# Semiparametric Mixtures of Regressions with Single-index for Model Based Clustering

Sijia Xiang[*]and Weixin Yao[†]

**Abstract**

In this article, we propose two classes of semiparametric mixture regression models with single-index for model based clustering. Unlike many semiparametric/nonparametric mixture regression models that can only be applied to low dimensional predictors, the new semiparametric models can easily incorporate high dimensional predictors into the nonparametric components. The proposed models are very general, and many of the recently proposed semiparametric/nonparametric mixture regression models are indeed special cases of the new models. Backfitting estimates and the corresponding modified EM algorithms are proposed to achieve optimal convergence rates for both parametric and nonparametric parts. We establish the identifiability results of the proposed two models and investigate the asymptotic properties of the proposed estimation procedures. Simulation studies are conducted to demonstrate the finite sample performance of the proposed models. An application of NBA data by new models reveals some new findings.

*Key words:* EM algorithm, Kernel regression, Mixture regression model, Model based clustering, Single-index models.

# 1   Introduction

Mixtures of regression models are commonly used as model based clustering methods to reveal the relationship among interested variables if the whole population is inhomogeneous and consists of several homogeneous subgroups. They have been widely used in many areas such as econometrics, biology, and epidemiology. For a general account of traditional parametric mixture models, please see, for example, Lindsay (1995), Böhning (1999), McLachlan and Peel (2000), and Frühwirth-Schnatter (2006). However, the traditional mixture of regression models requires strong parametric assumption: liner component regression functions, constant component variance, and constant component proportions. The fully parametric hierarchical mixtures of experts model (Jordan and Jacobs, 1994) has been proposed to allow the component proportions to depend on the covariates in machine learning. Recently, many semiparametric and nonparametric mixture regression models have been proposed to relax the parametric assumption of mixture regression models. See, for example, Young and Hunter (2010); Huang and Yao (2012); Cao and Yao (2012); Huang et al. (2013, 2014), among others. However, most of those existing semparametirc or nonparametric mixture regressions can only be applied for low dimensional predictors due to "curse of dimensionality". It will be desirable to be able to relax parametric assumptions of traditional mixtures of regression models when the dimension of predictors is high.

In this article, we propose a mixture of single-index models (MSIM) and a mixture of regression models with varying single-index proportions (MRSIP) to reduce the

dimension of high dimensional predictors before modeling them nonparametrically. Many existing popular models can be considered as special cases of the proposed two models. Huang et al. (2013) proposed the nonparametric mixture of regression models

$$Y|_{X=x} \sim \sum_{j=1}^{k} \pi_j(x)\phi(Y_i|m_j(x), \sigma_j^2(x)),$$

where $\pi_j(x), m_j(x)$, and $\sigma_j^2(x)$ are unknown smoothing functions, and $\phi(y|\mu, \sigma^2)$ is the normal density with mean $\mu$ and variance $\sigma^2$. Their proposed model can drastically reduce the modelling bias when the strong parametric assumption of traditional mixture of linear regression models does not hold. However, the above model is not applicable to high dimensional predictors due to the kernel estimation used for nonparametric parts. To solve the above problem, we propose a *mixture of single-index models*

$$Y|\boldsymbol{x} \sim \sum_{j=1}^{k} \pi_j(\boldsymbol{\alpha}^T \boldsymbol{x})\phi(Y_i|m_j(\boldsymbol{\alpha}^T \boldsymbol{x}), \sigma_j^2(\boldsymbol{\alpha}^T \boldsymbol{x})), \qquad (1.1)$$

in which the single index $\boldsymbol{\alpha}^T \boldsymbol{x}$ transfers the high dimensional nonparametric problem to a univariate nonparametric problem. When $k = 1$, model (1.1) reduces to a single index model (Ichimura, 1993; Hardle et al., 1993). If $\boldsymbol{x}$ is a scalar, then model (1.1) reduces to the nonparametric mixture of regression model proposed by Huang et al. (2013). Peng (2012) also applied the single index idea to the component means and variance and assumed that component proportions do not depend on the predictor $\boldsymbol{x}$. However, Peng (2012) did not give any theoretical properties of their proposed estimates.

Young and Hunter (2010) and Huang and Yao (2012) proposed a semiparametric

3

mixture of regression models

$$Y|_{X=x} \sim \sum_{j=1}^{k} \pi_j(\boldsymbol{x})\phi(Y_i|\boldsymbol{x}^T\boldsymbol{\beta}_j, \sigma_j^2),$$

where $\pi_j(\boldsymbol{x})$ is an unknown smoothing function, to combine nice properties of both nonparametric mixture regression models and traditional parametric mixture regression models. Their semiparametric mixture models assume that component proportions depend on covariates nonparametrically to reduce the modelling bias while component regression functions are still assumed to be linear to have better model interpretation. However, their estimation procedures cannot be applied if the dimension of predictors $\boldsymbol{x}$ is high due to kernel estimation used for $\pi_j(\boldsymbol{x})$. We propose a *mixture of regression models with varying single-index proportions*

$$Y|_{X=x} \sim \sum_{j=1}^{k} \pi_j(\boldsymbol{\alpha}^T\boldsymbol{x})\phi(Y_i|\boldsymbol{x}^T\boldsymbol{\beta}_j, \sigma_j^2), \tag{1.2}$$

which uses the idea of single index to model the nonparametric effect of predictors on component proportions, while allowing easy interpretation of linear component regression functions. When $k = 1$, model (1.2) reduces to the traditional linear regression model. If $\boldsymbol{x}$ is a scalar, then model (1.2) reduces to the semiparametric mixture models considered by Young and Hunter (2010) and Huang and Yao (2012). Modeling component proportions nonparametrically can reduce the modelling bias and better cluster the data when the traditional parametric assumptions of component proportions do not hold (Young and Hunter, 2010; Huang and Yao, 2012).

We prove the identifiability results of proposed two models under some mild conditions. We propose a modified EM algorithm by combining the ideas of backfitting algorithm, kernel estimation, and local likelihood to estimate global parameters and

nonparametric functions. In addition, the asymptotic properties of the proposed estimation procedures are also investigated. Simulation studies are conducted to demonstrate the finite sample performance of the proposed models. An application of NBA data by new models reveals some new interesting findings.

The rest of the paper is organized as follows. In Section 2, we introduce the MSIM and study its identifiability result. A one-step and a fully-iterated backfitting estimate are proposed, and their asymptotic properties are also studied. In Section 3, we introduce the MRSIP. The identifiability result and asymptotic properties of the proposed estimates are given. In Section 4 and Section 5, we use Monte Carlo studies and a real data example to demonstrate the finite sample performance of the proposed two models. A discussion section is given in Section 6 and we defer the technical conditions and proofs in the supplemental material.

