
Fitting large mixture models using stochastic component selection

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Traditional methods for unsupervised learning of finite mixture models require to
2 evaluate the likelihood of all components of the mixture. This becomes computa-
3 tionally prohibitive when the number of components is large, as it is, for example,
4 in the sum-product (transform) networks. As a remedy, we propose an approach
5 combining the expectation maximization and the Metropolis-Hastings algorithm
6 to evaluate only a small number of, stochastically sampled, components, thus
7 substantially reducing the computational cost. We put emphasis on generality of
8 our method, equipping it with the ability to train both shallow and deep mixture
9 models which involve complex, and possibly nonlinear, transformations. The
10 performance of our method is illustrated in a variety of synthetic and real-data
11 contexts, considering deep models, such as mixtures of normalizing flows and
12 sum-product (transform) networks.

13 1 Introduction

14 Finite mixture models [40] constitute a fundamental class of density estimation models. They
15 have been successfully applied in diverse fields, including bioinformatics [49], econometrics [10],
16 engineering [33], etc. A mixture model relies on a weighted sum of probability distributions—here
17 referred to as *components*—to cluster N unlabelled datapoints into K categories. The traditional
18 maximum likelihood techniques train the model by optimizing either (i) the marginal likelihood via
19 gradient-descent [50] or (ii) the evidence lower bound via variational methods [4], including the
20 expectation-maximization (EM) [13]. The dependence structure among approximate, variational,
21 distributions then ranges from the fully independent (mean-field) [25] to fully dependent [30]. The
22 sampling-based techniques target the posterior distribution using sequential Monte Carlo [9] or
23 Markov chain Monte Carlo [52], e.g. via the Gibbs [34] or Metropolis-Hastings sampling [38]. The
24 computational cost of these methods typically scales with $\mathcal{O}(TKND)$ operations, where N and K
25 are defined above, T is the number of iterations and D is the dimension of data.

26 Various methods to decrease the computational cost via any factor in $\mathcal{O}(TKND)$ have been proposed.
27 T can be lowered by proper initialization, e.g. the optimal seeding [5]; an efficient step-size schedule,
28 e.g. the line-search [58]; or increased estimation precision, e.g. the variance reduction [8]. N is often
29 reduced using the coreset methods, which approximate the original dataset by a weighted dataset
30 such that the exact and approximate marginal likelihoods are close. The weighted variants of the
31 variational [17, 59, 6] and sampling-based [39] methods then process the coresets. Reducing D relies
32 on the compression of data into smaller representations via random projections [53, 2], which is
33 achieved in two ways: (i) each data item is projected into an individual representation [11]; (ii) all
34 data items are projected into an overall representation, commonly referred to as sketch [28, 22].

35 Nevertheless, all the aforementioned techniques—including those with reduced computational cost—
36 evaluate all K components. This is very demanding for large models, and the problem is even more

Table 1: The computational features of various EM algorithms. We compare whether the methods (i) perform the computations with a reduced number of data (minibatching), (ii) update a lower number of statistics, (iii) make less evaluations of the conditional likelihood, and (iv) are suitable for training of deep models. Here, EM, SA, S, T, MC and MH stand for expectation-maximization, stochastic approximation, sparse, truncated, Monte Carlo and Metropolis-Hastings, respectively.

Feature / Algorithm	EM [13]	SAEM [44]	SSAEM [24]	TSAEM [18]	MCSAEM [1]	MHSAEM (ours)
$B < N$ datapoints	✗	✓	✓	✓	✓	✓
$M < K$ statistics	✗	✗	✓	✓	✓	✓
$M < K$ likelihoods	✗	✗	✗	✓	✗	✓
deep models	✗	✗	✗	✗	✗	✓

37 severe for mixtures involving intricate models, such as neural networks [21, 42], Gaussian processes
 38 [57], normalizing flows [48]; or deep mixtures, including sum-product (transform) networks [45, 47],
 39 deep Gaussian mixture models [55], etc. In spite of this, a little attention has been paid to the design
 40 of algorithms which does not evaluate all K components. The notable exceptions are the sparse EM
 41 algorithm [24] and the truncated variational EM algorithm [18], see Table 1 and Section 5 for details.
 42 Moreover, the methods are mostly tailored for a specific class of mixture models, e.g. the Gaussian
 43 mixture models.

44 In this paper, we make the following contributions:

- 45 • We propose an EM-based algorithm which relies on the MH sampler to stochastically evaluate less
 46 components in mixture models, substantially reducing the computational cost.
- 47 • We design our method to enable optimization of fairly generic EM objective functions, making it
 48 suitable for training of both shallow and deep mixture models.
- 49 • We apply our approach to Gaussian mixture models (GMMs) and their generalizations: sum-
 50 product-transform networks (SPTNs) and mixtures of real-valued non-volume preserving (real
 51 NVP) flows [15], reaching approximately $100\times$ speed-up compared to state-of-the-art methods.

52 2 Problem formulation

53 A finite mixture model characterizes the relation between an observed (known) variable, $x \in \mathcal{X} \subseteq \mathbb{R}^D$,
 54 and a latent (unknown) variable, $z \in \mathcal{Z} := \{1, \dots, K\}$, via the marginal (incomplete-data) likelihood
 55 in the following form:

$$p_\theta(x) = \sum_{k=1}^K p_{\eta_k}(x|z=k)p_{\pi_k}(z=k), \quad (1)$$

56 where $\theta := (\pi_1, \eta_1, \dots, \pi_K, \eta_K) \in \Theta$ are unknown parameters. Here, η_z are the parameters of the
 57 conditional likelihood, $p_{\eta_z}(x|z)$, and π_z is the weight which parameterizes the prior, $p_{\pi_z}(z) = \pi_z$,
 58 and satisfies $0 \leq \pi_k \leq 1$ for each $k \in \mathcal{Z}$ and $\sum_{k=1}^K \pi_k = 1$.

59 Given a set of independent and identically distributed data, $\mathbf{x} := (x_i)_{i=1}^N$, our goal is to learn the
 60 unknown parameters of the marginal log-likelihood,

$$\mathcal{L}(\theta) := \log p_\theta(\mathbf{x}) = \sum_{i=1}^N \log \sum_{k=1}^K p_{\eta_k}(x_i|z_i=k)p_{\pi_k}(z_i=k). \quad (2)$$

61 The marginalization in (2) is tractable for almost all forms of $p_{\eta_z}(x|z)$. Indeed, we consider $p_{\eta_z}(x|z)$
 62 to belong to an arbitrary family of η_z -differentiable probability distributions. However, we assume that
 63 K is high, making the marginalization in (2) computationally costly, thus rendering the optimization
 64 objective presumably intractable. Therefore, we want to design a computationally efficient algorithm,
 65 requiring only $M < K$ evaluations of $p_{\eta_z}(x|z)$ at each iteration.

