

# UC Riverside

## UC Riverside Previously Published Works

### Title

Outlier detection and robust mixture modeling using nonconvex penalized likelihood

### Permalink

<https://escholarship.org/uc/item/8vw5m7q9>

### Journal

Journal of Statistical Planning and Inference, 164(1)

### ISSN

0378-3758

### Authors

Yu, Chun

Chen, Kun

Yao, Weixin

### Publication Date

2015-09-01

### DOI

10.1016/j.jspi.2015.03.003

Peer reviewed

# Outlier Detection and Robust Mixture Modeling Using Nonconvex Penalized Likelihood

CHUN YU, \* KUN CHEN,<sup>†</sup> WEIXIN YAO,<sup>‡</sup>

## Abstract

Finite mixture models are widely used in a variety of statistical applications. However, the classical normal mixture model with maximum likelihood estimation is prone to the presence of only a few severe outliers. We propose a robust mixture modeling approach using a mean-shift formulation coupled with nonconvex sparsity-inducing penalization, to conduct simultaneous outlier detection and robust parameter estimation. An efficient iterative thresholding-embedded EM algorithm is developed to maximize the penalized log-likelihood. The efficacy of our proposed approach is demonstrated via simulation studies and a real application on Acidity data analysis.

**Key words:** EM algorithm; Mixture models; Outlier detection; Penalized likelihood.

## 1 Introduction

Nowadays finite mixture distributions are increasingly important in modeling a variety of random phenomena (see Everitt and Hand, 1981, Titterington, Smith and Markov, 1985,

---

\*Chun Yu is Instructor, School of Statistics, Jiangxi University of Finance and Economics, Nanchang 330013, China. Email: chuckyu0106@126.com.

<sup>†</sup>Kun Chen is Assistant Professor, Department of Statistics, University of Connecticut, Storrs, CT, 06269. Email: kun.chen@uconn.edu.

<sup>‡</sup>Weixin Yao is corresponding author, Associate Professor, Department of Statistics, University of California, Riverside, CA 92521. Email: weixin.yao@ucr.edu.

17 McLachlan and Basford, 1988, Lindsay, 1995, and Böhning, 1999). The  $m$ -component finite  
 18 normal mixture distribution has probability density

$$f(y; \boldsymbol{\theta}) = \sum_{i=1}^m \pi_i \phi(y; \mu_i, \sigma_i^2), \quad (1.1)$$

19 where  $\boldsymbol{\theta} = (\pi_1, \mu_1, \sigma_1; \dots; \pi_m, \mu_m, \sigma_m)^T$  collects all the unknown parameters,  $\phi(\cdot; \mu, \sigma^2)$  denotes  
 20 the density function of  $N(\mu, \sigma^2)$ , and  $\pi_j$  is the proportion of  $j$ th subpopulation with  $\sum_{j=1}^m \pi_j =$   
 21 1. Given observations  $(y_1, \dots, y_n)$  from model (1.1), the maximum likelihood estimator (MLE)  
 22 of  $\boldsymbol{\theta}$  is given by,

$$\hat{\boldsymbol{\theta}}_{\text{mle}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n \log \left\{ \sum_{j=1}^m \pi_j \phi(y_i; \mu_j, \sigma_j^2) \right\}, \quad (1.2)$$

23 which does not have an explicit form and is usually calculated by the EM algorithm (Dempster  
 24 et al. 1977).

25 The MLE based on the normality assumption possesses many desirable properties such as  
 26 asymptotic efficiency, however, the method is not robust and both parameter estimation and  
 27 inference can fail miserably in the presence of outliers. Our focus in this paper is hence on  
 28 robust estimation and accurate outlier detection in mixture model. For the simpler problem  
 29 of estimating of a single location, many robust methods have been proposed, including the M-  
 30 estimator (Huber, 1981), the least median of squares (LMS) estimator (Siegel 1982), the least  
 31 trimmed squares (LTS) estimator (Rousseeuw 1983), the S-estimates (Rousseeuw and Yohai  
 32 1984), the MM-estimator (Yohai 1987), and the weighted least squares estimator (REWLSE)  
 33 (Gervini and Yohai 2002). In contrast, there is much less research on robust estimation of  
 34 the mixture model, in part because it is not straightforward to replace the log-likelihood in  
 35 (1.2) by a robust criterion similar to M-estimation. Peel and McLachlan (2000) proposed a  
 36 robust mixture modeling using  $t$  distribution. Markatou (2000) and Qin and Priebe (2013)  
 37 proposed using a weighted likelihood for each data point to robustify the estimation procedure  
 38 for mixture models. Fujisawa and Eguchi (2005) proposed a robust estimation method in  
 39 normal mixture model using a modified likelihood function. Neykov et al. (2007) proposed  
 40 robust fitting of mixtures using the trimmed likelihood. Other related robust methods on

41 mixture models include Hennig (2002, 2003), Shen et al. (2004), Bai et al. (2012), Bashir and  
42 Carter (2012), Yao et al. (2014), and Song et al. (2014)

43 We propose a new robust mixture modelling approach based on a mean-shift model for-  
44 mulation coupled with penalization, which achieves simultaneous outlier detection and robust  
45 parameter estimation. A case-specific mean-shift parameter vector is added to the mean struc-  
46 ture of the mixture model, and it is assumed to be sparse for capturing the rare but possibly  
47 severe outlying effects caused by the potential outliers. When the mixture components are  
48 assumed to have equal variances, our method directly extends the robust linear regression ap-  
49 proaches proposed by She and Owen (2011) and Lee, MacEachern and Jung (2012). However,  
50 even in this case the optimization of the penalized mixture log-likelihood is not trivial, espe-  
51 cially for the SCAD penalty (Fan and Li, 2001). For the general case of unequal component  
52 variances, the variance heterogeneity of different components complicates the declaration and  
53 detection of the outliers, and we thus propose a general scale-free and case-specific mean-shift  
54 formulation to solve the general problem.

## 55 **2 Robust Mixture Model via Mean-Shift Penalization**

56 In this section, we introduce the proposed robust mixture modelling approach via mean-shift  
57 penalization (RMM). To focus on the main idea, we restrict our attention on the normal  
58 mixture model. The proposed approach can be readily extended to other mixture models,  
59 such as gamma mixture and logistic mixture. Due to the inherent difference between the case  
60 of equal component variances and the case of unequal component variances, we shall discuss  
61 two cases separately.

### 62 **2.1 RMM for Equal Component Variances**

63 Assume the mixture components have equal variances, i.e.,  $\sigma_1^2 = \dots = \sigma_m^2 = \sigma^2$ . The proposed  
64 robust mixture model with a mean-shift parameterization is to assume that the observations

65  $(y_1, \dots, y_n)$  come from the following mixture density

$$f(y_i; \boldsymbol{\theta}, \gamma_i) = \sum_{j=1}^m \pi_j \phi(y_i - \gamma_i; \mu_j, \sigma^2), \quad i = 1, \dots, n, \quad (2.1)$$

66 where  $\boldsymbol{\theta} = (\pi_1, \mu_1, \dots, \pi_m, \mu_m, \sigma)^T$ , and  $\gamma_i$  is the mean-shift parameter for the  $i$ th observation.  
 67 Apparently, without any constraints, the addition of the mean-shift parameters results in an  
 68 overly parameterized model. The key here is to assume that the vector  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_n)$  is  
 69 sparse, i.e.,  $\gamma_i$  is zero when the  $i$ th data point is a normal observation with any of the  $m$   
 70 mixture components, and only when the  $i$ th observation is an outlier,  $\gamma_i$  becomes nonzero to  
 71 capture the outlying effect. Therefore, the sparse estimation of  $\boldsymbol{\gamma}$  provides a direct way to  
 72 accommodate and identify outliers.

Due to the sparsity assumption of  $\boldsymbol{\gamma}$ , we propose to maximize the following penalized log-likelihood criterion to conduct model estimation and outlier detection,

$$pl_1(\boldsymbol{\theta}, \boldsymbol{\gamma}) = l_1(\boldsymbol{\theta}, \boldsymbol{\gamma}) - \sum_{i=1}^n \frac{1}{w_i} P_\lambda(|\gamma_i|) \quad (2.2)$$

73 where  $l_1(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \sum_{i=1}^n \log \left\{ \sum_{j=1}^m \pi_j \phi(y_i - \gamma_i; \mu_j, \sigma^2) \right\}$ ,  $w_i$ s are some prespecified weights,  $P_\lambda(\cdot)$   
 74 is some penalty function used to induce the sparsity in  $\boldsymbol{\gamma}$ , and  $\lambda$  is a tuning parameter con-  
 75 trolling the number of outliers, i.e., the number of nonzero  $\gamma_i$ . In practice,  $w_i$ s can be chosen  
 76 to reflect any available prior information about how likely that  $y_i$ s are outliers; to focus on the  
 77 main idea, here we mainly consider  $w_1 = w_2 = \dots = w_n = w$ , and discuss the choice of  $w$  for  
 78 different penalty functions.

