

# FEATURE SELECTION IN CLUSTERING PROBLEM

---

## TECHNICAL UPDATE

As a statistical tool, clustering analysis has been widely applied to a variety of scientific areas such as vector quantization, data mining, image processing, statistical data analysis, etc. In this paper the clustering problem is formulated as a combination of clustering and feature selection based on Gaussian Mixture Model, which is further optimized by Expectation-maximization Algorithm.

Often, Expectation maximization Algorithm provides a general solution for the parameter estimate in density mixture models. Nevertheless, it needs to pre assign an appropriate number of density components, that is, the number of clusters. Roughly, the mixture may overfit the data if too many components are utilized, whereas a mixture with too few components may not be flexible enough to approximate the true underlying model. Subsequently, the EM almost always leads to a poor estimate result if the number of components is misspecified. Unfortunately, from the practical viewpoint, it is hard or even impossible to know the exact cluster number in advance.

I here propose a method, which learns the model parameter via maximizing a weighted likelihood. Under a specific weight design, we then introduce Rival Penalized EM (RPEM) algorithm for density mixture clustering. The RPEM algorithm learns the parameter by making mixture components compete each other at each time step. Comparing Rival Penalized EM algorithm to EM algorithm, the RPEM fades out the redundant densities from a density mixture during parameter learning process. Hence, RPEM automatically selects an appropriate number of densities.

The convergence speed of the RPEM relies on the value of the learning rate. Often, by a rule of thumb, we arbitrarily set the learning rate at a small positive constant. If the value of learning rate is assigned too small, the algorithm will converge at a very slow speed. On the contrary, if it is too large, the algorithm may even oscillate. In general, it is a nontrivial task to assign an appropriate value to the learning rate, although we can pay extra efforts to make the learning rate dynamically change over time.

## OVERVIEW OF MAXIMUM WEIGHTED LIKELIHOOD (MWL) LEARNING FRAMEWORK

Suppose that an input  $x \in R^d$  comes from the following density mixture model:

$$P(\mathbf{x} | \Theta) = \sum_{j=1}^k \alpha_j p(\mathbf{x} | \theta_j), \quad \sum_{j=1}^k \alpha_j = 1, \quad (1)$$
$$\alpha_j > 0, \quad \forall 1 \leq j \leq k,$$

Where  $\Theta$  is the parameter set of  $\{\alpha_j, \theta_j\}_{j=1}^k$ . Furthermore,  $k$  is the number of components,  $\alpha_j$  is the mixture proportion of the  $j$ th component, and  $p(x | \theta_j)$  is a multivariate probability density function of  $x$  parameterized by  $\theta_j$ . As long as we know the value of  $\Theta$ , an input  $x$  can be classified into a certain cluster via its posterior probability:

$$h(j | \mathbf{x}, \Theta) = \frac{\alpha_j p(\mathbf{x} | \theta_j)}{P(\mathbf{x} | \Theta)} \quad (2)$$

Using the winner-take-all rule, that is, assigning an input  $x$  to cluster  $c$  if  $c = \arg\max_j h(j | x, \Theta)$ , or taking its soft version which assigns  $x$  to cluster  $j$  with the probability  $h(j | x, \Theta)$ . Therefore, how to estimate the parameter set  $\Theta$ , particularly without knowing the correct value of  $k$  in advance, is a key issue in density mixture clustering.

In the MWL learning framework, the parameter set is learned via maximizing the following Weighted Likelihood (WL) cost function:

$$l(\Theta) = \omega(\Theta; \mathbf{x}) + \nu(\Theta; \mathbf{x}) \quad (3)$$

$$\omega(\Theta; \mathbf{x}) = \int \sum_{j=1}^k g(j | \mathbf{x}, \Theta) \ln [\alpha_j p(\mathbf{x} | \theta_j)] dF(\mathbf{x}),$$
$$\nu(\Theta; \mathbf{x}) = - \int \sum_{j=1}^k g(j | \mathbf{x}, \Theta) \ln h(j | \mathbf{x}, \Theta) dF(\mathbf{x}), \quad (4)$$

Where  $g(j | x, \Theta)$ 's are the designable weights satisfying two conditions:

$$(1) \quad \sum_{j=1}^k g(j | x, \Theta) = 1$$

$$(2) \quad \text{For all, } j, g(j | x, \Theta) = 0, \text{ if } h(j | x, \Theta) = 0$$

Suppose that a set of  $N$  observations denoted as  $X = \{x_1, x_2, \dots, x_N\}$ , comes from the density mixture model in (1). The empirical WL function of (3) written as  $Q(\Theta; X)$ , can be given as :

$$Q(\Theta; X) = \omega(\Theta; X) + \nu(\Theta; X) \quad (5)$$

$$\begin{aligned} \omega(\Theta; X) &= \frac{1}{N} \sum_{t=1}^N \sum_{j=1}^k g(j | x_t, \Theta) \ln [\alpha_j p(x_t | \theta_j)], \\ \nu(\Theta; X) &= -\frac{1}{N} \sum_{t=1}^N \sum_{j=1}^k g(j | x_t, \Theta) \ln h(j | x_t, \Theta). \end{aligned} \quad (6)$$

Moreover, the weights  $g(j | x_t, \Theta)$ 's have been generally designed as:

$$g(j | x_t, \Theta) = (1 + \varepsilon_t) I(j | x_t, \Theta) - \varepsilon_t h(j | x_t, \Theta) \quad (7)$$

Subsequently, under a specific weight design, the papers [1,2] have presented the RPEM to learn  $\Theta$  via maximizing the WL function of (5) using a stochastic gradient ascent method, which is able to fade out the redundant densities gradually from a density mixture. Consequently, it can automatically select an appropriate number of density components in density mixture clustering. We summarize the main steps of the RPEM in following Algorithm. The experiments have shown the superior performance of the RPEM on automatic model selection, its convergence speed, however, relies on the value of learning rate.

#### ALGORITHM - RIVAL PENALIZED EXPECTATION MAXIMIZATION ALGORITHM

Initialization : Given a specific  $k(k \geq k^*, k^*$  is the true number of clusters), initialise the parameter  $\Theta$ .

Step 1: Given the current input  $x_t$  and the parameter estimate, written as  $\Theta^{(n)}$ , compute  $h(j | x_t, \Theta^{(n)})$ 's and  $g(j | x_t, \Theta^{(n)})$ 's via (2) and (7), respectively

Step 2: Given  $h(j | x_t, \Theta^{(n)})$ 's and  $g(j | x_t, \Theta^{(n)})$ 's, we update  $\Theta$  by  $\Theta^{(n+1)} = \Theta^{(n)} + \eta(\omega_t(\Theta; x_t)/\Theta)|_{\Theta^{(n)}}$ , where  $\eta$  is a small positive learning rate.

Step 3: Let  $n=n+1$  and go to step 1 for the next iteration until  $\Theta$  is converged.

Here I have proposed a MWL learning framework from the ML, through which a new RPEM algorithm has been proposed for density mixture clustering. The RPEM learns the density parameters by making mixture components compete each other at each time step. Not only are the associated parameters of the winning density component updated to adapt to an input, but

also all rivals' parameters are penalized with the strength proportional to the corresponding posterior density probabilities.

Compared to the EM algorithm, this intrinsic rival penalization mechanism enables the RPEM to automatically select an appropriate number of densities by fading out the redundant densities from a density mixture. Therefore using RPEM instead of EM may increase the accuracy and can be considered as a technical update in the model proposed in the paper.