

# SELF-PACED MIXTURE OF $T$ DISTRIBUTION MODEL

Yang Zhang<sup>†,\*</sup>, Qingtao Tang<sup>†,\*</sup>, Li Niu<sup>‡</sup>, Tao Dai<sup>†</sup>, Xi Xiao<sup>†</sup>, Shu-Tao Xia<sup>†</sup>

<sup>†</sup> Department of Computer Science and Technology, Tsinghua University, China

<sup>‡</sup> Department of Electrical and Computer Engineering, Rice University, U.S.

## 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.

**Index Terms**— Self-Paced Learning, Robustness, Gaussian Mixture Model, Mixture of  $t$  distribution

## 1. INTRODUCTION

Gaussian Mixture Model (GMM) [1] is a powerful probabilistic model for representing a population consisting of several subpopulations. Due to its satisfactory flexibility, good interpretability, and simple parameter learning, GMM has been widely used in many fields, including data mining, pattern recognition [2], machine learning, and statistical analysis.

In GMM, the basic model is  $p(\mathbf{x}) = \sum_{j=1}^g \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$ , where each Gaussian density  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$  is a component of the mixture and has its own mean  $\boldsymbol{\mu}_j$  and covariance  $\boldsymbol{\Sigma}_j$ , and the parameters  $\pi_j \geq 0$  are mixing coefficients satisfying  $\sum_{j=1}^g \pi_j = 1$ . However, due to the thin-tailed property of the Gaussian distribution, GMM may perform poorly on the data which contain a group or groups of observations with heavy tails or outliers [3].

In order to handle the data with heavy tails or outliers, heavy-tailed distributions have been introduced into mixture model. For instance, Laplace distribution has heavier tails

than Gaussian distribution, and  $t$  distribution is a robust generalization of Gaussian distribution [4–6]. In particular, mixture of Laplace distribution is proposed by [7] while Peel *et al.* [3] introduced mixture of  $t$  distribution (TMM) with each component assumed to be a  $t$  distribution. These mixture models alleviate the influence of outliers to a certain degree.

Unfortunately, the aforementioned approaches cannot achieve satisfactory performance when the outliers stay distant from the centroids [8]. The main reason is that these methods focus on alleviating the influence of outliers based on heavy-tailed distribution while the outliers are still included in the learning procedure. With that reason, these methods cannot remove the effect of outliers adequately. In fact, the performance of these models can be significantly compromised when the outliers are far from the centroids.

To further reduce the influence of outliers, we introduce the Self-Paced Learning (SPL) [9] mechanism into mixture model, and develop a novel method named Self-Paced Mixture of  $t$  distribution model (SPTMM). Analogous to human learning procedure (*e.g.*, a pupil is supposed to understand elementary algebra before he/she can move on to advanced algebra topic), SPL starts with the easiest samples and then gradually includes harder samples while keeping the outliers off the learning procedure. In the last few years, the effectiveness of SPL has been validated in many tasks, such as event detection [10], co-saliency detection [11] and mixture of regression [12]. In this paper, we apply SPL to address the outliers in mixture model by screening the outliers during the learning procedure, which has never been explored before.

In summary, this paper presents three major contributions:

- This is the first work of employing Self-Paced Learning (SPL) to mixture model<sup>1</sup>, 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.

<sup>1</sup>Here, we adopt the definition of mixture model in statistics, which corresponds to the mixture distribution [13].

\*Equal Contribution. (Corresponding author: Qingtao Tang.)

## 2. MIXTURE OF $T$ DISTRIBUTION

The  $t$  distribution is defined as follows.

**Definition 1.** 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

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

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  $\pi_j$  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.*,

$$\min_{\boldsymbol{\Psi}} - \sum_{i=1}^n \log \sum_{j=1}^g \pi_j t(\mathbf{x}_i; \nu_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j), \quad (2)$$

which can be solved by EM algorithm [3].

