

Breaking the gridlock in Mixture-of-Experts: Consistent and Efficient Algorithms

Ashok Vardhan Makkuva¹ Sewoong Oh² Sreeram Kannan³ Pramod Viswanath¹

Abstract

Mixture-of-Experts (MoE) is a widely popular model for ensemble learning and is a basic building block of highly successful modern neural networks as well as a component in Gated Recurrent Units (GRU) and Attention networks. However, present algorithms for learning MoE, including the EM algorithm and gradient descent, are known to get stuck in local optima. From a theoretical viewpoint, finding an efficient and provably consistent algorithm to learn the parameters remains a long standing open problem for more than two decades. In this paper, we introduce the first algorithm that learns the true parameters of a MoE model for a wide class of non-linearities with global consistency guarantees. While existing algorithms jointly or iteratively estimate the expert parameters and the gating parameters in the MoE, we propose a novel algorithm that breaks the deadlock and can directly estimate the expert parameters by sensing its echo in a carefully designed cross-moment tensor between the inputs and the output. Once the experts are known, the recovery of gating parameters still requires an EM algorithm; however, we show that the EM algorithm for this simplified problem, unlike the joint EM algorithm, converges to the true parameters. We empirically validate our algorithm on both the synthetic and real data sets in a variety of settings, and show superior performance to standard baselines.

1. Introduction

In this paper, we study a popular gated neural network architecture known as Mixture-of-Experts (MoE). MoE is a basic building block of highly successful modern neural networks like Gated Recurrent Units (GRU) and Attention networks. A key interesting feature of MoE is the presence of a gating mechanism that allows for specialization of experts in their respective domains. MoE allows for the underlying expert models to be simple while allowing to capture complex non-linear relations between the data. Ever since their inception more than two decades ago (Jacobs et al., 1991), they have been a subject of great research interest (Tresp, 2001; Collobert et al., 2002; Ng & Deisenroth, 2014; Theis & Bethge, 2015; Le et al., 2016; Gross et al., 2017; Sun et al., 2017; Wang et al., 2018) across multiple domains such as computer vision, natural language processing, speech recognition, finance, and forecasting.

The basic MoE model is the following: let $\mathbf{x} \in \mathbb{R}^d$ be the input feature vector and $y \in \mathbb{R}$ be the corresponding label. Then the discriminative model $P_{y|\mathbf{x}}$ for the k -mixture of experts (k -MoE) in the regression setting is:

$$\begin{aligned} P_{y|\mathbf{x}} &= \sum_{i=1}^k P_{i|\mathbf{x}} P_{y|\mathbf{x},i} \\ &= \sum_{i=1}^k \frac{e^{(\mathbf{w}_i^*, \mathbf{x})}}{\sum_j e^{(\mathbf{w}_j^*, \mathbf{x})}} \mathcal{N}(y|g(\langle \mathbf{a}_i^*, \mathbf{x} \rangle), \sigma^2). \end{aligned} \quad (1)$$

Figure 1 details the architecture for k -MoE.

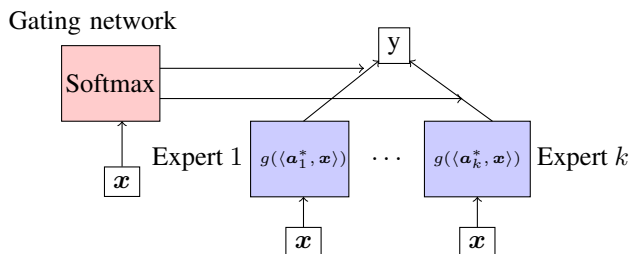


Figure 1: Architecture for k -MoE

¹Department of Electrical and Computer Engineering, Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, IL, USA ²Allen School of Computer Science & Engineering, University of Washington, Seattle, USA ³Department of Electrical Engineering, University of Washington, Seattle, USA. Correspondence to: Ashok Vardhan Makkuva <makkuva2@illinois.edu>.

The interpretation behind (1) is that for each input \mathbf{x} , the gating network chooses an expert based on the outcome

of a multinomial random variable $z \in [k]$, whose probability depends on \mathbf{x} in a parametric way, i.e. $z|\mathbf{x} \sim \text{softmax}(\langle \mathbf{w}_1^*, \mathbf{x} \rangle, \dots, \langle \mathbf{w}_k^*, \mathbf{x} \rangle)$. The chosen expert then generates the output y from a Gaussian distribution centred at a non-linear activation of \mathbf{x} , i.e. $g(\langle \mathbf{a}_z^*, \mathbf{x} \rangle)$, with variance σ^2 . We want to learn the expert parameters $\mathbf{a}_i^* \in \mathbb{R}^d$ (also referred to as the regressors) and the gating parameters $\mathbf{w}_i^* \in \mathbb{R}^d$, assuming we know the non-linear activation $g: \mathbb{R} \rightarrow \mathbb{R}$.

This problem of learning MoE has been a long standing open problem for more than two decades, even though it is a fundamental building block of several state-of-the-art gated neural network architectures. Gated neural networks such as GRUs and Sparsely-gated-MoEs have been widely successful in challenging tasks like machine translation (Chung et al., 2014; Shazeer et al., 2017; Vaswani et al., 2017). Parameters are typically learnt through (stochastic) gradient descent on a non-convex loss function. However, these methods do not possess any theoretical guarantees, even for the simplest gated neural network, which is the MoE.

On the other hand, existing guarantees for simpler models without gating units do not extend to MoEs. Consider the mixture of generalized linear models (M-GLMs) (Sedghi et al., 2014; Sun et al., 2014; Yi et al., 2016; Zhong et al., 2016), which is a strict simplification of the k -MoE model in (1), where $\mathbf{w}_i^* = 0$ for all $i \in \{1, \dots, k\}$. The learning in M-GLMs is usually done through a combination of spectral methods and greedy methods such as EM. A major limitation of these methods is that they rely critically on the fact that the mixing probability is a constant and hence they do not generalize to MoEs (see Section 2). In addition, the EM algorithm, which is the workhorse for learning in parametric mixture models, is prone to bad local minima (Sedghi et al., 2014; Balakrishnan et al., 2017; Zhong et al., 2016) (we independently verify this for MoEs in Section 4). These theoretical shortcomings and practical relevance of the MoE models lead to the following fundamental question:

Can we find an efficient and a consistent algorithm (with global initializations) that recovers the true parameters of the model with theoretical guarantees?

In this paper, we address this question precisely and make the following contributions:

1) First theoretical guarantees: We provide the first (poly-time) efficient algorithm that recovers the true parameters of a MoE model with global initializations (Theorem 1 and Theorem 2). We allow for a wide class of non-linearities which includes the popular choices of identity, sigmoid, and ReLU. To the best of our knowledge, ours is the first work to give global convergence guarantees for MoE.

