j = k \leq q \\ 0 & \text{elsewhere} \end{cases}$$

So , $\text{Tr}(\Sigma_y)$ achieves maximum when $\mathbf{B}^T = \mathbf{A}_q^T$
This ends the proof of the lemma . ■

Now , we observe that solving $\min_{\mathbf{H} \in \mathbb{R}^{n \times k}} \text{Tr}(\mathbf{K}(\mathbf{I}_n - \mathbf{H}\mathbf{H}^T))$ s.t. $\mathbf{H}^T \mathbf{H} = \mathbf{I}_k$ is equivalent to solving $\text{Tr}(\mathbf{K}) - \max_{\mathbf{H} \in \mathbb{R}^{n \times k}} \text{Tr}(\mathbf{K}\mathbf{H}\mathbf{H}^T)$ s.t. $\mathbf{H}^T \mathbf{H} = \mathbf{I}_k$. Now , by property of trace of a matrix , for two matrices \mathbf{A}, \mathbf{B} such that both \mathbf{AB} and \mathbf{BA} are defined , we have $\text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA})$. As \mathbf{K} is a fixed kernel matrix , it suffices to solve $\max_{\mathbf{H} \in \mathbb{R}^{n \times k}} \text{Tr}(\mathbf{K}\mathbf{H}\mathbf{H}^T)$ s.t. $\mathbf{H}^T \mathbf{H} = \mathbf{I}_k$. Taking $\mathbf{A} = \mathbf{KH}$ and $\mathbf{B} = \mathbf{H}^T$, we have to optimize $\max_{\mathbf{H} \in \mathbb{R}^{n \times k}} \text{Tr}(\mathbf{H}^T \mathbf{K} \mathbf{H})$ s.t. $\mathbf{H}^T \mathbf{H} = \mathbf{I}_k$.

Now , we shall make use of the lemma . Take $\mathbf{B} = \mathbf{H}$, as \mathbf{K} is the analogue of the variance- covariance matrix of the projected data-points in the feature space , it suffices to **obtain the optimal \mathbf{H} by taking the k eigenvectors that correspond to the k largest eigenvalues of \mathbf{K}** .

In the following , we shall discuss the algorithm for using **KKM** to obtain a clustering of the data-points .

Algorithm 3: kkmeans(X, \mathbf{K}, k)

Input: X, \mathbf{K}, k

Output: Z

- 1 Obtain the k eigenvectors that correspond to the k largest eigenvalues of $\mathbf{K} : v_1, \dots, v_k$ and let $\mathbf{M} = [v_1 : \dots : v_k]$.
 - 2 Normalize the rows of \mathbf{M} and set \mathbf{H} to be equal to this row-normalized \mathbf{M} .
 - 3 Run **kmeans**(\mathbf{H}, k) (Algorithm 2) and obtain Z .
-

4 Multiple Kernel K-means Clustering (MKKM)

In a multiple kernel clustering situation , as the name suggests , we have **multiple feature representations via a group of feature mappings** $\{\phi(\cdot)\}_{p=1}^m$. Each sample can be represented as

$$\phi_\mu(\mathbf{x}) = [\mu_1 \phi_1(\mathbf{x})^T, \dots, \mu_m \phi_m(\mathbf{x})^T]^T.$$

, where $\mu = [\mu_1, \dots, \mu_m]^T$. denotes the coefficients of each base kernel that we need to optimize during the learning procedure . We can calculate the kernel function over the above mapping as follows :

$$K_\mu(x_i, x_j) = \phi_\mu(x_i)^T \phi_\mu(x_j) = \sum_{p=1}^m \mu_p^2 K_p(x_i, x_j).$$

We follow a procedure similar to what we did in the earlier case . Here , we are replacing \mathbf{K} with \mathbf{K}_μ . We obtain the objective function for the **MKKM** as follows :

$$\min_{H \in \mathbb{R}^{n \times k}, \mu \in \mathbb{R}_+^m} \text{Tr}(K_\mu(\mathbf{I}_n - \mathbf{H}\mathbf{H}^T)) \quad s.t. \quad \mathbf{H}^T \mathbf{H} = \mathbf{I}_k, \mu^T \mathbf{1}_m = 1$$

We can solve this optimization problem as an **alternating minimization problem** by alternately updating \mathbf{H} and μ .