66 3 The EM algorithm

67 The maximum likelihood estimation seeks the parameters maximizing the marginal log-likelihood,
 68 $\theta^{ML} := \arg \max_{\theta \in \Theta} \mathcal{L}(\theta)$. The traditional EM algorithm [13] addresses this task indirectly, i.e. by
 69 optimizing the evidence lower bound (ELBO),

$$\mathcal{L}(\theta) \geq \mathcal{Q}(\theta) + \mathcal{H}(\hat{\theta}) := \text{ELBO}(\hat{\theta}), \quad (3)$$

70 where $\mathcal{H}(\hat{\theta}) := -\mathbb{E}_{p_{\hat{\theta}}(\mathbf{z}|\mathbf{x})}[\log p_{\hat{\theta}}(\mathbf{z}|\mathbf{x})]$ is the differential entropy at an estimate, $\hat{\theta} \in \Theta$, and

$$\mathcal{Q}(\theta) := \mathbb{E}_{p_{\hat{\theta}}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{z}, \mathbf{x})] = \sum_{i=1}^N \sum_{k=1}^K p_{\theta}(z_i = k | x_i) \log p_{\theta}(z_i = k, x_i) \quad (4)$$

71 is the EM objective function. Here, $p_{\theta}(\mathbf{z}, \mathbf{x})$ is the joint (complete-data) likelihood, and $p_{\theta}(\mathbf{z}|\mathbf{x})$ is
 72 the posterior distribution over the latent variables $\mathbf{z} := (z_i)_{i=1}^N$. Given an initial value, θ_0 , the EM
 73 algorithm produces a sequence of estimates, $(\theta_t)_{t=1}^T$, by alternating between the expectation (E) and
 74 maximization (M) steps,

$$\text{E-step: } \mathcal{Q}_{t-1}(\theta), \quad (5)$$

$$\text{M-step: } \theta_t := \arg \max_{\theta \in \Theta} \mathcal{Q}_{t-1}(\theta). \quad (6)$$

75 This sequence is guaranteed to monotonically tighten the ELBO, arriving at a local optimum of (2)
 76 under mild regularity assumptions [56].

77 The EM algorithm is computationally expensive, since (4) evaluates $p_{\theta}(z_i, x_i)$ for each $z_i \in \mathcal{Z}$ and
 78 $i \in (1, \dots, N)$. This has to be performed for all $t \in (1, \dots, T)$ in (5). Albeit the marginal factor,
 79 $p_{\pi_z}(z)$, is just the cheap categorical distribution, the conditional factor, $p_{\eta_z}(x|z)$, typically involves
 80 high-dimensional operations (e.g., the inversion of the full $D \times D$ -dimensional covariance matrices
 81 in the GMMs). Moreover, the M-step (6) is also expensive for large K . This holds despite that (6)
 82 can be reduced to closed-form updates of expected sufficient statistics for $p_{\eta_z}(x|z)$ belonging to the
 83 exponential family [44] (again, due to high D). All in all, the computational complexity of the EM
 84 algorithm scales with $\mathcal{O}(TDNK)$.

85 If (6) cannot be computed under a closed-form solution, one can resort to direct gradient-descent
 86 optimization of $\mathcal{Q}(\theta)$, where $\arg \max$ is replaced by one (or more) step(s) of a gradient descent
 87 technique. The EM algorithm is then referred to as the generalized EM algorithm [56].

88 4 The generalized MHSSEM algorithm

89 We design a version of the generalized EM algorithm suitable for scenarios where (4) can represent
 90 deep, discrete, latent variable models, thus being parameterized by possibly complex nonlinear
 91 transformations. We particularly focus on decreasing the the number of operations in the generalized
 92 EM algorithm from $\mathcal{O}(TDNK)$ to $\mathcal{O}(TDBM)$, where $B \ll N$ and $M \ll K$.

93 4.1 E-step

94 We reduce the cost of evaluating the EM objective function (4) by combining the minibatching (as
 95 used many times before) and the Monte Carlo sampling. Namely, the specific application of the latter
 96 to generic mixture models is the key contribution of this paper.

97 *Minibatching.* At each iteration, t , we compute the conditional expectation in (4) only for a subset—
 98 here referred to as a *minibatch*—of the original full dataset, i.e. $(x_i)_{i \in I}$. Here, I is a set of $B \ll N$
 99 indices, i , sampled uniformly without replacement from $(1, \dots, N)$. This substantially decreases the
 100 necessary computations compared to the full sweep over all N datapoints [23].

101 *Monte Carlo sampling.* For each $i \in I$, we want to draw $M \ll K$ random samples from $p_{\theta}(z_i|x_i)$ in
 102 order to obtain a Monte Carlo estimate of (4). The straightforward way to do this would be to draw
 103 the samples directly from $p_{\theta}(z_i|x_i)$. However, direct sampling from $p_{\theta}(z_i|x_i)$ does not lead to any
 104 substantial decrease in the number of operations. This is caused by the fact that even for a *single*
 105 sample of z_i , we have to first compute the normalizing factor, $p_{\theta}(x_i)$, to obtaining the posterior,

106 $p_\theta(z_i|x_i)$. This requires K expensive evaluations of $p_\theta(z_i, x_i)$, which is precisely what we want to
 107 avoid. Our approach is to resort to the Markov chain Monte Carlo (MCMC), which allows us to
 108 sample from $p_\theta(z_i|x_i)$, with the computational complexity decreasing to only a *single* evaluation of
 109 $p_\theta(z_i, x_i)$ per a *single* sample of z_i .

110 MCMC methods obviate the computation of the normalizing factor in $p_\theta(z_i|x_i)$ by simulating a
 111 Markov chain, $(z_{i,t})_{t=1}^T$, from a transition kernel, $z_{i,t} \sim P(z_{i,t-1}, \cdot)$, which leaves $p_\theta(z_i|x_i)$ as its
 112 unique stationary (invariant) distribution, starting from an initial value $z_{i,0}$. The specific form of P
 113 determines the structure of an MCMC method. We chose the Metropolis-Hastings (MH) sampler,
 114 which represents $P(z_{i,t-1}, z_{i,t})$ as follows: given $\bar{z}_i := z_{i,t-1}$, draw a sample from the *proposal*
 115 distribution $z_i \sim q(\cdot|\bar{z}_i)$, compute the acceptance ratio,

$$\alpha(\bar{z}_i, z_i) := \min \left\{ 1, \frac{p_{\eta_{z_i,t-1}}(x_i|z_i)\pi_{z_i,t-1}q(\bar{z}_i|z_i)}{p_{\eta_{\bar{z}_i,t-1}}(x_i|\bar{z}_i)\pi_{\bar{z}_i,t-1}q(z_i|\bar{z}_i)} \right\}, \quad (7)$$

116 and, if $u < \alpha(\bar{z}_i, z_i)$ —where u is drawn from a uniform distribution, $\text{Uniform}(0, 1)$ —accept the
 117 sample and set $z_{i,t} = z_i$; otherwise, set $z_{i,t} = \bar{z}_i$. For each $i \in I$ and $t \in (1, \dots, T)$, we repeat
 118 this process M times, construing a set $\mathbf{z}_{i,t} = (z_{i,t}^1, \dots, z_{i,t}^M)$. Therefore, at every current iteration,
 119 t , we continue to extend the chain from the point where we left at the previous iteration, $t - 1$, by
 120 taking $\bar{z}_i = z_{i,t-1}^M$. Under mild regularity assumptions [52], the chain passes the transition period
 121 (the burn-in phase), and the samples can then be used to approximate the conditional expectation in
 122 (4) as follows:

$$\hat{Q}_{t-1}(\theta) = \frac{1}{M} \sum_{i \in I} \sum_{z \in \mathbf{z}_{i,t}} \log p_{\eta_z}(x_i|z) \pi_z. \quad (8)$$

123 Note that, to ensure this approach is truly efficient, we have to draw only $M \ll K$ samples at each
 124 iteration, t ; otherwise, for $M \approx K$, we may rather compute the exact marginalization in (4), since it
 125 is tractable (but computationally costly).

