

# Minimum Wasserstein Distance Estimator Under Finite Location-Scale Mixtures



Qiong Zhang and Jiahua Chen

**Abstract** When a population exhibits heterogeneity, we often model it via a finite mixture: decompose it into several different but homogeneous subpopulations. Contemporary practice favors learning the mixtures by maximizing the likelihood for statistical efficiency and the convenient EM algorithm for numerical computation. Yet the maximum likelihood estimate (MLE) is not well defined for finite location-scale mixture in general. We hence investigate feasible alternatives to MLE such as minimum distance estimators. Recently, the Wasserstein distance has drawn increased attention in the machine learning community. It has intuitive geometric interpretation and is successfully employed in many new applications. Do we gain anything by learning finite location-scale mixtures via a minimum Wasserstein distance estimator (MWDE)? This chapter investigates this possibility in several respects. We find that the MWDE is consistent and derive a numerical solution under finite location-scale mixtures. We study its robustness against outliers and mild model mis-specifications. Our moderate scaled simulation study shows the MWDE suffers some efficiency loss against a penalized version of MLE in general without noticeable gain in robustness. We reaffirm the general superiority of the likelihood-based learning strategies even for the non-regular finite location-scale mixtures.

**Keywords** Finite mixture model · Location scale family · Minimum distance estimator · Penalized maximum likelihood estimator · Wasserstein distance.

## 1 Introduction

Let  $\mathcal{F} = \{f(\cdot|\boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta\}$  be a parametric distribution family with density function  $f(\cdot|\boldsymbol{\theta})$  with respect to some  $\sigma$ -finite measure. Denote by  $G = \sum_{k=1}^K w_k \{\boldsymbol{\theta}_k\}$  a

---

Q. Zhang · J. Chen (✉)

Department of Statistics, University of British Columbia, Vancouver, BC, Canada  
e-mail: [qiong.zhang@stat.ubc.ca](mailto:qiong.zhang@stat.ubc.ca); [jhchen@stat.ubc.ca](mailto:jhchen@stat.ubc.ca)

distribution assigning probability  $w_k$  on  $\theta_k \in \Theta$ . A distribution with the following density function:

$$f(x|G) = \int f(x|\theta) dG(\theta) = \sum_{k=1}^K w_k f(x|\theta_k)$$

is called a finite  $\mathcal{F}$  mixture. We call  $f(x|\theta)$  the subpopulation density function,  $\theta$  the subpopulation parameter, and  $w_k$  the mixing weight of the  $k$ th subpopulation. We use  $F(x|\theta)$  and  $F(x|G)$  for the cumulative distribution functions (CDF) of  $f(x|\theta)$  and  $f(x|G)$ , respectively. Let

$$\mathbb{G}_K = \left\{ G : G = \sum_{k=1}^K w_k \delta_{\theta_k}, 0 \leq w_k \leq 1, \sum_{k=1}^K w_k = 1, \theta_k \in \Theta \right\}$$

be a space of mixing distributions with at most  $K$  support points. A mixture distribution of (exactly) order  $K$  has its mixing distribution  $G$  being a member of  $\mathbb{G}_K - \mathbb{G}_{K-1}$ .

We study the problem of learning the mixing distribution  $G$  given a set of independent and identically distributed (IID) observations  $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$  from a mixture  $f(x|G)$ . Throughout the paper, we assume the order of  $G$  is known and  $\mathcal{F}$  is a known location-scale family. That is,

$$f(x|\theta) = \frac{1}{\sigma} f_0\left(\frac{x - \mu}{\sigma}\right)$$

for some probability density function  $f_0(x)$  with  $x \in \mathbb{R}$  with respect to Lebesgue measure where  $\theta = (\mu, \sigma)$  with  $\Theta = \{\mathbb{R} \times \mathbb{R}^+\}$ .

Finite mixture models provide a natural representation of heterogeneous population that is believed to be composed of several homogeneous subpopulations (Pearson 1894; Schork et al. 1996). They are also useful for approximating distributions with unknown shapes that are particularly relevant in image generation (Kolouri et al. 2018), image segmentation (Farnoosh & Zarpak 2008), object tracking (Santosh et al. 2013), and signal processing (Plataniotis & Hatzinak 2000).

In statistics, the most fundamental task is to learn the unknown parameters. In early days, the method of moments was the choice for its ease of computation (Pearson 1894) under finite mixture models. Nowadays, the maximum likelihood estimate (MLE) is the first choice due to its statistical efficiency and the availability of an easy-to-use EM algorithm. Under a finite location-scale mixture model, the log-likelihood function of  $G$  is given by

$$\ell_N(G|\mathcal{X}) = \sum_{n=1}^N \log f(x_n|G) = \sum_{n=1}^N \log \left\{ \sum_{k=1}^K \frac{w_k}{\sigma_k} f_0\left(\frac{x_n - \mu_k}{\sigma_k}\right) \right\}. \quad (1)$$

At an arbitrary mixing distribution  $G_\epsilon = 0.5\{(x_1, \epsilon)\} + 0.5\{(0, 1)\}$ , we have  $\ell_N(G_\epsilon|\mathcal{X}) \rightarrow \infty$  as  $\epsilon \rightarrow 0$ . Hence, the MLE of  $G$  is not well defined or is ill defined. Various remedies, such as penalized maximum likelihood estimate (pMLE), have been proposed to overcome this obstacle (Chen et al. 2008; Chen & Tan 2009). At the same time, MLE can be thought of a special minimum distance estimator. It minimizes a specific Kullback–Leibler divergence between the empirical distribution and the assumed model  $\mathcal{F}$ . Other divergences and distances have been investigated in the literature as in Choi (1969); Yakowitz (1969); Woodward et al. (1984); Clarke and Heathcote (1994); Cutler and Cordero-Brana (1996); Deely and Kruse (1968). Recently, the Wasserstein distance has drawn increased attention in machine learning community due to its intuitive interpretation and good geometric properties (Evans & Matsen 2012; Arjovsky et al. 2017). The Wasserstein distance-based estimator for learning finite mixture models is absent in the literature.

Are there any benefits to learn finite location-scale mixtures by the minimum Wasserstein distance estimator (MWDE)? This chapter answers this question from several angles. We find that the MWDE is consistent and derive a numerical solution under finite location-scale mixtures. We compare the robustness of the MWDE with pMLE in the presence of outliers and mild model mis-specifications. We conclude that the MWDE suffers some efficiency loss against pMLE in general without obvious gain in robustness. Through this chapter, we better understand the pros and cons of the MWDE under finite location-scale mixtures. We reaffirm the general superiority of the likelihood-based learning strategies even for the non-regular finite location-scale mixtures.

In the next section, we first introduce the Wasserstein distance and some of its properties. This is followed by a formal definition of the MWDE, a discussion of its existence, and consistency under finite location-scale mixtures. In Sect. 2.4, we give some algebraic results that are essential for computing 2-Wasserstein distance between the empirical distribution and the finite location-scale mixtures. We then develop a BFGS algorithm scheme for computing the MWDE of the mixing distribution. In addition, we briefly review the penalized likelihood approach and its numerical issues. In Sect. 3, we characterize the efficiency properties of the MWDE relative to pMLE in various circumstances via simulation. We also study their robustness when the data contains outliers, is contaminated, or when the model is mis-specified. We then apply both methods in an image segmentation example. We conclude the paper with a summary in Sect. 4.

## 2 Wasserstein Distance and the Minimum Distance Estimator

We introduce the Wasserstein distance and the minimum Wasserstein distance estimator in this section.

## 2.1 Wasserstein Distance

Wasserstein distance is a distance between probability measures. Let  $\Omega$  be a Polish space endowed with a ground distance  $D(\cdot, \cdot)$  and  $\mathcal{P}(\Omega)$  the space of Borel probability measures on  $\Omega$ . Let  $\eta \in \mathcal{P}(\Omega)$  be a probability measure. If for some  $p > 0$ ,

$$\int_{\Omega} D^p(x, x_0) \eta(dx) < \infty,$$

for some (and thus any)  $x_0 \in \Omega$ , we say  $\eta$  has finite  $p$ th moment. Denote by  $\mathcal{P}_p(\Omega) \subset \mathcal{P}(\Omega)$  the space of probability measures with finite  $p$ th moment. For any  $\eta, \nu \in \mathcal{P}(\Omega)$ , we use  $\Pi(\eta, \nu)$  to denote the space of the bivariate probability measures on  $\Omega \times \Omega$  whose marginals are  $\eta$  and  $\nu$ . Namely,

$$\Pi(\eta, \nu) = \left\{ \pi \in \mathcal{P}(\Omega^2) : \int_{\Omega} \pi(x, dy) = \eta(x), \int_{\Omega} \pi(dx, y) = \nu(y) \right\}.$$

The  $p$ -Wasserstein distance is defined as follows.

**Definition 1 ( $p$ -Wasserstein Distance)** For any  $\eta, \nu \in \mathcal{P}_p(\Omega)$  with  $p \geq 1$ , the  $p$ th Wasserstein distance between  $\eta$  and  $\nu$  is

$$W_p(\eta, \nu) = \left\{ \inf_{\pi \in \Pi(\eta, \nu)} \int_{\Omega^2} D^p(x, y) \pi(dx, dy) \right\}^{1/p}.$$

Suppose  $X$  and  $Y$  are two random variables whose distributions are  $F$  and  $G$  and induced probability measures are  $\eta$  and  $\nu$ . We regard the  $p$ -Wasserstein distance between  $\eta$  and  $\nu$  and also the distance between random variables or distributions:  $W_p(X, Y) = W_p(F, G) = W_p(\eta, \nu)$ .

The  $p$ -Wasserstein distance is a distance on  $\mathcal{P}_p(\Omega)$  as shown by Villani (2003, Theorem 7.3). For any  $\eta, \nu, \xi \in \mathcal{P}_p(\Omega)$ , it has the following properties:

1. Non-negativity:  $W_p(\eta, \nu) \geq 0$  and  $W_p(\eta, \nu) = 0$  if and only if  $\eta = \nu$ .
2. Symmetry:  $W_p(\eta, \nu) = W_p(\nu, \eta)$ .
3. Triangular inequality:  $W_p(\eta, \nu) \leq W_p(\eta, \xi) + W_p(\xi, \nu)$ .

The Wasserstein distance has many nice properties. Let us denote  $\eta_n \xrightarrow{d} \eta$  for convergence in distribution or measure. Villani (2003, Theorem 7.1.2) shows that it has the following properties:

**Property 1.** For any  $q \geq p \geq 1$ ,  $W_q(\eta, \nu) \geq W_p(\eta, \nu)$ .

**Property 2.**  $W_p(\eta_n, \eta) \rightarrow 0$  as  $n \rightarrow \infty$  if and only if both:

- (i)  $\eta_n \xrightarrow{d} \eta$ .
- (ii)  $\int D^p(x, x_0) \eta_n(dx) \rightarrow \int D^p(x, x_0) \eta(dx)$  for some (and thus any)  $x_0 \in \Omega$ .