## 3. SELF-PACED MIXTURE OF $T$ DISTRIBUTION

### 3.1. Introduction of Our model

Compared with GMM, TMM is able to mitigate the defect of outliers and thus is more robust to outliers, but it cannot remove the influence completely [8]. In fact, TMM is still prone to outliers, especially when the outliers are not close to the major data samples. To further reduce the influence of outliers, we introduce the Self-Paced Learning to TMM, leading to our Self-Paced Learning mixture of  $t$  distribution model (SPTMM). Specifically, based on (2), we introduce a latent binary variable  $v_i$  to indicate whether the sample  $\mathbf{x}_i$  is an outlier, and add a sparse regularizer of  $v_i$ :

$$E(\mathbf{v}; \boldsymbol{\Psi}, \lambda) = - \sum_{i=1}^n v_i \log \sum_{j=1}^g \pi_j t(\mathbf{x}_i; \nu_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) - \lambda \|\mathbf{v}\|_1, \quad (3)$$

where  $\lambda$  is a hyper-parameter,  $\mathbf{v} = (v_1, \dots, v_n)^T$  is a binary outlier indicator with  $v_i \in \{0, 1\}$ , and  $\|\mathbf{v}\|_1$  enforces the sparsity of  $\mathbf{v}$  since there exist only a few outliers in the training

samples. We expect  $v_i = 0$  if  $\mathbf{x}_i$  is an outlier and  $v_i = 1$  otherwise, so that only clean samples contribute to the learning procedure. Thus, we also refer to  $v_i$  as learning weight for  $\mathbf{x}_i$ . Note that  $\boldsymbol{\Psi}$  and  $\mathbf{v}$  are unknown and need to be learnt iteratively.

One interesting observation is that when  $\boldsymbol{\Psi}$  are fixed,  $\mathbf{v}$  can be easily learnt by using  $\lambda$  as threshold. Specifically, when  $-\log \sum_{j=1}^g \pi_j t(\mathbf{x}; \nu_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$  is larger than  $\lambda$ , the learnt  $v_i$  will be 0 (see Section 3.2.1 for technical derivation) and thus the corresponding  $\mathbf{x}_i$  will make no contribution when updating  $\boldsymbol{\Psi}$  next time. This meets our expectation because the outliers are usually far away from the mean of distributions with large negative log likelihood  $-\log \sum_{j=1}^g \pi_j t(\mathbf{x}; \nu_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$ .

### 3.2. EM algorithm for SPTMM

Inspired by Self-Paced Learning (SPL), we tend to start the learning procedure with the easiest samples and then gradually use harder samples while preventing outliers from being used. Hence, in our learning procedure, we initialize  $\lambda$  to a small value and enlarge  $\lambda$  gradually, until the learnt model parameters  $\boldsymbol{\Psi}$  remain unchanged, which is shown in Algorithm 1. With each fixed  $\lambda$ , we optimize (3) over  $\mathbf{v}$  and  $\boldsymbol{\Psi}$  iteratively until the objective in (3) converges. In particular, in each iteration with fixed  $\lambda$ , we update  $\mathbf{v}$  when fixing  $\boldsymbol{\Psi}$ , and then update  $\boldsymbol{\Psi}$  when fixing  $\mathbf{v}$ . In the following, we elaborate on the optimization *w.r.t.*  $\mathbf{v}$  and  $\boldsymbol{\Psi}$  in detail.

#### 3.2.1. Optimization of $\mathbf{v}$

When the model parameters  $\boldsymbol{\Psi}$  are fixed, we estimate  $\mathbf{v}$  by solving the following problem:

$$\min_{\mathbf{v} \in \{0,1\}} E(\mathbf{v}; \boldsymbol{\Psi}, \lambda) = \sum_{i=1}^n v_i \ell_i - \lambda \|\mathbf{v}\|_1, \quad (4)$$

where  $\ell_i = -\log \phi(\mathbf{x}_i; \boldsymbol{\Psi}) = -\log \sum_{j=1}^g \pi_j t(\mathbf{x}_i; \nu_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$ . Considering  $v_i \in \{0, 1\}$ , the problem in (4) can be written as

$$\min_{\mathbf{v} \in \{0,1\}} \sum_{i=1}^n v_i (\ell_i - \lambda). \quad (5)$$

It is obvious that the solution to (5) is

$$v_i = \begin{cases} 0 & \ell_i > \lambda, \\ 1 & \ell_i \leq \lambda. \end{cases} \quad (6)$$