## 2    Mixtures of Single-index Models

### 2.1    Model Definition and Identifiability

Assume that $\{(\boldsymbol{x}_i, Y_i), i = 1, ..., n\}$ is a random sample from the population $(\boldsymbol{x}, Y)$, where $\boldsymbol{x}$ is $p$-dimensional and $Y$ is univariate. Let $\mathcal{C}$ be a latent variable, and has a discrete distribution $P(\mathcal{C} = j|\boldsymbol{x}) = \pi_j(\boldsymbol{\alpha}^T\boldsymbol{x})$ for $j = 1, ..., k$. Conditional on $\mathcal{C} = j$ and $\boldsymbol{x}$, $Y$ follows a normal distribution with mean $m_j(\boldsymbol{\alpha}^T\boldsymbol{x})$ and variance $\sigma_j^2(\boldsymbol{\alpha}^T\boldsymbol{x})$. Without observing $\mathcal{C}$, the conditional distribution of $Y$ given $\boldsymbol{x}$ can be written as:

$$Y|\boldsymbol{x} \sim \sum_{j=1}^{k} \pi_j(\boldsymbol{\alpha}^T\boldsymbol{x})\phi(Y_i|m_j(\boldsymbol{\alpha}^T\boldsymbol{x}), \sigma_j^2(\boldsymbol{\alpha}^T\boldsymbol{x})).$$

The above model is the proposed mixture of single-index models. Throughout the paper, we assume that $k$ is fixed, and refer to model (1.1) as a finite semiparametric

5

mixture of regression models, since $\pi_j(\cdot)$, $m_j(\cdot)$ and $\sigma_j^2(\cdot)$ are all nonparametric. In the model (1.1), we use the same index $\boldsymbol{\alpha}$ for all components. But our proposed estimation procedure and asymptotic results can be easily extended to the cases where components have different index $\boldsymbol{\alpha}$.

Compared to Huang et al. (2013), the appeal of the proposed MSIM is that by using an index $\boldsymbol{\alpha}^T\boldsymbol{x}$, the so-called "curse of dimensionality" in fitting multivariate nonparametric regression functions is avoided. It is of dimension-reduction structure in the sense that, given the estimate of $\boldsymbol{\alpha}$, denoted by $\hat{\boldsymbol{\alpha}}$, we can use the univariate $\hat{\boldsymbol{\alpha}}^T\boldsymbol{x}$ as the covariate and simplify the model (1.1) to the nonparametric mixture regression model proposed by Huang et al. (2013). Therefore, model (1.1) is a reasonable compromise between fully parametric and fully nonparametric modeling.

Identifiability is a major concern for most mixture models. Some well known identifiability results of finite mixture models include: mixture of univariate normals is identifiable up to relabeling (Titterington et al., 1985) and finite mixture of regression models is identifiable up to relabeling provided that covariates have a certain level of variability (Henning, 2000). The following theorem establishes the identifiability result of the model (1.1) and its proof is given in the supplemental material.

**Theorem 2.1.** *Assume that*

1. *$\pi_j(z)$, $m_j(z)$, and $\sigma_j^2(z)$ are differentiable and not constant on the support of $\boldsymbol{\alpha}^T\boldsymbol{x}$, $j = 1, ..., k$;*

2. *The $\boldsymbol{x}$ are continuously distributed random variables that have a joint probability density function;*

3. *The support of $\boldsymbol{x}$ is not contained in any proper linear subspace of $\mathbb{R}^p$;*

4. *$\|\boldsymbol{\alpha}\| = 1$ and the first nonzero element of $\boldsymbol{\alpha}$ is positive;*

5. For any $1 \leq i \neq j \leq k$,

$$\sum_{l=0}^{1} \|m_i^{(l)}(z) - m_j^{(l)}(z)\|^2 + \sum_{l=0}^{1} \|\sigma_i^{(l)}(z) - \sigma_j^{(l)}(z)\|^2 \neq 0,$$

for any $z$ where $g^{(l)}$ is the $l$th derivative of $g$ and equal to $g$ if $l = 0$.

Then, model (1.1) is identifiable.

## 2.2  Estimation Procedure

In this subsection, we propose a one-step estimation procedure and a backfitting algorithm to estimate the nonparametric functions and the single index of the model (1.1).

Let $\ell^{*(1)}(\boldsymbol{\pi}, \boldsymbol{m}, \boldsymbol{\sigma}^2, \boldsymbol{\alpha})$ be the log-likelihood of the collected data $\{(\boldsymbol{x}_i, Y_i), i = 1, ..., n\}$ from the model (1.1). That is:

$$\ell^{*(1)}(\boldsymbol{\pi}, \boldsymbol{m}, \boldsymbol{\sigma}^2, \boldsymbol{\alpha}) = \sum_{i=1}^{n} \log\{\sum_{j=1}^{k} \pi_j(\boldsymbol{\alpha}^T \boldsymbol{x}_i)\phi(Y_i | m_j(\boldsymbol{\alpha}^T \boldsymbol{x}_i), \sigma_j^2(\boldsymbol{\alpha}^T \boldsymbol{x}_i))\}, \qquad (2.1)$$

where $\boldsymbol{\pi}(\cdot) = \{\pi_1(\cdot), ..., \pi_{k-1}(\cdot)\}^T$, $\boldsymbol{m}(\cdot) = \{m_1(\cdot), ..., m_k(\cdot)\}^T$, and $\boldsymbol{\sigma}^2(\cdot) = \{\sigma_1^2(\cdot), ..., \sigma_k^2(\cdot)\}^T$. Since $\boldsymbol{\pi}(\cdot)$, $\boldsymbol{m}(\cdot)$ and $\boldsymbol{\sigma}^2(\cdot)$ consist of nonparametric functions, (2.1) is not ready for maximization.

Note that for the model (1.1), the space spanned by the single index $\boldsymbol{\alpha}$ is in fact the central mean subspace of $Y|\boldsymbol{x}$ (Cook and Li, 2002) in the literature of sufficient dimension reduction. Therefore, we can employ existing sufficient dimension reduction methods to find an initial estimate of $\boldsymbol{\alpha}$. Please see, for example, Li (1991); Li, Zha, and Chiaromonte (2005); Wang and Xia (2008); Luo, Wang, and Tsai (2009); Ma and Zhu (2012a,b). In this article, we will simply employ sliced inverse regression (Li, 1991) to obtain an initial estimate of $\boldsymbol{\alpha}$, denoted by $\tilde{\boldsymbol{\alpha}}$.

Given the estimated single index $\tilde{\boldsymbol{\alpha}}$, the nonparametric functions $\boldsymbol{\pi}(z)$, $\boldsymbol{m}(z)$ and $\boldsymbol{\sigma}^2(z)$ can then be estimated by maximizing the following local log-likelihood function:

$$\ell_1^{(1)}(\boldsymbol{\pi}, \boldsymbol{m}, \boldsymbol{\sigma}^2) = \sum_{i=1}^{n} \log\{\sum_{j=1}^{k} \pi_j(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i)\phi(Y_i|m_j(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i), \sigma_j^2(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i))\}K_h(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i - z),$$

(2.2)

where $K_h(z) = \frac{1}{h}K(\frac{z}{h})$, $K(\cdot)$ is a kernel density function, and $h$ is a tuning parameter. Let $\hat{\boldsymbol{\pi}}(\cdot)$, $\hat{\boldsymbol{m}}(\cdot)$ and $\hat{\boldsymbol{\sigma}}^2(\cdot)$ be the estimates that maximize (2.2). The above estimates are the proposed *one-step estimate.*

We propose a modified EM-type algorithm to maximize $\ell_1^{(1)}$. In practice, we usually want to evaluate unknown functions at a set of grid points, which in this case, requires us to maximize local log-likelihood functions at a set of grid points. If we simply employ the EM algorithm separately for each grid point, the labels in the EM algorithm may change at different grid points, and we may not be able to get smoothed estimated curves (Huang and Yao, 2012). Therefore, we propose the following modified EM-type algorithm, which estimates the nonparametric functions simultaneously at a set of grid points, say $\{u_t, t = 1, ..., N\}$, and provides a unified label of each observation across all grid points.