Computing the Wasserstein distance involves a challenging optimization problem in general but has a simple solution under a special case. Suppose  $\Omega$  is the space of real numbers,  $D(x, y) = |x - y|$ , and  $F$  and  $G$  are univariate distributions. Let  $F^{-1}(t) := \inf\{x : F(x) \geq t\}$  and  $G^{-1}(t) := \inf\{x : G(x) \geq t\}$  for  $t \in [0, 1]$  be their quantile functions. We can easily compute the Wasserstein distance based on the following property:

**Property 3.**  $W_p(F, G) = \left\{ \int_0^1 |F^{-1}(t) - G^{-1}(t)|^p dt \right\}^{1/p}$ .

## 2.2 Minimum Wasserstein Distance Estimator

Let  $W_p(\cdot, \cdot)$  be the  $p$ -Wasserstein distance with ground distance  $D(x, y) = |x - y|$  for univariate random variables. Let  $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$  be a set of IID observations from finite location-scale mixture  $f(x|G)$  of order  $K$  and  $F_N(x) = N^{-1} \sum_{n=1}^N \mathbb{1}(x_n \leq x)$  be the empirical distribution. We introduce the MWDE of the mixing distribution  $G$  to be

$$\hat{G}_N^{\text{MWDE}} = \arg \inf_{G \in \mathbb{G}_K} W_p(F_N(\cdot), F(\cdot|G)) = \arg \inf_{G \in \mathbb{G}_K} W_p^p(F_N(\cdot), F(\cdot|G)). \quad (2)$$

As we pointed out earlier, the MLE is not well defined under finite location-scale mixtures. Is the MWDE well defined? We examine the existence or sensibility of the MWDE. We show that the MWDE exists when  $f_0(\cdot)$  satisfies certain conditions.

Assume that  $f_0(0) > 0$ ,  $f_0(x)$  is bounded, continuous, and has finite  $p$ th moment. Under these conditions, we can see

$$0 \leq W_p(F_N(\cdot), F(\cdot|G)) < \infty$$

for any  $G \in \mathbb{G}_K$ . When  $N \leq K$ , the solution to (2) merits special attention. Let  $G_\epsilon = \sum_{n=1}^N (1/N) \delta_{(x_n, \epsilon)}$  be a mixing distribution assigning probability  $1/N$  on  $\theta_n = (x_n, \epsilon)$ . When  $\epsilon \rightarrow 0$ , each subpopulation in the mixture  $f(x|G_\epsilon)$  degenerates to a point mass at  $x_n$ . Hence, as  $\epsilon \rightarrow 0$ ,

$$W_p(F_N(\cdot), F(\cdot|G_\epsilon)) \rightarrow 0.$$

Since none of  $G \in \mathbb{G}_K$  has zero distance from  $F_N(\cdot)$ , the MWDE does not exist unless we expand  $\mathbb{G}_K$  to include  $G_0 = \sum_{n=1}^N (1/N) \delta_{(x_n, 0)} = \lim G_\epsilon$ . To remove this technical artifact, in the MWDE definition, we expand the space of  $\sigma$  to  $[0, \infty)$ . We denote by  $F(\cdot|(\theta_0, 0))$  a distribution with point mass at  $x = \theta_0$ . With this expansion,  $G_0$  is the MWDE when  $N \leq K$ .

Let  $\delta = \inf\{W_p(F_N(\cdot), F(\cdot|G)) : G \in \mathbb{G}_K\}$ . Clearly,  $0 \leq \delta < \infty$ . By definition, there exists a sequence of mixing distributions  $G_m \in \mathbb{G}_K$  such that  $W_p(F_N(\cdot), F(\cdot|G_m)) \rightarrow \delta$  as  $m \rightarrow \infty$ . Suppose one mixing weight of  $G_m$  has limit 0. Removing this support point and rescaling, we get a new mixing distribution

sequence, and it still satisfies  $W_p(F_N(\cdot), F(\cdot|G_m)) \rightarrow \delta$ . For this reason, we assume that its mixing weights have non-zero limits by selecting converging subsequence if necessary to ensure the limits exist. Further, when the mixing weights of  $G_m$  assume their limiting values while keeping subpopulation parameters the same, we still have  $W_p(F_N(\cdot), F(\cdot|G_m)) \rightarrow \delta$  as  $m \rightarrow \infty$ . In the following discussion, we therefore discuss the sequence of mixing distributions whose mixing weights are fixed.

Suppose the first subpopulation of  $G_m$  has its scale parameter  $\sigma_1 \rightarrow \infty$  as  $m \rightarrow \infty$ . With the boundedness assumption on  $f_0(x)$ , the mass of this subpopulation will spread thinly over entire  $\mathbb{R}$  because  $\sigma_1^{-1} f_0((x - \mu_1)/\sigma_1) \rightarrow 0$  uniformly. For any fixed finite interval,  $[a, b]$ , this thinning makes

$$F(b|\theta_1) - F(a|\theta_1) \rightarrow 0$$

as  $m \rightarrow \infty$ . It implies that for any given  $t \in (0, 0.5)$ , we have

$$|F^{-1}(t|\theta_1)| + |F^{-1}(1-t|\theta_1)| \rightarrow \infty.$$

This further implies for any  $t \in (0, w_1/2)$ , we have

$$|F^{-1}(t|G_m)| + |F^{-1}(1-t|G_m)| \rightarrow \infty$$

as  $m \rightarrow \infty$ . In comparison, the empirical quantile satisfies  $x_{(1)} \leq F_N^{-1}(t) \leq x_{(N)}$  for any  $t$ . By Property 3 of  $W_p(\cdot, \cdot)$ , these lead to  $W_p(F_N(\cdot), F(\cdot|G_m)) \rightarrow \infty$  as  $m \rightarrow \infty$ . This contradicts the assumption  $W_p(F_N(\cdot), F(\cdot|G_m)) \rightarrow \delta$ . Hence,  $\sigma_1 \rightarrow \infty$  is not a possible scenario of  $G_m$  nor  $\sigma_k \rightarrow \infty$  for any  $k$ .

Can a subpopulation of  $G_m$  instead have its location parameter  $\mu \rightarrow \infty$ ? For definitiveness, let this subpopulation correspond to  $\theta_1$ . Note that at least  $w_1\{1 - F_0(0)\}$ -sized probability mass of  $F(x|G_m)$  is contained in the range  $[\mu_1, \infty)$ . Because of this, when  $\mu_1 \rightarrow \infty$ , we have  $F^{-1}(1-t|G_m) \rightarrow \infty$  for  $t = w_1\{1 - F_0(0)\}/2$ . Therefore,  $W_p(F_N(\cdot), F(\cdot|G_m)) \rightarrow \infty$  by Property 3. This contradicts  $W_p(F_N(\cdot), F(\cdot|G_m)) \rightarrow \delta < \infty$ . Hence,  $\mu_1 \rightarrow \infty$  is not a possible scenario of  $G_m$  either. For the same reason, we cannot have  $\mu_k \rightarrow \pm\infty$  for any  $k$ .

After ruling out  $\mu_k \rightarrow \pm\infty$  and  $\sigma_k \rightarrow \infty$ , we find  $G_m$  has a converging subsequence whose limit is a proper mixing distribution in  $\mathbb{G}_K$ . This limit is then an MWDE and the existence is verified.

The MWDE may not be unique, and the mixing distribution may lead to a mixture with degenerate subpopulations. We will show that the MWDE is consistent as the sample size goes to infinity. Thus, having degenerated subpopulations in the learned mixture is a mathematical artifact and also a sensible solution. In contrast, no matter how large the sample size becomes, there are always degenerated mixing distributions with unbounded likelihood values.

### 2.3 Consistency of MWDE

We consider the problem when  $\mathcal{X} = \{x_1, \dots, x_N\}$  are IID observations from a finite location-scale mixture of order  $K$ . The true mixing distribution is denoted as  $G^*$ . Assume that  $f_0(x)$  is bounded, continuous, and has finite  $p$ th moment. We say the location-scale mixture is identifiable if

$$F(x|G_1) = F(x|G_2)$$

for all  $x$  given  $G_1, G_2 \in \mathbb{G}_K$  implies  $G_1 = G_2$ . We allow subpopulation scale  $\sigma = 0$ . The most commonly used finite location-scale mixtures, such as the normal mixture, are well known to be identifiable (Teicher 1961). Holzmann et al. (2004) give a sufficient condition for the identifiability of general finite location-scale mixtures. Let  $\varphi(\cdot)$  be the characteristic function of  $f_0(t)$ . The finite location-scale mixture is identifiable if for any  $\sigma_1 > \sigma_2 > 0$ ,  $\lim_{t \rightarrow \infty} \varphi(\sigma_1 t) / \varphi(\sigma_2 t) = 0$ .

We consider the MWDE based on  $p$ -Wasserstein distance with ground distance  $D(x, y) = |x - y|$  for some  $p \geq 1$ . The MWDE under finite location-scale mixture model as defined in (2) is asymptotically consistent.

**Theorem 1** *With the same conditions on the finite location-scale mixture and same notations above, we have the following conclusions:*

1. For any sequence  $G_m \in \mathbb{G}_K$  and  $G^* \in \mathbb{G}_K$ ,  $W_p(F(\cdot|G_m), F(\cdot|G^*)) \rightarrow 0$  implies  $G_m \xrightarrow{d} G^*$  as  $m \rightarrow \infty$ .
2. The MWDE satisfies  $W_p(F(\cdot|G^*), F(\cdot|\hat{G}_N^{MWDE})) \rightarrow 0$  as  $N \rightarrow \infty$  almost surely.
3. The MWDE is consistent:  $W_p(\hat{G}_N^{MWDE}, G^*) \rightarrow 0$  as  $N \rightarrow \infty$  almost surely.

**Proof** We present these three conclusions in the current order that is easy to understand. For the sake of proof, a different order is better. For ease presentation, we write  $F^* = F(\cdot|G^*)$  and  $\hat{G} = \hat{G}_N^{MWDE}$  in this proof.

We first prove the second conclusion. By the triangular inequality and the definition of the minimum distance estimator, we have

$$W_p(F^*, F(\cdot|\hat{G}_N)) \leq W_p(F_N, F^*) + W_p(F_N, F(\cdot|\hat{G}_N)) \leq 2W_p(F_N, F^*).$$

Note that  $F_N$  is the empirical distribution and  $F^*$  is the true distribution; we have  $F_N(x) \rightarrow F^*(x)$  uniformly in  $x$  almost surely. At the same time, under the assumption that  $F_0(x)$  has finite  $p$ th moment,  $F^*(x)$  also has finite  $p$ th moment. The  $p$ th moment of  $F_N(x)$  converges to that of  $F^*(x)$  almost surely. Given the ground distance  $D(x, y) = |x - y|$ , the  $p$ th moment in Wasserstein distance sense is the usual moments in probability theory. By Property 2, we conclude  $W_p(F_N, F(\cdot|G^*)) \rightarrow 0$  as both conditions there are satisfied.