(i) **Optimizing \mathbf{H} given μ** . With μ fixed , we can obtain \mathbf{K}_μ and then apply **KKM** (Algorithm 3) with input : $\{ X, \mathbf{K}_\mu, k \}$. After the first

two steps of Algorithm 3 have been completed , we return the value of \mathbf{H} . This is the updated value of \mathbf{H} .

(ii) Optimizing μ given \mathbf{H} . In this step , we have to solve the following quadratic programming problem that is subject to linear constraints.

$$\min_{\mu \in R_+^m} \sum_{p=1}^m \mu_p^2 \text{Tr}(\mathbf{K}_p(\mathbf{I}_n - \mathbf{H}\mathbf{H}^T)) \quad s.t. \quad \mu^T \mathbf{1}_m = 1$$

We observe that this is a quadratic programming problem with non-negative constraints . This is a standard optimization problems that can be attacked with methods like gradient descent . Problems of this type have been discussed in [1] .

It has been seen that using a convex combination of the kernels , such as $\sum_{p=1}^m \mu_p \mathbf{K}_p$, often results in only kernel getting a nonzero weight and all others getting a zero weight . So , such a combination is not viable as it does not comply with the multi-view idea that is being deployed to ensure a better clustering of the data . Such situations have been discussed in [3] . Now , we shall write down the algorithm for using **MKKM** to obtain a clustering of the data points .

Algorithm 4: mkkmeans($X, \{\mathbf{K}_p\}_{p=1}^m, k$)

Input: $X, \{\mathbf{K}_p\}_{p=1}^m, k$

Output: Z

- 1 Initialize $\mu^{(0)} = \mathbf{1}_m/m$ and $t = 1$.
 - 2 Normalize the rows of \mathbf{M} and set \mathbf{H} to be equal to this row-normalized \mathbf{M} .
 - 3 Run **kmeans**(\mathbf{H}, k) (Algorithm 2) and obtain Z .
-

5 Multiple Kernel K Means with regularization

We should adequately consider the mutual influence of these kernels when updating kernel coefficients . So, the concept of *regularization* is introduced in [7] . To reduce the redundancy and enforce the diversity of the selected kernels, we need a regularization term that is able to characterise the correlation of each pair of kernels. To this end, we consider criterion $\mathcal{M}(\mathbf{K}_p, \mathbf{K}_q)$

to measure the correlation between K_p and K_q . A larger $\mathcal{M}(K_p, K_q)$ means high correlation between K_p and K_q , and a smaller one implies that their correlation is low. Therefore, a natural optimization criterion to prevent two highly correlated kernels from being selected can be defined as $\mu_p \mu_q \mathcal{M}(K_p, K_q)$. As observed, by minimizing this term, the risk of simultaneously assigning μ_p and μ_q with large weights can be greatly reduced. Also, this regularization increases the probability of jointly assigning μ_p and μ_q with larger weights as long as K_p and K_q are less correlated. As a consequence, this criterion is beneficial to promote the diversity of selected kernels.

We have the following regularization term :

$$\min_{\mu \in R_+^m} \sum_{p,q=1}^m \mu_p \mu_q \mathcal{M}(K_p, K_q) = \mu^T M \mu \quad s.t. \quad \mu^T \mathbf{1}_m = 1 \quad (5)$$

By integrating the matrix-induced regularization into the objective function of existing **MKKM**, we obtain the optimization problem of the proposed algorithm as follows ,

$$\begin{aligned} \min_{H \in R^{n \times k}, \mu \in R_+^m} & \quad \text{Tr}(K_\mu(\mathbf{I}_n - \mathbf{H}\mathbf{H}^T)) + \frac{\lambda}{2} \mu^T M \mu \\ s.t. & \quad \mathbf{H}^T \mathbf{H} = \mathbf{I}_k, \mu^T \mathbf{1}_m = 1 \end{aligned} \quad (6)$$

Here , λ is parameter to trade off the clustering cost function and the regularization term.

We have a two-step algorithm to solve the optimization problem in 6 alternatively.