Some commonly used sparsity-inducing penalty functions include the  $\ell_1$  penalty (Donoho and Johnstone, 1994a; Tibshirani, 1996, 1997)

$$P_\lambda(\gamma) = \lambda|\gamma|, \quad (2.3)$$

the  $\ell_0$  penalty (Antoniadis, 1997)

$$P_\lambda(\gamma) = \frac{\lambda^2}{2} I(\gamma \neq 0), \quad (2.4)$$

and the SCAD penalty (Fan and Li, 2001)

$$P_\lambda(\gamma) = \begin{cases} \lambda|\gamma| & \text{if } |\gamma| \leq \lambda, \\ -\left(\frac{\gamma^2 - 2a\lambda|\gamma| + \lambda^2}{2(a-1)}\right) & \text{if } \lambda < |\gamma| \leq a\lambda, \\ \frac{(a+1)\lambda^2}{2} & \text{if } |\gamma| > a\lambda, \end{cases} \quad (2.5)$$

79 where  $a$  is a constant usually set to be 3.7. In penalized estimation, each of the above penalty  
 80 forms corresponds to a thresholding rule, e.g.,  $\ell_1$  penalization corresponds to a soft-thresholding  
 81 rule, and  $\ell_0$  penalization corresponds to a hard-thresholding rule. It is also known that the  
 82 convex  $\ell_1$  penalization often leads to over-selection and substantial bias in estimation. In-  
 83 deed, as shown by She and Owen (2011) in the context of linear regression,  $\ell_1$  penalization in  
 84 mean-shift model has connections with M-estimation using Huber's loss and usually leads to  
 85 poor performance in outlier detection. Similar phenomenon is also observed in our extensive  
 86 numerical studies. However, if there are no high leverage outliers, the  $\ell_1$  penalty works well  
 87 and succeeds to detect the outliers, see for examples, Dalalyan and Keriven (2012); Dalalyan  
 88 and Chen (2012); Nguyen and Tran (2013).

89 We propose a thresholding embedded EM algorithm to maximize the objective function  
 90 (2.2). Let

$$z_{ij} = \begin{cases} 1 & \text{if the } i\text{th observation is from the } j\text{th component,} \\ 0 & \text{otherwise,} \end{cases}$$

91 and  $\mathbf{z}_i = (z_{i1}, \dots, z_{im})$ . The complete penalized log-likelihood function based on the complete  
 92 data  $\{(y_i, \mathbf{z}_i), i = 1, 2, \dots, n\}$  is

$$pl_1^c(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \sum_{i=1}^n \sum_{j=1}^m z_{ij} \log \{ \pi_j \phi(y_i - \gamma_i; \mu_j, \sigma^2) \} - \sum_{i=1}^n \frac{1}{w} P_\lambda(|\gamma_i|). \quad (2.6)$$

93 Based on the construction of the EM algorithm, in the E step, given the current estimate  
 94  $\boldsymbol{\theta}^{(k)}$  and  $\boldsymbol{\gamma}^{(k)}$  at the  $k$ th iteration, we need to find the conditional expectation of the complete  
 95 penalized log-likelihood function (2.6), i.e.,  $E\{p l_1^c(\boldsymbol{\theta}, \boldsymbol{\gamma}) \mid \boldsymbol{\theta}^{(k)}, \boldsymbol{\gamma}^{(k)}\}$ . The problem simplifies to  
 96 the calculation of  $E(z_{ij} \mid y_i; \boldsymbol{\theta}^{(k)}, \boldsymbol{\gamma}^{(k)})$ ,

$$p_{ij}^{(k+1)} = E(z_{ij} \mid y_i; \boldsymbol{\theta}^{(k)}, \boldsymbol{\gamma}^{(k)}) = \frac{\pi_j^{(k)} \phi(y_i - \gamma_i^{(k)}; \mu_j^{(k)}, \sigma^{2(k)})}{\sum_{j=1}^m \pi_j^{(k)} \phi(y_i - \gamma_i^{(k)}; \mu_j^{(k)}, \sigma^{2(k)})}.$$

In the M step, we then update  $(\boldsymbol{\theta}, \boldsymbol{\gamma})$  by maximizing  $E\{p l_1^c(\boldsymbol{\theta}, \boldsymbol{\gamma}) \mid \boldsymbol{\theta}^{(k)}, \boldsymbol{\gamma}^{(k)}\}$ . There is no explicit solution, except for the  $\pi_j$ s,

$$\pi_j^{(k+1)} = \frac{\sum_{i=1}^n p_{ij}^{(k+1)}}{n}.$$

We propose to alternately update  $\{\sigma, \mu_j, j = 1, \dots, m\}$  and  $\boldsymbol{\gamma}$  until convergence to get  $\{\mu_j^{(k+1)}, j = 1, \dots, m; \sigma^{(k+1)}, \boldsymbol{\gamma}^{(k+1)}\}$ . Given  $\boldsymbol{\gamma}$ ,  $\mu_j$ s and  $\sigma$  are updated by

$$\mu_j \leftarrow \frac{\sum_{i=1}^n p_{ij}^{(k+1)} (y_i - \gamma_i)}{\sum_{i=1}^n p_{ij}^{(k+1)}}, \quad \sigma^2 \leftarrow \frac{\sum_{j=1}^m \sum_{i=1}^n p_{ij}^{(k+1)} (y_i - \gamma_i - \mu_j)^2}{n}.$$

97 Given  $\mu_j$ s and  $\sigma$ ,  $\boldsymbol{\gamma}$  is updated by maximizing

$$\sum_{i=1}^n \sum_{j=1}^m p_{ij}^{(k+1)} \log \phi(y_i - \gamma_i; \mu_j, \sigma^2) - \sum_{i=1}^n \frac{1}{w} P_\lambda(|\gamma_i|),$$

98 which is equivalently to minimizing

$$\frac{1}{2} \left\{ \gamma_i - \sum_{j=1}^m p_{ij}^{(k+1)} (y_i - \mu_j) \right\}^2 + \frac{1}{w} \sigma^2 P_\lambda(|\gamma_i|), \quad (2.7)$$

99 separately for each  $\gamma_i$ ,  $i = 1, \dots, n$ . A detailed derivation is presented in the Appendix. For  
 100 either the  $\ell_1$  or the  $\ell_0$  penalty,  $w^{-1} \sigma^2 P_\lambda(|\gamma_i|) = \sigma P_{\lambda^*}(|\gamma_i|)$ , where  $\lambda^* = \frac{\sigma}{w} \lambda$ . Therefore, if  
 101  $\lambda$  is chosen data adaptively, we can simply set  $w = 1$  for these penalties. However, for the  
 102 SCAD penalty, such property does not hold and the solution may be affected nonlinearly by

103 the ratio  $\sigma^2/w$ . In order to mimic the unscaled SCAD and use the same  $a$  value as suggested  
 104 by Fan and Li (2001), we need to make sure that  $\sigma^2/w$  is close to 1. Therefore, we propose to  
 105 set  $w = \hat{\sigma}^2$  for SCAD penalty, where  $\hat{\sigma}^2$  is a robust estimate of  $\sigma^2$  such as the estimate from  
 106 the trimmed likelihood estimation (Neykov et al. 2007) or the estimator using the  $\ell_0$  penalty  
 107 assuming  $w = 1$ .

108 As shown in Proposition 1 below, when the  $\ell_1$  penalty is used, (2.7) is minimized by a soft  
 109 thresholding rule, and when the  $\ell_0$  penalty is used, (2.7) is minimized by a hard thresholding  
 110 rule. When the SCAD penalty is used, however, the problem is solved by a modified SCAD  
 111 thresholding rule, which is shown in Lemma 1.

112 **Proposition 1.** Let  $\xi_i = \sum_{j=1}^m p_{ij}^{(k+1)}(y_i - \mu_j)$ . Let  $w = 1$  in (2.7). When the penalty function  
 113 in (2.7) is the  $\ell_1$  penalty (2.8), the minimizer of (2.7) is given by

$$\hat{\gamma}_i = \Theta_{soft}(\xi_i; \lambda, \sigma) = \text{sgn}(\xi_i) (|\xi_i| - \sigma\lambda)_+, \quad (2.8)$$

114 where  $a_+ = \max(a, 0)$ . When the penalty function in (2.7) is the  $\ell_0$  penalty (2.9), the minimizer  
 115 of (2.7) is given by

$$\hat{\gamma}_i = \Theta_{hard}(\xi_i; \lambda, \sigma) = \xi_i I(|\xi_i| > \sigma\lambda), \quad (2.9)$$

116 where  $I(\cdot)$  denotes the indicator function.

117 **Lemma 1.** Let  $\xi_i = \sum_{j=1}^m p_{ij}^{(k+1)}(y_i - \mu_j)$ . Let  $w = \hat{\sigma}^2$  in (2.7), a robust estimator of  $\sigma^2$ .  
 118 When the penalty function in (2.7) is the SCAD penalty (2.5), the minimizer of (2.7) is given  
 119 by

120 1. when  $\sigma^2/\hat{\sigma}^2 < a - 1$ ,

$$\hat{\gamma}_i = \Theta_{scad}(\xi_i; \lambda, \sigma) = \begin{cases} \text{sgn}(\xi_i) \left( |\xi_i| - \frac{\sigma^2\lambda}{\hat{\sigma}^2} \right)_+, & \text{if } |\xi_i| \leq \lambda + \frac{\sigma^2\lambda}{\hat{\sigma}^2}, \\ \frac{\hat{\sigma}^2(a-1)\xi_i - \text{sgn}(\xi_i)a\lambda}{\hat{\sigma}^2(a-1)-1}, & \text{if } \lambda + \frac{\sigma^2\lambda}{\hat{\sigma}^2} < |\xi_i| \leq a\lambda, \\ \xi_i, & \text{if } |\xi_i| > a\lambda. \end{cases} \quad (2.10)$$



121 2. when  $a - 1 \leq \sigma^2/\hat{\sigma}^2 \leq a + 1$ ,

$$\hat{\gamma}_i = \Theta_{scad}(\xi_i; \lambda, \sigma) = \begin{cases} \text{sgn}(\xi_i) \left( |\xi_i| - \frac{\sigma^2 \lambda}{\hat{\sigma}^2} \right)_+, & \text{if } |\xi_i| \leq \frac{a+1+\frac{\sigma^2}{\hat{\sigma}^2}}{2} \lambda, \\ \xi_i, & \text{if } |\xi_i| > \frac{a+1+\frac{\sigma^2}{\hat{\sigma}^2}}{2} \lambda. \end{cases} \quad (2.11)$$

122 3. when  $\sigma^2/\hat{\sigma}^2 > a + 1$ ,

$$\hat{\gamma}_i = \Theta_{scad}(\xi_i; \lambda, \sigma) = \xi_i I(|\xi_i| > \sqrt{\frac{\sigma^2(a+1)}{\hat{\sigma}^2}} \lambda). \quad (2.12)$$

123 The detailed EM algorithm is summarized in Algorithm 1. For simplicity, we have used  
 124  $\Theta(\xi_i; \lambda, \sigma)$  to denote a general thresholding rule determined by the adopted penalty function,  
 125 e.g., the modified SCAD thresholding rule  $\Theta_{scad}(\xi_i; \lambda, \sigma)$  defined in Lemma 1. The convergence  
 126 property of the proposed algorithm is summarized in Theorem 2.1 below, which follows directly  
 127 from the property of the EM algorithm, and hence its proof is omitted.