Conclusion 3 is implied by Conclusions 1 and 2. With Conclusion 2 already established, we only need to prove Conclusion 1 to complete the whole proof. By Helly's lemma (Van der Vaart 2000, Lemma 2.5) again,  $G_m$  has a converging

subsequence though the limit can be a subprobability measure. Without loss of generality, we assume that  $G_m$  itself converges with limit  $\tilde{G}$ . If  $\tilde{G}$  is a subprobability measure, so would be  $F(\cdot|\tilde{G})$ . This will lead to

$$W_p(F(\cdot|G_m), F(\cdot|G^*)) \rightarrow W_p(F(\cdot|\tilde{G}), F(\cdot|G^*)) \neq 0,$$

which violates the theorem condition. If  $\tilde{G}$  is a proper distribution in  $\mathbb{G}_K$  and

$$W_p(F(\cdot|\tilde{G}), F(\cdot|G^*)) = 0,$$

then by identifiability condition, we have  $\tilde{G} = G^*$ . This implies  $G_m \rightarrow G^*$  and completes the proof.  $\square$

The multivariate normal mixture is another type of location-scale mixture. The above consistency result of MWDE can be easily extended to finite multivariate normal mixtures.

**Theorem 2** *Consider the problem when  $X = \{x_1, \dots, x_N\}$  are IID observations from a finite multivariate normal mixture distribution of order  $K$  and  $\hat{G}_N^{MWDE}$  is the minimum Wasserstein distance estimator defined by (2). Let the true mixing distribution be  $G^*$ . The MWDE is consistent:  $W_p(\hat{G}_N^{MWDE}, G^*) \rightarrow 0$  as  $N \rightarrow \infty$  almost surely.*

The rigorous proof is long though the conclusion is obvious. We offer a less formal proof based on several well-known probability theory results:

- (I) A multivariate random variable sequence  $Y_n$  converges in distribution to  $Y$  if and only if  $\mathbf{a}^\tau Y_n$  converges to  $\mathbf{a}^\tau Y$  for any unit vector  $\mathbf{a}$ .
- (II) If  $Y$  is multivariate normal if and only if  $\mathbf{a}^\tau Y$  is normal for all  $\mathbf{a}$ .
- (III) The normal distribution has finite moment of any order.

Let  $X_m$  be a random vector with distribution  $F(\cdot|G_m)$  for some  $G_m \in \mathbb{G}_K$ ,  $m = 0, 1, 2, \dots$ , in a general mixture model setting. Suppose as  $m \rightarrow \infty$ , with the notation we introduced previously

$$W_p(X_m, X_0) \rightarrow 0.$$

Then for any unit vector  $\mathbf{a}$ , based on property 2 of the Wasserstein distance and the result (I), we can see that

$$W_p(\mathbf{a}^\tau X_m, \mathbf{a}^\tau X_0) \rightarrow 0.$$

Next, we apply this result to normal mixture so that  $F(\cdot|G_m)$  becomes  $\Phi(\cdot|G_m)$  that stands for a finite multivariate normal mixture with mixing distribution  $G_m$ . In this case,  $X_m$  is a random vector with distribution  $\Phi(\cdot|G_m)$ . Let  $(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  be generic subpopulation parameters. We can see that the distribution of  $\mathbf{a}^\tau X_m$ ,  $\Phi_{\mathbf{a}}(\cdot|G_m)$  is a finite normal mixture with subpopulation parameters  $(\mathbf{a}^\tau \boldsymbol{\mu}_k, \mathbf{a}^\tau \boldsymbol{\Sigma}_k \mathbf{a})$ , and mixing



weighs the same as those of  $G_m$ . Let the mixing distributions after projection be  $G_{m,\mathbf{a}}$  and  $G_{0,\mathbf{a}}$ .

By the same argument in the proof of Theorem 1,

$$W_p(\Phi(\cdot|\hat{G}_N), \Phi(\cdot|G^*)) \rightarrow 0$$

almost surely as  $N \rightarrow \infty$ . This implies

$$W_p(\Phi_{\mathbf{a}}(\cdot|\hat{G}_N), \Phi_{\mathbf{a}}(\cdot|G^*)) \rightarrow 0$$

almost surely as  $N \rightarrow \infty$  for any  $\mathbf{a}$ . Hence, by Conclusion 1 of Theorem 1,  $\hat{G}_{N,\mathbf{a}} \xrightarrow{d} \hat{G}_{\mathbf{a}}^*$  almost surely for any unit vector  $\mathbf{a}$ . We therefore conclude the consistency result:  $\hat{G}_N \xrightarrow{d} \hat{G}^*$  almost surely.

## 2.4 Numerical Solution to MWDE

Both in applications and in simulation experiments, we need an effective way to compute the MWDE. We develop an algorithm that leverages the explicit form of the Wasserstein distance between two measures on  $\mathbb{R}$  for the numerical solution to the MWDE. The strategy works for any  $p$ -Wasserstein distance, but we only provide specifics for  $p = 2$  as it is the most widely used.

Let  $Y$  be a random variable with distribution  $f_0(\cdot)$ . Denote the mean and variance of  $Y$  by  $\mu_0 = \mathbb{E}(Y)$  and  $\sigma_0^2 = \text{Var}(Y)$ . Recall that  $G = \sum_{k=1}^K w_k\{(\mu_k, \sigma_k)\}$ . Let  $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(N)}$  be the order statistics,  $\bar{x}^2 = N^{-1} \sum_{n=1}^N x_n^2$ , and  $\xi_n = F^{-1}(n/N|G)$  be the  $(n/N)$ th quantile of the mixture for  $n = 0, 1, \dots, N$ . Let

$$T(x) = \int_{-\infty}^x t f_0(t) dt$$

and define

$$\begin{aligned} \Delta F_{nk} &= F_0\left(\frac{\xi_n - \mu_k}{\sigma_k}\right) - F_0\left(\frac{\xi_{n-1} - \mu_k}{\sigma_k}\right), \\ \Delta T_{nk} &= T\left(\frac{\xi_n - \mu_k}{\sigma_k}\right) - T\left(\frac{\xi_{n-1} - \mu_k}{\sigma_k}\right). \end{aligned}$$

When  $p = 2$ , we expand the squared  $W_2$  distance,  $\mathbb{W}_N$ , between the empirical distribution and  $F(\cdot|G)$  as follows:

$$\begin{aligned} \mathbb{W}_N(G) &= W_2^2(F_N(\cdot), F(\cdot|G)) \\ &= \int_0^1 \{F_N^{-1}(t) - F^{-1}(t|G)\}^2 dt \end{aligned}$$

$$\begin{aligned}
&= \overline{x^2} + \sum_{k=1}^K w_k \{ \mu_k^2 + \sigma_k^2 (\mu_0^2 + \sigma_0^2) + 2\mu_k \sigma_k \mu_0 \} \\
&\quad - 2 \sum_k w_k \{ \mu_k \sum_{n=1}^N x_{(n)} \Delta F_{nk} + \sigma_k \sum_{n=1}^N x_{(n)} \Delta T_{nk} \}.
\end{aligned}$$

The MWDE minimizes  $\mathbb{W}_N(G)$  with respect to  $G$ . The mixing weights and subpopulation-scale parameters in this optimization problem have natural constraints. We may replace the optimization problem with an unconstrained one by the following parameter transformation:

$$\begin{aligned}
\sigma_k &= \exp(\tau_k), \\
w_k &= \exp(t_k) / \left\{ \sum_{j=1}^K \exp(t_j) \right\}
\end{aligned}$$

for  $k = 1, 2, \dots, K$ . We may then minimize  $\mathbb{W}_N$  with respect to  $\{(\mu_k, \tau_k, t_k) : k = 1, 2, \dots, K\}$  over the unconstrained space  $\mathbb{R}^{3K}$ . Furthermore, we adopt the quasi-Newton BFGS algorithm (Nocedal & Wright 2006, Section 6.1). To use this algorithm, it is best to provide the gradients of  $\mathbb{W}_N(G)$ , which are given as follows:

$$\begin{aligned}
\frac{\partial}{\partial t_j} \mathbb{W}_N &= \sum_{k=1}^K \left\{ \frac{\partial w_k}{\partial t_j} \frac{\partial}{\partial w_k} \mathbb{W}_N \right\} = \sum_k w_j (\delta_{jk} - w_k) \frac{\partial}{\partial w_k} \mathbb{W}_N, \\
\frac{\partial}{\partial \mu_j} \mathbb{W}_N &= 2w_j \left\{ \mu_j + \sigma_j \mu_0 - \sum_{n=1}^N x_{(n)} \Delta F_{nj} \right\}, \\
\frac{\partial}{\partial \tau_j} \mathbb{W}_N &= 2w_j \left\{ \sigma_j (\mu_0^2 + \sigma_0^2) + \mu_j \mu_0 - \sum_{n=1}^N x_{(n)} \Delta T_{nj} \right\} \frac{\partial \sigma_j}{\partial \tau_j},
\end{aligned}$$

for  $j = 1, 2, \dots, K$ , where

$$\begin{aligned}
\frac{\partial}{\partial w_k} \mathbb{W}_N &= \{ \mu_k^2 + \sigma_k^2 (\mu_0^2 + \sigma_0^2) + 2\mu_k \sigma_k \mu_0 \} - 2 \sum_{n=1}^{N-1} \{ x_{(n+1)} - x_{(n)} \} \xi_n F(\xi_n | \mu_k, \sigma_k) \\
&\quad - 2 \{ \mu_k \sum_{n=1}^N x_{(n)} \Delta F_{nk} + \sigma_k \sum_{n=1}^N x_{(n)} \Delta T_{nk} \}.
\end{aligned}$$

Since  $\mathbb{W}_N(G)$  is non-convex, the algorithm may find a local minimum of  $\mathbb{W}_N(G)$  instead of a global minimum as required for MWDE. We use multiple initial values for the BFGS algorithm and regard the one with the lowest  $\mathbb{W}_N(G)$  value as the solution. We leave the algebraic details in the Appendix.

This algorithm involves computing the quantiles  $\xi_n$  and  $\Delta T_{nj}$ ; repeatedly that may lead to high computational cost. Since  $\xi_n$  is between  $\min_k F^{-1}(n/N|\theta_k)$  and  $\max_k F^{-1}(n/N|\theta_k)$ , it can be found efficiently via a bisection method. Fortunately,  $T(x)$  has simple analytical forms under two widely used location-scale mixtures that make the computation of  $\Delta T_{nj}$  efficient:

1. When  $f_0(t) = (2\pi)^{-1/2} \exp(-x^2/2)$ , which is the density function of the standard normal, we have  $tf_0(t) = -f_0'(t)$ . In this case, we find

$$T(x) = -f_0(x).$$

2. For a finite mixture of location-scale logistic distributions, we have

$$f_0(t) = \frac{\exp(-x)}{(1 + \exp(-x))^2}$$

and

$$T(x) = \int_{-\infty}^x tf_0(t)dt = \frac{x}{1 + \exp(-x)} - \log(1 + \exp(x)). \quad (3)$$

## 2.5 Penalized Maximum Likelihood Estimator

A well-investigated inference method under a finite mixture of location-scale families is the pMLE (Tanaka 2009; Chen et al. 2008). Chen et al. (2008) consider this approach for finite normal mixture models. They recommend the following penalized log-likelihood function:

$$p\ell_N(G|\mathcal{X}) = \ell_N(G|\mathcal{X}) - a_N \sum_k \left\{ s_x^2/\sigma_k^2 + \log \sigma_k^2 \right\}$$

for some positive  $a_N$  and sample variance  $s_x^2$ . The log-likelihood function is given in (1). They suggest us to learn the mixing distribution  $G$  via pMLE defined as

$$\hat{G}_N^{\text{pMLE}} = \arg \sup p\ell_N(G|\mathcal{X}).$$

The size of  $a_N$  controls the strength of the penalty, and a recommended value is  $N^{-1/2}$ . Regularizing the likelihood function via a penalty function fixes the problem caused by degenerated subpopulations (i.e., some  $\sigma_k = 0$ ). The pMLE is shown to be strongly consistent when the number of components has a known upper bound under the finite normal mixture model.