126 4.2 M-step

127 Assume for a moment that (6) with $Q_{t-1}(\theta)$ given by (8) has a closed-form solution, yielding an
 128 estimate of θ . Such an estimate would have a high variance, converging only for $M \rightarrow \infty$ and
 129 $T \rightarrow \infty$ [19]. The main reason is that the samples would not be reused over the iterations, t ,
 130 thus wasting computational resources. We consider that there is no closed-form solution of (6),
 131 and—to ensure that the samples (and thus computations) are recycled over the iterations—we use
 132 the stochastic approximation (SA) [51] to optimize (8). This is analogous to applying a stochastic
 133 gradient-descent method, $\theta_t = \theta_{t-1} + \gamma_t \nabla_\theta \hat{Q}_{t-1}(\theta)$, where γ_t is the step-size, satisfying the Robbins-
 134 Monro constraints, $\gamma_t \in [0, 1]$, $\sum_{t \geq 1} \gamma_t = \infty$, $\sum_{t \geq 1} \gamma_t^2 < \infty$, and ∇_θ is the gradient w.r.t. θ . In this
 135 way, the computations made in $\nabla_\theta \hat{Q}$ are accumulated via θ_t and reused over the iterations.

136 The parameters η_z have a different form based on a specific case of $p_{\eta_z}(x|z)$, whereas π_z is a
 137 permanent structure in (1). Therefore, without loss of generality, we split (6) into a generic part and a
 138 fixed part as follows:

$$\eta_{k,t} = \eta_{k,t-1} + \gamma_t \nabla_{\eta_k} \hat{Q}_{t-1}(\theta), \quad (9a)$$

$$\nu_{k,t} = \nu_{k,t-1} + \gamma_t \nabla_{\nu_k} \hat{Q}_{t-1}(\theta), \quad (9b)$$

139 where—to ensure that the probabilities, $(\pi_{k,t})_{k=1}^K$, satisfy the constraints (Section 2)—we transform
 140 $\nabla_{\pi_k} \hat{Q}$ via $\nu_k = \log \pi_k$ and optimize w.r.t. ν_k . Then, to obtain $(\pi_{k,t})_{k=1}^K$ from $\nu_t := (\nu_{k,t})_{k=1}^K$, we
 141 use the softmax function, i.e. $\pi_{k,t} := \text{softmax}(\nu_t)_k := \exp(\nu_{k,t}) / \sum_{l=1}^K \exp(\nu_{l,t})$.

142 Computing the gradients for all pairs of $(\nu_k, \eta_k)_{k=1}^K$ would be inefficient, especially since $\mathbf{z}_{i,t}$ contains
 143 only a small number of unique values of Z for $M \ll K$. Consequently, we compute $\nabla_{\eta_k} \hat{Q}$ and
 144 $\nabla_{\nu_k} \hat{Q}$ only for $k \in \text{unique}(\mathbf{z}_{i,t})$. We summarize the proposed approach in Algorithm 1.

145 4.3 Proposal distribution

146 The choice of the proposal distribution has a significant impact on the speed of convergence and the
 147 computational cost of the proposed algorithm. Here, we discuss various possible choices of $q(z_i|\bar{z}_i)$.

Algorithm 1 The generalized MHTSAEM algorithm

Input: $\theta_0, (\mathbf{z}_{i,0})_{i=1}^N, (\mathbf{x}_i)_{i=1}^N$ **Output:** $(\theta_t)_{t=1}^T$ **for** $t \in (1, \dots, T)$ or until convergence **do** form the set $I = (i_j)_{j=1}^B$ by sampling (without replacement) B indices $i \sim (1, \dots, N)$ **for** $i \in I$ **do** set \bar{z}_i as the last element of $\mathbf{z}_{i,t-1}$ **for** $j \in (1, \dots, M)$ **do** sample $z_i \sim q(z_i | \bar{z}_i)$ sample $u \sim \text{Uniform}(0, 1)$ compute $\alpha(\bar{z}_i, z_i)$ in (7) **if** $u < \alpha(\bar{z}_i, z_i)$ **then** set $z_{i,t}^j = z_i$ and $\bar{z}_i = z_i$ **else** set $z_{i,t}^j = \bar{z}_i$ **end if** **end for** set $\mathbf{z}_{i,t} = (z_{i,t}^1, \dots, z_{i,t}^M)$ **end for**

compute (8)

 compute (9) for $k \in \text{unique}(\mathbf{z}_{i,t})$ compute $\pi_{k,t} := \text{softmax}(\boldsymbol{\nu}_t)_k$ for $k \in Z$ **end for**

148 *Optimal proposal (O).* The optimal proposal distribution is $q(z_i | \bar{z}_i) := q(z_i) := p_\theta(z_i | x_i)$. This
149 ensures that the acceptance rate (7) is always $\alpha(\bar{z}_i, z_i) = 1$. However, the need to perform K
150 expensive evaluations of $p_\theta(z_i, x_i)$ before sampling from $p_\theta(z_i | x_i)$ is the reason we resorted to
151 the MH sampler in the first place. We consider this case only to set the upper limit on admissible
152 computational cost and to study the impact of sub-optimal proposal distributions.

153 *Uniform proposal (U).* The uniform distribution on the discrete interval from 1 to K , i.e. $q(z_i | \bar{z}_i) :=$
154 $q(z_i) := \text{Uniform}(1, K)$, is the simplest and computationally cheapest variant of the proposal
155 distribution. However, due to poor mixing properties, the algorithm may converge slowly for high K .

156 *Tabular proposal with forgetting (TF).* The key requirement to design a proposal distribution is to
157 restrict its computational complexity somewhere between that of the U and O proposals. One way to
158 satisfy this constraint is to use the Markov chain, $(z_{i,t})_{t=1}^T$, to learn a transition kernel, $p(z_i | \bar{z}_i)$, see,
159 e.g. [3]. Unfortunately, this would require us to store a table with K^2 entries for each $i \in (1, \dots, N)$,
160 which is very demanding even for moderate K and N . Therefore, we break the dependence in
161 the Markov chain and define: $q(z_i | \bar{z}_i) := q_{\boldsymbol{\alpha}_i}(z_i) := \mathcal{C}(\boldsymbol{\alpha}_i)$, where $\mathcal{C}(\boldsymbol{\alpha}_i) \propto \prod_{k=1}^K \alpha_{k,i}^{\mathbf{1}(z_i=k)}$ is the
162 categorical distribution with the weights $\boldsymbol{\alpha}_i := (\alpha_{1,i}, \dots, \alpha_{K,i})$. For $\mathcal{L}(\boldsymbol{\alpha}_i) := \sum_{\tau=1}^t \log q_{\boldsymbol{\alpha}_i}(z_{i,\tau})$,
163 we obtain an estimate of $\boldsymbol{\alpha}_i$ at iteration t as follows: $\boldsymbol{\alpha}_{i,t} := \arg \max_{\boldsymbol{\alpha}_i} \mathcal{L}(\boldsymbol{\alpha}_i) = \frac{\mathbf{n}_{i,t}}{t}$, with
164 $\mathbf{n}_{i,t} = \sum_{\tau=1}^t \mathbf{e}_{z_{i,\tau}}$, where \mathbf{e}_k is the standard basis vector (a one-hot vector) with one at k th position
165 and zeros otherwise. This can be further rewritten into a recursive form: $\mathbf{n}_{i,t} = \mathbf{n}_{i,t-1} + \mathbf{e}_{z_{i,t}}$ or,
166 using the Robbins-Monro step-size, $\mathbf{n}_{i,t} = (\mathbf{1} - \mathbf{e}_{z_{i,t}} \gamma_t) \odot \mathbf{n}_{i,t-1} + \gamma_t \mathbf{e}_{z_{i,t}}$, where $\mathbf{1}$ is the vector of
167 ones, and \odot is the Hadamard product. We refer to this case simply as “table with forgetting” (TF)
168 due to that it represents $N \times K$ table in the memory and γ_t is a forgetting factor.