128 **Theorem 2.1.** *Each iteration of E step and M step of Algorithm 1 monotonically non-decreases*  
 129 *the penalized log-likelihood (2.2), i.e.,  $pl_1(\boldsymbol{\theta}^{(k+1)}, \boldsymbol{\gamma}^{(k+1)}) \geq pl_1(\boldsymbol{\theta}^{(k)}, \boldsymbol{\gamma}^{(k)})$ , for all  $k \geq 0$ .*

## 130 2.2 RMM for Unequal Component Variances

131 When the component variances are unequal, the naive mean-shift model (2.1) can not be  
 132 directly applied, due to the scale difference in the mixture components. To illustrate further,  
 133 suppose the standard deviation in the first component is 1 and the standard deviation in the  
 134 second component is 4. If some weighted residual  $\xi_i$ , defined in Proposition 1, equals to 5, then  
 135 the  $i$ th observation is considered as an outlier if it is from the first component but should not be  
 136 regarded as an outlier if it belongs to the second component. This suggests that the declaration  
 137 of outliers in a mixture model shall take into account both the centers and the variabilities of  
 138 all the components, i.e., an observation is considered as an outlier in the mixture model only  
 139 if it is far away from all the component centers judged by their own component variabilities.

140 We propose a general scale-free mean-shift model to incorporate the information on com-

---

**Algorithm 1** Thresholding Embedded EM Algorithm for Equal Variances Case

---

Initialize  $\boldsymbol{\theta}^{(0)}$  and  $\boldsymbol{\gamma}^{(0)}$ . Set  $k \leftarrow 0$ .

**repeat**

E-Step: Compute the classification probabilities

$$p_{ij}^{(k+1)} = \mathbb{E}(z_{ij} | \mathbf{y}_i; \boldsymbol{\theta}^{(k)}) = \frac{\pi_j^{(k)} \phi(\mathbf{y}_i - \boldsymbol{\gamma}_i^{(k)}; \boldsymbol{\mu}_j^{(k)}, \sigma^{2(k)})}{\sum_{j=1}^m \pi_j^{(k)} \phi(\mathbf{y}_i - \boldsymbol{\gamma}_i^{(k)}; \boldsymbol{\mu}_j^{(k)}, \sigma^{2(k)})}.$$

M-Step: Update  $(\boldsymbol{\theta}, \boldsymbol{\gamma})$  by the following two steps:

1.  $\pi_j^{(k+1)} = \sum_{i=1}^n p_{ij}^{(k+1)} / n, j = 1, \dots, m$ .
2. Iterating the following steps until convergence to obtain  $\{\boldsymbol{\mu}_j^{(k+1)}, j = 1, \dots, m; \sigma^{2(k+1)}, \boldsymbol{\gamma}^{(k+1)}\}$ :

$$(2.a) \quad \boldsymbol{\gamma}_i \leftarrow \Theta(\boldsymbol{\xi}_i; \boldsymbol{\lambda}, \sigma), i = 1, \dots, n, \text{ where } \boldsymbol{\xi}_i = \sum_{j=1}^m p_{ij}^{(k+1)} (\mathbf{y}_i - \boldsymbol{\mu}_j),$$

$$(2.b) \quad \boldsymbol{\mu}_j \leftarrow \frac{\sum_{i=1}^n p_{ij}^{(k+1)} (\mathbf{y}_i - \boldsymbol{\gamma}_i)}{\sum_{i=1}^n p_{ij}^{(k+1)}}, j = 1, \dots, m,$$

$$(2.c) \quad \sigma^2 \leftarrow \frac{\sum_{j=1}^m \sum_{i=1}^n p_{ij}^{(k+1)} (\mathbf{y}_i - \boldsymbol{\gamma}_i - \boldsymbol{\mu}_j)^2}{n}.$$

$k \leftarrow k + 1$ .

**until** convergence.

---

141 ponent variability,

$$f(y_i; \boldsymbol{\theta}, \gamma_i) = \sum_{j=1}^m \pi_j \phi(y_i - \gamma_i \sigma_j; \mu_j, \sigma_j^2), \quad i = 1, \dots, n, \quad (2.13)$$

where with some abuse of notation,  $\boldsymbol{\theta}$  is redefined as  $\boldsymbol{\theta} = (\pi_1, \mu_1, \sigma_1, \dots, \pi_m, \mu_m, \sigma_m)^T$ . Given observations  $(y_1, y_2, \dots, y_n)$ , we estimate the parameters  $\boldsymbol{\theta}$  and  $\boldsymbol{\gamma}$  by maximizing the following penalized log-likelihood function:

$$pl_2(\boldsymbol{\theta}, \boldsymbol{\gamma}) = l_2(\boldsymbol{\theta}, \boldsymbol{\gamma}) - \sum_{i=1}^n \frac{1}{w_i} P_\lambda(|\gamma_i|), \quad (2.14)$$

142 where  $l_2(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \sum_{i=1}^n \log \left\{ \sum_{j=1}^m \pi_j \phi(y_i - \gamma_i \sigma_j; \mu_j, \sigma_j^2) \right\}$ . Since the  $\gamma_i$ s in (2.14) are scale free,  
 143 we can set  $w_1 = w_2 = \dots = w_n = 1$  when no prior information is available.

144 We again propose a thresholding embedded EM algorithm to maximize (2.14). As the  
 145 construction is similar to the case of equal variances, we omit the details of its derivation.  
 146 The proposed EM algorithm is presented in Algorithm 2, and here we shall briefly remark the  
 147 main changes. Unlike in the case of equal variances, the update of  $\sigma_j^2$  in (2.17), with other  
 148 parameters held fixed, does not have explicit solution in general and requires some numerical  
 149 algorithm to solve, e.g., the Newton-Raphson method; as the problem is one dimensional, the  
 150 computation remains very fast. In the case of unequal variances, the problem of updating  $\boldsymbol{\gamma}$ ,  
 151 with other parameters held fixed, is still separable in each  $\gamma_i$ , i.e., at the  $(k+1)$ th iteration,

$$\hat{\gamma}_i = \arg \min_{\gamma_i} \left\{ - \sum_{j=1}^m p_{ij}^{(k+1)} \log \phi(y_i - \gamma_i \sigma_j; \mu_j, \sigma_j^2) + P_\lambda(|\gamma_i|) \right\}.$$

It can be readily shown that the solution is given by simple thresholding rules. In particular, using the  $\ell_1$  penalty leads to  $\hat{\gamma}_i = \Theta_{\text{soft}}(\xi_i; \lambda, 1)$  and using the  $\ell_0$  penalty leads to  $\hat{\gamma}_i = \Theta_{\text{hard}}(\xi_i; \lambda, 1)$ , where  $\Theta_{\text{soft}}$  and  $\Theta_{\text{hard}}$  are defined in Proposition 1, and here in the case of unequal variance,  $\xi_i$  becomes

$$\xi_i = \sum_{j=1}^m \frac{p_{ij}^{(k+1)}}{\sigma_j} (y_i - \mu_j).$$

---

**Algorithm 2** Thresholding Embedded EM Algorithm for Unequal Variances Case

---

Initialize  $\boldsymbol{\theta}^{(0)}$  and  $\boldsymbol{\gamma}^{(0)}$ . Set  $k \leftarrow 0$ .

**repeat**

E-Step: Compute the classification probabilities

$$p_{ij}^{(k+1)} = E(z_{ij}|y_i; \boldsymbol{\theta}^{(k)}) = \frac{\pi_j^{(k)} \phi(y_i - \gamma_i^{(k)} \sigma_j^{(k)}; \mu_j^{(k)}, \sigma_j^{2(k)})}{\sum_{j=1}^m \pi_j^{(k)} \phi(y_i - \gamma_i^{(k)} \sigma_j^{(k)}; \mu_j^{(k)}, \sigma_j^{2(k)})}.$$

M-Step: Update  $(\boldsymbol{\theta}, \boldsymbol{\gamma})$  by the following two steps:

1.

$$\pi_j^{(k+1)} = \frac{\sum_{i=1}^n p_{ij}^{(k+1)}}{n}, j = 1, \dots, m.$$

2. Iterating the following steps until convergence to obtain  $\{\mu_j^{(k+1)}, \sigma_j^{2(k+1)}, j = 1, \dots, m, \boldsymbol{\gamma}^{(k+1)}\}$ :

$$(2.a) \quad \gamma_i \leftarrow \Theta(\xi_i; \lambda, 1), \text{ where } \xi_i = \sum_{j=1}^m p_{ij}^{(k+1)} (y_i - \mu_j) / \sigma_j, \quad (2.15)$$

$$(2.b) \quad \mu_j \leftarrow \frac{\sum_{i=1}^n p_{ij}^{(k+1)} (y_i - \gamma_i \sigma_j)}{\sum_{i=1}^n p_{ij}^{(k+1)}}, \quad (2.16)$$

$$(2.c) \quad \sigma_j^2 \leftarrow \arg \max_{\sigma_j} \sum_{i=1}^n p_{ij}^{(k+1)} \log \phi(y_i - \gamma_i \sigma_j; \mu_j, \sigma_j^2). \quad (2.17)$$

$k \leftarrow k + 1$ .