The penalized likelihood approach can be easily extended to a finite mixture of location-scale families. Let  $f_0(\cdot)$  be the density function in the location-scale family as before. We may replace the sample variance  $s_x^2$  in the penalty function by any scale-invariance statistic such as the sample inter-quartile range. This is applicable even if the variance of  $f_0(\cdot)$  is not finite.

We can use the EM algorithm for numerical computation. Let  $\mathbf{z}_n = (z_{n1}, \dots, z_{nK})$  be the membership vector of the  $n$ th observation. That is, the  $k$ th entry of  $\mathbf{z}_n$  is 1 when the response value  $x_n$  is an observation from the  $k$ th subpopulation and 0 otherwise. When the complete data  $\{(\mathbf{z}_n, x_n), n = 1, 2, \dots, N\}$  are available, the penalized complete data likelihood function of  $G$  is given by

$$p\ell_N^c(\mathcal{X}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \log \left\{ \frac{w_k}{\sigma_k} f_0 \left( \frac{x_n - \mu_k}{\sigma_k} \right) \right\} - a_N \sum_k \{s_x^2 / \sigma_k^2 + \log(\sigma_k^2)\}.$$

Given the observed data  $\mathcal{X}$  and proposed mixing distribution  $G^{(t)}$ , we have the conditional expectation

$$w_{nk}^{(t)} = \mathbb{E}(z_{nk} | \mathcal{X}, G^{(t)}) = \frac{w_k^{(t)} f(x_n | \mu_k^{(t)}, \sigma_k^{(t)})}{\sum_{j=1}^K w_j^{(t)} f(x_n | \mu_j^{(t)}, \sigma_j^{(t)})}.$$

After this E-step, we define

$$Q(G | G^{(t)}) = \sum_{n=1}^N \sum_{k=1}^K w_{nk}^{(t)} \log \left\{ \frac{w_k}{\sigma_k} f_0 \left( \frac{x_n - \mu_k}{\sigma_k} \right) \right\} - a_N \sum_k \{s_x^2 / \sigma_k^2 + \log(\sigma_k^2)\}.$$

Note that the subpopulation parameters are separated in  $Q(\cdot | \cdot)$ . The M-step is to maximize  $Q(G | G^{(t)})$  with respect to  $G$ . The solution is given by the mixing distribution  $G^{(t+1)}$  with mixing weights

$$w_k^{(t+1)} = N^{-1} \sum_{n=1}^N w_{nk}^{(t)}$$

and the subpopulation parameters

$$\theta_k^{(t+1)} = \arg \min_{\theta} \left\{ \sum_n w_{nk}^{(t)} \{\log \sigma - f(x_n | \theta)\} + a_N \{s_x^2 / \sigma^2 + \log \sigma^2\} \right\} \quad (4)$$

with the notational convention  $\theta = (\mu, \sigma)$ .

For general location-scale mixture, the M-step (4) may not have a closed-form solution, but it is merely a simple two-variable function. There are many effective algorithms in the literature to solve this optimization problem. The EM algorithm for pMLE increases the value of the penalized likelihood after each iteration. Hence, it should converge as long as the penalized likelihood function has an upper bound. We do not give a proof as it is a standard problem.

### 3 Experiments

We now study the performance of MWDE and pMLE under finite location-scale mixtures. We explore the potential advantages of the MWDE and quantify its efficiency loss, if any, by simulation experiments. Consider the following three location-scale families (Chen et al. 2020):

1. Normal distribution:  $f_0(x) = (2\pi)^{-1/2} \exp(-x^2/2)$ . Its mean and variance are given by  $\mu_0 = 0$  and  $\sigma_0^2 = 1$ .
2. Logistic distribution:  $f_0(x) = \exp(-x)/(1 + \exp(-x))^2$ . Its mean and variance are given by  $\mu_0 = 0$  and  $\sigma_0^2 = \pi^2/3$ .
3. Gumbel distribution (type I extreme-value distribution):  $f_0(x) = \exp(-x - \exp(-x))$ . Its mean and variance are given by  $\mu_0 = \gamma$  and  $\sigma_0^2 = \pi^2/6$ , where  $\gamma$  is the Euler constant.

We will also include a real-data example to compare the image segmentation result of using the MWDE and pMLE.

#### 3.1 Performance Measure

For vector-valued parameters, the commonly used performance metric of their estimators is the mean-squared error (MSE). A mixing distribution with finite and fixed support points can be regarded as a real-valued vector in theory. Yet the mean-squared errors of the mixing weights, the subpopulation means, and the subpopulation scales are not comparable in terms of the learned finite mixture. In this chapter, we use two performance metrics specific for finite mixture models. Let  $\hat{G}$  and  $G^*$  be the learned mixing distribution and the true mixing distribution. We use  $L_2$  distance between the learned mixture and the true mixture as the first performance metric. The  $L_2$  distance between two mixtures  $f(\cdot|G)$  and  $f(\cdot|\tilde{G})$  is defined to be

$$L_2(f(\cdot|G), f(\cdot|\tilde{G})) = \{\mathbf{w}^\tau S_{GG} \mathbf{w} - 2\mathbf{w}^\tau S_{G\tilde{G}} \tilde{\mathbf{w}} + \tilde{\mathbf{w}}^\tau S_{\tilde{G}\tilde{G}} \tilde{\mathbf{w}}\}^{1/2},$$

where  $S_{GG}$ ,  $S_{G\tilde{G}}$  and  $S_{\tilde{G}\tilde{G}}$  are three square matrices of size  $K \times K$  with their  $(n, m)$ th elements given by

$$\int f(x|\theta_n) f(x|\theta_m) dx, \quad \int f(x|\theta_n) f(x|\tilde{\theta}_m) dx, \quad \int f(x|\tilde{\theta}_n) f(x|\tilde{\theta}_m) dx.$$

Given an observed value  $x$  of a unit from the true mixture population, by Bayes' theorem, the most probable membership of this unit is given by

$$k^*(x) = \arg \max_k \{w_k^* f^*(x|\theta_k^*)\}.$$

Following the same rule, if  $\hat{G}$  is the learned mixing distribution, then the most likely membership of the unit with observed value  $x$  is

$$\hat{k}(x) = \arg \max_k \{\hat{w}_k f(x|\hat{\theta}_k)\}.$$

We cannot directly compare  $k^*(x)$  and  $\hat{k}(x)$  because the subpopulation themselves is not labeled. Instead, the adjusted Rand index (ARI) is a good performance metric for clustering accuracy. Suppose the observations in a dataset are divided into  $K$  clusters  $A_1, A_2, \dots, A_K$  by one approach, and  $K'$  clusters  $B_1, B_2, \dots, B_{K'}$  by another. Let  $N_i = \#(A_i)$ ,  $M_j = \#(B_j)$ ,  $N_{ij} = \#(A_i B_j)$  for  $i, j = 1, 2, \dots, K$ , where  $\#(A)$  is the number of units in set  $A$ . The ARI between these two clustering outcomes is defined to be

$$\text{ARI} = \frac{\sum_{i,j} \binom{N_{ij}}{2} - \binom{N}{2}^{-1} \sum_{i,j} \binom{N_i}{2} \binom{M_j}{2}}{\frac{1}{2} \sum_i \binom{N_i}{2} + \frac{1}{2} \sum_j \binom{M_j}{2} - \binom{N}{2}^{-1} \sum_{i,j} \binom{N_i}{2} \binom{M_j}{2}}.$$

When the two clustering approaches completely agree with each other, the ARI value is 1. When data are assigned to clusters randomly, the expected ARI value is 0. ARI values close to 1 indicate a high degree of agreement. We compute ARI based on clusters formed by  $k^*(x)$  and  $\hat{k}(x)$ .

For each simulation, we choose or generate a mixing distribution  $G^{*(r)}$  and then generate a random sample from mixture  $f(x|G^{*(r)})$ . This is repeated  $R$  times. Let  $\hat{G}^{(r)}$  be the learned  $G$  based on the  $r$ th dataset. We obtain the two performance metrics as follows:

1. Mean  $L_2$  distance:

$$\text{ML2} = R^{-1} \sum_{r=1}^R L_2(f(\cdot|\hat{G}^{(r)}), f(\cdot|G^{*(r)})).$$

2. Mean-adjusted Rand index:

$$\text{MARI} = R^{-1} \sum_{r=1}^R \text{ARI}(\hat{G}^{(r)}, G^{*(r)}).$$

The lower the ML2 and the higher the MARI, the better the estimator performs.

### 3.2 Performance Under Homogeneous Model

The homogeneous location-scale model is a special mixture model with a single subpopulation  $K = 1$ . Both MWDE and MLE are applicable for parameter estimation. There have been no studies of MWDE in this special case in the literature. It is therefore of interest to see how MWDE performs under this model.

Under three location-scale models given earlier, the MWDE has closed analytical forms. Using the same notation introduced, their analytical forms are as follows:

1. Normal distribution:

$$\hat{\mu}^{\text{MWDE}} = \bar{x}, \hat{\sigma}^{\text{MWDE}} = \sum_{n=1}^N x_{(n)} \{f_0(\xi_{n-1}) - f_0(\xi_n)\}.$$

2. Logistic distribution:

$$\hat{\mu}^{\text{MWDE}} = \bar{x}, \hat{\sigma}^{\text{MWDE}} = \frac{3}{\pi^2} \sum_{n=1}^N x_{(n)} \{T(\xi_n) - T(\xi_{n-1})\},$$

where  $T(x)$  is given in (3).

3. Gumbel distribution:

$$\hat{\mu}^{\text{MWDE}} = \{1 - \gamma r\}^{-1} \{\bar{x} - \gamma T\}, \hat{\sigma}^{\text{MWDE}} = T - r \hat{\mu}^{\text{MWDE}},$$

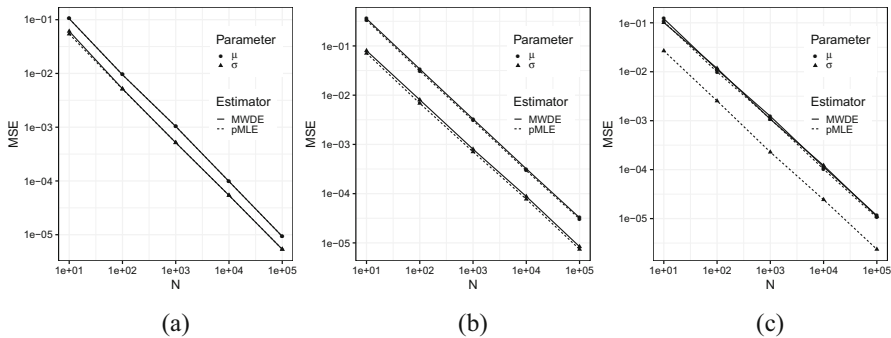
where

$$T = \{\gamma^2 + \pi^2/6\}^{-1} \sum_{n=1}^N x_{(n)} \int_{\xi_{n-1}}^{\xi_n} t f_0(t) dt$$

and  $r = \gamma/(\gamma^2 + \pi^2/6)$ .