169 5 Related work

170 *Stochastic approximation expectation-maximization.* The application of SA to prevent the evaluation
171 of all K components in mixture models has been overlooked for a long time. The reason is that the
172 original motivation to combine the EM algorithm with SA is to address the analytical intractability
173 of the expected value under $p_\theta(z|x)$ in (4), which is, however, almost always tractable for mixture
174 models. The intractability issue is addressed by either the Monte Carlo SAEM (MCSAEM) [12]
175 or the Markov chain Monte Carlo SAEM (MCMCSAEM) [31]. Applying the former approach to
176 mixture models would be inefficient, since it evaluates K joint distributions, $p_\theta(z, x)$, before drawing

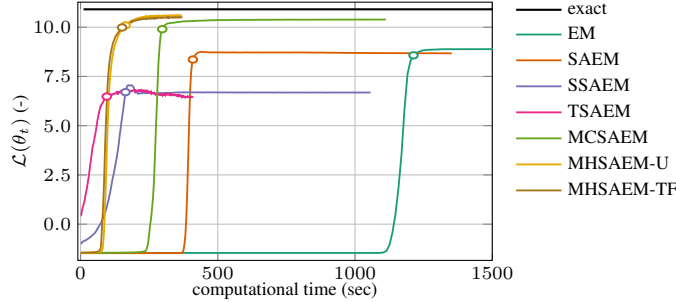


Figure 1: The training log-likelihood, $\mathcal{L}(\theta_t)$, versus the computational time (in seconds). Here, on the x-axis, the computational time at a current iteration, t , is obtained by accumulating the time from the previous iterations. \circ corresponds to $\mathcal{L}(\theta_{t_{95}})$, where t_{95} is the iteration of reaching 95% of $\max \mathcal{L}(\theta_t)$. The projection of \circ on the x-axis gives the time to reach $\mathcal{L}(\theta_{t_{95}})$. This experiment was performed with the following settings: $(D, K, N, \omega, B, M, T) = (10, 100, 10k, 0.1, 200, 2, 20k)$, see Section 6.1 for details. The results are averaged over five repetitions.

177 M samples from $p_\theta(z|x)$. Therefore, this method reduces only the computational cost of updating the
 178 *sufficient statistics*. This is addressed by the latter approach, where $M < K$ samples from a proposal
 179 distribution, $q(z|x)$, is used to calculate $p_\theta(z, x)$ and also the sufficient statistics. However, all these
 180 methods process all data at every iteration, providing only a limited advantage over the conventional
 181 EM algorithm. Minibatch versions of these techniques have recently been proposed [27, 32, 1].

182 All the above methods commonly assume $p_\theta(z, x)$ belonging to the exponential family. This provides
 183 a convenient, but limiting, property which allows (6) to be computed under a closed-form solution.
 184 The main contribution of our work is to release this restrictive assumption by admitting that $p_\theta(z, x)$
 185 (and thus \mathcal{Q}) is given by possibly complex and intractable transformations.

186 *Sparse and truncated variational techniques.* There is only a small body of methods explicitly
 187 reducing the number of evaluated components. Their common aspect is that they follow from the
 188 variational framework, where the exact posterior, $p_\theta(z|x)$, is approximated by a variational posterior,
 189 $q(z|x)$. This sparse, approximate, posterior is defined over a lower number of components, $M \ll K$,
 190 such that only the important components are selected, relying on relaxation of the hard EM algorithm
 191 from taking a single $M = 1$ assignment [26] to taking multiple $M \ll K$ assignments. The sparse
 192 SAEM (SSAEM) algorithm [24] selects the components by a quick partial sorting of the posterior
 193 probabilities, $p_\theta(z|x)$. Again, this requires K evaluations of $p_\theta(z, x)$ before the sorting, thus only
 194 reducing the amount of updated statistics. Similarly, the truncated SAEM (TSAEM) algorithm [18]
 195 selects $M < K$ cluster-to-cluster and $\bar{M} < K$ cluster-to-datapoint minimal Euclidean distances,
 196 preventing the problem in the SSAEM algorithm. However, all these distances are evaluated for all
 197 components in a pairwise manner, leading to K^2 -computational complexity, which makes the saving
 198 dubious. Similarly as before, these methods assume $p_\theta(z, x)$ to belong to the exponential family.

199 We summarize the distinguishing features of the above discussed methods in Table 1.

200 6 Experiments

201 To demonstrate the key features of our algorithm—its low computational complexity, competitive
 202 learning performance, and generality—we use it below to train: (i) GMMs on synthetic datasets, and
 203 (ii) SPTNs [47] and (iii) mixtures of real NVP flows [48] on real datasets. All experiments have been
 204 performed on a Slurm cluster equipped with Intel Xeon Scalable Gold 6146 with 384GB of RAM.

205 6.1 Gaussian mixture models

206 Consider the special case of a data-generating distribution given by (1), with the components taking
 207 the form of the multivariate Gaussian distribution, $p_{\eta_z}(x|z) = \mathcal{N}(x; \mu_z, \Sigma_z)$, where μ_z is the mean
 208 value and Σ_z is the covariance matrix. The difficulty of learning GMMs heavily depends on the
 209 degree of interaction among all mixture components, hence having the ability to generate synthetic
 210 datasets with arbitrary overlap characteristics between all pairs of components is crucial for systematic

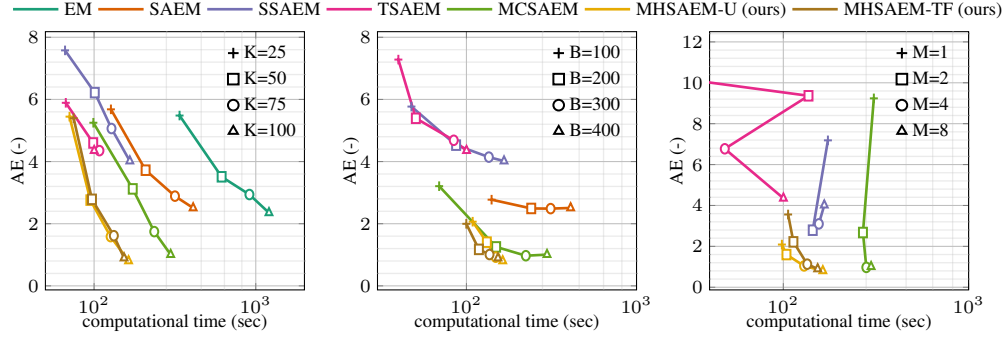


Figure 2: The absolute error, $AE = |\mathcal{L}(\theta_{t_{95}}) - \mathcal{L}(\theta)|$, versus the computational time (in seconds). All experiments use the following settings: $(D, K, N, \omega, B, M, T) = (10, 100, 10k, 0.1, 200, 2, 20k)$, where the number of components, K , (left), the batchsize, B , (middle) and the number of samples, M , (right) change for different values denoted by (+, \square , \circ , \triangle). At each of these points (marks), we perform an experiment as illustrated in Figure 1, find $\mathcal{L}(\theta_{t_{95}})$ to compute the AE, and record the time corresponding to t_{95} . The results are averaged over five repetitions.