**Algorithm 2.1.** *Modified EM-type algorithm to maximize (2.2) given the single index estimate $\tilde{\boldsymbol{\alpha}}$.*

**E-step:** *Calculate the expectations of component labels based on estimates from $l^{th}$ iteration:*

$$p_{ij}^{(l+1)} = \frac{\pi_j^{(l)}(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i)\phi(Y_i|m_j^{(l)}(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i), \sigma_j^{2(l)}(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i))}{\sum_{j=1}^{k} \pi_j^{(l)}(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i)\phi(Y_i|m_j^{(l)}(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i), \sigma_j^{2(l)}(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i))},$$

(2.3)

*where $i = 1, \ldots, n, j = 1, \ldots, k$.*

8

**M-step:** *Update the estimates*

$$\pi_j^{(l+1)}(z) = \frac{\sum_{i=1}^n p_{ij}^{(l+1)} K_h(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i - z)}{\sum_{i=1}^n K_h(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i - z)}, \tag{2.4}$$

$$m_j^{(l+1)}(z) = \frac{\sum_{i=1}^n p_{ij}^{(l+1)} Y_i K_h(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i - z)}{\sum_{i=1}^n p_{ij}^{(l+1)} K_h(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i - z)}, \tag{2.5}$$

$$\sigma_j^{2(l+1)}(z) = \frac{\sum_{i=1}^n p_{ij}^{(l+1)} (Y_i - m_j^{(l+1)}(z))^2 K_h(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i - z)}{\sum_{i=1}^n p_{ij}^{(l+1)} K_h(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i - z)}, \tag{2.6}$$

*for $z \in \{u_t, t = 1, ..., N\}$ and $j = 1, \ldots, k$. We then update $\pi_j^{(l+1)}(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i)$, $m_j^{(l+1)}(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i)$ and $\sigma_j^{2(l+1)}(\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i)$, $i = 1, ..., n$, by linear interpolating $\pi_j^{(l+1)}(u_t)$, $m_j^{(l+1)}(u_t)$ and $\sigma_j^{2(l+1)}(u_t)$, $t = 1, ..., N$, respectively.*

Note that in the M-step, the nonparametric functions are estimated simultaneously at a set of grid points, and therefore, the classification probabilities in the the E-step can be estimated globally to avoid the label switching problem (Yao and Lindsay, 2009). If the sample size $n$ is not too large, one can also take all $\{\tilde{\boldsymbol{\alpha}}^T \boldsymbol{x}_i, i = 1, \ldots, n\}$ as grid points for $z$ in the M-step.

The initial estimate $\tilde{\boldsymbol{\alpha}}$ by SIR does not make use of the mixture information and thus is not efficient. Given one step estimate $\hat{\boldsymbol{\pi}}(\cdot)$, $\hat{\boldsymbol{m}}(\cdot)$ and $\hat{\boldsymbol{\sigma}}^2(\cdot)$, we can further improve the estimate of $\boldsymbol{\alpha}$ by maximizing

$$\ell_2^{(1)}(\boldsymbol{\alpha}) = \sum_{i=1}^n \log\{\sum_{j=1}^k \hat{\pi}_j(\boldsymbol{\alpha}^T \boldsymbol{x}_i)\phi(Y_i|\hat{m}_j(\boldsymbol{\alpha}^T \boldsymbol{x}_i), \hat{\sigma}_j^2(\boldsymbol{\alpha}^T \boldsymbol{x}_i))\}, \tag{2.7}$$

with respect to $\boldsymbol{\alpha}$. The proposed *fully iterative backfitting estimator* of $\boldsymbol{\alpha}$, denoted by $\hat{\boldsymbol{\alpha}}$, iterates the above two steps until convergence.

**Algorithm 2.2.** *Fully iterative backfitting estimator (FIB)*

**Step 1:** *Apply sliced inverse regression (SIR) to obtain an initial estimate of the*

9

*single index parameter $\boldsymbol{\alpha}$, denoted by $\tilde{\boldsymbol{\alpha}}$.*

**Step 2:** *Given $\tilde{\boldsymbol{\alpha}}$, apply the modified EM-algorithm (2.3)—(2.6) to maximize $\ell_1^{(1)}$ in (2.2) to obtain the estimates $\hat{\boldsymbol{\pi}}(\cdot)$, $\hat{\boldsymbol{m}}(\cdot)$, and $\hat{\boldsymbol{\sigma}}^2(\cdot)$.*

**Step 3:** *Given $\hat{\boldsymbol{\pi}}(\cdot)$, $\hat{\boldsymbol{m}}(\cdot)$, and $\hat{\boldsymbol{\sigma}}^2(\cdot)$ from Step 2, update the estimate of $\boldsymbol{\alpha}$ by maximizing $\ell_2^{(1)}$ in (2.7).*

**Step 4:** *Iterate Steps 2 - 3 until convergence.*

## 2.3  Asymptotic Properties

The asymptotic properties of the proposed estimates are investigated below. Let $\boldsymbol{\theta}(z) = (\boldsymbol{\pi}^T(z), \boldsymbol{m}^T(z), (\boldsymbol{\sigma}^2)^T(z))^T$. Define

$$\ell(\boldsymbol{\theta}(z), y) = \log \sum_{j=1}^{k} \pi_j(z)\phi\{y | m_j(z), \sigma_j^2(z)\},$$

$$q_1(z) = \frac{\partial \ell(\boldsymbol{\theta}(z), y)}{\partial \theta},$$

$$q_2(z) = \frac{\partial^2 \ell(\boldsymbol{\theta}(z), y)}{\partial \theta \partial \theta^T},$$

$$\mathcal{I}_\theta^{(1)}(z) = -E[q_2(Z) | Z = z],$$

$$\Lambda_1(u|z) = E[q_1(z) | Z = u].$$

Under further conditions defined in the supplemental material, the asymptotic properties of the one-step estimates $\hat{\boldsymbol{\pi}}(\cdot)$, $\hat{\boldsymbol{m}}(\cdot)$, and $\hat{\boldsymbol{\sigma}}^2(\cdot)$ are given in the following theorem.

**Theorem 2.2.** *Assume that conditions (C1)-(C7) in the supplemental material hold.*

*Then, as $n \to \infty$, $h \to 0$ and $nh \to \infty$, we have*

$$\sqrt{nh}\{\hat{\boldsymbol{\theta}}(z) - \boldsymbol{\theta}(z) - \mathcal{B}_1 + o_p(h^2)\} \xrightarrow{D} N\{0, \nu_0 f^{-1}(z)\mathcal{I}_\theta^{(1)}(z)\}, \qquad (2.8)$$

*where*

$$\mathcal{B}_1(z) = \mathcal{I}_\theta^{(1)-1} \left\{ \frac{f'(z)\Lambda_1'(z|z)}{f(z)} + \frac{1}{2}\Lambda_1''(z|z) \right\} \kappa_2 h^2,$$

*with $f(\cdot)$ the marginal density function of $\boldsymbol{\alpha}^T \boldsymbol{x}$, $\kappa_l = \int t^l K(t)dt$ and $\nu_l = \int t^l K^2(t)dt$.*

Note that the asymptotic variance of $\hat{\boldsymbol{\theta}}(z)$ is the same as those given in Huang et al. (2013). Thus, the nonparametric functions can be estimated with the same accuracy as it would have if the single index $\boldsymbol{\alpha}^T \boldsymbol{x}$ were known. This is expected since the single index $\boldsymbol{\alpha}$ can be estimated at a root $n$ convergence rate which is much faster than $\hat{\boldsymbol{\theta}}(z)$. In addition, note that the one-step estimates of $\boldsymbol{\theta}(z)$ have the same asymptotic variance (up to the first order) as the full iterative backfitting algorithm but with much less computations. Our simulation results in Section 4 further confirm this result.