The MLEs under the logistic and Gumbel distributions do not have an easy-to-use analytical form. We employ a numerical optimization program to solve for MLE. We generate samples of sizes between  $N = 10$  and  $N = 100,000$  with  $R = 1000$  repetitions. Under the homogeneous model, it is most convenient to compute the MSE of the location and scale parameters separately. Due to the invariance property, we generate data from distributions with  $\mu = 0$  and  $\sigma = 1$ . The simulation results are summarized as plots in Fig. 1. Both the x and y axes in these plots are in logarithm scale. For both MLE and MWDE, their log-MSE and  $\log(N)$  values are close to the straight lines with slope  $-1$ . This phenomenon indicates that both estimators have the expected convergence rates  $O(N^{-1/2})$  as the sample size  $N \rightarrow \infty$ .

The performances of the estimators for the location parameter and scale parameter are different. For the location parameter under all three models, the lines formed by MLE and MWDE are nearly indistinguishable though the MLE is always below



**Fig. 1** The MSEs of the MWDE and MLE for location and scale parameters versus sample size  $N$  under different homogeneous models. **(a)** Normal. **(b)** Logistic. **(c)** Gumbel

the MWDE. For the scale parameter  $\sigma$ , the MLE is also more efficient than the MWDE, but the difference is negligible under the normal and logistic models. Under the Gumbel model, the MWDE is less efficient.

In summary, using MWDE under a homogeneous model may not be preferred but may be acceptable under the normal and logistic models. We do not investigate the performance of MWDE under Gumbel mixture due to its efficiency loss under the homogeneous model. With these observations, we move to its performance under finite location-scale mixtures.

### 3.3 Efficiency and Robustness Under Finite Location-Scale Mixtures

We next study the efficiency and robustness of the MWDE for learning finite location-scale mixtures. Since the MLE is not well defined, we compare the performance of MWDE with the pMLE (Chen & Tan 2009) instead. We compare their performances when the mixture model is correctly specified, when the data is contaminated, or when the model is mildly mis-specified.

#### 3.3.1 Efficiency

A widely employed two-component mixture model (Cutler & Cordero-Brana 1996; Zhu 2016) has a density function in the following form:

$$f(x|G) = pf(x|0, a) + (1 - p)f(x|b, 1) \tag{5}$$

for some density function  $f(\cdot|\theta)$  from a location-scale family. Namely, we have  $K = 2$  is known, the mixing weights be  $w_1 = p, w_2 = 1 - p$ , and subpopulation



parameters be  $\theta_1 = (0, a)$  and  $\theta_2 = (b, 1)$ . By choosing different combinations of  $p$ ,  $a$ , and  $b$ , we obtain mixtures with different properties. Due to the invariance property, we need to consider only the case where one of the location parameters is 0 and one of the scale parameters is 1.

We generate samples from  $f(x|G)$  according to the following scheme: generate an observation  $Y$  from distribution with density function  $f_0(x)$ , and let

$$X = \begin{cases} aY, & \text{with probability } p; \\ Y + b, & \text{otherwise.} \end{cases} \quad (6)$$

We can easily see that the distribution of  $X$  is  $f(x|G)$  specified earlier.

The level of difficulty to precisely estimate the mixture largely depends on the degree of overlap between the subpopulations. Let

$$o_{j|i} = \mathbb{P}(w_i f(X|\mu_i, \sigma_i) < w_j f(X|\mu_j, \sigma_j) | X \sim f(x|\mu_i, \sigma_i)).$$

This is the probability of a unit from subpopulation  $i$  misclassified as a unit in subpopulation  $j$  by the maximum posterior rule. The degree of overlap between the  $i$ th and  $j$ th subpopulations is therefore

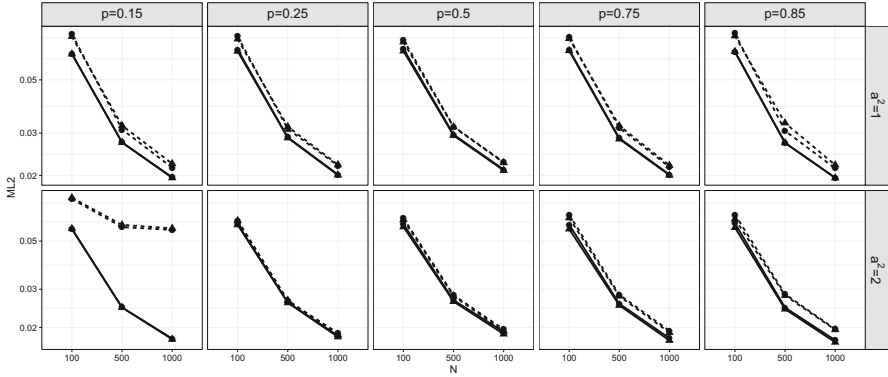
$$o_{ij} = o_{j|i} + o_{i|j}. \quad (7)$$

We employ the following  $a$ ,  $b$ , and  $p$  values in our simulation experiments:

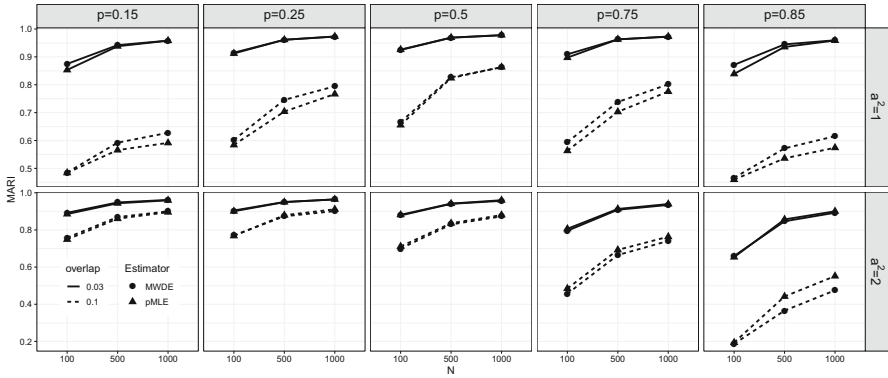
1. Mixing proportion  $p = 0.15, 0.25, 0.5, 0.75, 0.85$ .
2. Scale of the first subpopulation  $a^2 = 1, 2$ .
3. Location parameter  $b$  values such that  $o_{12} = 0.03, 0.1$ .

The combination of these choices leads to 24 mixtures with various shapes. The sample size  $N$  in our experiments is chosen to be 100, 500, and 1000, respectively.

We obtain the average  $L_2$  distance (ML2) and adjusted Rand index (MARI) based on  $R = 1000$  repetitions on data generated from normal and logistic mixture distributions as specified by (6). Figures 2 and 3, respectively, contain plots of ML2 and MARI of the WMDE and pMLE estimators against sample size  $N$  under these two models. We can see that when the sample size increases, ML2 of both estimators decreases and MARI of both estimators increases, supporting the theory that both WMDE and pMLE are consistent. Under the normal mixture, these two estimators have nearly equal  $L_2$  distances. The MWDE slightly outperforms pMLE in terms of the MARI, when the degree of overlap is large ( $o_{12} = 0.1$ ) and the two subpopulations have both equal scale and highly unbalanced weights. Under logistic mixture, as shown in plots (a) and (b) of Fig. 3, the pMLE always outperforms the MWDE in terms of the  $L_2$  distance. In terms of the MARI, the MWDE is better when the scale parameters are equal and weights are highly unbalanced. When the scale parameters are different, the pMLE is better than MWDE when  $p > 0.5$  and worse than MWDE when  $p < 0.5$ .



(a)



(b)

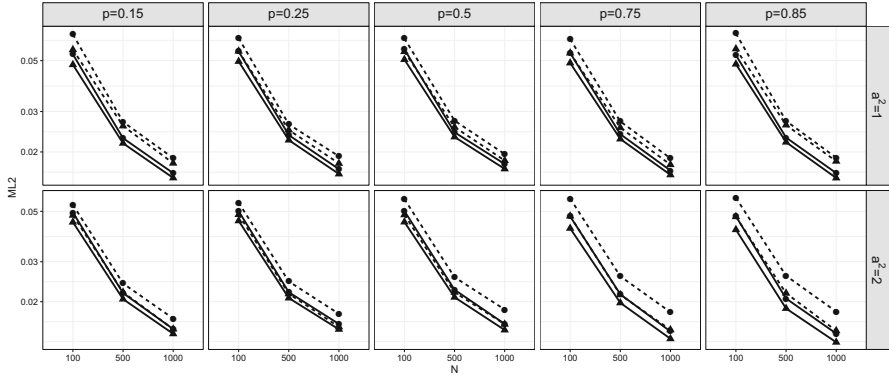
**Fig. 2** Performances of pMLE and MWDE under 2-component normal mixture. (a)  $L_2$  distance. (b) Adjusted Rand index

We next investigate the performance of the MWDE and pMLE for learning 3-component normal mixtures. We come up with 8 such distributions with different configurations. The three subpopulations have the same or different weights and same or different scale parameter values. They lead to different degrees of overlap as defined by

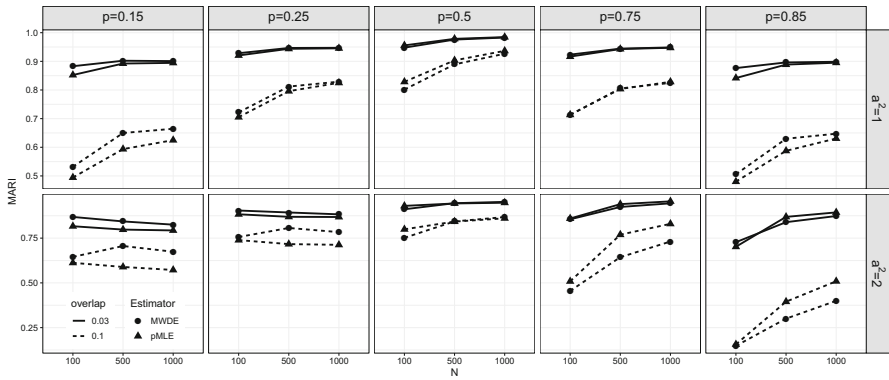
$$\text{MeanOmega} = \text{mean}_{1 \leq i < j \leq 3} \{o_{ij}\},$$

where  $o_{ij}$  is the degree of overlap between subpopulations  $i$  and  $j$  in (7). See Table 1 for detailed parameter values.