2) Algorithmic innovations: Existing algorithms jointly or iteratively estimate the expert parameters and the gating

parameters in the MoE and can get stuck in local minima. In this paper, we propose a novel algorithm that breaks the gridlock and can directly estimate the expert parameters by sensing its echo in a cross-moment tensor between the inputs and the output (Algorithm 1 and Algorithm 2). Once the experts are known, the recovery of gating parameters still requires an EM algorithm; however, we show that the EM algorithm for this simplified problem, unlike the joint EM algorithm, converges to the true parameters. The proofs of global convergence of EM as well as the design of the cross-moment tensor are of independent mathematical interest.

3) Novel transformations: In this paper, we introduce the novel notion of ‘‘Cubic and Quadratic Transform (CQT)’’. These are polynomial transformations on the output labels tailored to specific non-linear activation functions and the noise variance. The key utility of these transforms is to equip MoEs with a supersymmetric tensor structure in a principled way (Theorem 1).

Related work. While there is a huge literature on MoEs ((Yuksel et al., 2012; Masoudnia & Ebrahimpour, 2014) are detailed surveys), there are relatively few works on its learning guarantees. (Jordan & Xu, 1995) is the first work to analyze the local convergence of joint-EM for both the gating and the expert parameters. As noted earlier, however, EM is prone to bad local minima. In contrast, our algorithms have *global convergence* guarantees. It is important to note that even for the simpler problem of mixtures of Gaussians, it is known that EM gets stuck in local minima, whenever number of mixtures, k , is at least 3 (Jin et al., 2016), whereas we can handle $2k - 1 < d$ with global convergence.

The simplified versions of MoE, M-GLMs, are widely studied in the literature. The key techniques for parameter inference in M-GLMs include EM algorithm, spectral methods, convex relaxations, and their variants. (Yi et al., 2014; Balakrishnan et al., 2017) prove convergence of EM for 2-mixtures of linear regressions; in contrast, we handle $k \geq 2$ mixtures for a wide class of non-linearities and provide global convergence. (Sedghi et al., 2014) construct a 3rd-order supersymmetric tensor containing the regressors as its rank-1 components. However, this approach fails to generalize for MoE. (Zhong et al., 2016) use a similar tensor construction followed by EM to learn the parameters; however, they can only handle linear noiseless mixtures and no gating parameters. In contrast, our algorithms can handle non-linearities and the gating parameters. (Chen et al., 2014) use a convex objective to learn the regressors for a special setting of 2-mixtures of linear regressions. Similar to earlier approaches, this relaxation too does not generalize to $k > 2$.

Notation. In this paper, we denote Euclidean vectors by bold face lowercase letters \mathbf{a}, \mathbf{b} , etc., and scalars by plain lowercase letters y, z , etc. We use $\mathcal{N}(y|\mu, \sigma^2)$ either to denote the density or the distribution of a Gaussian random

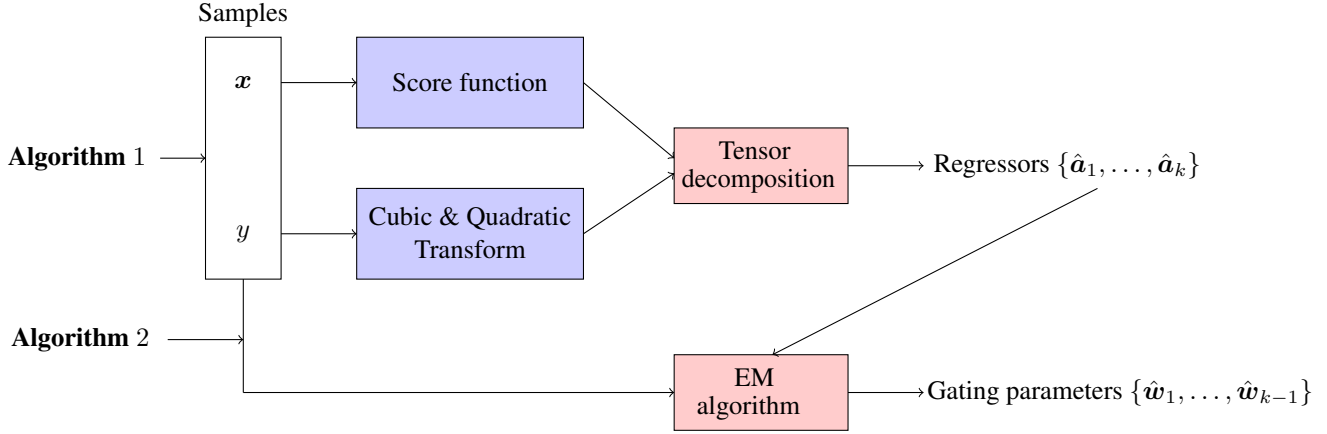


Figure 2: Algorithm to learn the MoE parameters. **Algorithm 1**: First we take non-linear transformations on the samples (\mathbf{x}_i, y_i) to compute the tensors $\mathcal{T}_2, \mathcal{T}_3$. Spectral decomposition on $\mathcal{T}_2, \mathcal{T}_3$ recovers the regressors. **Algorithm 2**: EM uses the learnt regressors and samples to learn the gating parameters with random initializations

variable y with mean μ and variance σ^2 , depending on the context. $[d] \triangleq \{1, \dots, d\}$. $\text{Perm}[d]$ denotes the set of all permutations on $[d]$. We use \otimes to denote the tensor outer product of vectors in \mathbb{R}^d . $\mathbf{x}^{\otimes 3}$ denotes $\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}$, where $(\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x})_{ijk} = x_i x_j x_k$. $\text{sym}(\mathbf{x} \otimes \mathbf{y} \otimes \mathbf{z})$ denotes the symmetrized version of $\mathbf{x} \otimes \mathbf{y} \otimes \mathbf{z}$, i.e. $\text{sym}(\mathbf{x} \otimes \mathbf{y} \otimes \mathbf{z})_{ijk} = \sum_{\sigma \in \text{Perm}[d]} x_{\sigma(i)} y_{\sigma(j)} z_{\sigma(k)}$. $\mathbf{e}_i, i \in [d]$ denotes the standard basis vectors for \mathbb{R}^d . Through out the paper, we assume that $\mathbf{w}_k^* = 0$, without loss of generality.

2. Algorithms

In this section, we present our algorithms to learn the regression and gating parameters *separately*. Figure 2 summarizes our algorithm. First we take a moment to highlight the issues of the existing approaches.