The next theorem gives the asymptotic results of the $\hat{\boldsymbol{\alpha}}$ given by full iterative backfitting algorithm.

**Theorem 2.3.** *Assume that conditions (C1)-(C8) in the supplemental material hold. Then, as $n \to \infty$, $nh^4 \to 0$, and $nh^2/\log(1/h) \to \infty$,*

$$\sqrt{n}(\hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}) \xrightarrow{D} N(0, \boldsymbol{Q}_1^{-1}), \qquad (2.9)$$

*where*

$$\boldsymbol{Q}_1 = E\left[ \{\boldsymbol{x}\boldsymbol{\theta}'(Z)\}q_2(Z)\{\boldsymbol{x}\boldsymbol{\theta}'(Z)\}^T - \boldsymbol{x}\boldsymbol{\theta}'(Z)q_2(Z)\mathcal{I}_\theta^{(1)-1}(Z)E\{q_2(Z)[\boldsymbol{x}\boldsymbol{\theta}'(Z)]^T|Z\} \right].$$

# 3 Mixtures of Regression Models with Varying Single-Index Proportions

## 3.1 Model Definition and Identifiability

The MRSIP assumes that $P(\mathcal{C} = j|\boldsymbol{x}) = \pi_j(\boldsymbol{\alpha}^T\boldsymbol{x})$ for $j = 1, ..., k$, and conditional on $\mathcal{C} = j$ and $\boldsymbol{x}$, $Y$ follows a normal distribution with mean $\boldsymbol{x}^T\boldsymbol{\beta}_j$ and variance $\sigma_j^2$. That is,

$$Y|\boldsymbol{x} \sim \sum_{j=1}^{k} \pi_j(\boldsymbol{\alpha}^T\boldsymbol{x})N(\boldsymbol{x}^T\boldsymbol{\beta}_j, \sigma_j^2).$$

Since $\pi_j(\cdot)$'s are nonparametric, model (1.2) is also a finite semiparametric mixture of regression models. The linear component regression functions $\boldsymbol{x}^T\boldsymbol{\beta}_j$ enjoy simple interpretation, while nonparametric functions $\pi_j(\boldsymbol{\alpha}^T\boldsymbol{x})$ can incorporate the effects of predictors on component proportions more flexibly to reduce the modeling bias. See Young and Hunter (2010); Huang et al. (2013) for more information. We first prove the identifiability result of the model (1.2) in the following theorem and its proof is given in the supplemental material.

**Theorem 3.1.** *Assume that*

1. *$\pi_j(z) > 0$ are differentiable and not constant on the support of $\boldsymbol{\alpha}^T\boldsymbol{x}$, $j = 1, ..., k$;*

2. *The component of $\boldsymbol{x}$ are continuously distributed random variables that have a joint probability density function;*

3. *The support of $\boldsymbol{x}$ contains an open set in $\mathbb{R}^p$ and is not contained in any proper linear subspace of $\mathbb{R}^p$;*

4. *$\|\boldsymbol{\alpha}\| = 1$ and the first nonzero element of $\boldsymbol{\alpha}$ is positive;*

5. $(\boldsymbol{\beta}_j, \sigma_j^2)$, $j = 1, ..., k$, are distinct pairs.

Then, model (1.2) is identifiable.

## 3.2    Estimation Procedure

The log-likelihood of the collected data for the model (1.2) is:

$$\ell^{*(2)}(\boldsymbol{\pi}, \boldsymbol{\sigma}^2, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \sum_{i=1}^{n} \log\{\sum_{j=1}^{k} \pi_j(\boldsymbol{\alpha}^T \boldsymbol{x}_i)\phi(Y_i|\boldsymbol{x}_i^T\boldsymbol{\beta}_j, \sigma_j^2)\}, \qquad (3.1)$$

where $\boldsymbol{\pi}(\cdot) = \{\pi_1(\cdot), ..., \pi_{k-1}(\cdot)\}^T$, $\boldsymbol{\sigma}^2 = \{\sigma_1^2, ..., \sigma_k^2\}^T$, and $\boldsymbol{\beta} = \{\boldsymbol{\beta}_1, ..., \boldsymbol{\beta}_k\}^T$. Since $\boldsymbol{\pi}(\cdot)$ consists of nonparametric functions, (3.1) is not ready for maximization. We propose a backfitting algorithm to iterate between estimating the parameters $(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\sigma}^2)$ and the nonparametric functions $\boldsymbol{\pi}(\cdot)$.

Given the estimates of $(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\sigma}^2)$, say $(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\sigma}}^2)$, then $\boldsymbol{\pi}(\cdot)$ can be estimated locally by maximizing the following local log-likelihood function:

$$\ell_1^{(2)}(\boldsymbol{\pi}) = \sum_{i=1}^{n} \log\{\sum_{j=1}^{k} \pi_j(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_i)\phi(Y_i|\boldsymbol{x}_i^T\hat{\boldsymbol{\beta}}_j, \hat{\sigma}_j^2)\}K_h(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_i - z). \qquad (3.2)$$

Let $\hat{\boldsymbol{\pi}}(\cdot)$ be the estimate that maximizes (3.2). We can then further update the estimate of $(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\sigma}^2)$ by maximizing

$$\ell_2^{(2)}(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\sigma}^2) = \sum_{i=1}^{n} \log\{\sum_{j=1}^{k} \hat{\pi}_j(\boldsymbol{\alpha}^T \boldsymbol{x}_i)\phi(Y_i|\boldsymbol{x}_i^T\boldsymbol{\beta}_j, \sigma_j^2)\}. \qquad (3.3)$$

The backfitting algorithm by iterating the above two steps can be summarized as follows.

**Algorithm 3.1.** *Backfitting algorithm to estimate the model (1.2).*

**Step 1:** *Obtain an initial estimate of $(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\sigma}^2)$.*

**Step 2:** *Given $(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\sigma}}^2)$, use the following modified EM-type algorithm to maximize $\ell_1^{(2)}$ in (3.2).*

*__E-step:__ Calculate the expectations of component labels based on estimates from $l^{th}$ iteration:*

$$p_{ij}^{(l+1)} = \frac{\pi_j^{(l)}(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_i)\phi(Y_i|\boldsymbol{x}_i^T \hat{\boldsymbol{\beta}}_j, \hat{\sigma}_j^2)}{\sum_{j=1}^{k} \pi_j^{(l)}(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_i)\phi(Y_i|\boldsymbol{x}_i^T \hat{\boldsymbol{\beta}}_j, \hat{\sigma}_j^2)}, \tag{3.4}$$

*where $i = 1, \ldots, n, j = 1, ..., k$. __M-step:__ Update the estimate*

$$\pi_j^{(l+1)}(z) = \frac{\sum_{i=1}^{n} p_{ij}^{(l+1)} K_h(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_i - z)}{\sum_{i=1}^{n} K_h(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_i - z)} \tag{3.5}$$

*for $z \in \{u_t, t = 1, ..., N\}$. We then update $\pi_j^{(l+1)}(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_i)$, $i = 1, ..., n$ by linear interpolating $\pi_j^{(l+1)}(u_t)$, $t = 1, ..., N$.*

**Step 3:** *Given $\hat{\boldsymbol{\pi}}(\cdot)$ from Step 2, update $(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\sigma}}^2)$ by maximizing (3.3). We propose to iterate between updating $\boldsymbol{\alpha}$ and $(\boldsymbol{\beta}, \boldsymbol{\sigma})$.*

*Step 3.1: Given $\hat{\boldsymbol{\alpha}}$, update $(\boldsymbol{\beta}, \boldsymbol{\sigma}^2)$.*