211 evaluation of performance of learning algorithms [43]. Traditional techniques usually define overlap
 212 (or separation) of components only in terms of their mean vectors and maximum eigenvalues of the
 213 covariance matrices, not accounting for their rotation and mixing weights (see [36] for a detailed
 214 treatment of the problem). We therefore use a more objective measure of the clustering complexity
 215 defined by the total probability of misclassification [41], which allows to generate data with a
 216 user-defined degree of maximum pairwise overlap, ω .

217 *Experiment settings:* We generate the parameters of (1), and the corresponding dataset, uniquely for a
 218 given quadruple (D, K, N, ω) . Therefore, the parameters of the generative model are known and we
 219 can measure and display the convergence of the training log-likelihood, $\mathcal{L}(\theta_t)$, compared to the exact
 220 log-likelihood, $\mathcal{L}(\theta)$, for $t = (1, \dots, T)$. We are further interested in the absolute error between the
 221 training log-likelihood at the iteration of reaching 95% of its maximum value, t_{95} , and the exact
 222 log-likelihood, i.e. $AE = |\mathcal{L}(\theta_{t_{95}}) - \mathcal{L}(\theta)|$.

223 We also measure the computational time until reaching t_{95} . We have used 95% of the maximum
 224 value instead of the maximum value to prevent cases, where the model oscillate around target value,
 225 making the estimate of convergence time very noisy (for example MCSAEM in Figure 1).

226 *Algorithms:* The GMMs belong to the exponential family of probability distributions. This allows us
 227 to find a closed-form, recursive, solution of (6), relying on a Robbins-Monro type of the step-size
 228 sequence, $(\gamma_t)_{t=1}^T$, [7, 44]. In this setting, we compare our MHAEM algorithm with a number of
 229 related methods in Table 1. Note we use the acronyms U and TF to specify the proposal distribution of
 230 the MHAEM algorithm (Section 4.3). However, we do not use the O-proposal, since the MHAEM-
 231 O algorithm is equivalent to the MCSAEM algorithm. All the SA-variants in Table 1 use a minibatch
 232 of size B . The key quantity to reduce the number of evaluated components and/or sufficient statistics
 233 in the SSAEM, TSAEM, MCSAEM and MHAEM algorithms is collectively denoted by M (Section
 234 5). Note that we always keep $M = \bar{M}$ in the TSAEM algorithm (see Figure 1 and 2 for concrete
 235 numbers). We use the step-size given by $\gamma_t = 1$ for $t = 1, \dots, 50$ and $\gamma_t = 0.05$ otherwise. In
 236 this section, to counteract the issue of attaining poor local optima, we equip *all* algorithms with the
 237 anti-annealing schedule $(\beta_t)_{t=1}^T$, starting with $\beta_1 = 0.1$, reaching $\beta_{2/3T} = 1.2$, and decreasing back
 238 to $\beta_T = 1.0$, see [43] for details. The initial estimates of: (i) μ_k are uniformly drawn from the unit
 239 hyper-cube, (ii) Σ_k are fixed to unit diagonal matrix, and (iii) π_k are uniformly drawn from the unit
 240 interval (followed by normalization).

241 *Results:* Figure 1 shows that the EM [13] and SAEM [44] algorithms take the longest time to
 242 converge, attaining a poor local optima. On the other hand, the MCSAEM [1] and MHAEM (U
 243 and TF) algorithms achieve $\mathcal{L}(\theta_{t_{95}})$ closest to the likelihood $\mathcal{L}(\theta)$ of the true model. Moreover, both
 244 MHAEM algorithms reach this value in the shortest time compared to all the other methods. The
 245 SSAEM [24] and TSAEM [18] algorithms are comparable in terms of the computational time, but
 246 they both provide the lowest $\mathcal{L}(\theta_{t_{95}})$. In Figure 2, we investigate sensitivity of fitting the model to

247 increasing values of K , B and M by measuring the time and the likelihood again. In all the cases,
248 the proposed MHSAEM algorithms achieve the lowest AE in the shortest time.

249 SSAEM and TSAEM algorithms failed to converge for $M > 2$ and for $K > 50$ respectively. We
250 believe this is caused by selecting only M maximal probabilities in the SSAEM (or distances in the
251 TSAEM) algorithm (Section 5), which prevents certain, but not a negligible number of, components
252 from being updated, thus providing only a crude approximation of $p_\theta(z|x)$. The results then suffer
253 from substantial variational gap to the exact log-likelihood (Figure 1). On the contrary, MH sampler
254 provides samples which consistently approximate $p_\theta(z|x)$ despite evaluating much lower number of
255 components in each step.

256 6.2 Sum-product transform networks

257 The sum product networks (SPNs) are a deep learning extension of finite mixture models. They can
258 be interpreted as a mixture of trees [60], where each tree corresponds to a component. Therefore,
259 they can be cast into the form of (1), but the number of components grows exponentially with their
260 depth. In this section, we use recently proposed SPTNs which introduce additional transformation
261 nodes to provide better expressiveness than the SPNs (SPTNs effectively generalize SPNs and flow
262 models into one large family of models).

263 *Experimental settings:* We use 19 real datasets from the UCI database [16, 37, 35, 54], preprocessed
264 in the same way as in [46]. For each experiment, we randomly split the data into 64%, 16% and 20%
265 for training, validation and testing, respectively. We calculate the average log-likelihood on the test
266 set and measure again the time to reach 95% of the maximal training log-likelihood, $\mathcal{L}(\theta_{t_{95}})$.

267 To evaluate various (possibly shallow and/or deep) architectures of SPTNs, we fit each dataset with
268 all the following combinations of hyper-parameters¹: $s \in (8, 32, 128)$, $b \in (2, 4, 6, 8)$, $l \in (2, 3, 4)$,
269 where s is the number of children of each sum node, b is the number of partitions of each product
270 node, and l is the number of layers (one layer contains sum and product nodes). The number of
271 components of the SPTN, after its conversion into (1), is given as follows: $K = s^l$. Note that the
272 maximum number of components for the investigated parameters of the SPTN is 268,435,456. To
273 reduce the space of possible architectures, we restrict ourselves only to (i) the leaf nodes given by
274 $\mathcal{N}(0, \mathbf{I})$; (ii) affine transformations fixed to the singular value decomposition, choosing the the Givens
275 parameterization for the unitary matrices [47]; and (iii) no sharing of any type of nodes [47].