Figure 4 contains plots of the  $L_2$  and MARI values of two estimators. It is seen that the pMLE consistently outperforms MWDE in terms of  $L_2$  but the difference



(a)

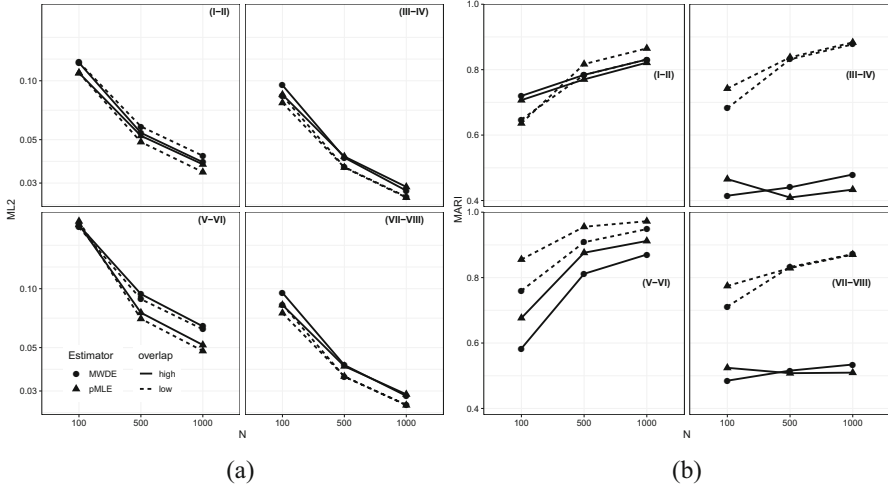


(b)

**Fig. 3** Performances of pMLE and MWDE under 2-component logistic mixture. (a)  $L_2$  distance. (b) Adjusted Rand index

**Table 1** Parameter values of 3-component normal mixtures

	MeanOmega	$w_1$	$w_2$	$w_3$	$\mu_1$	$\mu_2$	$\mu_3$	$\sigma_1$	$\sigma_2$	$\sigma_3$
I	0.288 (low)	0.4	0.5	0.1	-2	0	1	0.3	2	0.4
II	0.367 (high)	0.4	0.5	0.1	-2	0	1	0.3	1	0.4
III	0.097 (low)	0.3	0.5	0.2	-3	0	3	1	1	1
IV	0.249 (high)	0.3	0.5	0.2	-2	0	2	1	1	1
V	0.148 (low)	1/3	1/3	1/3	-1	0	1	1.5	0.1	0.5
VI	0.267 (high)	1/3	1/3	1/3	-0.5	0	0.5	1.5	0.1	0.5
VII	0.091 (low)	1/3	1/3	1/3	-3	0	3	1	1	1
VIII	0.226 (high)	1/3	1/3	1/3	-2	0	2	1	1	1



**Fig. 4** Performances of pMLE and MWDE under 3-component normal mixture. (a)  $L_2$  distance. (b) Adjusted Rand index

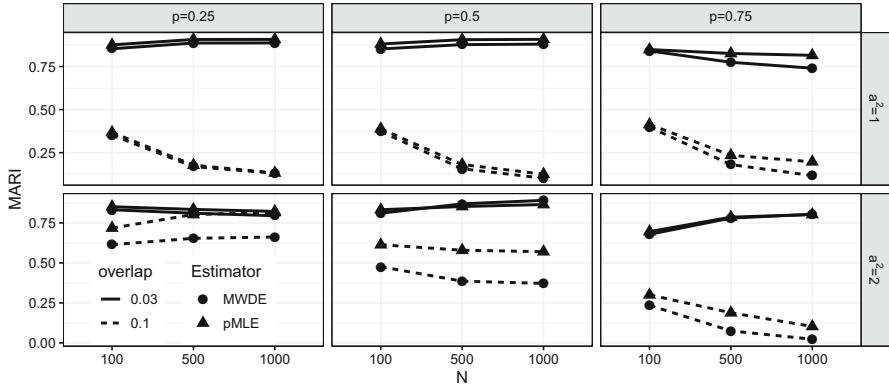
is small. The performances of the MWDE and pMLE are mixed in terms of MARI and the differences are small. The pMLE is clearly better under I and II.

### 3.3.2 Robustness

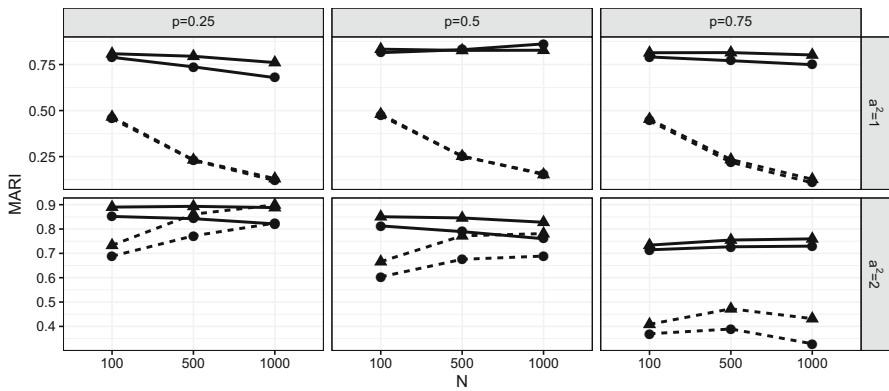
Robustness is another important property of estimators. Sample mean is the most efficient unbiased estimator of the population mean in terms of variance under normality or some other well-known parametric models. However, the value of the sample mean changes dramatically even if the dataset contains merely a single extreme value. Sample median offers a respectable alternative and still has high efficiency across a broader range of parametric models.

In the context of learning finite location-scale mixture models, both pMLE and MWDE rely on a parametric distribution family assumption through  $f_0(x)$ . How important is to have  $f_0(x)$  correctly specified? We shed some light into this problem by simulation experiments in this section. We learn finite normal mixtures assuming  $K = 2$  but generate data from the following distributions:

1. Mixture with outliers:  $(1 - \alpha)\{p\phi(x|0, a) + (1 - p)\phi(x|b, 1)\} + \alpha\phi(x|8, 1)$  with  $\alpha = 0.01$  and  $\phi(x|\mu, \sigma) = \exp(-(x - \mu)^2/2\sigma^2)/\sqrt{2\pi}\sigma^2$ .
2. Mixture contaminated:  $(1 - \alpha)\{p\phi(x|0, a) + (1 - p)\phi(x|b, 1)\} + \alpha\phi(x|b/2, 7)$  with  $\alpha = 0.01$ .
3. Mixture mis-specified I:  $p f_0(x|0, a) + (1 - p) f_0(x|b, 1)$  with  $f_0(x)$  being Student-t with 4 degrees of freedom.



(a)



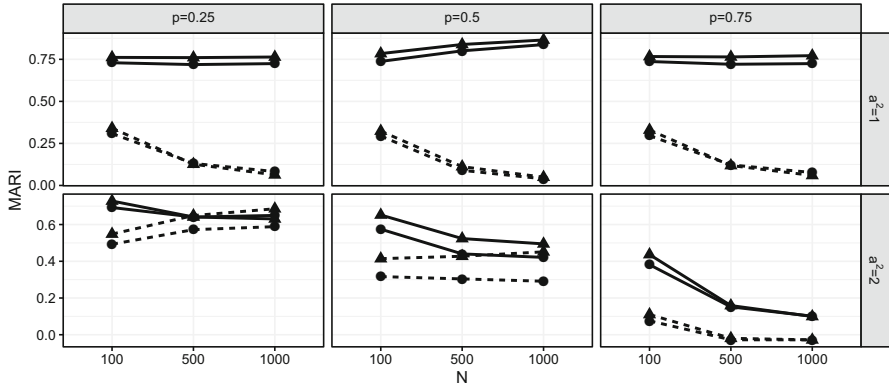
(b)

**Fig. 5** Adjusted Rand index based on pMLE and MWDE when data contains outliers or is contaminated. (a) Mixture with outliers. (b) Mixture contaminated

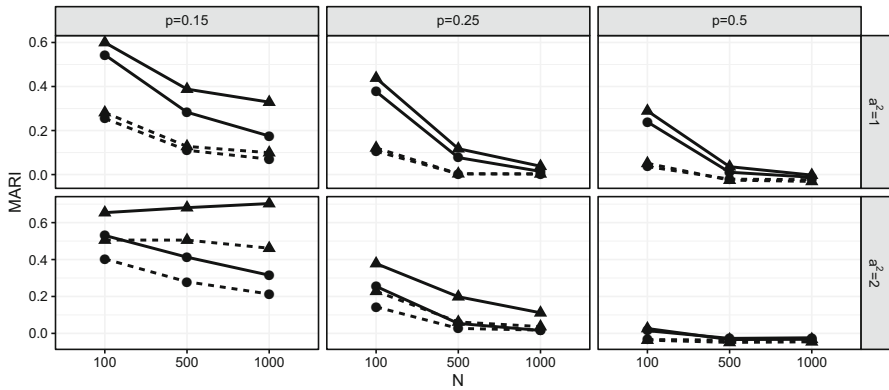
4. Mixture mis-specified II:  $p f_1(x|0, a) + (1 - p) f_2(x|b, 1)$  with  $f_1(x)$  and  $f_2(x)$  being Student-t with 2 and 4 degrees of freedom.

In every case, we use the combinations of the  $a$ ,  $b$ , and  $p$  values the same as before. We regard  $(1 - \alpha)\{p\phi(x|0, a) + (1 - p)\phi(x|b, 1)\}$  as the true distribution in all cases and computed the MARI accordingly.

We obtain the MARI values based on  $R = 1000$  repetitions with sample sizes  $N = 100, 500, \text{ and } 1000$ , see Figs. 5 and 6. We see that when the degree of overlap is low, MWDE and pMLE have similar performances. When the subpopulation variance is larger ( $a^2 = 2$ ), the performance of pMLE is generally better. In general, we conclude that pMLE is preferred.



(a)



(b)

**Fig. 6** Adjusted Rand index based on pMLE and MWDE when subpopulation distributions are mis-specified. (a) Mixture mis-specified I. (b) Mixture mis-specified II

Statistical inference usually becomes more accurate when the sample size increases. This is not the case in this simulation experiment. We can see that MARI often decreases (becomes less accurate) when the sample size increases. This is not caused by simulation error. When the model is mis-specified, the learned model does not converge to the “true model” as  $N \rightarrow \infty$ . Hence, the inference does not necessarily improve. The moral of this simulation study is that the MWDE is not more robust than the pMLE, against our intuition.

### 3.4 Image Segmentation

Image segmentation aims to partition an image into regions, each with a reasonably homogeneous visual appearance or corresponds to objects or parts of objects (Bishop 2006, Chapter 9). In this section, we perform image segmentation with finite normal mixtures, a common practice in the machine learning community.

Each pixel in an image is represented by three numbers within the range of  $[0, 1]$  that corresponds to the intensities of the Red, Green, and Blue (RGB) channels. Since the intensity values are always between 0 and 1, unlike the common practice in the literature, we feel obliged to transform the intensity values to ensure the normal mixture model fits better. Let  $y = \Phi^{-1}((x + 1/N)/(1 + 2/N))$  with  $x$  being the intensity and  $N$  the total number of pixels in the image. We then learn a two-component normal mixture on  $y$  values from each channel. Namely, we learn three normal mixtures on red, green, and blue channels, respectively.