It is worth noting that the training samples  $\mathbf{x}_i$  with larger  $\ell_i$  are more likely to be outliers. To explain more,  $\ell_i$  is an increasing function of  $(\mathbf{x}_i - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j)$  (see (1)), which is the Mahalanobis distance between sample  $\mathbf{x}_i$  and component centroid  $\boldsymbol{\mu}_j$ . Therefore,  $\ell_i$  can be used to measure the distance between sample  $\mathbf{x}_i$  and component centroids. As outliers are far from the component centroids, the training sample  $\mathbf{x}_i$  with larger  $\ell_i$  is more likely to be an outlier. According to (6), the

learning weight  $v_i$  of outliers would be set to 0, which can allow us to use only clean samples for learning model parameters  $\Psi$ .

### 3.2.2. Optimization of $\Psi$

Fixing  $\mathbf{v}$ , we estimate  $\Psi$  by minimizing the negative log marginal likelihood (3):

$$\min_{\Psi} E(\Psi; \lambda, \mathbf{v}) \Leftrightarrow \min_{\Psi} \sum_{i=1}^n v_i \ell_i - \lambda \|\mathbf{v}\|_1 \Leftrightarrow \min_{\Psi} \sum_{i=1}^n v_i \ell_i. \quad (7)$$

As  $v_i \in \{0, 1\}$ , (7) is the same optimization problem as (2). Thus, a similar EM algorithm can be applied to solve (7). For parameters  $\Psi = \{\boldsymbol{\pi}, \boldsymbol{\nu}, \boldsymbol{\mu}, \boldsymbol{\Sigma}\}$ , we set the degrees of freedom  $\nu_j$  ( $j = 1, 2, \dots, g$ ) to 4 following the practice in [14] and update  $\{\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}\}$  using the following equations:

$$\pi_j^* = \frac{1}{n} \sum_{i=1}^n v_i \hat{\tau}_{ij}^*, \quad (8)$$

$$\boldsymbol{\mu}_j^* = \frac{\sum_{i=1}^n v_i \hat{\tau}_{ij}^* \hat{u}_{ij}^* \mathbf{x}_i}{\sum_{i=1}^n v_i \hat{\tau}_{ij}^* \hat{u}_{ij}^*}, \quad (9)$$

$$\boldsymbol{\Sigma}_j^* = \frac{\sum_{i=1}^n v_i \hat{\tau}_{ij}^* \hat{u}_{ij}^* (\mathbf{x}_i - \boldsymbol{\mu}_j^*)(\mathbf{x}_i - \boldsymbol{\mu}_j^*)^T}{\sum_{i=1}^n v_i \hat{\tau}_{ij}^*}, \quad (10)$$

where

$$\hat{\tau}_{ij}^* = \frac{\pi_i t(\mathbf{x}_i; \nu_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{j=1}^g \pi_j t(\mathbf{x}_i; \nu_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}, \quad (11)$$

$$\hat{u}_{ij}^* = \frac{\nu_j + p}{\nu_j + (\mathbf{x}_i - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j)}. \quad (12)$$

The definitions of  $\hat{\tau}$  and  $\hat{u}$  are given by Peel *et al.* [3]. Note that the above two sets of variables  $\{\pi_j^*, \boldsymbol{\mu}_j^*, \boldsymbol{\Sigma}_j^*\}$  and  $\{\hat{\tau}_{ij}^*, \hat{u}_{ij}^*\}$  are updated in an iterative fashion based on EM algorithm. More precisely, we update  $\{\hat{\tau}_{ij}^*, \hat{u}_{ij}^*\}$  in E-step and update  $\{\pi_j^*, \boldsymbol{\mu}_j^*, \boldsymbol{\Sigma}_j^*\}$  in M-step.

### 3.2.3. Summary of the Algorithm

The whole optimization algorithm is summarized in Algorithm 1, which has an outer loop and an inner loop. In the outer loop, we start with small  $\lambda$  and increase  $\lambda$  by  $\lambda \leftarrow a\lambda$  in each iteration until the model parameters  $\Psi$  remain unchanged in several steps, which follows the strategy of SPL. In the inner loop, we update the binary outlier indicator  $\mathbf{v}$  and the model parameters  $\Psi$  iteratively until the objective in (3) converges.