276 *Algorithms:* We evaluate only on the MHSAEM-U algorithm—due to its favourable computational
277 complexity and simplicity—and compare it with the stochastic gradient-descent (SGD) algorithm,
278 which is routinely used to train SP(T)Ns [45, 47]. In this case, SGD in each iteration performs
279 computations over all subtrees of the network, whereas the MHSAEM-U algorithm computes with
280 only $M = 1$ subtrees, thus we should observe speed-up of the computations. In our implementation,
281 both these methods perform optimization of their respective objective functions—the log-likelihood
282 (2) for SGD and the EM objective (8) for MHSAEM-U—via the use of the automatic differentiation
283 and the ADAM optimizer [29], using $B = 100$ and $T = 20000$.

284 *Results:* Since each dataset might benefit from a different architecture, Table 6.2 shows the test
285 log-likelihood of the architectures selected according to the best likelihood measured on the validation
286 set and the corresponding speed-up. The test log-likelihoods reveal that the MHSAEM-U algorithm
287 outperforms the SGD algorithm on 10 out of 19 datasets, which was not originally the goal, but
288 the added stochasticity helps to escape poor local minima. The speed-up demonstrates lower
289 computational complexity of the MHSAEM-U algorithm on 17 out of 19 datasets, which was the
290 main goal. The magic-telescope and wine datasets show approximately $102\times$ and $75\times$ speed-up,
291 respectively, while on very small datasets (pima-indians and iris), the SGD is faster due to
292 effective implementation. In the supplementary material, we present Table 3, exhibiting the same
293 trends on a fixed architecture.

294 6.3 Mixtures of real NVP flows

295 We consider another class of mixture models (1), where each component $p_{\eta_z}(x|z)$ is transformed by
296 the flow model—real NVP [15]. These transformations are parameterized via deep neural networks,
297 allowing for flexible adjustment of the learning capacity of each component.

¹We have set a hard limit to train a single model to 24h, which is default on our Slurm cluster.

Table 2: The speed-up and test log-likelihood, $\mathcal{L}^{\text{test}}$, for the SGD and MHTSAEM-U algorithms. The test log-likelihood (higher is better) is computed for the best model, with the corresponding K , which is selected based on the validation log-likelihood. The speed-up is computed as the ratio of MHTSAEM-U to SGD, i.e. their time to reach 95% of the training log-likelihood. The results are averaged over five repetitions. Then, the higher test log-likelihood is highlighted with bold blue, and no speed-up is highlighted with red. The average rank is computed as the standard competition (“1224”) ranking [14] on each dataset (lower is better).

dataset	Sum-product transform networks					Mixtures of real NVP flows				
	SGD		MHTSAEM-U			SGD		MHTSAEM-U		
	speed-up	$\mathcal{L}^{\text{test}}$	K	$\mathcal{L}^{\text{test}}$	K	speed-up	$\mathcal{L}^{\text{test}}$	K	$\mathcal{L}^{\text{test}}$	K
breast-cancer-wisconsin	4.66	-4.66	64	1.43	1024	0.63	-99.85	32	-39.31	128
cardiotocography	10.55	59.52	512	31.04	1024	9.85	54.34	32	56.08	128
magic-telescope	102.53	-3.65	512	-5.03	1024	3.74	-3.97	8	-4.22	8
pendigits	4.89	0.88	1024	-4.86	16384	4.17	1.46	8	0.48	8
pima-indians	0.37	-8.54	64	-7.62	64	1.35	-20.09	128	-16.33	128
wall-following-robot	3.43	1.84	1024	-11.3	16384	22.21	-14.26	128	-17.56	128
waveform-1	4.35	-26.14	64	-23.91	1024	3.72	-34.12	8	-33.42	8
waveform-2	4.82	-26.21	64	-23.91	1024	4.12	-34.15	8	-33.64	8
yeast	20.57	10.26	512	5.18	1024	14.49	6.61	128	9.59	128
ecoli	1.86	-5.5	64	-0.22	1024	2.15	-11.37	128	-10.64	128
ionosphere	1.88	-20.27	64	-5.93	512	2.74	-87.01	128	-42.75	128
iris	0.23	-10.65	64	-1.49	16384	3.28	-16.34	128	-9.21	32
page-blocks	12.18	12.21	512	6.84	1024	44.95	17.13	128	17.94	32
parkinsons	1.46	-21.85	64	0.5	512	3.09	-566.58	128	-33.31	32
sonar	2.96	-95.39	512	-69.29	64	2.52	-622.2	128	-88.81	128
statlog-segment	1.44	47.35	512	26.53	16384	38.49	35.84	128	42.04	32
statlog-vehicle	2.97	-4.25	64	-5.45	1024	6.78	-31.34	32	-26.43	128
wine	75.42	-25.99	1024	-13.27	1024	2.05	-171.58	128	-25.57	128
rank		1.56		1.44			1.83		1.17	

298 *Experimental settings:* We use the same experimental settings and evaluation metrics as in Section
299 6.2. We apply the mixture model on all datasets, changing the number of components as follows:
300 $K \in (8, 32, 128)$. Each real NVP-based component in the mixture model has (i) the translation
301 function parameterized via multi-layer perceptron with a single hidden layer of dimension 10, using
302 the rectified linear activation function; and (ii) the scale function parameterized via the same network
303 except with the hyperbolic tangent activation function. We do not use the batch normalization [15] and
304 we stack two layers of the translation-scale transformation (we have used implementation from [20]).

305 *Algorithms:* The algorithms and their settings are the same as those in Section 6.2.

306 *Results:* The experimental results are presented in right part of Table 6.2. They are similar to those
307 obtained in the previous section. In terms of the test log-likelihood, the MHTSAEM-U algorithm
308 outperforms the SGD algorithm on all but three datasets, and it provides a substantial speed-up on all
309 datasets except one. The test likelihood of models with the real NVP flows is most of the time worse
310 than that of SPTNs with the affine transformations. As explained in the supplementary, this is due to
311 the overfitting, which has been observed in [47].

312 7 Conclusion

313 This paper has presented a method to decrease computational complexity of fitting mixture models,
314 including their generalizations, such as sum-product-(transform) networks and mixtures of flow
315 models. The speed-up is achieved by evaluating and updating only a single component (per iteration),
316 where the Metropolis-Hasting algorithm ensures sampling of components from a proper posterior. An
317 experimental comparison on all three classes of models mentioned above confirmed the theoretical
318 expectations. The method significantly speeds-up the fitting time and, importantly, without sacrificing
319 the quality of the fit. In fact, the likelihood was better than that of the models fitted by the EM
320 algorithm or the SGD algorithm in more than 50% of cases. We attribute this to higher stochasticity,
321 which helps to escape from poor local minima.

322 In the experiments, the proposed method has used a uniform proposal distribution in the MH sampler.
323 Despite outperforming the alternative methods, we conjecture that this limits the speed of convergence.
324 Therefore, we believe that there is still a room for improvement in the implementation. We plan to
325 address these issues in future work.