*__E-step:__ Calculate the classification probabilities:*

$$p_{ij}^{(l+1)} = \frac{\hat{\pi}_j(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_i)\phi(Y_i|\boldsymbol{x}_i^T \boldsymbol{\beta}_j^{(l)}, \sigma_j^{2(l)})}{\sum_{j=1}^{k} \hat{\pi}_j(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_i)\phi(Y_i|\boldsymbol{x}_i^T \boldsymbol{\beta}_j^{(l)}, \sigma_j^{2(l)})}, \quad j = 1, ..., k. \tag{3.6}$$

*__M-step:__ Update $\boldsymbol{\beta}$ and $\boldsymbol{\sigma}^2$:*

$$\boldsymbol{\beta}_j^{(l+1)} = (\boldsymbol{S}^T \boldsymbol{R}_j^{(l+1)} \boldsymbol{S})^{-1} \boldsymbol{S}^T \boldsymbol{R}_j^{(l+1)} \boldsymbol{y}, \tag{3.7}$$

$$\sigma_j^{2(l+1)} = \frac{\sum_{i=1}^{n} p_{ij}^{(l+1)}(Y_i - \boldsymbol{x}_i^T \boldsymbol{\beta}_j^{(l+1)})^2}{\sum_{i=1}^{n} p_{ij}^{(l+1)}}, \tag{3.8}$$

14

*where $j = 1, ..., k$, $\boldsymbol{R}_j^{(l+1)} = diag\{p_{ij}^{(l+1)}, ..., p_{nj}^{(l+1)}\}$, and $\boldsymbol{S} = (\boldsymbol{x}_1, ..., \boldsymbol{x}_n)^T$.*

*Step 3.2: Given $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\sigma}}^2)$, update $\boldsymbol{\alpha}$ by maximizing the following log-likelihood*

$$\ell_3^{(2)}(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \log\{\sum_{j=1}^{k} \hat{\pi}_j(\boldsymbol{\alpha}^T \boldsymbol{x}_i) \phi(Y_i | \boldsymbol{x}_i^T \hat{\boldsymbol{\beta}}_j, \hat{\sigma}_j^2)\}.$$

*Step 3.3: Iterate Steps 3.1-3.2 until convergence.*

**Step 4:** *Iterate Steps 2-3 until convergence.*

There are many ways to obtain an initial estimate of $(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\sigma}^2)$. In our numerical studies, we get an initial estimate of $(\boldsymbol{\beta}, \boldsymbol{\sigma}^2)$ by fitting traditional mixtures of linear regression models. Using resulting hard-clustering results as new response variable, we apply SIR to get an initial estimate of $\boldsymbol{\alpha}$.

## 3.3 Asymptotic Properties

Let $(\hat{\boldsymbol{\pi}}(z), \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\sigma}}^2)$ be the resulting estimate of backfitting Algorithm 3.1. In this section, we investigate their asymptotic properties. Let $\boldsymbol{\eta} = (\boldsymbol{\beta}^T, (\boldsymbol{\sigma}^2)^T)^T$ and $\boldsymbol{\lambda} = (\boldsymbol{\alpha}^T, \boldsymbol{\eta}^T)^T$. Define

$$\ell(\boldsymbol{\pi}(z), \boldsymbol{\lambda}, \boldsymbol{x}, y) = \log \sum_{j=1}^{k} \pi_j(z) \phi\{y | \boldsymbol{x}^T \boldsymbol{\beta}_j, \sigma_j^2\},$$

$$q_\pi(z) = \frac{\partial \ell(\boldsymbol{\pi}(z), \lambda, x, y)}{\partial \boldsymbol{\pi}},$$

$$q_{\pi\pi}(z) = \frac{\partial^2 \ell(\boldsymbol{\pi}(z), \lambda, x, y)}{\partial \boldsymbol{\pi} \partial \boldsymbol{\pi}^T}.$$

Similarly, define $q_\lambda$, $q_{\lambda\lambda}$, and $q_{\pi\eta}$. Denote $\mathcal{I}_\pi^{(2)}(z) = -E[q_{\pi\pi}(Z)|Z = z]$ and $\Lambda_2(u|z) = E[q_\pi(z)|Z = u]$.

Under some regularity conditions, the asymptotic properties of $\hat{\boldsymbol{\pi}}(z)$ are given in the following theorem and its proof is given in the supplemental material.

**Theorem 3.2.** *Assume that conditions (C1)-(C4) and (C9)-(C11) in the supplemental material hold. Then, as $n \to \infty$, $h \to 0$ and $nh \to \infty$, we have*

$$\sqrt{nh}\{\hat{\boldsymbol{\pi}}(z) - \boldsymbol{\pi}(z) - \mathcal{B}_2(z) + o_p(h^2)\} \xrightarrow{D} N\{0, \nu_0 f^{-1}(z)\mathcal{I}_\pi^{(2)}(z)\}, \qquad (3.9)$$

*where*

$$\mathcal{B}_2(z) = \mathcal{I}_\pi^{(2)-1}\left\{\frac{f'(z)\Lambda_2'(z|z)}{f(z)} + \frac{1}{2}\Lambda_2''(z|z)\right\}\kappa_2 h^2.$$

The asymptotic property of the parametric estimate $\hat{\boldsymbol{\lambda}}$ is given in the following theorem and its proof is given in the supplemental material.

**Theorem 3.3.** *Assume that conditions (C1)-(C4) and (C9)-(C12) in the supplemental material hold. Then, as $n \to \infty$, $nh^4 \to 0$, and $nh^2/\log(1/h) \to \infty$,*

$$\sqrt{n}(\hat{\boldsymbol{\lambda}} - \boldsymbol{\lambda}) \xrightarrow{D} N(0, \boldsymbol{Q}_2^{-1}),$$

*where,*

$$\boldsymbol{Q}_2 = E\left[q_{\pi\pi}(Z)\begin{pmatrix}\boldsymbol{x}\boldsymbol{\pi}'(Z)\\ \boldsymbol{I}\end{pmatrix}\left\{\begin{pmatrix}\boldsymbol{x}\boldsymbol{\pi}'(Z)\\ \boldsymbol{I}\end{pmatrix} - \begin{pmatrix}\mathcal{I}_\pi^{(2)-1}(Z)E\{q_{\pi\pi}(Z)(\boldsymbol{x}\boldsymbol{\pi}'(Z))^T|Z\}\\ \mathcal{I}_\pi^{(2)-1}(Z)E\{q_{\pi\eta}(Z)|Z\}\end{pmatrix}\right\}^T\right].$$

# 4   Simulation Studies

In this section, we conduct simulation studies to test the performance of the proposed models and estimation procedures.

The performance of the estimates of the mean functions $m_j(\cdot)$'s in the model (1.1)

16

is measured by the square root of the average square errors (RASE)

$$RASE_m^2 = N^{-1} \sum_{j=1}^{k} \sum_{t=1}^{N} [\hat{m}_j(u_t) - m_j(u_t)]^2.$$

In our simulation, we set $N = 100$. Similarly, we can define the $RASE$ for variance functions $\sigma_j^2(\cdot)$'s and proportion functions $\pi_j(\cdot)$'s, denoted by $RASE_{\sigma^2}$ and $RASE_{\pi}$, respectively.