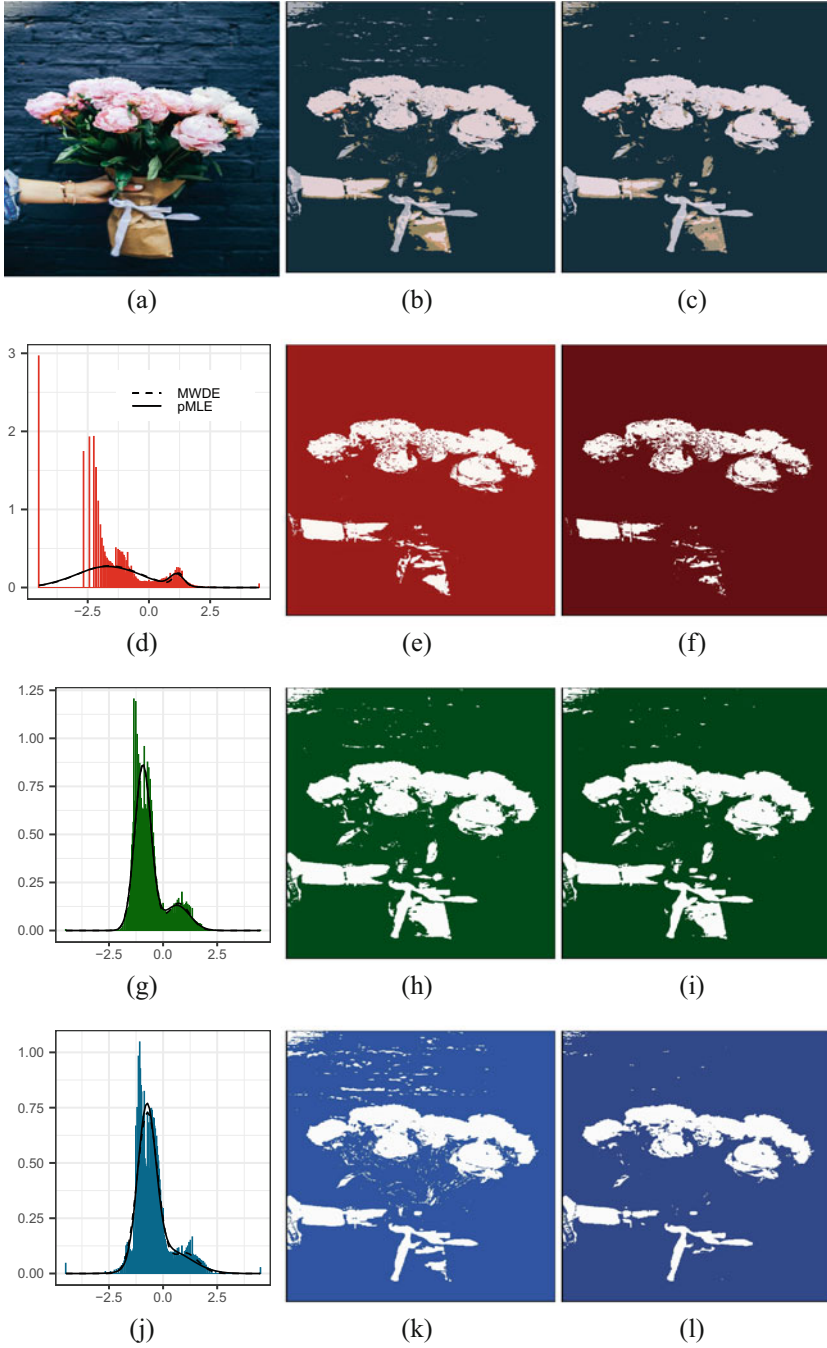
We use the maximum posterior probability rule to assign each pixel to one of two clusters. We then form an image segment by pixels assigned to the same cluster. We visualize the segregated images channel by channel by re-drawing the image with the original intensity value replaced by the average intensity of the pixels assigned to the specific cluster.

The segregated images depend heavily on the fitted mixture distributions. We compare the segregated images obtained by the normal mixtures learned via the pMLE and MWDE. We retrieved an image from Pexels<sup>1</sup> as shown in Fig. 7a. Clark (2015) resized the original high-resolution image to  $433 \times 650$  grids using Lanczos filter. We learn a normal mixture of order  $K = 2$  for each channel based on resized datasets and evaluated its utility of segregating the foreground and the background.

We present the specifications of the learned mixing distributions by pMLE and MWDE in Table 2. Plots (d), (g), and (j) in Fig. 7 are histograms of the transformed intensity values of RGB channels, together with the mixture densities learned via pMLE and MWDE. The corresponding segmented images are shown as plots (e), (h), and (k) for pMLE and (f), (i), and (l) for MWDE. The estimated parameter values and the fitted density on the red and green channels based on these two approaches are very similar. For the blue channel, the fitted densities and the segmentation results are very similar although the estimated parameter values of the second component are quite different. Both approaches can produce images with meaningful structures segregating foreground from background.

There are two clusters in each of 3 channels leading to 8 refined clusters. We may paint each pixel with the average RGB intensity triplet according to these 8 refined clusters. The re-created images via pMLE and MWDE, respectively, are shown in (b) and (c). We note these two images are very similar, showing that both learning strategies are effective.

<sup>1</sup> <https://www.pinterest.se/pin/761952830692007143/>.



**Fig. 7** Flower image and its segmentation outcomes. (a), (b) and (c): original image; aggregated images based on segmentation outcomes via pMLE and MWDE. (d), (g) and (j): histograms of pixel intensity of Red, Green, and Blue channels together with the fitted mixtures. (e), (h) and (k): segregated images via PMLE in RGB channels. (f), (i) and (l): segregated images via MWDE in RGB channels



**Table 2** Estimated mixing distributions of the flower image by pMLE and MWDE.

Channel	Estimator	$w_1$	$w_2$	$\mu_1$	$\mu_2$	$\sigma_1$	$\sigma_2$
Red	pMLE	0.896	0.104	-1.668	1.139	1.321	0.277
	MWDE	0.915	0.085	-1.617	1.220	1.316	0.213
Green	pMLE	0.804	0.196	-0.935	0.637	0.373	0.595
	MWDE	0.819	0.181	-0.926	0.724	0.378	0.510
Blue	pMLE	0.735	0.265	-0.753	0.268	0.414	1.034
	MWDE	0.862	0.138	-0.722	1.019	0.473	0.592

## 4 Conclusion

The MWDE provides another approach for learning finite location-scale mixtures. We have shown the MWDE is well defined and consistent. Our moderate scaled simulation study shows it suffers some efficiency loss against a penalized version of MLE in general without noticeable gain in robustness. We gain the knowledge on the benefits and drawbacks of the MWDE under finite location-scale mixtures. We reaffirm the general superiority of the likelihood-based learning strategies even for non-regular models.

**Acknowledgments** The authors would like to thank Richard Schonberg for proofreading the manuscript.

## Appendix

### Numerically Friendly Expression of $W_2(F_N, F(\cdot|G))$

To learn the finite mixture distribution through MWDE, we must compute

$$\mathbb{W}_N(G) = W_2^2(F_N(\cdot), F(\cdot|G)) = \int_0^1 \{F_N^{-1}(t) - F^{-1}(t|G)\}^2 dt$$

for finite location-scale mixture

$$F(\cdot|G) = \sum_{k=1}^K \pi_k F(\cdot|\theta_k) = \sum_{k=1}^K \pi_k \sigma_k^{-1} F_0((x - \mu_k)/\sigma_k).$$

We write  $\mathbb{E}_k(\cdot)$  as expectation under distribution  $F(\cdot|\theta_k)$ . For instance,

$$\mathbb{E}_k\{X^2\} = \mu_k^2 + \sigma_k^2(\mu_0^2 + \sigma_0^2) + 2\mu_k\sigma_k\mu_0.$$

Let  $I_n = ((n-1)/N, n/N]$  for  $n = 1, 2, \dots, N$  so that  $F_N^{-1}(t) = x_{(n)}$  when  $t \in I_n$ , where  $x_{(n)}$  is the  $n$ th order statistic. For ease of notation, we write  $x_{(n)}$  as  $x_n$ . Over this interval, we have

$$\int_{I_n} \{F_N^{-1}(t) - F^{-1}(t|G)\}^2 dt = \int_{I_n} [x_n^2 - 2x_n F^{-1}(t|G) + \{F^{-1}(t|G)\}^2] dt. \quad (8)$$

The integration of the first term in (8), after summing over  $n$ , is given by

$$\sum_{n=1}^N \int_{I_n} x_n^2 dt = N^{-1} \sum_n x_n^2 = \bar{x}^2.$$

The integration of the third term in (8) is

$$\sum_{n=1}^N \int_{I_n} \{F^{-1}(t|G)\}^2 dt = \int_{-\infty}^{\infty} x^2 f(x|G) dx = \sum_{k=1}^K w_k \mathbb{E}_k\{X^2\}.$$

Let  $\xi_0 = -\infty$ ,  $\xi_{N+1} = \infty$ , and  $\xi_n = F^{-1}(n/N|G)$  for  $n = 1, \dots, N$ . Denote

$$\Delta F_{nk} = F(\xi_n|\theta_k) - F(\xi_{n-1}|\theta_k)$$

and

$$T(x) = \int_{-\infty}^x t f_0(t) dt, \quad \Delta T_{nk} = T((\xi_n - \mu_k)/\sigma_k) - T((\xi_{n-1} - \mu_k)/\sigma_k).$$

Then

$$\begin{aligned} \int_{I_n} F^{-1}(t|G) dt &= \sum_k w_k \int_{\xi_{n-1}}^{\xi_n} x f(x|\mu_k, \sigma_k) dx \\ &= \sum_k w_k \{\mu_k \Delta F_{nk} + \sigma_k \Delta T_{nk}\}. \end{aligned}$$

These lead to numerically convenient expression

$$\mathbb{W}_N(G) = \bar{x}^2 + \sum_k w_k \mathbb{E}_k\{X^2\} - 2 \sum_k w_k \{\mu_k \Delta F_{nk} + \sigma_k \Delta T_{nk}\}.$$

To most effectively use BFGS algorithm, it is best to provide gradients of the objective function. Here are some numerically friendly expressions of some partial derivatives.