326 **8 Broader impact statement**

327 The presented method decreases the computational complexity of fitting large (and deep) mixture
328 models, which leads to five to hundred time speed-up depending on a size of the problem (although
329 negative exceptions occurs). We believe this line of research, which we want to continue, to have
330 important benefits. First, it is directly related to decrease in energy consumption and in production of
331 CO₂ (we expect similar rates as the speedup). Second, it has a positive effect on financial aspects of
332 deploying (and experimenting with) mixture models. Third, it decreases the hardware requirements,
333 as in all experiments presented above the model was fitted on a single-core.

334 **References**

- 335 [1] S. Allasonnière and J. Chevallier. A new class of stochastic EM algorithms: Escaping
336 local maxima and handling intractable sampling. *Computational Statistics & Data Analysis*,
337 159:107159, 2021.
- 338 [2] S. Ayesha, M. K. Hanif, and R. Talib. Overview and comparative study of dimensionality
339 reduction techniques for high dimensional data. *Information Fusion*, 59:44–58, 2020.
- 340 [3] D. S. Bai. Efficient estimation of transition probabilities in a Markov chain. *The Annals of*
341 *Statistics*, pages 1305–1317, 1975.
- 342 [4] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians.
343 *Journal of the American statistical Association*, 112(518):859–877, 2017.
- 344 [5] J. Blömer and K. Bujna. Adaptive seeding for Gaussian mixture models. In *Pacific-asia*
345 *conference on knowledge discovery and data mining*, pages 296–308. Springer, 2016.
- 346 [6] T. Campbell and B. Beronov. Sparse variational inference: Bayesian coresets from scratch.
347 *arXiv preprint arXiv:1906.03329*, 2019.
- 348 [7] O. Cappé and E. Moulines. On-line expectation–maximization algorithm for latent data models.
349 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(3):593–613,
350 2009.
- 351 [8] J. Chen, J. Zhu, Y. Teh, and T. Zhang. Stochastic expectation maximization with variance
352 reduction. *Advances in Neural Information Processing Systems*, page 7967, 2018.
- 353 [9] N. Chopin. A sequential particle filter method for static models. *Biometrika*, 89(3):539–552,
354 2002.
- 355 [10] G. Compiani and Y. Kitamura. Using mixtures in econometric models: A brief review and some
356 new results. *The Econometrics Journal*, 19(3):C95–C127, 2016.
- 357 [11] S. Dasgupta. Learning mixtures of Gaussians. In *40th Annual Symposium on Foundations of*
358 *Computer Science (Cat. No. 99CB37039)*, pages 634–644. IEEE, 1999.
- 359 [12] B. Delyon, M. Lavielle, E. Moulines, et al. Convergence of a stochastic approximation version
360 of the EM algorithm. *The Annals of Statistics*, 27(1):94–128, 1999.
- 361 [13] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data
362 via the EM algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*,
363 39(1):1–22, 1977.
- 364 [14] J. Demšar. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine*
365 *Learning Research*, 7:1–30, 2006.
- 366 [15] L. Dinh, J. Sohl-Dickstein, and S. Bengio. Density estimation using real NVP. In *5th Interna-*
367 *tional Conference on Learning Representations, ICLR 2017*, 2017.
- 368 [16] D. Dua and C. Graff. UCI machine learning repository, 2017.

- 369 [17] D. Feldman, M. Faulkner, and A. Krause. Scalable training of mixture models via coresets. In
370 *Proceedings of the 24th International Conference on Neural Information Processing Systems*,
371 pages 2142–2150, 2011.
- 372 [18] D. Forster and J. Lücke. Can clustering scale sublinearly with its clusters? A variational EM
373 acceleration of GMMs and k-means. In *International Conference on Artificial Intelligence and*
374 *Statistics*, pages 124–132. PMLR, 2018.
- 375 [19] G. Fort, E. Moulines, et al. Convergence of the Monte Carlo expectation maximization for
376 curved exponential families. *Annals of Statistics*, 31(4):1220–1259, 2003.
- 377 [20] J. Franců. Continuousflows.jl. <https://github.com/janfrancu/ContinuousFlows.jl>,
378 2020.
- 379 [21] K. Greff, S. van Steenkiste, and J. Schmidhuber. Neural expectation maximization. In *Proceed-*
380 *ings of the 31st International Conference on Neural Information Processing Systems*, pages
381 6694–6704, 2017.
- 382 [22] R. Gribonval, A. Chatalic, N. Keriven, V. Schellekens, L. Jacques, and P. Schniter. Sketching
383 datasets for large-scale learning (long version). *arXiv preprint arXiv:2008.01839*, 2020.
- 384 [23] M. D. Hoffman, D. M. Blei, C. Wang, and J. Paisley. Stochastic variational inference. *Journal*
385 *of Machine Learning Research*, 14(5), 2013.
- 386 [24] M. C. Hughes and E. B. Sudderth. Fast learning of clusters and topics via sparse posteriors.
387 *arXiv preprint arXiv:1609.07521*, 2016.
- 388 [25] K. Humphreys and D. Titterton. Approximate Bayesian inference for simple mixtures. In
389 *COMPSTAT*, pages 331–336. Springer, 2000.
- 390 [26] B.-H. Juang and L. R. Rabiner. The segmental K-means algorithm for estimating parameters
391 of hidden Markov models. *IEEE Transactions on acoustics, speech, and signal Processing*,
392 38(9):1639–1641, 1990.
- 393 [27] B. Karimi, M. Lavielle, and É. Moulines. On the convergence properties of the mini-batch EM
394 and MCEM algorithms, 2019.
- 395 [28] N. Keriven, A. Bourrier, R. Gribonval, and P. Pérez. Sketching for large-scale learning of
396 mixture models. *Information and Inference: A Journal of the IMA*, 7(3):447–508, 2018.
- 397 [29] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint*
398 *arXiv:1412.6980*, 2014.
- 399 [30] A. Kucukelbir, D. Tran, R. Ranganath, A. Gelman, and D. M. Blei. Automatic differentiation
400 variational inference. *The Journal of Machine Learning Research*, 18(1):430–474, 2017.
- 401 [31] E. Kuhn and M. Lavielle. Coupling a stochastic approximation version of EM with an MCMC
402 procedure. *ESAIM: Probability and Statistics*, 8:115–131, 2004.
- 403 [32] E. Kuhn, C. Matias, and T. Rebafka. Properties of the stochastic approximation EM algorithm
404 with mini-batch sampling. *Statistics and Computing*, 30(6):1725–1739, 2020.
- 405 [33] A. Lagrange, M. Fauvel, and M. Grizonnet. Large-scale feature selection with Gaussian mixture
406 models for the classification of high dimensional remote sensing images. *IEEE Transactions on*
407 *Computational Imaging*, 3(2):230–242, 2017.
- 408 [34] M. Lavine and M. West. A Bayesian method for classification and discrimination. *Canadian*
409 *Journal of Statistics*, 20(4):451–461, 1992.
- 410 [35] M. Little, P. McSharry, S. Roberts, D. Costello, and I. Moroz. Exploiting nonlinear recurrence
411 and fractal scaling properties for voice disorder detection. *Nature Precedings*, pages 1–1, 2007.
- 412 [36] R. Maitra and V. Melnykov. Simulating data to study performance of finite mixture modeling
413 and clustering algorithms. *Journal of Computational and Graphical Statistics*, 19(2):354–376,
414 2010.