For illustration purposes, we suppose that $k = 2$ in (1). We assume without loss of generality that $\mathbf{w}_k^* = \mathbf{w}_2^* = 0$ and denote $\mathbf{w}_1^* = \mathbf{w}^*$. Thus the 2-MoE model is given by $P_{y|\mathbf{x}}$:

$$\frac{e^{\langle \mathbf{w}^*, \mathbf{x} \rangle} \mathcal{N}(y|g(\langle \mathbf{a}_1^*, \mathbf{x} \rangle), \sigma^2)}{1 + e^{\langle \mathbf{w}^*, \mathbf{x} \rangle}} + \frac{\mathcal{N}(y|g(\langle \mathbf{a}_2^*, \mathbf{x} \rangle), \sigma^2)}{1 + e^{\langle \mathbf{w}^*, \mathbf{x} \rangle}} \quad (2)$$

Issues with traditional tensor methods. In the far simplified setting of the absence of the gating parameter, i.e. $\mathbf{w}^* = 0 \in \mathbb{R}^d$, we see that 2-MoE reduces to 2-uniform mixture of GLMs. In this case, for $\mathbf{x} \sim \mathcal{N}(0, I_d)$, the standard approach is to construct a 3rd-order tensor \mathcal{T} by regressing the output y on the score transformation $\mathcal{S}_3(\mathbf{x}) \triangleq \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} - \sum_{i \in [d]} \text{sym}(\mathbf{x} \otimes \mathbf{e}_i \otimes \mathbf{e}_i)$, i.e.

$$\begin{aligned} \mathcal{T} \triangleq \mathbb{E}[y \cdot \mathcal{S}_3(\mathbf{x})] &= \frac{1}{2} \mathbb{E}[g'''(\langle \mathbf{a}_1^*, \mathbf{x} \rangle)] \cdot (\mathbf{a}_1^*)^{\otimes 3} \\ &+ \frac{1}{2} \mathbb{E}[g'''(\langle \mathbf{a}_2^*, \mathbf{x} \rangle)] \cdot (\mathbf{a}_2^*)^{\otimes 3}. \end{aligned} \quad (3)$$

Here the second equality follows from the generalized Stein's lemma that $\mathbb{E}[f(\mathbf{x}) \cdot \mathcal{S}_3(\mathbf{x})] = \mathbb{E}[\nabla_{\mathbf{x}}^{(3)} f(\mathbf{x})]$ under some regularity conditions on $f: \mathbb{R}^d \mapsto \mathbb{R}$ (see Lemma 2 in Appendix A). Then the regressors can be learned through spectral decomposition on \mathcal{T} , where the uniqueness of decomposition follows from (Kruskal, 1977). If we apply a similar technique for 2-MoE in (2), we obtain that

$$\begin{aligned} \mathbb{E}[y \cdot \mathcal{S}_3(\mathbf{x})] &= \sum_{i=1,2} \alpha_i (\mathbf{a}_i^*)^{\otimes 3} + \beta_i \text{sym}(\mathbf{a}_i^* \otimes \mathbf{a}_i^* \otimes \mathbf{w}^*) \\ &+ \gamma_i \text{sym}(\mathbf{a}_i^* \otimes \mathbf{w}^* \otimes \mathbf{w}^*) + \delta (\mathbf{w}^*)^{\otimes 3}, \end{aligned} \quad (4)$$

where $\alpha_i, \beta_i, \gamma_i, \delta$ are some scalar constants depending on the parameters $\mathbf{a}_1^*, \mathbf{a}_2^*, \mathbf{w}^*$ and g (see Appendix D.1 for the proof). Thus (4) reveals that traditional spectral methods do not yield a supersymmetric tensor of the desired parameters for MoEs. In fact, (4) contains all the 3rd-order rank-1 terms formed by $\mathbf{a}_1^*, \mathbf{a}_2^*$ and \mathbf{w}^* . Hence we cannot recover these parameters uniquely. Note that the inherent coupling between the regressors $\mathbf{a}_1^*, \mathbf{a}_2^*$ and the gating parameter \mathbf{w}^* in (2) manifests as a cross tensor in (4). This coupling serves as a key limitation for the traditional methods which critically rely on the fact that the mixing probability $p = \frac{1}{2}$ in (4) is a constant. In fact, we recover (3) by letting $\mathbf{w}^* = 0$ in (4).

Issues with EM algorithm. EM algorithm is the workhorse for parameter learning in both the k -MoE and HME models (Jordan & Jacobs, 1994). However, it is well known that EM is prone to spurious minima and existing theoretical results only establish local convergence for the regressors and the gating parameters. Indeed, our numerical experiments in Section 4.3 verify this fact. Figure 3b and Figure 3c highlight that joint-EM often gets stuck in bad local minima.

2.1. The proposed algorithm for learning MoE

In order to tackle these challenges, we take a different route and propose to estimate the regressors and gating parameters *separately*. To gain intuition about our approach, let us consider 2-MoE model in (2) with $\sigma = 0$ and linear g . Then we have that y either equals $\langle \mathbf{a}_1^*, \mathbf{x} \rangle$ with probability $\sigma(\langle \mathbf{w}^*, \mathbf{x} \rangle)$ or equals $\langle \mathbf{a}_2^*, \mathbf{x} \rangle$ with probability $1 - \sigma(\langle \mathbf{w}^*, \mathbf{x} \rangle)$, where $\sigma(\cdot)$ is the sigmoid function. If we exactly know \mathbf{w}^* , we can recover \mathbf{a}_1^* and \mathbf{a}_2^* by solving a simple linear regression problem since we can recover the true latent variable $z \in \{1, 2\}$ with high probability. Similarly, if we know \mathbf{a}_1^* and \mathbf{a}_2^* , it is easy to see that we can recover \mathbf{w}^* by solving a binary linear classification problem. Thus knowing either the regressors or the gating parameters makes the estimation of other parameters easier. However, how do we first obtain one set of parameters without any knowledge about the other?

Our approach precisely addresses this question and breaks the *grid lock*. We show that we can extract the regressors \mathbf{a}_1^* and \mathbf{a}_2^* without knowing \mathbf{w}^* at all, just using the samples. Although we explain our approach with two mixtures, all claims are made precise for general k in Theorems 1 and 2, and the algorithms are written for general k as well in Algorithms 1 and 2.

STEP 1: ESTIMATION OF REGRESSORS

To learn the regressors, we first pre-process $\mathbf{x} \sim \mathcal{N}(0, I_d)$ using the score transformations \mathcal{S}_3 and \mathcal{S}_2 , i.e.

$$\mathcal{S}_3(\mathbf{x}) \triangleq \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} - \sum_{i \in [d]} \text{sym}(\mathbf{x} \otimes \mathbf{e}_i \otimes \mathbf{e}_i), \quad (5)$$

$$\mathcal{S}_2(\mathbf{x}) \triangleq \mathbf{x} \otimes \mathbf{x} - I. \quad (6)$$

These score functions can be viewed as higher-order feature extractors from the inputs. As we have seen in (3), these transformations suffice to learn the parameters in M-GLMs. However this approach fails in the context of MoE, as highlighted in (4). Can we still construct a supersymmetric tensor for MoE?

