

SELF-PACED MIXTURE OF t DISTRIBUTION MODEL

Yang Zhang^{†,*}, Qingtao Tang^{†,*}, Li Niu[‡], Tao Dai[†], Xi Xiao[†], Shu-Tao Xia[†]

[†] Department of Computer Science and Technology, Tsinghua University, China

[‡] Department of Electrical and Computer Engineering, Rice University, U.S.

*Equal Contribution

Email: yangzhan16@mails.tsinghua.edu.cn



Abstract

Gaussian mixture model (GMM) is a powerful probabilistic model for representing the probability distribution of observations in the population. However, the fitness of Gaussian mixture model can be significantly degraded when the data contain a certain amount of outliers. Although there are certain variants of GMM (e.g., mixture of Laplace, mixture of t distribution) attempting to handle outliers, none of them can sufficiently mitigate the effect of outliers if the outliers are far from the centroids. Aiming to remove the effect of outliers further, this paper introduces a Self-Paced Learning mechanism into mixture of t distribution, which leads to Self-Paced Mixture of t Distribution Model (SPTMM). We derive an Expectation-Maximization based algorithm to train SPTMM and show SPTMM is able to screen the outliers. To demonstrate the effectiveness of SPTMM, we apply the model to density estimation and clustering. Finally, the results indicate that SPTMM outperforms other methods, especially on the data with outliers.

Main Contributions

- This is the first work of employing Self-Paced Learning (SPL) to mixture model, with the aim to effectively remove the influence of outliers.
- We propose our SPTMM method which integrates SPL with TMM, and develop an EM based algorithm to solve the corresponding optimization problem.
- Extensive experiments demonstrate the superiority of our SPTMM method for density estimation and clustering.

Related Works

Mixture of t Distribution

The t distribution is defined as follows. A p -dim random vector $\mathbf{x} \in \mathbb{R}^p$ follows the p -variate t distribution with degrees of freedom $\nu \in \mathbb{R}_+$, mean $\boldsymbol{\mu} \in \mathbb{R}^n$, and correlation matrix $\boldsymbol{\Sigma} \in \Pi(p)$ if its joint probability density function (PDF) is given by

$$t(\mathbf{x}|\nu, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{\Gamma[(\nu+p)/2]}{\Gamma(\nu/2)\nu^{p/2}\pi^{p/2}|\boldsymbol{\Sigma}|^{1/2}} \cdot \left[1 + \frac{1}{\nu}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]^{-\frac{\nu+p}{2}}$$

The mixture of t distribution model (TMM) is a linear superposition of g -component t distribution, i.e.,

$$\phi(\mathbf{x}; \boldsymbol{\Psi}) = \sum_{j=1}^g \pi_j t(\mathbf{x}; \nu_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j),$$

where π_j is the mixing coefficient of the j -th component and $\boldsymbol{\Psi} = \{\boldsymbol{\pi}, \boldsymbol{\nu}, \boldsymbol{\mu}, \boldsymbol{\Sigma}\}$, in which $\boldsymbol{\pi} = (\pi_1, \pi_2, \dots, \pi_g)^T$, $\boldsymbol{\nu} = (\nu_1, \nu_2, \dots, \nu_g)^T$, $\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_g)$, and $\boldsymbol{\Sigma} = (\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2, \dots, \boldsymbol{\Sigma}_g)$.

Given the dataset $\mathcal{D} = \{\mathbf{x}_i | i = 1, 2, \dots, n\}$, where $\mathbf{x}_i \in \mathbb{R}^p$ denotes a p -dim sample, the model parameters of TMM $\boldsymbol{\Psi}$ can be estimated by minimum the negative log likelihood, i.e.,

$$\min_{\boldsymbol{\Psi}} - \sum_{i=1}^n \log \sum_{j=1}^g \pi_j t(\mathbf{x}_i; \nu_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j), \quad (1)$$

which can be solved by EM algorithm.

Self-paced Learning

The objective function can be written as

$$\sum_{i=1}^n v_i \ell_i + f(\mathbf{v}, \lambda)$$

where v_i is learning weight of \mathbf{x}_i , ℓ_i is negative log likelihood of \mathbf{x}_i , $f(\mathbf{v}, \lambda)$ is a regularizer, if

- $f(\mathbf{v}, \lambda)$ is convex with respect to $v_i \in [0, 1]$;
- $v^*(\lambda, \ell)$ is monotonically decreasing with respect to ℓ , and it holds that $\log_{\ell} \ell \rightarrow 0) v^*(\lambda, \ell) = 1, \log_{\ell} \ell \rightarrow \inf) v^*(\lambda, \ell) = 0$;
- $v^*(\lambda, \ell)$ is monotonically increasing with respect to λ , and it holds that $\log_{\lambda} \lambda \rightarrow 0) v^*(\lambda, \ell) = 0, \log_{\lambda} \lambda \rightarrow \inf) v^*(\lambda, \ell) = 1$, where $v^*(\lambda, \ell) = \arg \min_{v \in [0, 1]} v \ell + f(\mathbf{v}, \lambda)$.

The Proposed Model

Objective Function

The objective function is given by

$$E(\mathbf{v}; \boldsymbol{\Psi}, \lambda) = - \sum_{i=1}^n v_i \log \sum_{j=1}^g \pi_j t(\mathbf{x}_i; \nu_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) - \lambda \|\mathbf{v}\|_1$$

where $\boldsymbol{\Psi} = \{\boldsymbol{\pi}, \boldsymbol{\nu}, \boldsymbol{\mu}, \boldsymbol{\Sigma}\}$; λ is a hyper-parameter, as a threshold when fixing $\boldsymbol{\Psi}$; v is the learning weight, also the outlier indicator with $v \in \{0, 1\}$; $\lambda \|\mathbf{v}\|_1$ enforces the sparsity of v since there exist only a few outliers in the training samples.