- 415 [37] O. L. Mangasarian and W. H. Wolberg. Cancer diagnosis via linear programming. Technical
416 report, University of Wisconsin-Madison Department of Computer Sciences, 1990.
- 417 [38] J.-M. Marin, K. Mengersen, and C. P. Robert. Bayesian modelling and inference on mixtures of
418 distributions. *Handbook of statistics*, 25:459–507, 2005.
- 419 [39] C. A. McGrory, D. C. Ahfock, J. A. Horsley, and C. L. Alston. Weighted Gibbs sampling for
420 mixture modelling of massive datasets via coresets. *Stat*, 3(1):291–299, 2014.
- 421 [40] G. J. McLachlan, S. X. Lee, and S. I. Rathnayake. Finite mixture models. *Annual review of*
422 *statistics and its application*, 6:355–378, 2019.
- 423 [41] V. Melnykov, W.-C. Chen, and R. Maitra. MixSim: An R package for simulating data to study
424 performance of clustering algorithms. *Journal of Statistical Software*, 51(12):1, 2012.
- 425 [42] T. Monnier, T. Groueix, and M. Aubry. Deep transformation-invariant clustering. In *Conference*
426 *on Neural Information Processing Systems (NeurIPS 2020)*, 2020.
- 427 [43] I. Naim and D. Gildea. Convergence of the EM algorithm for Gaussian mixtures with unbalanced
428 mixing coefficients. In *Proceedings of the 29th International Conference on International*
429 *Conference on Machine Learning*, pages 1427–1431, 2012.
- 430 [44] H. D. Nguyen, F. Forbes, and G. J. McLachlan. Mini-batch learning of exponential family finite
431 mixture models. *Statistics and Computing*, pages 1–18, 2020.
- 432 [45] R. Peharz, A. Vergari, K. Stelzner, A. Molina, X. Shao, M. Trapp, K. Kersting, and Z. Ghahra-
433 mani. Random sum-product networks: A simple and effective approach to probabilistic deep
434 learning. In *Uncertainty in Artificial Intelligence*, pages 334–344. PMLR, 2020.
- 435 [46] T. Pevný. Loda: Lightweight on-line detector of anomalies. *Machine Learning*, 102(2):275–304,
436 2016.
- 437 [47] T. Pevný, V. Šmídl, M. Trapp, O. Poláček, and T. Oberhuber. Sum-product-transform networks:
438 Exploiting symmetries using invertible transformations. *arXiv preprint arXiv:2005.01297*,
439 2020.
- 440 [48] G. G. Pires and M. A. Figueiredo. Variational mixture of normalizing flows. *arXiv preprint*
441 *arXiv:2009.00585*, 2020.
- 442 [49] A. Rau, C. Maugis-Rabusseau, M.-L. Martin-Magniette, and G. Celeux. Co-expression analysis
443 of high-throughput transcriptome sequencing data with Poisson mixture models. *Bioinformatics*,
444 31(9):1420–1427, 2015.
- 445 [50] R. A. Redner and H. F. Walker. Mixture densities, maximum likelihood and the EM algorithm.
446 *SIAM review*, 26(2):195–239, 1984.
- 447 [51] H. Robbins and S. Monro. A stochastic approximation method. *The annals of mathematical*
448 *statistics*, pages 400–407, 1951.
- 449 [52] C. Robert and G. Casella. *Monte Carlo statistical methods*. Springer Science & Business Media,
450 2013.
- 451 [53] W. Sibli, P. Kuntz, and F. Meyer. A review on dimensionality reduction for multi-label
452 classification. *IEEE Transactions on Knowledge and Data Engineering*, 2019.
- 453 [54] J. P. Siebert. Vehicle recognition using rule based methods. 1987.
- 454 [55] C. Viroli and G. J. McLachlan. Deep gaussian mixture models. *Statistics and Computing*,
455 29(1):43–51, 2019.
- 456 [56] C. F. J. Wu. On the convergence properties of the EM algorithm. *The Annals of statistics*, pages
457 95–103, 1983.
- 458 [57] D. Wu and J. Ma. An effective EM algorithm for mixtures of Gaussian processes via the MCMC
459 sampling and approximation. *Neurocomputing*, 331:366–374, 2019.

- 460 [58] W. Xiang, A. Karfoul, C. Yang, H. Shu, and R. L. B. Jeannès. An exact line search scheme to
461 accelerate the EM algorithm: Application to Gaussian mixture models identification. *Journal of*
462 *computational science*, 41:101073, 2020.
- 463 [59] M. Zhang, Y. Fu, K. M. Bennett, and T. Wu. Computational efficient variational Bayesian Gaus-
464 sian mixture models via coresnet. In *2016 International Conference on Computer, Information*
465 *and Telecommunication Systems (CITS)*, pages 1–5. IEEE, 2016.
- 466 [60] H. Zhao, P. Poupart, and G. Gordon. A unified approach for learning the parameters of sum-
467 product networks. In *Proceedings of the 30th International Conference on Neural Information*
468 *Processing Systems*, pages 433–441, 2016.

469 **Checklist**

- 470 1. For all authors...
- 471 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
472 contributions and scope? [Yes]
- 473 (b) Did you describe the limitations of your work? [Yes] Our main contribution is compu-
474 tational speedup. Cases where it was not achieved are highlighted in the experimental
475 section.
- 476 (c) Did you discuss any potential negative societal impacts of your work? [No] We do not
477 foresee any potential negative impact.
- 478 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
479 them? [Yes]
- 480 2. If you are including theoretical results...
- 481 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 482 (b) Did you include complete proofs of all theoretical results? [N/A]
- 483 3. If you ran experiments...
- 484 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
485 imental results (either in the supplemental material or as a URL)? [Yes] The code is
486 available in a github repository. All dataset are public from the UCI database.
- 487 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
488 were chosen)? [Yes] See Section 6.
- 489 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
490 ments multiple times)? [No] We report only average of Monte Carlo repetitions, the
491 error bars were too small to have any visual impact in the reported logarithmic scale.
- 492 (d) Did you include the total amount of compute and the type of resources used (e.g., type
493 of GPUs, internal cluster, or cloud provider)? [Yes]
- 494 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 495 (a) If your work uses existing assets, did you cite the creators? [Yes] We use 20 datasets
496 from UCI, we cite the required papers for each dataset, mostly the UCI database and
497 few additional publications.
- 498 (b) Did you mention the license of the assets? [No] The data are publically available, we
499 comply with the requirement on citing appropriate publications.
- 500 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 501 (d) Did you discuss whether and how consent was obtained from people whose data you're
502 using/curating? [N/A]
- 503 (e) Did you discuss whether the data you are using/curating contains personally identifiable
504 information or offensive content? [N/A]
- 505 5. If you used crowdsourcing or conducted research with human subjects...
- 506 (a) Did you include the full text of instructions given to participants and screenshots, if
507 applicable? [N/A]
- 508 (b) Did you describe any potential participant risks, with links to Institutional Review
509 Board (IRB) approvals, if applicable? [N/A]
- 510 (c) Did you include the estimated hourly wage paid to participants and the total amount
511 spent on participant compensation? [N/A]