*Example 1:* We conduct a simulation for a two-component MSIM:

$$\pi_1(z) = 0.5 + 0.3\sin(\pi z) \text{ and } \pi_2(z) = 1 - \pi_1(z),$$

$$m_1(z) = 3 - \sin(2\pi z/\sqrt{3}) \text{ and } m_2(z) = \cos(\sqrt{3}\pi z),$$

$$\sigma_1(z) = 0.7 + \sin(3\pi z)/15 \text{ and } \sigma_2(z) = 0.3 + \cos(1.3\pi z)/10.$$

where $z_i = \boldsymbol{\alpha}^T \boldsymbol{x}_i$, $\boldsymbol{x}_i$ are trivariate with independent uniform (0,1) components, and the direction parameter is $\boldsymbol{\alpha} = (1,1,1)/\sqrt{3}$. The sample sizes $n = 200$, $n = 400$, and $n = 800$ are conducted over 500 repetitions. To estimate $\boldsymbol{\alpha}$, we use sliced inverse regression (SIR) and the fully iterative backfitting estimate (FIB). To estimate the nonparametric functions, we apply the one-step estimate (OS) and FIB. For FIB, we use both true value (T) and SIR (S) as the initial values.

We first select a proper bandwidth for estimating $\boldsymbol{\pi}(\cdot)$, $\boldsymbol{m}(\cdot)$ and $\boldsymbol{\sigma}^2(\cdot)$. Based on Theorem 2.2, one can calculate theoretical optimal bandwidth by minimizing asymptotic mean squared errors. However, the theoretical optimal bandwidth depends on many unknown quantities, which are not easy to estimate in practice. In our examples, we propose to use the following cross-validation (CV) method to choose the bandwidth. Let $\mathscr{D}$ be the full data set, and divide $\mathscr{D}$ into a training set $\mathscr{R}_l$ and a test set $\mathscr{T}_l$. That is, $\mathscr{R}_l \cup \mathscr{T}_l = \mathscr{D}$ for $l = 1, ..., L$. We use the training set $\mathscr{R}_l$ to obtain the estimates $\{\hat{\boldsymbol{\pi}}(\cdot), \hat{\boldsymbol{m}}(\cdot), \hat{\boldsymbol{\sigma}}^2(\cdot), \hat{\boldsymbol{\alpha}}\}$. We then evaluate $\boldsymbol{\pi}(\cdot)$, $\boldsymbol{m}(\cdot)$ and $\boldsymbol{\sigma}^2(\cdot)$ for the

test data set. For each $(\boldsymbol{x}_t, y_t) \in \mathscr{T}_l$, we calculate the classification probability as

$$\hat{p}_{tj} = \frac{\hat{\pi}_j(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_t) \phi(y_t | \hat{m}_j(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_t), \hat{\sigma}_j^2(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_t))}{\sum_{j=1}^k \hat{\pi}_j(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_t) \phi(y_t | \hat{m}_j(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_t), \hat{\sigma}_j^2(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_t))}, \tag{4.1}$$

for $j = 1, ..., k$. We consider the regular $CV$, which is defined by

$$CV(h) = \sum_{l=1}^L \sum_{t \in \mathscr{T}_l} (y_t - \hat{y}_t)^2,$$

where $\hat{y}_t = \sum_{j=1}^k \hat{p}_{tj} \hat{m}_j(\hat{\boldsymbol{\alpha}}^T \boldsymbol{x}_t)$. We also implemented the likelihood based cross validation to choose the bandwidth and the results are similar but with more computations.

We set $L = 10$ and randomly partition the data. We repeat the procedure 30 times, and take the average of the selected bandwidth as the optimal bandwidth, denoted by $\hat{h}$. In the simulation, we consider three different bandwidths, $\hat{h} \times n^{-2/15}$, $\hat{h}$ and $1.5\hat{h}$, which correspond to the under-smoothing, appropriate smoothing and over-smoothing condition, respectively.

Table 1 reports the MSEs of $\hat{\boldsymbol{\alpha}}$ (true value times 100) and Table 2 contains the mean and standard deviation of $RASE_\pi$, $RASE_m$, and $RASE_{\sigma^2}$. Based on Table 1, we can see that the proposed fully iterative backfitting estimates (FIB) give much better results than SIR, which is reasonable since FIB makes use of mixture information while SIR does not. Based on Table 2, we can see that OS provides close estimates to FIB, although FIB generally provides slightly smaller RASEs than OS for finite sample size. This verified the theoretical results stated in Section 2.3.

In addition, from Tables 1 and 2, we can see that the proposed bandwidth selection procedure based on cross validation works reasonably well since the appropriate bandwidths chosen by CV usually provide the estimate that is or is close to the best

one. Furthermore, FIB(S) provides similar results to FIB(T). Therefore, SIR provides good initial values for the proposed fully iterative estimates.

Table 1: MSE of $\hat{\boldsymbol{\alpha}}$ (true value times 100) for Example 1.

| | | SIR | FIB(T) | | | FIB(S) | | |
|---|---|---|---|---|---|---|---|---|
| | | | $h = 0.054$ | $h = 0.109$ | $h = 0.164$ | $h = 0.054$ | $h = 0.109$ | $h = 0.164$ |
| | $\alpha_1$ | 0.881 | 0.099 | 0.126 | 0.128 | 0.287 | 0.130 | 0.147 |
| $n = 200$ | $\alpha_2$ | 0.829 | 0.113 | 0.144 | 0.124 | 0.324 | 0.144 | 0.137 |
| | $\alpha_3$ | 1.066 | 0.110 | 0.152 | 0.137 | 0.388 | 0.154 | 0.167 |
| | | | $h = 0.045$ | $h = 0.100$ | $h = 0.149$ | $h = 0.045$ | $h = 0.100$ | $h = 0.149$ |
| | $\alpha_1$ | 0.435 | 0.066 | 0.046 | 0.046 | 0.125 | 0.050 | 0.045 |
| $n = 400$ | $\alpha_2$ | 0.447 | 0.063 | 0.054 | 0.051 | 0.121 | 0.055 | 0.052 |
| | $\alpha_3$ | 0.411 | 0.062 | 0.052 | 0.052 | 0.123 | 0.053 | 0.052 |
| | | | $h = 0.037$ | $h = 0.091$ | $h = 0.137$ | $h = 0.037$ | $h = 0.091$ | $h = 0.137$ |
| | $\alpha_1$ | 0.215 | 0.047 | 0.022 | 0.029 | 0.063 | 0.035 | 0.024 |
| $n = 800$ | $\alpha_2$ | 0.256 | 0.034 | 0.035 | 0.040 | 0.044 | 0.029 | 0.027 |
| | $\alpha_3$ | 0.226 | 0.065 | 0.031 | 0.058 | 0.062 | 0.050 | 0.030 |

*Example 2:* We conduct a simulation for a two-component MRSIP:

$$\pi_1(z) = 0.5 - 0.35 \sin(\pi z) \text{ and } \pi_2(z) = 1 - \pi_1(z),$$

$$m_1(\boldsymbol{x}) = 1 + 3x_2 \text{ and } m_2(\boldsymbol{x}) = -1 + 2x_1 + 3x_3,$$

$$\sigma_1^2 = 0.7 \text{ and } \sigma_2^2 = 0.6,$$

where $m_1(\boldsymbol{x})$ and $m_2(\boldsymbol{x})$ are the regression functions for the first and second components, respectively. Therefore, $\boldsymbol{\beta}_1 = (1, 0, 3, 0)$ and $\boldsymbol{\beta}_2 = (-1, 2, 0, 3)$. $\boldsymbol{x}_i$ are trivariate with independent uniform $(0,1)$ components, and the single index parameter is $\boldsymbol{\alpha} = (1, 1, 1)/\sqrt{3}$. MRSIP with true value (T) and SIR (S) as initial values are used to fit the data, and the results are compared to the traditional mixture of linear regression models (MixLinReg). The bandwidth for MRSIP is chosen based on the cross validation similar to Example 1.