To answer this question in a principled way, we introduce the notion of ‘‘Cubic and Quadratic Transform (CQT)’’ for the labels, i.e.

$$\mathcal{P}_3(y) \triangleq y^3 + \alpha y^2 + \beta y, \quad \mathcal{P}_2(y) \triangleq y^2 + \gamma y.$$

The coefficients (α, β, γ) in these polynomial transforms are obtained by solving a linear system of equations (see Appendix C). For the special case of $g = \text{linear}$, we obtain $\mathcal{P}_3(y) = y^3 - 3(1 + \sigma^2)y$ and $\mathcal{P}_2(y) = y^2$. These special transformations are specific to the choice of non-linearity g and the noise variance σ . The key intuition behind the design of these transforms is that we can nullify the cross moments and obtain supersymmetric tensor in (3) if we

regress $\mathcal{P}_3(y)$ instead of y , for properly chosen constants α and β . This is made mathematically precise in Theorem 1. A similar argument holds for $\mathcal{P}_2(y)$ too. In addition, the choice of these polynomials is unique in the sense that any other polynomial transformations fail to yield the desired tensor structure. Using these transforms, we construct two special tensors $\hat{\mathcal{T}}_3 \in (\mathbb{R}^d)^{\otimes 3}$ and $\hat{\mathcal{T}}_2 \in (\mathbb{R}^d)^{\otimes 2}$. Later we use the robust tensor power method (Anandkumar et al., 2014) on these tensors to learn the regressors. Algorithm 1 details our learning procedure. Theorem 1 establishes the theoretical justification for our algorithm.

Algorithm 1 Learning the regressors

- 1: **Input:** Samples $(\mathbf{x}_i, y_i), i \in [n]$
 - 2: Compute $\hat{\mathcal{T}}_3 = \frac{1}{n} \sum_{i=1}^n \mathcal{P}_3(y_i) \cdot \mathcal{S}_3(\mathbf{x}_i)$ and $\hat{\mathcal{T}}_2 = \frac{1}{n} \sum_{i=1}^n \mathcal{P}_2(y_i) \cdot \mathcal{S}_2(\mathbf{x}_i)$
 - 3: $\hat{\mathbf{a}}_1, \dots, \hat{\mathbf{a}}_k = \text{Rank-}k \text{ tensor decomposition on } \hat{\mathcal{T}}_3 \text{ using } \hat{\mathcal{T}}_2$
-

STEP 2: ESTIMATION OF GATING PARAMETERS

To gain intuition for estimating the gating parameters, let $g = \text{linear}$ in (2) for simplicity. Moreover, assume that we know both \mathbf{a}_1^* and \mathbf{a}_2^* . Then taking conditional expectation on y , we obtain from (2) that

$$\begin{aligned} \mathbb{E}[y|\mathbf{x}] &= f(\langle \mathbf{w}^*, \mathbf{x} \rangle) \cdot \langle \mathbf{a}_1^*, \mathbf{x} \rangle + (1 - f(\langle \mathbf{w}^*, \mathbf{x} \rangle)) \cdot \langle \mathbf{a}_2^*, \mathbf{x} \rangle, \\ &= \langle \mathbf{a}_2^*, \mathbf{x} \rangle + f(\langle \mathbf{w}^*, \mathbf{x} \rangle) \cdot \langle \mathbf{a}_1^* - \mathbf{a}_2^*, \mathbf{x} \rangle, \end{aligned} \quad (7)$$

where f is the sigmoid function. Thus,

$$\mathbb{E} \left[\frac{y - \langle \mathbf{a}_2^*, \mathbf{x} \rangle}{\langle \mathbf{a}_1^* - \mathbf{a}_2^*, \mathbf{x} \rangle} \middle| \mathbf{x} \right] = \frac{\mathbb{E}[y|\mathbf{x}] - \langle \mathbf{a}_2^*, \mathbf{x} \rangle}{\langle \mathbf{a}_1^* - \mathbf{a}_2^*, \mathbf{x} \rangle} = f(\langle \mathbf{w}^*, \mathbf{x} \rangle).$$

Note that since \mathbf{x} is Gaussian, $\langle \mathbf{a}_1^* - \mathbf{a}_2^*, \mathbf{x} \rangle$ is non-zero with probability 1. Hence, to recover \mathbf{w}^* , in view of Stein’s lemma, we may write

$$\begin{aligned} \mathbb{E} \left[\left(\frac{y - \langle \mathbf{a}_2^*, \mathbf{x} \rangle}{\langle \mathbf{a}_1^* - \mathbf{a}_2^*, \mathbf{x} \rangle} \right) \cdot \mathbf{x} \right] &\stackrel{\mathbf{x}}{=} \mathbb{E} [f(\langle \mathbf{w}^*, \mathbf{x} \rangle) \cdot \mathbf{x}] \\ &= \mathbb{E} [f'(\langle \mathbf{w}^*, \mathbf{x} \rangle)] \cdot \mathbf{w}^* \\ &= \mathbb{E}_{Z \sim \mathcal{N}(0,1)} f'(\| \mathbf{w}^* \| Z) \cdot \mathbf{w}^* \\ &\propto \mathbf{w}^*. \end{aligned}$$

However, it turns out that the above chain of equalities does not hold. Surprisingly, the first equality, which essentially is the law of iterated expectations, is not valid in this case as $\frac{y - \langle \mathbf{a}_2^*, \mathbf{x} \rangle}{\langle \mathbf{a}_1^* - \mathbf{a}_2^*, \mathbf{x} \rangle}$ is not integrable since it is a mixture of two Cauchy distributions, as proved in Appendix D.4. Thus the above analysis highlights the difficulty of learning the gating parameters even in the simplest setting of two linear mixtures. Can we still learn \mathbf{w}^* using method of moments (MoM)? In Theorem 3, we precisely address this question

and show that we can still provably recover the gating parameters using MoM, by designing clever transformations on the data to infer the parameters of a Cauchy mixture distribution.

While Theorem 3 highlights that gating parameters can be learnt using the method of moments for 2-MoE, we still need a principled approach to learn these parameters for a more generic setting of k -MoE. Recall that the traditional joint-EM algorithm randomly initializes both the regressors and the gating parameters and updates them iteratively. Figure 3b and Figure 3c highlight that this procedure is prone to spurious minima. Can we still learn the gating parameters with *global initializations*? To address this question, we utilize the regressors learnt from Algorithm 1. In particular, we use EM algorithm to update *only* the gating parameters, while fixing the regressors $\hat{\mathbf{a}}_1, \dots, \hat{\mathbf{a}}_k$. We show in Theorem 2 that, with *global/random* initializations, this variant of EM algorithm learns the true parameters. To the best of our knowledge, this is the first global convergence result for EM for $k > 2$ mixtures. This motivates the following algorithm ($\varepsilon > 0$ is some error tolerance):

Algorithm 2 Learning the gating parameter

- 1: **Input:** Samples $(\mathbf{x}_i, y_i), i \in [n]$ and regressors $\hat{\mathbf{a}}_1, \dots, \hat{\mathbf{a}}_k$ from Algorithm 1
 - 2: $t \leftarrow 0$
 - 3: Initialize \mathbf{w}_0 uniformly randomly in its domain Ω
 - 4: **while** (Estimation error $< \varepsilon$) **do**
 - 5: Compute the posterior $p_{\mathbf{w}_t}^{(i)}$ according to (9) for each $j \in [k]$ and $i \in [n]$
 - 6: Compute $Q(\mathbf{w}|\mathbf{w}_t)$ according to (8) using empirical expectation
 - 7: Set $\mathbf{w}_{t+1} = \operatorname{argmax}_{\mathbf{w} \in \Omega} Q(\mathbf{w}|\mathbf{w}_t)$
 - 8: $t \leftarrow t + 1$
 - 9: Estimation error = $\|\mathbf{w}_t - \mathbf{w}_{t-1}\|$
 - 10: **end while**
-

