

# KLIC

*Stephen Coleman*

*25/11/2019*

## Kernel learning integrative clustering

### Elevator pitch

Takes independent clusterings of the individuals (i.e. local clusterings) and combines these for a global clustering. Allows different local clusterings to contribute with different strengths to the final clustering. No specifications on the type of model used to create the original clusterings.

### Intro

Kernel Learning Integrative Clustering (KLIC) is a **sequential analysis** method of integrative clustering. This in comparison to **post-processing** or **joint** methods. Are these labels a real thing?

KLIC is an extension of Cluster-of-Cluster Analysis (COCA). If one has  $L$  different datasets of measurements for the same  $N$  individuals to which one applies independent clustering methods, COCA then turns the similarity matrices that result from this into a global clustering by combining the matrices in a method similar to Consensus Clustering (the original paper). KLIC extends this by allowing different similarity matrices to have different weights in how they contribute to the global clustering. In short, KLIC applies multiple kernel  $k$ -means clustering to similarity matrices generated for individual datasets.

To understand KLIC one must understand the following:

- COCA;
- the kernel trick;
- $k$ -means clustering; and
- multiple kernel  $k$ -means clustering.

### The kernel trick

This is a computational trick to avoid operations. It aims to do an analysis in a high-dimensional space while only considering calculations in the original space.

#### Definition: Positive definite kernel

or simply a kernel,  $\delta$ , is a symmetric map:

$$\delta : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

which for all  $x_1, \dots, x_N \in \mathcal{X}$ , the matrix  $\Delta$  defined by entries  $\Delta_{ij} = \delta(x_i, x_j)$ , is positive semi-definite.

#### Definition: Kernel matrix

or **Gram matrix**,  $\Delta$ , is the positive semi-definite matrix defined by a kernel  $\delta$  applied to data  $\mathcal{X} = (x_1, \dots, x_N)$  with entries  $\Delta_{ij} = \delta(x_i, x_j)$ .

**Definition: Feature map**

For each kernel  $\delta$  there exists a **feature map**  $\phi(\cdot)$ , which maps the original data  $\mathcal{X} = (x_1, \dots, x_N)$  to some new feature space taking values in some inner product of  $\mathcal{X}$  defined by:

$$\delta(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle.$$

Thus if one is interested in working in some feature space that can be defined in terms of inner products, one may use kernels to avoid computations in analysing the data in said space.

 **$k$ -means clustering**

$k$ -means clustering assigns  $N$  points to  $K$  different clusters. These clusters are defined by the points assigned to them - they are the centroid of the assigned points. The aim of this algorithm is to minimise the sum of all squared distance between the points and their assigned centroid.

Let  $m_k$  denote the  $k$ th centroid. Then if  $\mathbf{Z}$  is the  $N \times K$  classification matrix, with

$$z_{ik} = \begin{cases} 1 & \text{if point } x_i \text{ is assigned to cluster } k, \\ 0 & \text{else.} \end{cases}$$

In classic  $k$ -means clustering, each point can only be assigned to one cluster. We use the following notation:

$$\begin{aligned} \sum_{k=1}^K z_{ik} &= 1 \text{ for all } i \in \{1, \dots, N\} \\ N_k &= \sum_{i=1}^N z_{ik} \\ m_k &= \frac{1}{N_k} \sum_{i=1}^N z_{ik} x_i \text{ for all } k \in \{1, \dots, K\} \end{aligned}$$

We can consider this an optimisation problem.

$$\underset{\mathbf{Z}}{\operatorname{argmin}} \sum_{i=1}^N \sum_{k=1}^K z_{ik} \|x_i - m_k\|_2^2$$

If we redefine this problem in some feature space defined by the feature map  $\phi(\cdot)$ , we can use the kernel trick in this context. First we write the problem in terms of the feature map:

$$\begin{aligned} \underset{\mathbf{Z}}{\operatorname{argmin}} \sum_{i=1}^N \sum_{k=1}^K z_{ik} \|\phi(x_i) - m_k^*\|_2^2 \\ m_k^* = \frac{1}{N_k} \sum_{i=1}^N z_{ik} \phi(x_i) \end{aligned}$$

If we define the  $K \times K$  matrix  $\mathbf{L}$  with  $(k, k)$ th entries of  $N_k$  and 0's elsewhere and the gram matrix  $\mathbf{\Delta}$  as the matrix with  $(i, j)$ th entries  $\delta(x_i, x_j)$ , then, according to Gönen and Margolin (2014), this can be rephrased as a trace maximisation problem:

$$\underset{\mathbf{Z}}{\operatorname{argmax}} \operatorname{tr}(\mathbf{L}^{\frac{1}{2}} \mathbf{Z}^T \mathbf{\Delta} \mathbf{Z} \mathbf{L}^{\frac{1}{2}} - \mathbf{\Delta})$$

This is subject to the constraints:

$$\begin{aligned} \mathbf{Z} \mathbf{1}_K &= \mathbf{1}_N \\ z_{ik} &\in \{0, 1\} \end{aligned}$$

The binary variables,  $z_{ik}$ , make this problem very difficult to solve (Gönen and Margolin 2014). Relaxing this constraint and letting  $\mathbf{Z} \mathbf{L}^{\frac{1}{2}} = \mathbf{H}$  and letting  $\mathbf{H}$  take arbitrary real values constrained by  $\mathbf{H}^T \mathbf{H} = \mathbf{I}_k$  makes the optimisation problem more feasible:

$$\underset{\mathbf{H}}{\operatorname{argmax}} \operatorname{tr}(\mathbf{H}^T \mathbf{\Delta} \mathbf{H} - \mathbf{\Delta}).$$

One can solve this by performing Kernel-PCA on the Gram matrix  $\mathbf{\Delta}$  and setting  $\mathbf{H}$  to the  $k$  largest eigevalues. To fianlly acquire a clustering solution, one can normalise all rows of  $\mathbf{H}$  to be on the unit sphere and then implement  $k$ -means clustering on this normalised matrix.

## References

Gönen, Mehmet, and Adam A Margolin. 2014. “Localized Data Fusion for Kernel K-Means Clustering with Application to Cancer Biology.” In *Advances in Neural Information Processing Systems*, 1305–13.