Table 2: Mean and Standard Deviation of RASEs for Example 1.

| | OS | FIB(T) | | | FIB(S) | | |
|---|---|---|---|---|---|---|---|
| n=200 | $h = 0.125$ | $h = 0.054$ | $h = 0.109$ | $h = 0.164$ | $h = 0.054$ | $h = 0.109$ | $h = 0.164$ |
| $\pi$ | 0.044(0.017) | 0.057(0.015) | 0.043(0.016) | 0.049(0.017) | 0.058(0.015) | 0.043(0.016) | 0.049(0.017) |
| $\mu$ | 0.227(0.063) | 0.181(0.098) | 0.176(0.046) | 0.287(0.056) | 0.178(0.086) | 0.177(0.051) | 0.288(0.059) |
| $\sigma^2$ | 0.197(0.084) | 0.175(0.169) | 0.163(0.081) | 0.246(0.071) | 0.162(0.131) | 0.164(0.095) | 0.247(0.080) |
| n=400 | $h = 0.108$ | $h = 0.045$ | $h = 0.100$ | $h = 0.149$ | $h = 0.045$ | $h = 0.100$ | $h = 0.149$ |
| $\pi$ | 0.023(0.008) | 0.032(0.008) | 0.023(0.008) | 0.027(0.009) | 0.032(0.008) | 0.023(0.008) | 0.027(0.009) |
| $\mu$ | 0.118(0.022) | 0.093(0.045) | 0.100(0.022) | 0.169(0.020) | 0.094(0.046) | 0.100(0.022) | 0.169(0.020) |
| $\sigma^2$ | 0.104(0.035) | 0.089(0.077) | 0.093(0.045) | 0.143(0.028) | 0.089(0.077) | 0.093(0.045) | 0.143(0.028) |
| n=800 | $h = 0.094$ | $h = 0.037$ | $h = 0.091$ | $h = 0.137$ | $h = 0.037$ | $h = 0.091$ | $h = 0.137$ |
| $\pi$ | 0.013(0.004) | 0.017(0.003) | 0.012(0.004) | 0.016(0.004) | 0.017(0.003) | 0.012(0.004) | 0.016(0.004) |
| $\mu$ | 0.062(0.010) | 0.050(0.023) | 0.056(0.010) | 0.102(0.011) | 0.050(0.023) | 0.056(0.010) | 0.101(0.010) |
| $\sigma^2$ | 0.055(0.015) | 0.049(0.046) | 0.052(0.015) | 0.086(0.010) | 0.049(0.046) | 0.051(0.012) | 0.085(0.010) |

Table 3 reports the MSEs of parameter estimates, and Table 4 contains the MSEs of $\hat{\boldsymbol{\alpha}}$ and the average of $RASE_\pi$. From Table 3, we can see that MRSIP works comparable to MixLinReg when the sample size is small, and outperforms MixLinReg when sample size is large (such as $n = 400$ or $800$). By reducing the modelling bias of component proportions, MRSIP is able to better classify observations into two components and thus provide better component regression parameters. Based on Table 4, it is clear that MRSIP provides better estimates of component proportions than MixLinReg since the constant assumption of component proportions by MixLinReg is violated. From both tables, we can see that MRSIP(S) provides similar results to MRSIP(T), which demonstrates that SIR provides good initial values for MRSIP.

Table 3: The MSEs of parameters (true value times 100) for Example 2.

| | | $\beta_{10}$ | $\beta_{11}$ | $\beta_{12}$ | $\beta_{13}$ | $\beta_{20}$ | $\beta_{21}$ | $\beta_{22}$ | $\beta_{23}$ | $\sigma_1^2$ | $\sigma_2^2$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $n = 200$ | MRSIP(S) | 46.37 | 32.78 | 34.73 | 37.61 | 11.19 | 16.55 | 15.05 | 16.36 | 4.649 | 1.754 |
| | MRSIP(T) | 51.91 | 33.62 | 39.01 | 37.25 | 11.10 | 16.56 | 15.07 | 16.04 | 4.584 | 1.649 |
| $h = 0.131$ | MixLinReg | 50.87 | 33.67 | 42.53 | 34.68 | 12.03 | 12.66 | 18.84 | 12.30 | 4.250 | 1.265 |
| $n = 400$ | MRSIP(S) | 13.83 | 11.89 | 14.19 | 11.47 | 5.541 | 6.332 | 6.767 | 7.165 | 1.631 | 0.721 |
| | MRSIP(T) | 14.79 | 12.49 | 14.84 | 11.59 | 5.513 | 6.254 | 6.632 | 6.926 | 1.672 | 0.675 |
| $h = 0.103$ | MixLinReg | 29.03 | 14.97 | 29.46 | 15.72 | 8.045 | 5.967 | 12.46 | 6.269 | 1.864 | 0.626 |
| $n = 800$ | MRSIP(S) | 6.324 | 4.491 | 6.150 | 4.736 | 2.365 | 2.973 | 2.773 | 3.584 | 0.669 | 0.334 |
| | MRSIP(T) | 6.788 | 4.614 | 6.820 | 4.922 | 2.301 | 2.829 | 2.718 | 3.348 | 0.691 | 0.307 |
| $h = 0.080$ | MixLinReg | 21.89 | 6.866 | 21.84 | 8.223 | 5.413 | 3.163 | 8.775 | 3.640 | 0.848 | 0.352 |

# 5   Real Data Example

We illustrate the proposed methodology by an analysis of "The effectiveness of National Basketball Association guards". There are many ways to measure the (statistical) performance of guards in the National Basket Association (NBA). Of interest is how the height of the player (Height), minutes per game (MPG) and free throw percentage (FTP) affect points per game (PPM) (Chatterjee et al., 1995).

The data set contains some descriptive statistics for all 105 guards for the 1992-1993 season. Since players playing very few minutes are quite different from those who play a sizable part of the season, we only look at those players playing 10 or more minutes per game and appearing in 10 or more games. In addition, Michael Jordan is an outlier, so we also omit him from our data analysis. These exclude 10 players (Chatterjee et al., 1995). We divide each variable by its corresponding standard deviation, so that they have comparable numerical scales. An optimal bandwidth is selected at 0.344 by CV procedure. Figure 1(a) contains the estimated mean functions and hard-clustering results, denoted by dots and squares, respectively. The 95% confidence interval for $\hat{\boldsymbol{\alpha}}$ based on MSIM are (0.134,0.541), (0.715,0.949) and

Table 4: The MSEs of single index parameter and the average of $\text{RASE}_\pi$ (true value times 100) for Example 2.

|  |  | $\alpha_1$ | $\alpha_2$ | $\alpha_3$ | $\text{RASE}_\pi$ |
|---|---|---|---|---|---|
| $n = 200$ | MRSIP(S) | 5.709 | 19.30 | 5.996 | 18.87 |
|  | MRSIP(T) | 4.984 | 9.449 | 4.896 | 17.86 |
| $h = 0.131$ | MixLinReg | - | - | - | 28.98 |
| $n = 400$ | MRSIP(S) | 2.682 | 6.968 | 3.029 | 13.74 |
|  | MRSIP(T) | 2.113 | 3.019 | 1.902 | 12.98 |
| $h = 0.103$ | MixLinReg | - | - | - | 28.23 |
| $n = 800$ | MRSIP(S) | 0.980 | 2.527 | 1.585 | 10.35 |
|  | MRSIP(T) | 0.892 | 0.979 | 0.969 | 9.960 |
| $h = 0.080$ | MixLinReg | - | - | - | 28.04 |

(0.202,0.679). Therefore, MPG is the most influential factor on PPM. This might be partly explained by that coaches tend to let good players with higher PPM play longer minutes per game (i.e., higher MPG).