3. Theoretical analysis

In this section, we provide the theoretical guarantees for our algorithms in the population setting. We first formally state our assumptions and justify the rationale behind them:

1. \mathbf{x} follows standard Gaussian distribution, i.e. $\mathbf{x} \sim \mathcal{N}(0, I_d)$.
2. $\|\mathbf{a}_i^*\|_2 = 1$ for $i \in [k]$ and $\|\mathbf{w}_i^*\|_2 \leq R$ for $i \in [k-1]$, with some $R > 0$.
3. $\mathbf{a}_i^*, i \in [k]$ are linearly independent and \mathbf{w}_i^* is orthogonal to $\operatorname{span}\{\mathbf{a}_1^*, \dots, \mathbf{a}_k^*\}$ for $i \in [k-1]$.
4. The non-linearity $g: \mathbb{R} \rightarrow \mathbb{R}$ is (α, β, γ) -valid, which

we define in Appendix C. For example, this class includes $g = \text{linear}$, sigmoid and ReLU.

Remark. We note that the Gaussianity of the input distribution and norm constraints on the parameters are standard assumptions in the learning of neural networks literature (Janzamin et al., 2015; Li & Yuan, 2017; Ge et al., 2017; Zhong et al., 2017; Du et al., 2017; Safran & Shamir, 2017) and also that of M-GLMs (Sedghi et al., 2014; Yi et al., 2016; Zhong et al., 2016; Balakrishnan et al., 2017). An interpretation behind Assumption 3 is that if we think of \mathbf{x} as a high-dimensional feature vector, distinct sub-features of \mathbf{x} are used to perform the two distinct tasks of classification (using \mathbf{w}_i^* 's) and regression (using \mathbf{a}_i^* 's). We note that we need the above assumptions only for the technical analysis. In Section 4.1 and Section 4.2, we empirically verify that our algorithms work well in practice even under the relaxation of these assumptions. Thus we believe that the assumptions are merely technical artifacts.

We are now ready to state our results.

Theorem 1 (Recovery of regression parameters). *Let (\mathbf{x}, y) be generated according to the true model (1). Under the above assumptions, we have that*

$$\mathcal{T}_2 \triangleq \mathbb{E}[\mathcal{P}_2(y) \cdot \mathcal{S}_2(\mathbf{x})] = \sum_{i=1}^k c_g' \mathbb{E}[P_{i|\mathbf{x}}] \cdot \mathbf{a}_i^* \otimes \mathbf{a}_i^*,$$

$$\mathcal{T}_3 \triangleq \mathbb{E}[\mathcal{P}_3(y) \cdot \mathcal{S}_3(\mathbf{x})] = \sum_{i=1}^k c_{g,\sigma} \mathbb{E}[P_{i|\mathbf{x}}] \cdot \mathbf{a}_i^* \otimes \mathbf{a}_i^* \otimes \mathbf{a}_i^*,$$

where c_g' and $c_{g,\sigma}$ are two non-zero constants depending on g and σ . Hence the regressors \mathbf{a}_i^* 's can be learnt through tensor decomposition on \mathcal{T}_2 and \mathcal{T}_3 .

Proof. (Sketch) To highlight the central ideas behind the proof, first let $g = \text{linear}$. From (1) we get that

$$\mathbb{E}[y|\mathbf{x}] = \sum_{i \in [k]} p_i^*(\mathbf{x}) \langle \mathbf{a}_i^*, \mathbf{x} \rangle,$$

where $p_i^*(\mathbf{x}) \triangleq P_{i|\mathbf{x}}$ for $i \in [k]$. Taking the cross moment of y with $\mathcal{S}_3(\mathbf{x})$ and using Lemma 2 we obtain that

$$\begin{aligned} \mathbb{E}[y \cdot \mathcal{S}_3(\mathbf{x})] &= \sum_{i \in [k]} \mathbb{E}[p_i^*(\mathbf{x}) \langle \mathbf{a}_i^*, \mathbf{x} \rangle \cdot \mathcal{S}_3(\mathbf{x})] \\ &= \sum_{i \in [k]} \mathbb{E}[\nabla_{\mathbf{x}}^{(3)}(p_i^*(\mathbf{x}) \langle \mathbf{a}_i^*, \mathbf{x} \rangle)]. \end{aligned}$$

Notice that had $p_i^*(\mathbf{x})$ been a constant in the above equation, we would obtain a supersymmetric tensor easily as is the case with M-GLMs. However, $\mathbb{E}[\nabla_{\mathbf{x}}^{(3)}(p_i^*(\mathbf{x}) \langle \mathbf{a}_i^*, \mathbf{x} \rangle)]$ now contains all the third-order rank-1 terms involving the tensor product of $\mathbf{w}_1^*, \dots, \mathbf{w}_{k-1}^*$ and \mathbf{a}_i^* for any fixed i . Our key

insight is that this issue can be avoided if we cleverly transform y . In particular, we consider a cubic transformation $\mathcal{P}_3(y) = y^3 - 3y(1 + \sigma^2)$ and obtain that

$$\mathbb{E}[\mathcal{P}_3(y)|\mathbf{x}] = \sum_{i \in [k]} p_i^*(\mathbf{x}) (\langle \mathbf{a}_i^*, \mathbf{x} \rangle^3 - 3\langle \mathbf{a}_i^*, \mathbf{x} \rangle)$$

Now it turns out that after using the orthogonality of \mathbf{w}_i^* and \mathbf{a}_i^* , and the fact $\mathbb{E}[p(Z)] = \mathbb{E}[p'(Z)] = \mathbb{E}[p''(Z)] = 0$ for 3rd-Hermite polynomial $p(z) = z^3 - 3z$ and $Z \sim \mathcal{N}(0, 1)$, we can nullify the cross-moments between \mathbf{w}_i^* 's and \mathbf{a}_i^* 's to obtain that

$$\mathbb{E}[\mathcal{P}_3(y) \cdot \mathcal{S}_3(\mathbf{x})] = 6 \sum_{i \in [k]} \mathbb{E}[p_i^*(\mathbf{x})] \cdot (\mathbf{a}_i^*)^{\otimes 3}.$$

Similarly, we can show that $\mathbb{E}[\mathcal{P}_2(y) \cdot \mathcal{S}_2(\mathbf{x})] = 2 \sum_{i \in [k]} \mathbb{E}[p_i^*(\mathbf{x})] \cdot (\mathbf{a}_i^*)^{\otimes 2}$. For a general non-linearity $g: \mathbb{R} \rightarrow \mathbb{R}$, we can similarly design cubic and quadratic polynomials $\mathcal{P}_3 = y^3 + \alpha y^2 + \beta y$ and $\mathcal{P}_2 = y^2 + \gamma y$ such that we can still construct supersymmetric tensors involving the regressors. In order to obtain the unique set of coefficients (α, β, γ) , we need to solve a linear system of equations, which we describe in Appendix C. \square