- (i) Optimizing \mathbf{H} with fixed μ . Given μ , the optimization in 6 w.r.t \mathbf{H} is a standard kernel k-means clustering problem, and the \mathbf{H} can be obtained by solving 4 with given K_μ ;
- (ii) Optimizing μ with fixed \mathbf{H} . Given \mathbf{H} , the optimization in 6 w.r.t μ is a quadratic programming with linear constraints. In specific, we can obtain μ by solving the following problem :

$$\min_{\mu \in R_+^m} \frac{1}{2} \mu^T (2Z + \lambda M) \mu \quad s.t. \quad \mathbf{H}^T \mathbf{H} = \mathbf{I}_k, \mu^T \mathbf{1}_m = 1 \quad (7)$$

where $\mathbf{Z} = \text{diag}([Tr(\mathbf{K}_1(\mathbf{I}_n - \mathbf{H}\mathbf{H}^T)), \dots, Tr(\mathbf{K}_m(\mathbf{I}_n - \mathbf{H}\mathbf{H}^T))])$

Algorithm 5: mkkmeans-reg($\{\mathbf{K}_p\}_{p=1}^m, k, \lambda, \epsilon_0$)

Input: $\{\mathbf{K}_p\}_{p=1}^m, k, \lambda, \epsilon_0$

Output: \mathbf{H}, μ

- 1 Initialize $\mu^{(0)} = \mathbf{1}_m/m$ and $t = 1$.
 - 2 **repeat**
 - 3 $\mathbf{K}_\mu^{(t)} = \sum_{p=1}^m (\mu_p^{(t-1)})^2 \mathbf{K}_p$
 - 4 Update $\mathbf{H}^{(t)}$ by solving (4) with given \mathbf{K}_μ^t .
 - 5 Update $\mu^{(t)}$ by solving (7) with given \mathbf{H}^t .
 - 6 $t = t + 1$
 - 7 **until** $(obj^{(t-1)} - obj^{(t)}) / obj^{(t)} \leq \epsilon_0$
-

We have to design \mathbf{M} appropriate for a given clustering task. For example, $\mathbf{M}(\mathbf{K}_p, \mathbf{K}_q)$ could be defined according to some commonly used criteria such as Kullback-Leibler (KL) divergence [5] , maximum mean discrepancy (Gretton et al. 2006) [6] and Hilbert-Schmidt independence criteria (HSIC), to name just a few. One choice for this is $\mathbf{M}(\mathbf{K}_p, \mathbf{K}_q) = Tr(\mathbf{K}_p^T \mathbf{K}_q)$. Designing proper \mathbf{M} to satisfy various requirements of clustering tasks is interesting and worth exploring .

In this report , we have studied various unsupervised multi-view clustering methods that are used to group data points .

References

- [1] Can Li A Conjugate Gradient Type Method for the Nonnegative Constraints Optimization Problems In *Journal of Applied Mathematics* Volume 2013, Article ID 986317
- [2] Greg Hamerly and Charles Elkan. Learning the K in K-Means . In *Advances in Neural Information Processing Systems* vol 17 ,2004
- [3] Bernhard Scholkopf , Alexander Smols , Klaus-Robert Muller. Non-linear component analysis as a kernel eigenvalue problem. In *Neural Comput.* 10(5):12991319.

- [4] Shapiro, S.S . How to test normality and other distributional assumptions In *The ASQC basic references in quality control: statistical techniques 3*, pp. 178.
- [5] F. Topsøe . Some inequalities for information divergence and related measures of discrimination In *IEEE Transactions on Information Theory archive* Volume 46 Issue 4, July 2000 Pages 1602-1609
- [6] A. Gretton , K.M. Borgwardt , M.J. Rasch, B. Scholkopf , A. Smola . A Kernel Two-Sample Test In *ournal of Machine Learning Research* Vol 13 (2012) pg 723-773
- [7] Xinwang Liu, Yong Dou, Jianping Yin, Lei Wang, En Zhu Multiple Kernel k-Means Clustering with Matrix-Induced Regularization In *Thirtieth AAAI Conference on Artificial Intelligence (AAAI-16)*
- [8] I.T. Jolliffe . Principal Component Analysis , 2nd ed. , Springer (2002)
- [9] Gonen, M., and Margolin, A. A. . Localized data fusion for kernel k-means clustering with application to cancer biology. In *NIPS (2104)* , 13051313.
- [10] Lu, Y.; Wang, L.; Lu, J.; Yang, J.; and Shen, C. 2014. Multiple kernel clustering based on centered kernel alignment. In *Pattern Recognition 47(11):3656–3664* (2014) .
- [11] Cai, X.; Nie, F.; and Huang, H. 2013 . Multi-view k-means clustering on big data . In *IJCAI* .
- [12] Du, L.; Zhou, P.; Shi, L.; Wang, H.; Fan, M.; Wang, W.; and Shen, Y.-D. 2015 . Robust multiple kernel k-means clustering using l_{21} norm. In *IJCAI* , 34763482 .