**until** convergence

---

As the  $\gamma_i$ s become scale free, the thresholding rule for solving SCAD becomes much simpler, and it is given by (2.10) when setting the quantity  $\sigma^2/\hat{\sigma}^2 = 1$ , i.e.,

$$\hat{\gamma}_i = \Theta_{SCAD}(\xi_i; \lambda, 1) = \begin{cases} \text{sgn}(\xi_i)(|\xi_i| - \lambda)_+, & \text{if } |\xi_i| \leq 2\lambda, \\ \frac{(a-1)\xi_i - \text{sgn}(\xi_i)a\lambda}{a-2}, & \text{if } 2\lambda < |\xi_i| \leq a\lambda, \\ \xi_i, & \text{if } |\xi_i| > a\lambda. \end{cases}$$

152 Similar to Theorem 2.1, it is easy to check that the monotone non-decreasing property  
 153 remains hold for Algorithm 2. We note that in both algorithms, we have used an iterative  
 154 algorithm aiming to fully maximize the expected complete log-likelihood under penalization.  
 155 It can be seen that in this blockwise coordinate descent algorithm, each loop of (2.a) – (2.c)  
 156 monotonically non-decreases the objective function. Therefore, an alternative strategy is to  
 157 run (2.a) – (2.c) only a few times or even just once in each M-step; the resulting generalized  
 158 EM algorithm continues to possess the desirable convergence property. Based on our limited  
 159 experience, however, this method generally does not lead to significant saving in computation,  
 160 because the iterations in the M-step only involve simple operations and partially solving M-  
 161 step may slow down the overall convergence. Nevertheless, it is worthwhile to point out this  
 162 strategy, as it can be potentially useful when more complicated penalization methods are  
 163 required.

### 164 **2.3 Tuning Parameter Selection**

165 When using robust estimation or outlier detection methods, it is usually required to choose a  
 166 “threshold” value, e.g., the percentage of observations to eliminate, or the cutoff to declare ex-  
 167 treme residuals. In our method, selecting “threshold” becomes the tuning parameter selection  
 168 problem in penalized regression (2.2) and (2.14). As such, many well-developed methodologies  
 169 including cross validation and information criterion based approaches are all applicable, and  
 170 the turning parameter  $\lambda$  can be selected in an objective way, based on predictive power of  
 171 the model or the balance between model goodness of fit and complexity. Here, we provide  
 172 a data adaptive way to select  $\lambda$  based on a Bayesian information criterion (BIC), due to its

173 computation efficiency and proven superior performance on variable selection,

$$\text{BIC}(\lambda) = -l_j^*(\lambda) + \log(n)\text{df}(\lambda), \quad (2.18)$$

174 where  $j = 1$  or  $2$ ,  $l_j^*(\lambda) = l_j(\hat{\boldsymbol{\theta}}(\lambda), \hat{\boldsymbol{\gamma}}(\lambda))$  is the mixture log-likelihood evaluated at the estimator  
175  $(\hat{\boldsymbol{\theta}}(\lambda), \hat{\boldsymbol{\gamma}}(\lambda))$  obtained by maximizing the penalized likelihood criterion (2.2) or (2.14) with  $\lambda$   
176 being the tuning parameter, and  $\text{df}(\lambda)$  is the model degrees of freedom which is estimated by  
177 the sum of the number of nonzero  $\gamma$  values and the number of mixture component parameters.  
178 In practice, the optimal tuning parameter  $\lambda$  is chosen by minimizing  $\text{BIC}(\lambda)$  over a grid of  
179 100  $\lambda$  values, equally spaced on the log scale between  $\lambda_{\min}$  and  $\lambda_{\max}$ , where  $\lambda_{\max}$  is some  
180 large value of  $\lambda$  resulting in all zero values in  $\hat{\boldsymbol{\gamma}}$ , corresponding to the case of no outlier, and  
181  $\lambda_{\min}$  is some small value of  $\lambda$  resulting in roughly 40% nonzero values in  $\hat{\boldsymbol{\gamma}}$ , since in reality  
182 the proportion of outliers is usually quite small. The models with various  $\lambda$  values are fitted  
183 sequentially. The previous solution is used as the initial value for fitting the next model to  
184 speed up the computation. As such, our proposed method is able to search conveniently over  
185 a whole spectrum of possible models.

186 In mixture model, it is a foremost task to determine the number of mixture component  
187  $m$ . The problem may be resolved based on prior knowledge of the underlying data generation  
188 process. In many applications where no prior information is available, we suggest to conduct  
189 the penalized mixture model analysis with a few plausible  $m$  values, and use the proposed BIC  
190 criterion to select both the number of component  $m$  and the amount of penalization  $\lambda$ .

## 191 **3 Simulation**

### 192 **3.1 Setups**

193 We conduct simulation studies to investigate the effectiveness of the proposed approach and  
194 compare it with several existing methods. We consider both the case of equal variances and the  
195 case of unequal variances. In each setup to be elaborated below, we first generate independent  
196 observations from a normal mixture distribution; a few outliers are then created by adding

197 random mean-shift to some of the observations. The sample size is set to  $n = 200$ , and we  
198 consider two proportions of outliers, i.e.,  $p_{\mathcal{O}} = 5\%$  and  $p_{\mathcal{O}} = 10\%$ . The number of replicates is  
199 200 for each simulation setting.

200 Example 1: The samples  $(y_1, y_2, \dots, y_n)$  are generated from model (2.1) with  $\pi_1 = 0.3$ ,  $\mu_1 =$   
201  $0$ ,  $\pi_2 = 0.7$ ,  $\mu_2 = 8$ , and  $\sigma = 1$ . That is, the size of the first component  $n_1$  is generated  
202 from a binomial distribution with  $n_1 \sim \text{Bin}(n, p = 0.3)$ , and consequently the size of the  
203 second component is given by  $n_2 = n - n_1$ . To create  $100p_{\mathcal{O}}\%$  outliers, we randomly  
204 choose  $3np_{\mathcal{O}}/10$  many observations from component 1, and each of them is added a  
205 random mean shift  $\gamma \sim \text{Unif}([-5, -7])$ . Similarly  $7np_{\mathcal{O}}/10$  outliers are created by adding  
206 random mean shift  $\gamma \sim \text{Unif}([5, 7])$  to observations from component 2.

207 Example 2: The samples  $(y_1, y_2, \dots, y_n)$  are generated from model (2.13) with  $\pi_1 = 0.3$ ,  $\mu_1$   
208  $= 0$ ,  $\sigma_1 = 1$ ,  $\pi_2 = 0.7$ ,  $\mu_2 = 8$ , and  $\sigma_2 = 2$ . All other settings are the same as in  
209 Example 1, except that when generating outliers, we add an amount  $\text{Unif}([-5\sigma_1, -7\sigma_1])$   
210 to observations from component 1 and  $\text{Unif}([5\sigma_2, 7\sigma_2])$  to observations from component  
211 2.

212 In the above simulation examples, the majority of data points form two well-separated  
213 clusters. There are very few extreme observations (5% or 10%), which are far away from both  
214 the cluster centers. As such, it is appropriate to model these anomaly observations as outliers  
215 in a two-component mixture model.

## 216 3.2 Methods and Evaluation Metrics

217 We use our proposed RMM approaches with several different penalty forms including  $\ell_0$ ,  $\ell_1$  and  
218 SCAD penalties, denoted as Soft, Hard and SCAD, respectively. For each penalty, our approach  
219 efficiently produces a solution path with varying numbers of outliers. The optimal solution is  
220 selected by the BIC criterion. To investigate the performance of BIC and to better understand  
221 the true potential of each penalization method, we also report an “oracle” estimator, which is  
222 defined as the solution having the best outlier detection performance along the fitted solution

223 path. When there are multiple such solutions on the path, we choose the one gives the best  
 224 parameter estimates. These oracle estimators are denoted as  $\text{Soft}_{\mathcal{O}}$ ,  $\text{Hard}_{\mathcal{O}}$  and  $\text{SCAD}_{\mathcal{O}}$ . We  
 225 note that the oracle estimators rely on the knowledge of the true parameter values, and thus  
 226 they are not feasible to compute in practice. Nevertheless, as we shall see below, they provide  
 227 interesting information about the behaviors of different penalty forms. We also compare our  
 228 RMM approach to the nonrobust maximum likelihood estimation method (MLE) and the  
 229 robust trimmed likelihood estimation method (TLE) proposed by Neykov et al. (2007), with  
 230 the percentage of trimmed data  $\alpha$  set to either 0.05 ( $\text{TLE}_{0.05}$ ) or 0.10 ( $\text{TLE}_{0.1}$ ). TLE methods  
 231 require a cutoff value  $\eta$  to identify extreme residuals; following Gervini and Yohai (2002), we  
 232 set  $\eta = 2.5$ .

233 To evaluate the outlier detection performance, we report (1) the proportion of masking  
 234 (M%), i.e., the fraction of undetected outliers, (2) the proportion of swapping (S%), i.e., the  
 235 fraction of good points labeled as outliers, and (3) the joint detection rate (JD%), i.e., the  
 236 proportion of simulations with 0 masking. Ideally,  $\text{M}\% \approx 0\%$ ,  $\text{S}\% \approx 0\%$  and  $\text{JD}\% \approx 100\%$ .  
 237 To evaluate the performance of parameter estimation, we report both the mean squared errors  
 238 (MSE) and the robust median squared errors (MeSE) of the parameter estimates.