Optimization of \mathbf{v}

Fixing $\boldsymbol{\Psi}$, we estimate \mathbf{v} by solving

$$\min_{\mathbf{v} \in \{0, 1\}} E(\mathbf{v}; \boldsymbol{\Psi}, \lambda) = \sum_{i=1}^n v_i \ell_i - \lambda \|\mathbf{v}\|_1$$

Considering $v_i \in \{0, 1\}$, the problem can be written as

$$\min_{\mathbf{v} \in \{0, 1\}} \sum_{i=1}^n v_i (\ell_i - \lambda)$$

It is obvious that the solution is

$$v_i = \begin{cases} 0 & \ell_i > \lambda, \\ 1 & \ell_i \leq \lambda. \end{cases}$$

Optimization of $\boldsymbol{\Psi}$

Fixing \mathbf{v} , we estimate $\boldsymbol{\Psi}$ by solving

$$\min_{\boldsymbol{\Psi}} E(\boldsymbol{\Psi}; \lambda, \mathbf{v}) \Leftrightarrow \min_{\boldsymbol{\Psi}} \sum_{i=1}^n v_i \ell_i - \lambda \|\mathbf{v}\|_1 \Leftrightarrow \min_{\boldsymbol{\Psi}} \sum_{i=1}^n v_i \ell_i$$

Use EM algorithm to optimize $\boldsymbol{\Psi}$ iteratively.

Summary of the algorithm

Input: dataset $\mathcal{D} = \{\mathbf{x}_i, i = 1, 2, \dots, n\}$, learning pace $a > 1$, number of components

g

Output: $\boldsymbol{\Psi}_j (j = 1, 2, \dots, g)$

Initialize $\boldsymbol{\Psi}$ by the result of k -means.

Initialize λ to the median of $\ell_i, i = 1, 2, \dots, n$.

while $\Delta \boldsymbol{\Psi} \simeq 0$ **do**

while not converged do

 Update $\mathbf{v} = \arg \max_{\mathbf{v} \in \{0, 1\}} E(\mathbf{v}; \boldsymbol{\Psi}, \lambda)$.

while not converged do

E-step

 Update $\hat{\tau}_{ij}$ and \hat{u}_{ij} , which are intermediate variables in TMM.

M-step

 Update $\boldsymbol{\Psi}$.

end

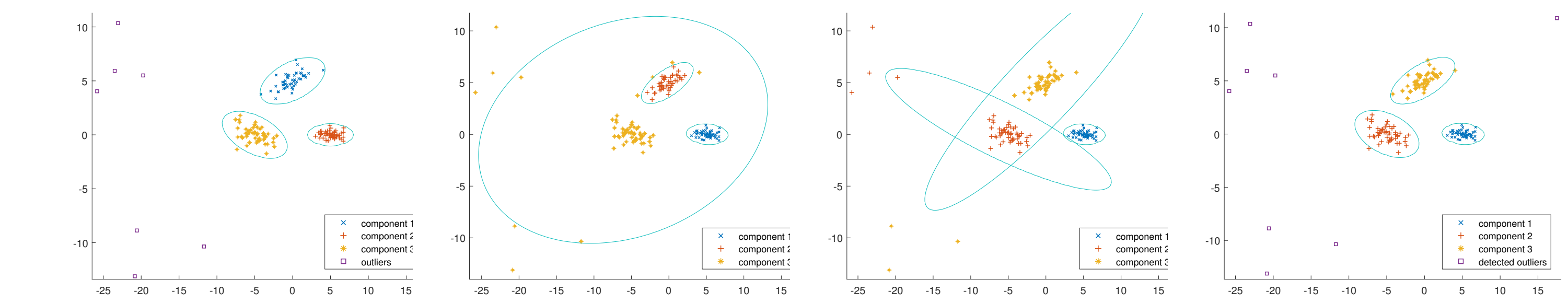
end

$\lambda \leftarrow a\lambda$.

end

Experiments

On the synthetic dataset



Synthetic dataset

GMM

TMM

SPTMM

On real dataset

	Clean				Noisy				Clean				Noisy			
	MSE	W/B	DB	Dunn	MSE	W/B	DB	Dunn	MSE	W/B	DB	Dunn	MSE	W/B	DB	Dunn
k -means	2.80	0.28	0.59	2.35	3.55	0.69	0.72	2.06	3.44	2.23	0.69	2.66	3.64	2.50	0.72	2.61
GMM	2.87	0.28	0.60	2.37	4.67	1.16	1.81	0.63	3.54	2.48	0.59	3.34	3.75	2.79	1.14	1.35
TMM	2.93	0.29	0.61	2.37	4.66	0.81	1.41	0.89	3.62	2.74	0.57	3.33	3.80	3.02	0.58	3.46
SPTMM	2.55	0.29	0.61	2.37	2.94	0.53	0.34	2.30	0.98	0.99	0.43	3.80	1.67	1.05	0.43	3.66

Seeds dataset

Thyroid dataset

Conclusions

In this paper, we depicted a novel model SPTMM which integrates the Self-Paced Learning mechanism into mixture of t distribution, in order to improve the mixture models' ability of handling outliers. Given the model, we developed an EM based algorithm that can solve the optimization problem in SPTMM efficiently. In addition to the mathematical justification, the experiments also display the value of the model. The results demonstrated that SPTMM clearly outperforms k -means, GMM and TMM for estimating the covariance matrix in the distributions. With respect to clustering, SPTMM is shown to be the best performer in most cases, in particular for the data with outliers. In the future, we would like to assess if SPTMM can be improved to perform better in a clean environment.