Once we obtain \mathcal{T}_2 and \mathcal{T}_3 , the recovery guarantees for the regressors \mathbf{a}_i^* follow from the standard tensor decomposition guarantees, for example, Theorem 4.3 and Theorem 5 of (Anandkumar et al., 2014). We assume that the learnt regressors \mathbf{a}_i are such that $\max_{i \in [k]} \|\mathbf{a}_i - \mathbf{a}_i^*\|_2 = \sigma^2 \varepsilon$ for some $\varepsilon > 0$. Now we present our theoretical results for global convergence of EM. First we briefly recall the algorithm. Let Ω denote the domain of our gating parameters, defined as

$$\Omega = \{\mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_{k-1}) : \|\mathbf{w}_i\|_2 \leq R, \forall i \in [k-1]\}.$$

Then the population EM for the mixture of experts consists of the following two steps:

- **E-step:** Using the current estimate \mathbf{w}_t to compute the function $Q(\cdot|\mathbf{w}_t)$,
- **M-step:** $\mathbf{w}_{t+1} = \operatorname{argmax}_{\mathbf{w} \in \Omega} Q(\mathbf{w}|\mathbf{w}_t)$,

where the function $Q(\cdot|\mathbf{w}_t)$ is the expected log-likelihood of the complete data distribution with respect to current posterior distribution. Mathematically,

$$\begin{aligned} Q(\mathbf{w}|\mathbf{w}_t) &\triangleq \mathbb{E}_{(\mathbf{x}, y)} \mathbb{E}_{P_{z|\mathbf{x}, y, \mathbf{w}_t}} [\log P_{\mathbf{w}}(\mathbf{x}, z, y)] \\ &= \mathbb{E}_{(\mathbf{x}, y)} \mathbb{E}_{P_{z|\mathbf{x}, y, \mathbf{w}_t}} [\log P(\mathbf{x}) P_{\mathbf{w}}(z|\mathbf{x}) P(y|\mathbf{x}, z)] \\ &= \mathbb{E}_{(\mathbf{x}, y)} \mathbb{E}_{P_{z|\mathbf{x}, y, \mathbf{w}_t}} [\log P_{\mathbf{w}}(z|\mathbf{x})] + \text{const.} \\ &= \mathbb{E} \left[\sum_{i \in [k-1]} p_{\mathbf{w}_t}^{(i)}(\mathbf{w}_i^\top \mathbf{x}) - \left(1 + \sum_{i \in [k-1]} e^{\mathbf{w}_i^\top \mathbf{x}}\right) \right] \\ &\quad + \text{const.} \end{aligned} \quad (8)$$

where const refers to terms not depending on \mathbf{w} , $P_{\mathbf{w}}(z = i|\mathbf{x}) = \exp(\mathbf{w}_i^\top \mathbf{x}) / \sum_j \exp(\mathbf{w}_j^\top \mathbf{x})$ and $p_{\mathbf{w}_t}^{(i)} \triangleq \mathbb{P}[z = i|\mathbf{x}, y, \mathbf{w}_t]$ corresponds to the posterior probability for the i^{th} expert, given by

$$\begin{aligned} p_{\mathbf{w}_t}^{(i)} &= \frac{p_{i,t}(\mathbf{x}) \mathcal{N}(y|g(\mathbf{a}_i^\top \mathbf{x}), \sigma^2)}{\sum_{j \in [k]} p_{j,t}(\mathbf{x}) \mathcal{N}(y|g(\mathbf{a}_j^\top \mathbf{x}), \sigma^2)}, \quad (9) \\ p_{i,t}(\mathbf{x}) &= \frac{e^{(\mathbf{w}_t)_i^\top \mathbf{x}}}{1 + \sum_{j \in [k-1]} e^{(\mathbf{w}_t)_j^\top \mathbf{x}}}. \end{aligned}$$

In (8), the expectation is with respect to the true distribution of (\mathbf{x}, y) , given by (1). Thus the EM can be viewed as a deterministic procedure which maps $\mathbf{w}_t \mapsto M(\mathbf{w}_t)$ where

$$M(\mathbf{w}) = \operatorname{argmax}_{\mathbf{w}' \in \Omega} Q(\mathbf{w}'|\mathbf{w}).$$

When the estimated regressors \mathbf{a}_i equal the true parameters \mathbf{a}_i^* , it follows from the self-consistency property of the EM that the true parameter \mathbf{w}^* is a fixed-point for the EM operator M , i.e. $M(\mathbf{w}^*) = \mathbf{w}^*$ (McLachlan & Krishnan, 2007). However, this does not guarantee that EM converges to \mathbf{w}^* . In the following theorem, we show that even when the regressors are known approximately, EM algorithm converges to the true gating parameters at a geometric rate upto an additive error, under *global* initializations. For the error metric, we define $\|\mathbf{w} - \mathbf{w}'\| \triangleq \max_{i \in [k-1]} \|\mathbf{w}_i - \mathbf{w}'_i\|_2$ for any $\mathbf{w}, \mathbf{w}' \in \Omega$. We assume that $R = 1$ for simplicity. (Our results extend straightforwardly to general R).

Theorem 2. *Let $\varepsilon > 0$ be such that $\max_i \|\mathbf{a}_i - \mathbf{a}_i^*\|_2 = \sigma^2 \varepsilon$. There exists a constant $\sigma_0 > 0$ such that whenever $0 < \sigma < \sigma_0$, for any random initialization $\mathbf{w}_0 \in \Omega$, the population-level EM updates on the gating parameter $\{\mathbf{w}\}_{t \geq 0}$ converge almost geometrically to the true parameter \mathbf{w}^* upto an additive error, i.e.*

$$\|\mathbf{w}_t - \mathbf{w}^*\| \leq (\kappa_\sigma)^t \|\mathbf{w}_0 - \mathbf{w}^*\| + \kappa \varepsilon \sum_{i=0}^{t-1} \kappa_\sigma^i,$$

where κ_σ, κ are dimension-independent constant depending on g and σ such that $\kappa_\sigma \xrightarrow{\sigma \rightarrow 0} 0$ and $\kappa \leq (k-1) \frac{\sqrt{6(2+\sigma^2)}}{2}$ for $g = \text{linear, sigmoid and ReLU}$.

Proof. (Sketch) One can show that the $Q(\cdot|\mathbf{w}_t)$ defined in (8) is a strongly concave function. Moreover, if we let $\varepsilon = 0$ and $\mathbf{w}_t = \mathbf{w}^*$, we have from the self-consistency of EM that $\operatorname{argmax} Q(\cdot|\mathbf{w}^*) = \mathbf{w}^*$. Thus if we can show that the functions are $Q(\cdot|\mathbf{w}_t)$ and $Q(\cdot|\mathbf{w}^*)$ ‘‘sufficiently close’’ whenever \mathbf{w}_t and \mathbf{w}^* are close, we can use the EM convergence analysis tools from (Balakrishnan et al., 2017) to show that their corresponding maximizers also stay close upto a scaling factor determined by κ_σ above. Then it follows that the EM updates converge geometrically. \square

Remark. In the M-step of the EM algorithm, the next iterate is chosen so that the function $Q(\cdot|\mathbf{w}_t)$ is maximized. Instead we can perform an ascent step in the direction of the gradient of $Q(\cdot|\mathbf{w}_t)$ to produce the next iterate, i.e. $\mathbf{w}_{t+1} = \Pi_{\Omega}(\mathbf{w}_t + \alpha \nabla Q(\mathbf{w}_t|\mathbf{w}_t))$, where $\Pi_{\Omega}(\cdot)$ is the projection operator. This variant of EM algorithm is known as *Gradient EM*. In Appendix G, we show that Gradient EM also enjoys similar convergence guarantees.