239 A very important usage of mixture model is for clustering. From the fitted mixture model,  
 240 the Bayes classification rule assigns the  $i$ th observation to cluster  $j$  such that  $j = \arg \max_k p_{ik}$ ,  
 241 where  $p_{ik}$ ,  $k = 1, \dots, m$ , are the set of cluster probabilities for the  $i$ th observation directly  
 242 produced from the EM algorithm. We thus compute the average misclassification rate (Mis%)  
 243 to evaluate the clustering performance of each method. We note that for mixture models, there  
 244 are well-known label switching issues (Celeux, et al., 2000; Stephens, 2000; Yao and Lindsay,  
 245 2009; Yao, 2012a, 2012b). Roughly speaking, the mixture likelihood function is invariant to the  
 246 permutation of the component labels, so that the component parameters are not identifiable  
 247 marginally since they are exchangeable. As a consequence, the estimation results from different  
 248 simulation runs are not directly comparable, as the mixture components in each simulation run  
 249 can be labeled arbitrarily. In our examples, the component labels in each simulation are aligned  
 250 to the reference label of the true parameter values, i.e., the labels are chosen by minimizing  
 251 the distance from the resulting parameter estimates to the true parameter values.



### 252 3.3 Results

253 The simulation results are summarized in Tables 1 to 4. Not surprisingly, MLE fails in all the  
254 cases. This demonstrates that robust mixture modeling is indeed needed in the presence of  
255 rare but severe outliers.

256 In case of equal variances, both Hard and SCAD perform very well, and their performance  
257 on outlier detection is very similar to their oracle counterparts. While the Soft method per-  
258 forms well in outlier detection when  $p_{\mathcal{O}} = 5\%$ , its performance becomes much worse when  
259  $p_{\mathcal{O}} = 10\%$  mainly due to masking. The observed nonrobustness of Soft is consistent with  
260 the results in She and Owen (2011). In terms of parameter estimation, Hard and  $\text{Hard}_{\mathcal{O}}$  per-  
261 form the best among the RMM methods. On the other hand,  $\text{SCAD}_{\mathcal{O}}$  performs better than  
262  $\text{Soft}_{\mathcal{O}}$  and they are slightly outperformed by SCAD and Soft, respectively. This interesting  
263 phenomenon reveals some important behaviors of the penalty functions. When using the  $\ell_0$   
264 penalty, the effect of an outlier is completely captured by its estimated mean-shift parameter  
265 whose magnitude is not penalized, so once an observation is detected as an outlier, i.e., its  
266 mean-shift parameter is estimated to be nonzero, it does not affect parameter estimation any  
267 more. However, when using  $\ell_1$  type penalty, due to its inherit shrinkage effects on the mean-  
268 shift parameters, the model tries to accommodate the effects of severe outliers in estimation.  
269 Even if an observation is detected as an outlier with a nonzero mean-shift, it may still partially  
270 affects parameter estimation as the magnitude of the mean-shift parameter is shrunk towards  
271 zero. As a consequence, the oracle estimator which has the best outlier detection performance  
272 does not necessarily leads to the best estimation. Since the SCAD penalty can be regarded as  
273 a hybrid between  $\ell_0$  and  $\ell_1$ , it exhibits behaviors that are characteristics of both of  $\ell_0$  and  $\ell_1$ .  
274 Further examination of the simulation results reveals that  $\text{Soft}_{\mathcal{O}}$  ( $\text{SCAD}_{\mathcal{O}}$ ) tends to require  
275 a stronger penalty than the Soft (SCAD) estimator in order to reduce false positives, which  
276 induces heavier shrinkage of  $\gamma$ , and consequently the former is distorted more by the outliers  
277 than the latter. The TLE method leads to satisfactory results when the trimming proportion  
278 is correctly specified. It loses efficiency when the trimming proportion is too large and fails  
279 to be robust when the trimming proportion is too small. Our RMM methods can achieve

280 comparable performance to the oracle TLE that assumes the correct trimming proportion.

281 In case of unequal variances, the behaviors of the RMM estimators and their oracle counter-  
282 parts are similar to those in the case of equal variances. Hard still performs the best among all  
283 feasible estimators in both outlier detection and parameter estimation. SCAD and Soft work  
284 satisfactorily when  $p_{\mathcal{O}} = 5\%$ . However, when  $p_{\mathcal{O}} = 10\%$ , the two methods may fail to detect  
285 outliers and their average masking rates become 18.72% and 55.67%, respectively. Again, this  
286 can be explained by the shrinkage effects on the mean-shift parameters induced by the penalty  
287 forms. Nevertheless, SCAD is affected much less and thus performs much better in parameter  
288 estimation than Soft.

289 We have investigated the problem of selecting the number of mixture components using  
290 the proposed BIC criterion. In Example 2 with unequal variances and  $p_{\mathcal{O}} = 5\%$ , we use the  
291 RMM method to fit models with 2, 3, and 4 mixture components. The two-component model  
292 is selected 100%, 98% and 63% of the time when using Hard, SCAD and Soft, respectively,  
293 based on 200 simulated datasets. The results are similar using Example 1 and/or  $p_{\mathcal{O}} = 10\%$ .  
294 These results again suggest that RMM works well with nonconvex penalty forms. In Table 5,  
295 we compare the average computation times. As expected, RMM tends to be slightly slower  
296 than TLE and MLE, mainly because the M-step has to be solved by an iterative procedure. In  
297 general, the computation time of RMM increases slightly as the proportion of outliers increases,  
298 and the case of unequal variances needs slightly longer time to compute than the case of  
299 equal variances. Nevertheless, the proposed RMM method remains to be very computationally  
300 efficient and the speed can be further improved with more careful implementation. (A user-  
301 friendly R package for RMM will be made available to the public).

302 In summary, our RMM approach using nonconvex penalization, together with the proposed  
303 BIC criterion, achieves the dual goal of accuracy outlier detection and robust parameter esti-  
304 mation. In practice, the proportion of extreme outliers is usually very small in mixture model  
305 setup, and we suggest to use either the  $\ell_0$  or the SCAD penalty. Other nonconvex penalty  
306 forms such as the minimax concave penalty (MCP) (Zhang, 2010) can also be used.

## 307 4 Acidity Data Analysis

308 We apply the proposed robust procedure to Acidity dataset (Crawford, 1994; Crawford et al.,  
309 1992). The observations are the logarithms of an acidity index measured in a sample of 155  
310 lakes in north-central Wisconsin. More details on the data can be found in Crawford (1994),  
311 Crawford et al. (1992), and Richardson and Green (1997). Figure 1 shows the histogram of  
312 the observed acidity indices.

313 Following Richardson and Green (1997), we fit the data by a three-component normal  
314 mixture model with equal variances, using both the traditional MLE method and the proposed  
315 RMM approach with  $\ell_0$  penalty. The tuning parameter in RMM is selected by BIC. Table  
316 6 shows the parameter estimates. In the original data, there does not appear to be outliers,  
317 and the proposed RMM approach results in very similar parameter estimates to that of the  
318 traditional MLE. This shows that RMM does not lead to efficiency loss when there is no outlier,  
319 and its performance is as good as that of MLE.

320 Following Peel and McLachlan (2000), to examine the effects of outliers, we add one outlier  
321 ( $y = 12$ ) to the original data. While RMM is not influenced by the outlier and gives similar  
322 parameter estimates to the case of no outliers, MLE leads to very different parameter estimates.  
323 Note the first and second components are estimated to have the same mean based on MLE,  
324 thus the model essentially has only two components. We then add three identical outliers  
325 ( $y = 12$ ) to the data. As expected, RMM still provides similar estimates as before. However,  
326 MLE fits a new component to the outliers and gives drastically different estimates comparing  
327 to the case of no outliers. In fact, in both cases, RMM successfully detects the added extreme  
328 observations as outliers, so that the parameter estimation remains unaffected. This example  
329 shows that our proposed RMM method provides a stable and robust way for fitting mixture  
330 models, especially in the presence of severe outliers.

## 331 5 Discussion

332 We have developed a robust mixture modelling approach under the penalized estimation frame-  
333 work. Our robust method with nonconvex penalization is capable of conducting simultaneous

334 outlier detection and robust parameter estimation. The method has comparable performance  
335 to TLE that uses an oracle trimming proportion. However, our method can efficiently produce  
336 a solution path consisting of solutions with varying number of outliers, so that the proportion  
337 of outliers and the accommodation of them in estimation can both be efficiently determined  
338 data adaptively.

339 There are many directions for future research. It is pressing to investigate the theoretical  
340 properties of the proposed RMM approach, e.g., the selection consistency of outlier detection.  
341 As RMM is formulated as a penalized estimation problem, the many established results on  
342 penalized variable selection may shed light on this problem; see. e.g., Khalili (2007) and  
343 Stadler (2010). Our proposed general scaled-dependent outlier detection model shares similar  
344 idea with the reparameterized model proposed by Stadler (2010), and our model can be written  
345 as a penalized mixture regression problem. However, their approach for establishing the oracle  
346 properties of the penalized estimator is not directly applicable to our problem, as in our case  
347 the design matrix associated with the mean-shift parameters becomes a fixed identity matrix of  
348 dimension  $n$ . We have mainly focused on normal mixture model in this paper, but the method  
349 can be readily extended to other mixture models, such as mixtures of binomial and mixtures  
350 of Poisson. It would also be interesting to extend the method to multivariate mixture models  
351 and mixture regression models.

## 352 **Acknowledgements**

353 We thank the two referees and the Associate Editor, whose comments and suggestions have  
354 helped us to improve the paper significantly. Yao's research is supported by NSF grant DMS-  
355 1461677.