## 4. EXPERIMENTS

In this section, we evaluate our SPTMM model on two tasks: density estimation and clustering. For density estimation, we

---

### Algorithm 1: Self-paced Mixture of $t$ distribution

---

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

**Output:**  $\Psi_j$  ( $j = 1, 2, \dots, g$ )

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

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

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

**while** *not converged* **do**

Update  $\mathbf{v} = \arg \max_{\mathbf{v} \in \{0,1\}} E(\mathbf{v}; \Psi, \lambda)$  by (6).

**while** *not converged* **do**

**E-step**

Update  $\hat{\tau}_{ij}$  and  $\hat{u}_{ij}$  by (11) and (12), respectively.

**M-step**

Update  $\Psi$  by (8), (9) and (10).

**end**

**end**

$\lambda \leftarrow a\lambda$ .

**end**

---

use a synthetic dataset, and compare our SPTMM method with GMM and TMM. For clustering, we use real-world datasets and additionally include  $k$ -means algorithm as a baseline besides GMM and TMM.

### 4.1. Experimental results in density estimation

This section is devoted to comparing GMM, TMM and SPTMM *w.r.t.* estimating the density of a synthetic dataset.

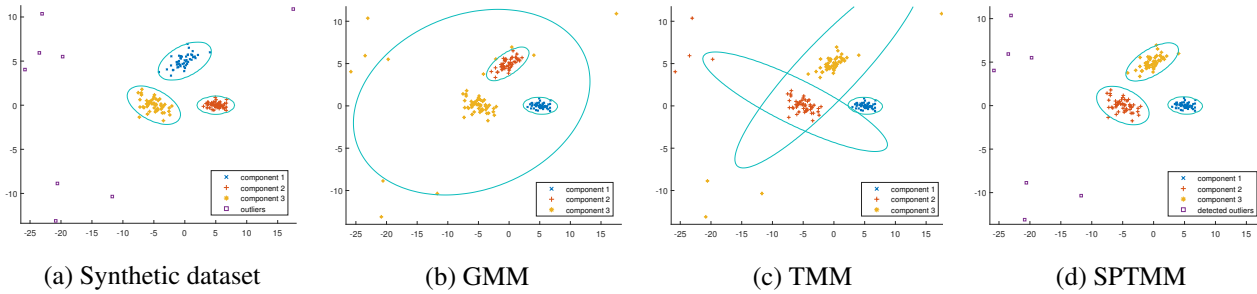
Firstly, a synthetic dataset is generated with a 3-component GMM. To evaluate the robustness of GMM, TMM and SPTMM to outliers, we add 5% outliers<sup>2</sup> to the synthetic dataset. The synthetic dataset is shown in **Fig. 1** (a).

The experimental results of GMM, TMM and SPTMM on the synthetic dataset are illustrated in **Fig. 1** (b) (c) (d), respectively. It is obvious that the covariance matrix of GMM is severely influenced by the outliers. Although TMM reduces the influence to a degree, the covariance of TMM is still affected heavily. In contrast, our SPTMM model gives much better estimation for the covariance matrix, since all the identified outliers are not used in the fitting procedure.

### 4.2. Experimental results in clustering

In this section, we evaluate the performance of our SPTMM model and other baselines, including  $k$ -means, GMM and TMM, for clustering on 4 real-world datasets, *i.e.*, Bezdekiris, Iris, Seeds and Thyroid [16]. Besides, to validate the robustness of these algorithms, we add 5% outliers to each real dataset. The real datasets without and with 5% outliers are referred to as Clean and Noisy respectively in Table 1.

<sup>2</sup>In this paper, we adopt the definition of outliers in [15], *i.e.*,  $[Q_1 - \alpha(Q_3 - Q_1), Q_3 + \alpha(Q_3 - Q_1)]$ .



**Fig. 1:** (a) is Synthetic dataset with outliers. (b), (c) and (d) are fitting results of GMM, TMM, SPTMM, respectively. The colored lines represent 99% confidence ellipses.