**Lemma 1** Let  $\delta_{jk} = 1$  when  $j = k$  and  $\delta_{jk} = 0$  when  $j \neq k$ . For  $n = 1, \dots, N$  and  $j = 1, 2, \dots, K$ , we have

$$\begin{aligned}\frac{\partial}{\partial w_j} F(\xi_n|\theta_k) &= f(\xi_n|\theta_k) \frac{\partial \xi_n}{\partial w_j}, \\ \frac{\partial}{\partial \mu_j} F(\xi_n|\theta_k) &= f(\xi_n|\theta_k) \left( \frac{\partial \xi_n}{\partial \mu_j} - \delta_{jk} \right), \\ \frac{\partial}{\partial \sigma_j} F(\xi_n|\theta_k) &= f(\xi_n|\theta_k) \left( \frac{\partial \xi_n}{\partial \sigma_j} - \left\{ \frac{\xi_n - \mu_k}{\sigma_k} \right\} \delta_{jk} \right),\end{aligned}$$

and

$$\begin{aligned}\frac{\partial}{\partial w_j} T \left( \frac{\xi_n - \mu_k}{\sigma_k} \right) &= f(\xi_n|\theta_k) \left( \frac{\xi_n - \mu_k}{\sigma_k} \right) \frac{\partial \xi_i}{\partial w_j}, \\ \frac{\partial}{\partial \mu_j} T \left( \frac{\xi_n - \mu_k}{\sigma_k} \right) &= f(\xi_n|\theta_k) \left( \frac{\xi_n - \mu_k}{\sigma_k} \right) \left( \frac{\partial \xi_n}{\partial \mu_j} - \delta_{jk} \right), \\ \frac{\partial}{\partial \sigma_j} T \left( \frac{\xi_n - \mu_k}{\sigma_k} \right) &= f(\xi_n|\theta_k) \left( \frac{\xi_n - \mu_k}{\sigma_k} \right) \left\{ \frac{\partial \xi_i}{\partial \sigma_j} - \left( \frac{\xi_n - \mu_k}{\sigma_k} \right) \delta_{jk} \right\}.\end{aligned}$$

Furthermore, we have

$$\begin{aligned}\frac{\partial \xi_n}{\partial \mu_k} &= \frac{w_k f(\xi_i|\theta_k)}{f(\xi_n|G)}, \\ \frac{\partial \xi_n}{\partial \sigma_k} &= \frac{w_k f(\xi_n|\theta_k)}{f(\xi_i|G)} \left( \frac{\xi_n - \mu_k}{\sigma_k} \right), \\ \frac{\partial \xi_n}{\partial w_k} &= -\frac{F(\xi_n|\theta_k)}{f(\xi_n|G)}.\end{aligned}$$

Based on this lemma, it is seen that

$$\begin{aligned}\frac{\partial}{\partial \mu_j} \mathbb{W}_N &= 2w_j(\mu_j + \sigma_j \mu_0) - 2w_j \sum_{n=1}^N x_{(n)} \Delta F_{nj} \\ &\quad - 2 \sum_{n=1}^N \sum_k w_k \mu_k x_{(n)} \left\{ \frac{\partial F_0(\xi_n|\theta_k)}{\partial \mu_j} - \frac{\partial F_0(\xi_{n-1}|\theta_k)}{\partial \mu_j} \right\} \\ &\quad - 2 \sum_{n=1}^N \sum_k w_k \sigma_k x_{(n)} \frac{\partial}{\partial \mu_j} \left\{ T \left( \frac{\xi_n - \mu_k}{\sigma_k} \right) - T \left( \frac{\xi_{n-1} - \mu_k}{\sigma_k} \right) \right\}\end{aligned}$$

with  $F_0(\xi_0|\theta_k) = 0$ ,  $F_0(\xi_{N+1}|\theta_k) = 1$ ,  $T\left(\frac{\xi_0 - \mu_k}{\sigma_k}\right) = 0$ , and  $T\left(\frac{\xi_{N+1} - \mu_k}{\sigma_k}\right) = \int_{-\infty}^{\infty} t f_0(t) dt$  is a constant that does not depend on any parameters. Substituting the partial derivatives in Lemma 1, we then get

$$\frac{\partial}{\partial \mu_j} \mathbb{W}_N = 2w_j(\mu_j + \sigma_j \mu_0) - 2w_j \sum_{n=1}^N x_{(n)} \Delta F_{nj}$$

$$\begin{aligned}
& -2 \sum_{n=1}^{N-1} x_{(n)} \xi_n \sum_k w_k f(\xi_n | \mu_k, \sigma_k) \left( \frac{\partial \xi_n}{\partial \mu_j} - \delta_{jk} \right) \\
& + 2 \sum_{n=1}^{N-1} x_{(n)} \xi_{n-1} \sum_k w_k f(\xi_{n-1} | \mu_k, \sigma_k) \left( \frac{\partial \xi_{n-1}}{\partial \mu_j} - \delta_{jk} \right) \\
& = 2w_j \left\{ \mu_j + \sigma_j \mu_0 - \sum_{n=1}^N x_{(n)} \Delta F_{nj} \right\}.
\end{aligned}$$

Similarly, we have

$$\begin{aligned}
\frac{\partial}{\partial \sigma_j} \mathbb{W}_N &= 2w_j \{ \sigma_j (\mu_0^2 + \sigma_0^2) + \mu_j \mu_0 - \sum_{n=1}^N x_{(n)} \Delta \mu_{nj} \}, \\
\frac{\partial}{\partial w_k} \mathbb{W}_N &= \{ \mu_k^2 + \sigma_k^2 (\mu_0^2 + \sigma_0^2) + 2\mu_k \sigma_k \mu_0 \} - 2 \sum_{n=1}^{N-1} \{ x_{(n+1)} - x_{(n)} \} \xi_i F(\xi_n | \theta_k) \\
& \quad - 2 \left\{ \mu_k \sum_{n=1}^N x_{(n)} \Delta F_{nk} + \sigma_k \sum_{n=1}^N x_{(n)} \Delta T_{nk} \right\}.
\end{aligned}$$

Computing the quantiles of the mixture distribution  $F(\cdot|G)$  for each  $G$  is one of the most demanding tasks. The property stated in the following lemma allows us to develop a bi-section algorithm.

**Lemma 2** Let  $F(x|G) = \sum_{k=1}^K F(x|\mu_k, \sigma_k)$  be a  $K$ -component mixture, and  $\xi(t) = F^{-1}(t|G)$  and  $\xi_k(t) = F^{-1}(t|\theta_k)$ , respectively, the  $t$ -quantile of the mixture and its  $k$ th subpopulation. For any  $t \in (0, 1)$ ,

$$\min_k \xi_k(t) \leq \xi(t) \leq \max_k \xi_k(t). \quad (9)$$

**Proof** Since  $F(x|\theta)$  has a continuous CDF, we must have  $F(\xi_k(t)|\theta_k) = t$ . By the monotonicity of the CDF  $F(\cdot|\theta_k)$ , we have

$$F(\min_k \xi_k(t)|\theta_k) \leq F(\xi_k(t)|\theta_k) \leq F(\max_k \xi_k(t)|\theta_k).$$

Multiplying by  $w_k$  and summing over  $k$  lead to

$$F(\min_k \xi_k(t)|G) \leq t \leq F(\max_k \xi_k(t)|G).$$

This implies (9) and completes the proof.  $\square$

In view of this lemma, we can easily find the quantiles of  $F(\cdot|\theta_k)$  to form an interval containing the targeting quantile of  $F(\cdot|G)$ . We can quickly find  $F^{-1}(t|G)$  value through a bi-section algorithm.

## References

- Arjovsky, M., Chintala, S., & Bottou, L. (2017). Wasserstein GAN. Preprint. arXiv:1701.07875.
- Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.
- Chen, J., & Tan, X. (2009). Inference for multivariate normal mixtures. *Journal of Multivariate Analysis*, 100(7), 1367–1383.
- Chen, J., Tan, X., & Zhang, R. (2008). Inference for normal mixtures in mean and variance. *Statistica Sinica*, 18(2), 443–465.
- Chen, J., Li, P., & Liu, G. (2020). Homogeneity testing under finite location-scale mixtures. *Canadian Journal of Statistics*, 48(4), 670–684.
- Choi, K. (1969). Estimators for the parameters of a finite mixture of distributions. *Annals of the Institute of Statistical Mathematics*, 21(1), 107–116.
- Clark, A. (2015). Pillow (PIL Fork) documentation.
- Clarke, B., & Heathcote, C. (1994). Robust estimation of k-component univariate normal mixtures. *Annals of the Institute of Statistical Mathematics*, 46(1), 83–93.
- Cutler, A., & Cordero-Brana, O. I. (1996). Minimum Hellinger distance estimation for finite mixture models. *Journal of the American Statistical Association*, 91(436), 1716–1723.
- Deely, J., & Kruse, R. (1968). Construction of sequences estimating the mixing distribution. *The Annals of Mathematical Statistics*, 39(1), 286–288.
- Evans, S. N., & Matsen, F. A. (2012). The phylogenetic Kantorovich–Rubinstein metric for environmental sequence samples. *Journal of the Royal Statistical Society: Series B (Methodological)*, 74(3), 569–592.
- Farnoosh, R., & Zarpak, B. (2008). Image segmentation using Gaussian mixture model. *IUST International Journal of Engineering Science*, 19(1–2), 29–32.
- Holzmann, H., Munk, A., & Stratmann, B. (2004). Identifiability of finite mixtures-with applications to circular distributions. *Sankhyā: The Indian Journal of Statistics*, 66(3), 440–449.
- Kolouri, S., Rohde, G. K., & Hoffmann, H. (2018). Sliced Wasserstein distance for learning Gaussian mixture models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 3427–3436).
- Nocedal, J., & Wright, S. (2006). *Numerical optimization*. Springer Science & Business Media.
- Pearson, K. (1894). Contributions to the mathematical theory of evolution. *Philosophical Transactions of the Royal Society of London. A*, 185(326-330), 71–110.
- Plataniotis, K. N., & Hatzinak, D. (2000). Gaussian mixtures and their applications to signal processing. In S. Stergiopoulos (Ed.), *Advanced signal processing handbook: theory and implementation for radar, sonar, and medical imaging real time systems* (vol. 25, chapter 3, pp. 3-1–3-35, 1st edn). Boca Raton: CRC Press.
- Santosh, D. H. H., Venkatesh, P., Poornesh, P., Rao, L. N., & Kumar, N. A. (2013). Tracking multiple moving objects using Gaussian mixture model. *International Journal of Soft Computing and Engineering (IJSCE)*, 3(2), 114–119.
- Schork, N. J., Allison, D. B., & Thiel, B. (1996). Mixture distributions in human genetics research. *Statistical Methods in Medical Research*, 5(2), 155–178.
- Tanaka, K. (2009). Strong consistency of the maximum likelihood estimator for finite mixtures of location–scale distributions when penalty is imposed on the ratios of the scale parameters. *Scandinavian Journal of Statistics*, 36(1), 171–184.

- Teicher, H. (1961). Identifiability of mixtures. *The Annals of Mathematical Statistics*, 32(1), 244–248.
- Van der Vaart, A. W. (2000). *Asymptotic statistics* (vol. 3). Cambridge University Press.
- Villani, C. (2003). *Topics in optimal transportation* (vol. 58). American Mathematical Society.
- Woodward, W. A., Parr, W. C., Schucany, W. R., & Lindsey, H. (1984). A comparison of minimum distance and maximum likelihood estimation of a mixture proportion. *Journal of the American Statistical Association*, 79(387), 590–598.
- Yakowitz, S. (1969). A consistent estimator for the identification of finite mixtures. *The Annals of Mathematical Statistics*, 40(5), 1728–1735.
- Zhu, D. (2016). A two-component mixture model for density estimation and classification. *Journal of Interdisciplinary Mathematics*, 19(2), 311–319.