356 **Appendix**

357 **Derivation of Equation (2.7)**

358 The estimate of  $\gamma$  is obtained by maximizing

$$\sum_{i=1}^n \sum_{j=1}^m p_{ij}^{(k+1)} \log \phi(y_i - \gamma_i; \mu_j, \sigma^2) - \sum_{i=1}^n \frac{1}{w} P_\lambda(|\gamma_i|).$$

359 The problem is separable in each  $\gamma_i$ , and thus each  $\gamma_i$  can be updated by minimizing

$$- \sum_{j=1}^m p_{ij}^{(k+1)} \log \phi(y_i - \gamma_i; \mu_j, \sigma^2) + \frac{1}{w} P_\lambda(|\gamma_i|).$$

360 Using the form of the normal density, the solution has the following form,

$$\hat{\gamma}_i = \arg \min_{\gamma_i} \sum_{j=1}^m p_{ij} \left\{ \frac{1}{2} \log(\sigma^2) + \frac{(y_i - \gamma_i - \mu_j)^2}{2\sigma^2} \right\} + \frac{1}{w} P_\lambda(|\gamma_i|).$$

Note that  $\sum_{j=1}^m p_{ij} \log(\sigma^2)$  does not depend on  $\gamma$ , and

$$\sum_{j=1}^m p_{ij} \frac{(y_i - \gamma_i - \mu_j)^2}{2\sigma^2} = \frac{1}{2\sigma^2} \left[ \left\{ \gamma_i - \sum_{j=1}^m p_{ij}(y_i - \mu_j) \right\}^2 + \text{const} \right].$$

361 It follows that

$$\hat{\gamma}_i = \arg \min_{\gamma_i} \frac{1}{2\sigma^2} \left[ \left\{ \gamma_i - \sum_{j=1}^m p_{ij}(y_i - \mu_j) \right\}^2 \right] + \frac{1}{w} P_\lambda(|\gamma_i|).$$

362 **Proof of Lemma 1**

363 The penalized least squares has the following form:

$$g(\gamma) = \frac{1}{2}(\xi - \gamma)^2 + \frac{\sigma^2}{\hat{\sigma}^2} P_\lambda(\gamma) \tag{5.1}$$

364 where  $\xi = \{\sum_{j=1}^m p_{ij}(y_i - \mu_j)\} / (\sum_{j=1}^m p_{ij})$ . For simplicity, we have omitted the subscripts in  
 365  $\gamma_i$  and  $\xi_i$ . The first derivative of  $g(\gamma)$  with respect to  $\gamma$  is

$$g'(\gamma) = \gamma - \xi + \text{sgn}(\gamma) \frac{\sigma^2}{\hat{\sigma}^2} P'_\lambda(\gamma).$$

366 We first discuss some possible solutions of (5.1) in three cases.

367 Case 1: when  $|\gamma| \leq \lambda$ , the problem becomes an  $\ell_1$  penalized problem, and the solution, if  
 368 feasible, is given by  $\hat{\gamma}_1 = \text{sgn}(\xi) (|\xi| - \sigma^2\lambda/\hat{\sigma}^2)_+$ .

Case 2: when  $\lambda < |\gamma| \leq a\lambda$ ,  $g''(\gamma) = 1 - \sigma^2/\hat{\sigma}^2/(a-1)$ . The second derivative is positive if  
 $\sigma^2/\hat{\sigma}^2 < a-1$ . The solution, if feasible, is given by

$$\hat{\gamma}_2 = \frac{\frac{\hat{\sigma}^2}{\sigma^2}(a-1)\xi - \text{sgn}(\xi)a\lambda}{\frac{\hat{\sigma}^2}{\sigma^2}(a-1) - 1}.$$

369 Case 3: when  $|\gamma| > a\lambda$ ,  $g''(\gamma) = 1$ . The solution, if feasible, is given by  $\hat{\gamma}_3 = \xi$ .

370 The above three cases indicate that the solution depends on the value  $\sigma^2/\hat{\sigma}^2$  and  $\xi$ . Since  
 371 equation (5.1) is symmetric about  $\xi$  and  $\Theta(-\xi; \lambda) = -\Theta(\xi; \lambda)$ , we shall only discuss the case  
 372  $\xi \geq 0$ .

373 We now derive the solution  $\hat{\gamma}$  in the following scenarios.

374 Scenario 1:  $\sigma^2/\hat{\sigma}^2 < a-1$ .

- 375 1. When  $\xi > a\lambda$ ,  $\gamma$  satisfies Case 3. Then  $\hat{\gamma} = \hat{\gamma}_3$ .
- 376 2. When  $\lambda + \sigma^2\lambda/\hat{\sigma}^2 < \xi \leq a\lambda$ ,  $\gamma$  satisfies Case 2. Then  $\hat{\gamma} = \hat{\gamma}_2$ .
- 377 3. When  $\xi \leq \lambda + \sigma^2\lambda/\hat{\sigma}^2$ ,  $\gamma$  satisfies Case 1. Then  $\hat{\gamma} = \hat{\gamma}_1$ .

378 Scenario 2:  $a-1 \leq \sigma^2/\hat{\sigma}^2 \leq a+1$ . Case 2 is not feasible.

- 379 1. When  $\xi \leq a\lambda$ , based on Case 1,  $\hat{\gamma} = \hat{\gamma}_1$ .
- 380 2. When  $a\lambda \leq \xi \leq \lambda + \sigma^2\lambda/\hat{\sigma}^2$ . As  $|\hat{\gamma}_1| \leq \lambda$  and  $|\hat{\gamma}_3| \geq a\lambda$ , they are both possible solutions.  
 381 Define  $d = g(\hat{\gamma}_1) - g(\hat{\gamma}_3)$ . Then  $\hat{\gamma} = \hat{\gamma}_3$  if  $d > 0$  and  $\hat{\gamma} = \hat{\gamma}_1$  if  $d < 0$ . It can be verified

382 that  $d > 0$  if  $\xi > \frac{a+1+\frac{\sigma^2}{\hat{\sigma}^2}}{2}\lambda$ , and  $d < 0$  if  $\xi < \frac{a+1+\frac{\sigma^2}{\hat{\sigma}^2}}{2}\lambda$ . When  $\xi = \frac{a+1+\frac{\sigma^2}{\hat{\sigma}^2}}{2}\lambda$ , both  $\hat{\gamma}_1$  and  
383  $\hat{\gamma}_3$  are minimizers; in (2.11) we have taken  $\hat{\gamma} = \hat{\gamma}_1$ .

384 3. When  $\xi > \lambda + \sigma^2\lambda/\hat{\sigma}^2$ , then  $\xi > a\lambda$ . Based on Case 3,  $\hat{\gamma} = \xi$ .

385 Scenario 3:  $\sigma^2/\hat{\sigma}^2 > a + 1$ . Case 2 is not feasible.

386 1. When  $\xi > \sigma^2\lambda/\hat{\sigma}^2$ , it is easy to see that  $\hat{\gamma} = \xi$ .

387 2. When  $0 \leq \xi \leq \sigma^2\lambda/\hat{\sigma}^2$ ,  $\hat{\gamma}_1 = 0$  and  $d = g(\hat{\gamma}_1) - g(\hat{\gamma}_3) = \xi^2/2 - \sigma^2(a+1)\lambda^2/(2\hat{\sigma}^2)$ . It  
388 follows that  $d > 0$  if  $\xi > \sqrt{\frac{\sigma^2(a+1)}{\hat{\sigma}^2}}\lambda$ ,  $d < 0$  if  $\xi < \sqrt{\frac{\sigma^2(a+1)}{\hat{\sigma}^2}}\lambda$ . When  $\xi = \sqrt{\frac{\sigma^2(a+1)}{\hat{\sigma}^2}}\lambda$ ,  
389 both  $\hat{\gamma}_1 = 0$  and  $\hat{\gamma}_3 = \xi$  are minimizers; in (2.12) we have taken  $\hat{\gamma} = \hat{\gamma}_1 = 0$ .

390 Combining the three scenarios leads to the modified SCAD thresholding rule in Lemma 1. We  
391 note that in practice, as  $\sigma^2/\hat{\sigma}^2$  is close to one, Scenarios 2 and 3 are highly unlikely to occur.

## 392 References

393 Antoniadis, A. (1997). Wavelets in Statistics: A Review (with discussion). *Journal of the*  
394 *Italian Statistical Association*, 6, 97-144.

395 Bai, X., Yao, W., and Boyer, J. E. (2012). Robust fitting of mixture regression models. *Com-*  
396 *putational Statistics and Data Analysis*, 56, 2347-2359.

397 Bashir, S. and Carter, E. (2012). Robust mixture of linear regression models. *Communications*  
398 *in Statistics-Theory and Methods*, 41, 3371-3388.

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

401 Celeux, G., Hurn, M., and Robert, C. P. (2000). Computational and inferential difficulties  
402 with mixture posterior distributions. *Journal of the American Statistical Association*, 95,  
403 957-970.

404 Crawford, S. L., Degroot, M. H., Kadane, J. B., and Small, M. J. (1992). Modeling  
405 lake-chemistry distributions-approximate Bayesian methods for estimating a finite-mixture  
406 model. *Technometrics*, 34, 441-453.

407 Crawford, S. L. (1994). An application of the Laplace method to finite mixture distributions.  
408 *Journal of the American Statistical Association*, 89, 259-267.

409 Dalalyan, A. S. and Keriven, R.(2012). Robust estimation for an inverse problem arising in  
410 multiview geometry. *Journal of Mathematical Imaging and Vision*, 43, 10-23.

411 Dalalyan, A. S. and Chen, Y. (2012). Fused sparsity and robust estimation for linear models  
412 with unknown variance. *In Advances in Neural Information Processing Systems*, 25: NIPS,  
413 1268-1276

414 Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete  
415 data via the EM algorithm (with discussion). *Journal of Royal Statistical Society, Ser B.*,  
416 39, 1-38.

417 Donoho, D. L. and Johnstone, I. M. (1994a), Ideal Spatial Adaptation by Wavelet shrinkage,  
418 *Biometrika*, 81, 425-455.

419 Everitt, B. S. and Hand D. J. (1981). Finite Mixture Distributions. Chapman and Hall, London

420 Fan, J. and Li, R. (2001). Variable Selection via Nonconcave Penalized Likelihood and its  
421 Oracle Properties. *Journal of the American Statistical Association*, 96, 1348-1360.