(a) Bezdekiris									(b) Iris								
Clean					Noisy				Clean					Noisy			
	MSE	W/B	DB	Dunn	MSE	W/B	DB	Dunn	MSE	W/B	DB	Dunn	MSE	W/B	DB	Dunn	
<i>k</i> -means	<b>0.93</b>	<b>0.30</b>	0.62	2.03	2.48	0.75	0.53	1.73	<b>0.93</b>	<b>0.31</b>	0.57	2.03	2.48	0.72	0.54	1.76	
GMM	1.12	0.39	0.74	1.88	2.53	0.78	0.84	0.93	1.13	0.40	0.74	1.88	2.95	0.98	1.12	1.03	
TMM	1.43	0.61	<b>0.29</b>	1.94	2.66	0.88	0.59	1.67	1.45	0.62	<b>0.28</b>	1.95	2.92	1.07	0.73	1.18	
SPTMM	1.51	0.65	<b>0.29</b>	<b>2.78</b>	<b>1.51</b>	<b>0.61</b>	<b>0.26</b>	<b>3.69</b>	1.45	0.62	<b>0.28</b>	<b>2.68</b>	<b>1.45</b>	<b>0.62</b>	<b>0.28</b>	<b>2.89</b>	

(c) Seeds									(d) Thyroid								
Clean					Noisy				Clean					Noisy			
	MSE	W/B	DB	Dunn	MSE	W/B	DB	Dunn	MSE	W/B	DB	Dunn	MSE	W/B	DB	Dunn	
<i>k</i> -means	2.80	<b>0.28</b>	<b>0.59</b>	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	<b>0.28</b>	0.60	<b>2.37</b>	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	<b>2.37</b>	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	<b>2.55</b>	0.29	0.61	<b>2.37</b>	<b>2.94</b>	<b>0.53</b>	<b>0.34</b>	<b>2.30</b>	<b>0.98</b>	<b>0.99</b>	<b>0.43</b>	<b>3.80</b>	<b>1.67</b>	<b>1.05</b>	<b>0.43</b>	<b>3.66</b>	

**Table 1:** MSE, W/B, DBI and Dunn of the clustering results on real datasets.

We evaluate the clustering performance of all these algorithms with internal clustering evaluation metrics [17], including clustering mean squared error (MSE) [18],  $\frac{WSS}{BSS}$  (WSS: within-cluster sum of squares; BSS: between-cluster sum of squares) [19], Davies-Bouldin index (DB) [20] and Dunn index (Dunn) [21]. Note that for MSE, W/B, DB, smaller value indicates better performance while for Dunn, larger value indicates better performance.

Clustering results reported in Table 1, show that on the Clean Bezdekiris, Iris and Seeds datasets, our SPTMM method performs best with half of the evaluation metrics and on the Clean Thyroid dataset, SPTMM achieves the best performance with all the evaluation metrics. Besides, we can observe that on all Noisy datasets, SPTMM significantly outperforms all the baselines with all the evaluation metrics. This attributes to that SPTMM can eliminate most outliers from the fitting procedure.

## 5. CONCLUSION

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.

## 6. ACKNOWLEDGMENTS

We thank Doctor Weipeng Huang, Insight Centre for Data Analytics, University College Dublin, for comments that greatly improved the manuscript. This work is supported by the National Natural Science Foundation of China under grant Nos. 61371078, 61771273, and the R&D Program of Shenzhen under grant Nos. JCYJ20140509172959977, JSGG20150512162853495, ZDSYS20140509172959989, JCYJ20160331184440545.