To evaluate the prediction performance of the proposed models and compared them to linear regression model and mixture of linear regression models, we used $d$-fold cross-validation with $d$=5, 10, and also Monte-Carlo cross-validation (MCCV) (Shao, 1993). In MCCV, the data were partitioned 500 times into disjoint training subsets (with size $n - d$) and test subsets (with size $d$). The mean squared prediction error evaluated at the test data sets over 500 replications are reported as boxplots in Figure 1(b). Apparently, the MSIM and the MRSIP have superior prediction power than the linear regression model or the mixture of linear regression models, and MSIM is more favorable than the MRSIP for this data set. The two groups of guards our new models found might be explained by the difference between shooting guards and passing guards.
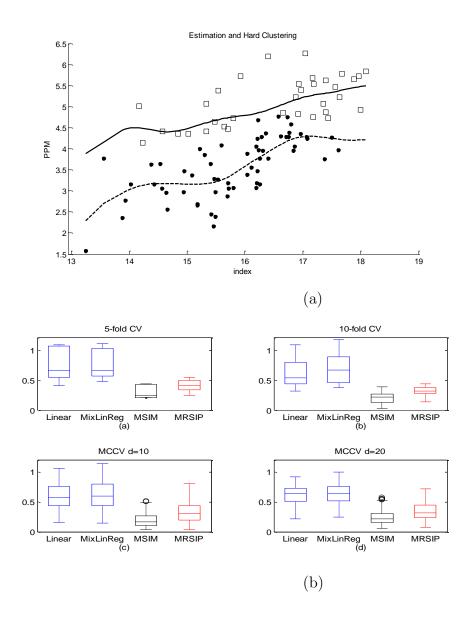
(a)



(b)

Figure 1: NBA data: (a) Estimated mean functions and a hard-clustering result; (b) Prediction accuracy: 5-fold CV; 10-fold CV; MCCV d=10; MCCV d=20.

# 6    Discussion

In this paper, we propose two finite semiparametric mixture of regression models and provide the modified EM algorithms to estimate them. We establish the identifiability results of the new models and investigate the asymptotic properties of the proposed

estimation procedures. Throughout the article, we assume that the number of components is known and fixed, but it requires more research to select the number of components for the proposed semiparametric mixture models. It will be interesting to know whether the recently proposed EM test (Chen and Li, 2009; Li and Chen, 2010) can be extended to the proposed semiparametric mixture models. In addition, it is also interesting to build some formal model selection procedure to compare different semiparametric mixture models. In the real data application, we use the cross-validation criteria to compare different models. When the models are nested, one might use generalized likelihood ratio statistic proposed by Fan et al. (2001) to test any parametric assumption for the semiparametric models. Furthermore, the assumption of fixed dimension of predictors can be relaxed and the proposed models can be extended to the cases where the dimension of predictors $p$ also diverges with the sample size $n$. This might be done by using the idea of penalized local likelihood if the sparsity assumption is added on the predictors.

# References

Böhning, D.(1999), *Computer-Assisted Analysis of Mixtures and Applications*, Boca Raton, FL: Chapman and Hall/CRC.

Cao, J. and Yao, W. (2012). Semiparametric mixture of binomial regression with a degenerate component. *Statistica Sinica*, 22, 27-46.

Chatterjee, S., Handcock, M.S. and Simmonoff, J.S. (1995). *A casebook for a first course in statistics and data analysis.* John Wiley & Sons, Inc.

Chen, J. and Li, P. (2009). Hypothesis test for normal mixture models: The EM approach. *The Annals of Statistics*, 37, 2523-2542.

Cook, R. D. and Li, B. (2002). Dimension reduction for conditional mean in regression. *Annals of Statistics*, 30, 455-474.

Fan, J., Zhang, C. and Zhang, J. (2001). Generalized likelihood ratio statistics and Wilks phenomenon. *The Annals of Statistics*, 29, 153-193.

Frühwirth-Schnatter, S. (2006), *Finite Mixture and Markov Switching Models*, Springer, New York.

Goldfeld, S.M. and Quandt, R.E. (1973). A Markov model for switching regressions. *Journal of Econometrics*, 1, 3-6.

Härdle, W., Hall, P. and Ichimura, H. (1993). Optimal smoothing in single-index models. *The Annals of Statistics*, 21, 157-178.

Henning, C. (2000). Identifiability of models for clusterwise linear regression. *Journal of Classification*, 17, 273-296.

Huang, M., Li, R. and Wang, S. (2013). Nonparametric mixture of regression models. *Journal of the American Statistical Association*, 108, 929-941.

Huang, M., Li, R., and Wang, H., and Yao, W. (2014). Estimating Mixture of Gaussian Processes by Kernel Smoothing. *Journal of Business and Economics Statistics*, 32, 259-270.

Huang, M. and Yao, W. (2012). Mixture of regression models with varying mixing proportions: a semiparametric approach. *Journal of the American Statistical Association*, 107, 711-724.

Ichimura, H. (1993). Semiparametric least squares (SLS) and weighted SLS estimation of single-index models. *Journal of Econometrics*, 58, 71-120.

Jordan, M. I. and Jacobs, R. A. (1994). Hierarchical mixtures of experts and the EM algorithm. *Neural Computation.* 6, 181-214.

Li, K. (1991). Sliced inverse regression for dimension reduction. *Journal of the American Statistical Association*, 86(414), 316-327.

Li, P. and Chen, J. (2010). Testing the order of a finite mixture. *Journal of the American Statistical Association*, 105, 1084-1092.

Li, B., Zha, H. and Chiaromonte, F. (2005). Contour regression: a general approach to dimension reduction. *Annals of Statistics*, 33, 1580-1616.

Ma, Y. and Zhu, L. (2012a). A semiparametric approach to dimension reduction. *Journal of the American Statistical Association*, 107(497), 168-179.

Ma, Y. and Zhu, L. (2012). Efficient estimation in sufficient dimension reduction. *Annals of Statistics.*

Lindsay, B. G., (1995), *Mixture Models: Theory, Geometry, and Applications*, NSF-CBMS Regional Conference Series in Probability and Statistics v 5, Hayward, CA: Institure of Mathematical Statistics.

Luo, R., Wang, H., and Tsai, C. L. (2009). Contour projected dimension reduction. *Annals of Statistics*, 37, 3743-3778.

McLachlan, G. J. and Peel, D. (2000), *Finite Mixture Models*, New York: Wiley.

Shao, J. (1993). Linear models selection by cross-validation. *Journal of the American Statistical Association*, 88, 486-494.

Titterington, D., Smith, A., and Makov, U. (1985). Statistical analysis of finite mixture distribution. Wiley.

Wang, H. and Xia, Y. (2008). Sliced regression for dimension reduction. *Journal of the American Statistical Association*, 103, 811-821.

Xiang, S. and Yao, W(2015). Semiparametric mixtures of nonparametric regressions. *Computational Statistics & Data Analysis*, submitted for publication.

Young, D.S. and Hunter, D.R.(2010). Mixtures of regressions with predictors dependent mixing proportions. *Computational Statistics and Data Analysis*, 54, 2253-2266.

Yao, W. and Lindsay, B. G.(2009). Bayesian mixture labeling by highest posterior density. *Journal of American Statistical Association*, 104, 758-767.

Zeng, P. (2012). Finite Mixture of Heteroscedastic Single-Index Models. *Open Journal of Statistics*, 2, 12-20.