422 Fujisawa, H. and Eguchi S. (2005). Robust estimation in the normal mixture model. *Journal*  
423 *of Statistical Planning and Inference*, 1-23.

424 Gervini, D. and Yohai, V. J. (2002), A Class of Robust and Fully Efficient Regression Estima-  
425 tors. *The Annals of Statistics*, 30, 583-616.

426 Goldfeld, S. M. and Quandt, R. E. (1973). A Markov model for switching regression. *Journal*  
427 *of Econometrics*, 1, 3-15.



- 428 Hennig, C. (2000). Identifiability of models for clusterwise linear regression. *Journal of Clas-*  
429 *sification*. 17, 273-296.
- 430 Hennig, C. (2002). Fixed point clusters for linear regression: computation and comparison.  
431 *Journal of Classification*, 19, 249-276
- 432 Hennig, C. (2003). Clusters, outliers, and regression: Fixed point clusters. *Journal of Multi-*  
433 *variate Analysis*, 86, 183-212.
- 434 Huber, P.J. (1981), Robust Statistics. *New York: John Wiley and Sons*.
- 435 Khalili, A., Chen, J.H. (2007) Variable Selection in Finite Mixture of Regression Models.  
436 *Journal of the American Statistical Association*, 102, 1025–1038.
- 437 Lee, Y., MacEachern, S. N., and Jung, Y. (2012), Regularization of Case-Specific Parameters  
438 for Robustness and Efficiency. *Statistical Science*, 27, 350-372.
- 439 Lindsay, B. G. (1995). Mixture Models: Theory, Geometry and Applications, *NSF-CBMS*  
440 *Regional Conference Series in Probability and Statistics*, Vol. 5. Institute of mathematical  
441 Statistics and the American Statistical Association, Alexandria, VA.
- 442 Markatou, M. (2000). Mixture models, robustness, and the weighted likelihood methodology.  
443 *Biometrics*, 56, 483-486.
- 444 McLachlan, G. J. and Basford, K. E. (1988). Mixture Models: Inference and Applications to  
445 Clustering. Marcel Dekker, New York.
- 446 McLachlan, G. J. and Peel, D. (2000). *Finite Mixture Models*. New York: Wiley.
- 447 Neykov, N., Filzmoser, P., Dimova, R., and Neytchev, P. (2007). Robust fitting of mixtures  
448 using the trimmed likelihood estimator. *Computational Statistics and Data Analysis*, 52,  
449 299-308.
- 450 Nguyen, N. H. and Tran, T. D. (2013). Robust Lasso With Missing and Grossly Corrupted  
451 Observations. *IEEE Transactions on Information Theory*, 59, 2036-2058.

- 452 Peel, D. and McLachlan, G. J. (2000). Robust mixture modelling using the t distribution  
453 *Statistics and Computing*, 10, 339-348.
- 454 Richardson, S. and Green, P. J. (1997). On Bayesian analysis of mixtures with an unknown  
455 number of components (with discussion). *Journal of The Royal Statistical Society Series, B*  
456 , 59, 731-792.
- 457 Rousseeuw, P.J.(1983), Multivariate Estimation with High Breakdown Point. Resaerch Report  
458 No. 192, Center for Statistics and Operations research, VUB Brussels.
- 459 Rousseeuw, P.J. and Yohai, V. J. (1984). Robust Regression by Means of S-estimators. *Robust*  
460 *and Nonlinear Time series*, J. Franke, W. Härdle and R. D. Martin (eds.), Lectures Notes in  
461 Statistics 26, 256-272, New York: Springer.
- 462 She, Y. and Owen, A. (2011). Outlier Detection Using Nonconvex Penalized Regression. *Jour-*  
463 *nal of the American Statistical Association*, 106(494), 626-639.
- 464 Shen, H., Yang, J., and Wang, S. (2004). Outlier detecting in fuzzy switching regression models.  
465 *Artificial Intelligence: Methodology, Systems, and Applications Lecture Notes in Computer*  
466 *Science*, 2004, Vol. 3192/2004, 208-215.
- 467 Siegel, A.F. (1982), Robust Regression Using Repeated Medians. *Biometrika*, 69, 242-244.
- 468 Song, W., Yao, W., and Xing Y. (2014). Robust mixture regression model fitting by laplace  
469 distribution. *Computational Statistics and Data Analysis*, 71, 128-137.
- 470 Stadler, N., Buhlmann, P., and van de Geer, S. (2010).  $\ell_1$ -penalization for mixture regression  
471 models. *Test*, 19(2), 209–256.
- 472 Stephens, M. (2000). Dealing with label switching in mixture models. *Journal of Royal Statis-*  
473 *tical Society, Ser B.*, 62, 795-809.
- 474 Tibshirani, R. J. (1996). Regression Shrinkage and Selection via the LASSO. *Journal of The*  
475 *Royal Statistical Society Series, B*58, 267-288.

476 Tibshirani, R. J. (1996). The LASSO Method for Variable Selection in the Cox Model. *Statistics*  
477 *in Medicine*, 16, 385-395.

478 Titterton, D. M., Smith, A. F. M., and Makov, U. E. (1985). Statistical Analysis of Finite  
479 Mixture Distribution. Wiley, New York.

480 Yao, W. (2012a). Model based labeling for mixture models. *Statistics and Computing*, 22,  
481 337-347.

482 Yao, W. (2012b). Bayesian mixture labeling and clustering. *Communications in Statistics -*  
483 *Theory and Methods*, 41, 403-421.

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

486 Yao, W., Wei, Y., and Yu, C. (2014). Robust Mixture Regression Using T-Distribution. *Com-*  
487 *putational Statistics and Data Analysis*, 71, 116-127.

488 Yohai, V. J. (1987), High Breakdown-point and High Efficiency Robust Estimates for Regres-  
489 sion. *The Annals of Statistics*, 15, 642-656.

490 Zhang, C.-H. (2010). Nearly unbiased Variable Selection Under Minimax Concave Penalty.  
491 *Ann. Statist.* 38 (2), 894–942.

Table 1: Simulation results for the case of equal variances with  $n = 200$  and  $p_{\mathcal{O}} = 5\%$ .

|                  | Hard   | Hard $_{\mathcal{O}}$ | SCAD   | SCAD $_{\mathcal{O}}$ | Soft   | Soft $_{\mathcal{O}}$ | TLE $_{0.05}$ | TLE $_{0.10}$ | MLE    |
|------------------|--------|-----------------------|--------|-----------------------|--------|-----------------------|---------------|---------------|--------|
| M%               | 0.00   | 0.00                  | 0.00   | 0.00                  | 0.00   | 0.00                  | 0.06          | 0.06          | –      |
| S%               | 0.27   | 0.02                  | 0.99   | 0.03                  | 0.42   | 0.03                  | 1.04          | 3.34          | –      |
| JD%              | 100.00 | 100.00                | 100.00 | 100.00                | 100.00 | 100.00                | 99.44         | 99.44         | –      |
| Mis%             | 0.26   | 0.02                  | 0.94   | 0.02                  | 0.40   | 0.03                  | 0.07          | 5.01          | 15.53  |
| MeSE( $\pi$ )    | 0.001  | 0.001                 | 0.001  | 0.001                 | 0.001  | 0.001                 | 0.001         | 0.002         | 0.002  |
| MSE( $\pi$ )     | 0.002  | 0.002                 | 0.002  | 0.002                 | 0.002  | 0.002                 | 0.002         | 0.003         | 0.030  |
| MeSE( $\mu$ )    | 0.018  | 0.017                 | 0.035  | 0.052                 | 0.055  | 0.065                 | 0.017         | 0.031         | 0.293  |
| MSE( $\mu$ )     | 0.023  | 0.022                 | 0.041  | 0.063                 | 0.061  | 0.071                 | 0.022         | 0.038         | 11.150 |
| MeSE( $\sigma$ ) | 0.009  | 0.010                 | 0.067  | 0.191                 | 0.176  | 0.231                 | 0.008         | 0.064         | 0.952  |
| MSE( $\sigma$ )  | 0.016  | 0.016                 | 0.088  | 0.191                 | 0.198  | 0.242                 | 0.012         | 0.071         | 25.478 |

Table 2: Simulation results for the case of equal variances with  $n = 200$  and  $p_{\mathcal{O}} = 10\%$ .

|                  | Hard   | Hard $_{\mathcal{O}}$ | SCAD   | SCAD $_{\mathcal{O}}$ | Soft  | Soft $_{\mathcal{O}}$ | TLE $_{0.05}$ | TLE $_{0.10}$ | MLE    |
|------------------|--------|-----------------------|--------|-----------------------|-------|-----------------------|---------------|---------------|--------|
| M%               | 0.00   | 0.00                  | 0.00   | 0.00                  | 12.11 | 0.00                  | 24.53         | 0.00          | –      |
| S%               | 0.32   | 0.04                  | 2.89   | 0.04                  | 0.80  | 0.04                  | 0.19          | 1.19          | –      |
| JD%              | 100.00 | 100.00                | 100.00 | 100.00                | 72.78 | 100.00                | 2.78          | 100.00        | –      |
| Mis%             | 0.29   | 0.05                  | 2.61   | 0.03                  | 1.93  | 0.04                  | 5.94          | 0.09          | 22.28  |
| MeSE( $\pi$ )    | 0.001  | 0.001                 | 0.001  | 0.001                 | 0.001 | 0.001                 | 0.004         | 0.001         | 0.003  |
| MSE( $\pi$ )     | 0.002  | 0.002                 | 0.002  | 0.002                 | 0.002 | 0.002                 | 0.009         | 0.002         | 0.053  |
| MeSE( $\mu$ )    | 0.020  | 0.019                 | 0.061  | 0.183                 | 0.171 | 0.212                 | 0.840         | 0.019         | 0.918  |
| MSE( $\mu$ )     | 0.023  | 0.024                 | 0.066  | 0.209                 | 0.230 | 0.231                 | 1.093         | 0.023         | 14.125 |
| MeSE( $\sigma$ ) | 0.012  | 0.010                 | 0.120  | 0.700                 | 0.590 | 0.815                 | 9.164         | 0.010         | 2.648  |
| MSE( $\sigma$ )  | 0.016  | 0.014                 | 0.139  | 0.698                 | 0.742 | 0.809                 | 6.345         | 0.012         | 12.599 |