## 7. REFERENCES

- [1] Geoffrey J Mclachlan and Kaye E Basford, “Mixture models. inference and applications to clustering,” *Applied Statistics*, vol. 38, no. 2, 1988.
- [2] Sebastian Thrun, Dieter Fox, Wolfram Burgard, and Frank Dellaert, “Robust monte carlo localization for mobile robots,” *Artificial Intelligence*, vol. 128, no. 1-2, pp. 99–141, 2001.
- [3] David Peel and Geoffrey J Mclachlan, “Robust mixture modelling using the t distribution,” *Statistics and Computing*, vol. 10, no. 4, pp. 339–348, 2000.
- [4] Qingtao Tang, Li Niu, Yisen Wang, Tao Dai, Wangpeng An, Jianfei Cai, and Shu-Tao Xia, “Student-t process regression with student-t likelihood,” in *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, 2017, pp. 2822–2828.
- [5] Qingtao Tang, Tao Dai, Li Niu, Yisen Wang, Shu-Tao Xia, and Jianfei Cai, “Robust survey aggregation with student-t distribution and sparse representation,” in *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, 2017, pp. 2829–2835.
- [6] Qingtao Tang, Yisen Wang, and Shu-Tao Xia, “Student-t process regression with dependent student-t noise,” in *ECAI 2016 - 22nd European Conference on Artificial Intelligence, 29 August-2 September 2016, The Hague, The Netherlands - Including Prestigious Applications of Artificial Intelligence (PAIS 2016)*, 2016, pp. 82–89.
- [7] Debjani Bhowmick, Ac Davison, Darlene R Goldstein, and Yann Ruffieux, “A laplace mixture model for identification of differential expression in microarray experiments,” *Biostatistics*, vol. 7, no. 4, pp. 630, 2006.
- [8] G. J. Mclachlan, S K Ng, and R. W. Bean, “Robust cluster analysis via mixture models,” *Austrian Journal of Statistics*, vol. 35, pp. 157–174, 2006.
- [9] M. Pawan Kumar, Benjamin Packer, and Daphne Koller, “Self-paced learning for latent variable models,” in *International Conference on Neural Information Processing Systems*, 2010, pp. 1189–1197.
- [10] Jiang Lu, Deyu Meng, Teruko Mitamura, and Alexander G. Hauptmann, “Easy samples first: Self-paced reranking for zero-example multimedia search,” in *ACM International Conference on Multimedia*, 2014, pp. 547–556.
- [11] D. Zhang, D. Meng, and J. Han, “Co-saliency detection via a self-paced multiple-instance learning framework,” *IEEE Transactions on Pattern Analysis & Machine Intelligence*, vol. 39, no. 5, pp. 865–878, 2016.
- [12] Longfei Han, Dingwen Zhang, Dong Huang, Xiaojun Chang, Jun Ren, Senlin Luo, and Junwei Han, “Self-paced mixture of regressions,” in *Twenty-Sixth International Joint Conference on Artificial Intelligence*, 2017, pp. 1816–1822.
- [13] Bruce G. Lindsay, *Mixture Models: Theory, Geometry, and Applications*, Institute of Mathematical Statistics, American Statistical Association, 1995.
- [14] Jarno Vanhatalo, Pasi Jylänki, and Aki Vehtari, “Gaussian process regression with student- $t$  likelihood,” in *Annual Conference on Neural Information Processing Systems*, 2009, pp. 1910–1918.
- [15] J W Tuckey, *Exploratory Data Analysis*, Addison-Wesley Pub. Co., 1977.
- [16] M. Lichman, “UCI machine learning repository,” 2013.
- [17] Yanchi Liu, Zhongmou Li, Hui Xiong, Xuedong Gao, and Junjie Wu, “Understanding of internal clustering validation measures,” in *IEEE International Conference on Data Mining*, 2011, pp. 911–916.
- [18] Pang Ning Tan, Michael Steinbach, and Vipin Kumar, *Introduction to Data Mining, (First Edition)*, Addison-Wesley Longman Publishing Co., Inc., 2005.
- [19] Hans-Peter Kriegel, Erich Schubert, and Arthur Zimek, “The (black) art of runtime evaluation: Are we comparing algorithms or implementations?,” *Knowledge and Information Systems*, vol. 52, no. 2, pp. 341–378, Aug 2017.
- [20] David L. Davies and Donald W. Bouldin, *A Cluster Separation Measure*, IEEE Computer Society, 1979.
- [21] J. C. Dunn, “A fuzzy relative of the isodata process and its use in detecting compact well-separated clusters,” *Journal of Cybernetics*, vol. 3, no. 3, pp. 32–57, 1974.