MoM to learn gating parameters. In Theorem 2, we proved that EM algorithm provably recovers the true gating parameters for any $k \geq 2$ mixtures. In this section, we show that for the special case of $k = 2$, we can learn \mathbf{w}^* (upto the unit direction) using an alternative procedure involving MoM. First we define

$$\text{Ratio}(\mathbf{x}, y) \triangleq \frac{y - \langle \mathbf{a}_2, \mathbf{x} \rangle}{\langle \mathbf{a}_1 - \mathbf{a}_2, \mathbf{x} \rangle} \quad (10)$$

The following theorem establishes that the the CDF of the random variable $\text{Ratio}(\mathbf{x}, y)$, when regressed on input \mathbf{x} , is proportional to \mathbf{w}^* .

Theorem 3. *Suppose that $(\mathbf{a}_1, \mathbf{a}_2) = (\mathbf{a}_1^*, \mathbf{a}_2^*)$. Then we have that*

$$\mathbb{E}[\mathbb{1}\{\text{Ratio}(\mathbf{x}, y) \leq 0.5\} \cdot \mathbf{x}] = \alpha \mathbf{w}^*,$$

where $\alpha \in \mathbb{R}$ is a scalar given by $\alpha = \mathbb{E}[f'(\langle \mathbf{w}^*, \mathbf{x} \rangle) \cdot (1 - 2\Phi(\frac{|\langle \mathbf{a}_1 - \mathbf{a}_2, \mathbf{x} \rangle|}{2\sigma}))]$.

Proof. (Sketch) We first show that $\text{Ratio}(\mathbf{x}, y)$ is a mixture of Cauchy distributions. Then we show that $\mathbb{E}[\mathbb{1}\{\text{Ratio}(\mathbf{x}, y) \leq z\}|\mathbf{x}] = \mathbb{P}[\text{Ratio} \leq z|\mathbf{x}] = f(\mathbf{w}^\top \mathbf{x})\Phi\left(\frac{(z-1)|\Delta_x|}{\sigma}\right) + (1 - f(\mathbf{w}^\top \mathbf{x}))\Phi\left(\frac{z|\Delta_x|}{\sigma}\right)$ where $\Delta_x = (\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}$. Then our result follows from taking the first moment of the indicator random variable with \mathbf{x} and Stein’s lemma. \square

4. Experiments

In this section, we empirically validate our algorithm in various settings and compare its performance to that of EM on both synthetic and real world datasets¹. In both the scenarios, we found that our algorithm consistently outperforms the existing approaches. For the tensor decomposition in our Algorithm 1, we use the Orth-ALS package by (Sharan & Valiant, 2017). In all the synthetic experiments, we first draw the regressors $\{\mathbf{a}_i^*\}_{i=1}^k$ i.i.d uniformly from the unit sphere \mathbb{S}^{d-1} . The input distribution $P_{\mathbf{x}}$ and the generation of \mathbf{w}_i^* ’s are detailed for each experiment. Then the labels y_i are generated according to the true k -MoE model in (1) for linear activation. Additional experiments in this setting with non-linear activations are detailed in Appendix H.1. Experiments with real world data are provided in Section 4.4.

¹Codes are available at this repository [MoE codes](#).

4.1. Non-gaussian inputs

In this section we let the input distribution to be mixtures of Gaussians (GMM). We let $k = 2, d = 10$ and $\sigma = 0.1$. The gating parameter $\mathbf{w}^* \in \mathbb{R}^{10}$ is uniformly chosen from the unit sphere \mathbb{S}^9 . To generate the input features, we first randomly draw $\mu_1, \mu_2 \in \mathbb{S}^9$, and generate n i.i.d. samples $\mathbf{x}_i \sim p\mathcal{N}(\mu_1, I_d) + (1-p)\mathcal{N}(\mu_2, I_d)$, where $p \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$. Here $n = 2000$. Since \mathbf{x} is a 2-GMM, its score functions $\mathcal{S}_3(\mathbf{x}), \mathcal{S}_2(\mathbf{x})$ are computed using the densities of Gaussian mixtures (Janzamin et al., 2014). To gauge the performance of our algorithm, we measure the correlation of our learned parameters $\mathbf{a}_1, \mathbf{a}_2$ and \mathbf{w} with the ground truth, i.e.

$$\text{Regressor Fit}(\mathbf{a}_1, \mathbf{a}_2) = \max_{\pi} \min_{i \in \{1,2\}} |\langle \mathbf{a}_{\pi(i)}, \mathbf{a}_i^* \rangle|, \quad (11)$$

where $\pi : \{1, 2\} \rightarrow \{1, 2\}$ is a permutation. Similarly, for the gating parameter, we define

$$\text{Gating Fit}(\mathbf{w}) = |\langle \mathbf{w}, \mathbf{w}^* \rangle|. \quad (12)$$

Here we assume that all the parameters are unit-normalized. The closer the values of fit are to 1, the closer the learnt parameters are to the ground truth. As shown in Table 1, our algorithms are able to learn the ground truth very accurately in a variety of settings, as indicated by the measured fit. This highlights the fact that our algorithms are robust to the input distributions.

4.2. Non-orthogonal parameters

In this section we verify that our algorithms still work well in practice even under the relaxation of Assumption 3. For the experiments, we consider the similar setting as before with $k = 2, d = 10, \sigma = 0.1$ and the gating parameter \mathbf{w}^* is drawn uniformly from \mathbb{S}^9 without the orthogonality restriction. We let $\mathbf{x}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, I_d)$. We choose $n = 2000$. We use RegressorFit and GatingFit defined in (11) and (12) respectively, as our performance metrics. From Table 2, we can see that the performance of our algorithms is almost the same across both the settings. In both the scenarios, our fit is consistently greater than 0.9.

In Figure 3a, we plotted $\text{GatingFit}(\mathbf{w}_t)$ vs. the number of iterations t , as \mathbf{w}_t is updated according to Algorithm 2, over 10 independent trials. We observe that the learned parameters converge to the true parameters in less than 5 iterations.