Table 3: Simulation results for the case of unequal variances with  $n = 200$  and  $p_{\mathcal{O}} = 5\%$ .

|                  | Hard   | Hard $_{\mathcal{O}}$ | SCAD   | SCAD $_{\mathcal{O}}$ | Soft   | Soft $_{\mathcal{O}}$ | TLE $_{0.05}$ | TLE $_{0.10}$ | MLE     |
|------------------|--------|-----------------------|--------|-----------------------|--------|-----------------------|---------------|---------------|---------|
| M%               | 0.00   | 0.00                  | 0.00   | 0.00                  | 0.00   | 0.00                  | 0.94          | 0.06          | –       |
| S%               | 0.13   | 0.04                  | 1.12   | 0.23                  | 1.32   | 0.29                  | 0.73          | 3.12          | –       |
| JD%              | 100.00 | 100.00                | 100.00 | 100.00                | 100.00 | 100.00                | 93.89         | 99.44         | –       |
| Mis%             | 0.51   | 0.44                  | 1.48   | 1.35                  | 2.24   | 1.87                  | 3.88          | 6.22          | 44.82   |
| MeSE( $\pi$ )    | 0.001  | 0.001                 | 0.001  | 0.003                 | 0.004  | 0.006                 | 0.001         | 0.001         | 0.024   |
| MSE( $\pi$ )     | 0.002  | 0.002                 | 0.002  | 0.005                 | 0.004  | 0.007                 | 0.008         | 0.002         | 0.148   |
| MeSE( $\mu$ )    | 0.038  | 0.042                 | 0.051  | 0.081                 | 0.063  | 0.087                 | 0.042         | 0.056         | 77.214  |
| MSE( $\mu$ )     | 0.048  | 0.051                 | 0.068  | 0.115                 | 0.080  | 0.134                 | 3.060         | 0.073         | 141.426 |
| MeSE( $\sigma$ ) | 0.022  | 0.019                 | 0.149  | 0.730                 | 1.121  | 2.133                 | 0.026         | 0.112         | 7.711   |
| MSE( $\sigma$ )  | 0.028  | 0.024                 | 0.177  | 1.474                 | 1.121  | 2.345                 | 0.172         | 0.121         | 10.154  |

Table 4: Simulation results for the case of unequal variances with  $n = 200$  and  $p_{\mathcal{O}} = 10\%$ .

|                  | Hard  | Hard $_{\mathcal{O}}$ | SCAD  | SCAD $_{\mathcal{O}}$ | Soft   | Soft $_{\mathcal{O}}$ | TLE $_{0.05}$ | TLE $_{0.10}$ | MLE     |
|------------------|-------|-----------------------|-------|-----------------------|--------|-----------------------|---------------|---------------|---------|
| M%               | 0.08  | 0.00                  | 18.72 | 1.70                  | 55.67  | 1.90                  | 24.44         | 1.11          | –       |
| S%               | 0.10  | 0.07                  | 2.49  | 0.83                  | 0.20   | 0.94                  | 0.06          | 0.77          | –       |
| JD%              | 98.33 | 100.00                | 66.67 | 68.67                 | 5.56   | 65.33                 | 1.11          | 83.89         | –       |
| Mis%             | 0.46  | 0.42                  | 6.14  | 4.35                  | 11.48  | 4.82                  | 23.96         | 7.65          | 47.99   |
| MeSE( $\pi$ )    | 0.001 | 0.002                 | 0.002 | 0.019                 | 0.030  | 0.023                 | 0.024         | 0.002         | 0.112   |
| MSE( $\pi$ )     | 0.002 | 0.003                 | 0.008 | 0.019                 | 0.032  | 0.025                 | 0.066         | 0.049         | 0.168   |
| MeSE( $\mu$ )    | 0.036 | 0.037                 | 0.095 | 0.165                 | 0.212  | 0.193                 | 10.861        | 0.044         | 79.288  |
| MSE( $\mu$ )     | 0.044 | 0.046                 | 0.136 | 0.222                 | 0.265  | 0.239                 | 17.001        | 21.439        | 193.846 |
| MeSE( $\sigma$ ) | 0.029 | 0.024                 | 0.613 | 7.306                 | 11.553 | 7.734                 | 11.059        | 0.028         | 13.128  |
| MSE( $\sigma$ )  | 0.035 | 0.033                 | 3.416 | 6.088                 | 11.396 | 7.482                 | 10.261        | 0.917         | 16.203  |

Table 5: Comparison of average computation times in seconds. To make fair comparison, each reported time is the average computation time per each tuning parameter and simulated dataset.

| Example | $p_{\mathcal{O}}$ | Hard  | SCAD  | Soft  | TLE <sub>0.05</sub> | TLE <sub>0.1</sub> | MLE   |
|---------|-------------------|-------|-------|-------|---------------------|--------------------|-------|
| 1       | 5%                | 0.039 | 0.041 | 0.042 | 0.041               | 0.042              | 0.016 |
| 1       | 10%               | 0.043 | 0.043 | 0.046 | 0.089               | 0.045              | 0.025 |
| 2       | 5%                | 0.081 | 0.128 | 0.166 | 0.083               | 0.076              | 0.008 |
| 2       | 10%               | 0.084 | 0.112 | 0.201 | 0.179               | 0.088              | 0.007 |

Table 6: Parameter estimation in Acidity data analysis.

|      | #outlier | $\pi_1$ | $\pi_2$ | $\pi_3$ | $\mu_1$ | $\mu_2$ | $\mu_3$ | $\sigma$ |
|------|----------|---------|---------|---------|---------|---------|---------|----------|
| MLE  | 0        | 0.589   | 0.138   | 0.273   | 4.320   | 5.682   | 6.504   | 0.365    |
|      | 1        | 0.327   | 0.324   | 0.349   | 4.455   | 4.455   | 6.448   | 0.687    |
|      | 3        | 0.503   | 0.478   | 0.019   | 5.105   | 5.105   | 12.00   | 1.028    |
| Hard | 0        | 0.588   | 0.157   | 0.255   | 4.333   | 5.720   | 6.545   | 0.336    |
|      | 1        | 0.591   | 0.157   | 0.252   | 4.333   | 5.723   | 6.548   | 0.334    |
|      | 3        | 0.597   | 0.157   | 0.246   | 4.333   | 5.729   | 6.553   | 0.331    |

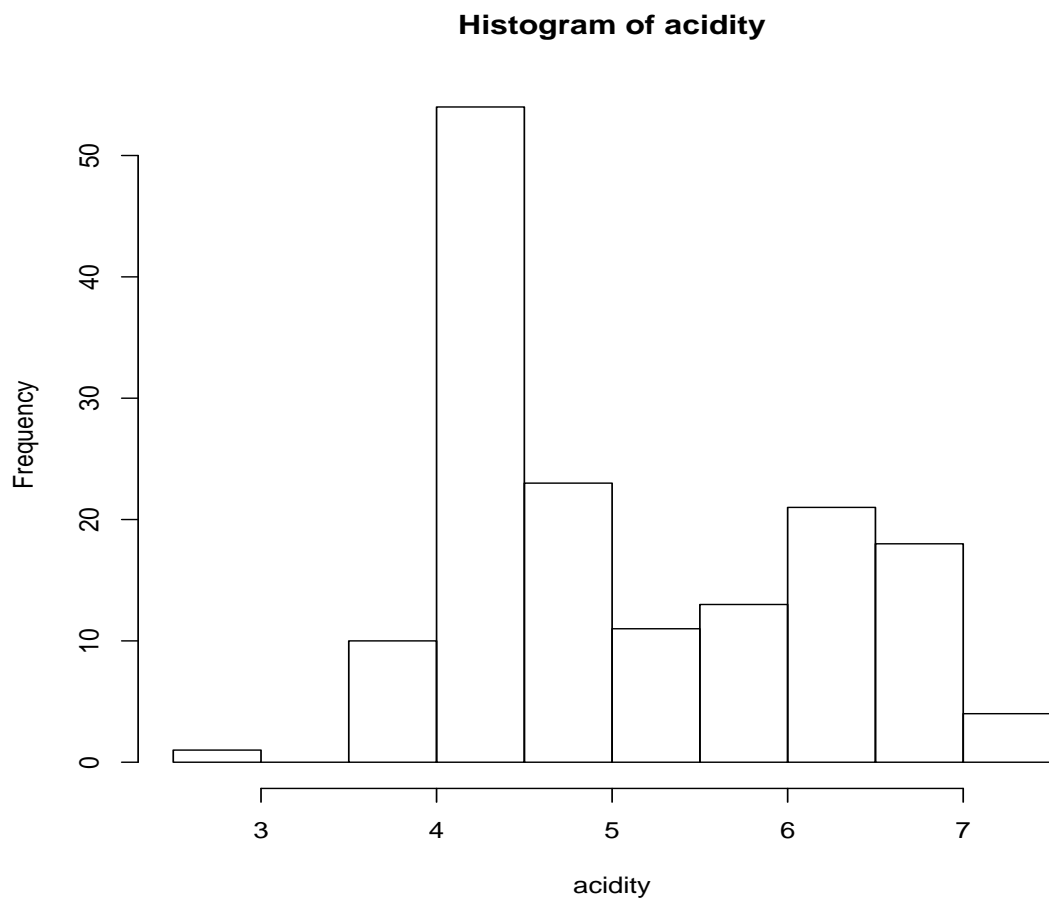


Figure 1: Histogram for Acidity data