4.3. Comparison to joint-EM

Here we compare the performance of our algorithm with that of the joint-EM. We let the number of mixture components be $k = 3$ and $k = 4$. We let $\mathbf{x} \sim \mathcal{N}(0, I_d)$ and the gating parameters are drawn uniformly from \mathbb{S}^9 . If $\mathbf{A} = [\mathbf{a}_1 \dots \mathbf{a}_k]$

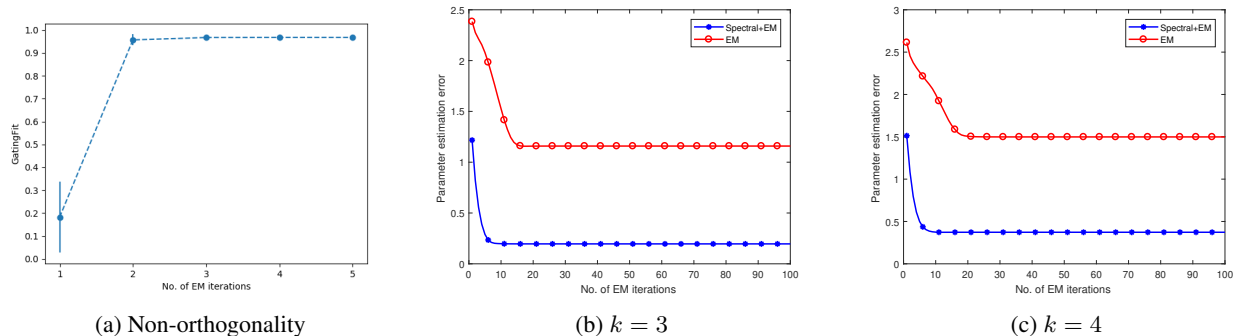


Figure 3: (a): GatingFit for our algorithm under non-orthogonality setting. (b),(c): Estimation error $\mathcal{E}(\mathbf{A}, \mathbf{W})$ of our algorithm vs. joint-EM algorithm. Our algorithm is significantly better than the joint-EM under random initializations.

Table 1: Fit of our learned parameters for non-Gaussian inputs

	$p = 0.1$	$p = 0.3$	$p = 0.5$	$p = 0.7$	$p = 0.9$
Regressor Fit	0.93 ± 0.06	0.94 ± 0.02	0.92 ± 0.04	0.92 ± 0.02	0.91 ± 0.06
Gating Fit	0.9 ± 0.1	0.97 ± 0.01	0.93 ± 0.04	0.96 ± 0.03	0.97 ± 0.01

Table 2: Performance of our algorithm under orthogonal and non-orthogonal settings

	Regressor Fit	Gating Fit
Non-orthogonal	0.9 ± 0.08	0.96 ± 0.02
Orthogonal	0.93 ± 0.03	0.96 ± 0.03

and $\mathbf{W} = [\mathbf{w}_1 \dots \mathbf{w}_{k-1} \ 0]$ denote the estimated expert and gating parameters respectively, our evaluation metric is \mathcal{E} , the Frobenious norm of the parameter error accounting for the best possible permutation $\pi : [k] \rightarrow [k]$, i.e. $\mathcal{E}(\mathbf{A}, \mathbf{W}) = \inf_{\pi} \|\mathbf{A} - \mathbf{A}_{\pi}^*\|_F + \|\mathbf{W} - \mathbf{W}_{\pi}^*\|_F$, where $\mathbf{A}_{\pi}^* = [\mathbf{a}_{\pi(1)}^* \dots \mathbf{a}_{\pi(k)}^*]$ denotes the permuted regression parameter matrix and similarly for \mathbf{W}_{π}^* . In Figure 3b and Figure 3c, we compare the performance of our algorithm with the joint-EM algorithm for $n = 8000, d = 10, \sigma = 0.5$. The plotted estimation error $\mathcal{E}(\mathbf{A}, \mathbf{W})$ is averaged for 10 trials. It is clear that our algorithm is able to recover the true parameters thus resulting in much smaller parameter error than the joint-EM which often gets stuck in local optima. In addition, our algorithm is able to learn these parameters in very few iterations, often less than 10 iterations. We also find that our algorithm consistently outperforms the joint-EM for different choices of non-linearities, number of samples, number of mixtures, etc. (details provided in Appendix H). Note that the above error metric $\mathcal{E}(\mathbf{A}, \mathbf{W})$ is close to zero if and only if Regressor Fit and Gating Fit is close to one.

4.4. Real data

To highlight the generalizability of our algorithm, in Appendix H.2 of the supplement, we compare the performance

of our algorithm to that of the standard approaches on a variety of real world datasets. Results from these experiments highlight the fact that in the real world scenario, where the underlying data is not generated according to a MoE model, our approach still learns a superior set of parameters as opposed to the existing algorithms. This fact is reflected in the lowest prediction errors obtained by our algorithm.

5. Discussion

In this paper we provided the first provable and globally consistent algorithm that can learn the true parameters of a MoE model. We believe that ideas from (Sedghi et al., 2014) can be naturally extended for the finite sample complexity analysis of the tensor decomposition to learn the regressors and similarly, techniques from (Balakrishnan et al., 2017) can be extended to the finite sample EM convergence analysis for the gating parameters. While we have focused here on parameter recovery, however, there are no statistical bounds on output prediction error when the data is not generated from the model. MoE models are known to be capable of fitting general functions, and getting statistical guarantees on learning in such regimes is an interesting direction for future work.

Acknowledgements

This work is partly supported by NSF grants 1927712 and 1815535, NSF awards CNS-1718270, 1651236, 1703403, and the Army Research Office under grant W911NF1810332.

References

- Anandkumar, A., Ge, R., Hsu, D., Kakade, S. M., and Telgarsky, M. Tensor decompositions for learning latent variable models. *J. Mach. Learn. Res.*, 15(1):2773–2832, January 2014. ISSN 1532-4435.
- Balakrishnan, S., Wainwright, M. J., and Yu, B. Statistical guarantees for the EM algorithm: From population to sample-based analysis. *The Annals of Statistics*, 45(1): 77–120, 2017.
- Brooks, T., Pope, D., and Marcolini, A. Airfoil self-noise and prediction. Technical report, NASA, 1989. URL <https://archive.ics.uci.edu/ml/datasets/Airfoil+Self-Noise>.
- Chen, Y., Yi, X., and Caramanis, C. A convex formulation for mixed regression with two components: Minimax optimal rates. In *Conference on Learning Theory*, pp. 560–604, 2014.
- Chung, J., Gülçehre, Ç., Cho, K., and Bengio, Y. Empirical evaluation of gated recurrent neural networks on sequence modeling. abs/1412.3555, 2014.
- Collobert, R., Bengio, S., and Bengio, Y. A parallel mixture of SVMs for very large scale problems. *Neural Computing*, 2002.
- Du, S. S., Lee, J. D., Tian, Y., Póczos, B., and Singh, A. Gradient descent learns one-hidden-layer cnn: Don’t be afraid of spurious local minima. *arXiv preprint arXiv:1712.00779*, 2017.
- Ge, R., Lee, J. D., and Ma, T. Learning one-hidden-layer neural networks with landscape design. *arXiv preprint arXiv:1711.00501*, 2017.
- Gross, S., Szlam, A., et al. Hard mixtures of experts for large scale weakly supervised vision. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pp. 5085–5093. IEEE, 2017.
- Jacobs, R. A., Jordan, M. I., Nowlan, S. J., and Hinton, G. E. Adaptive mixtures of local experts. *Neural Computation*, 1991.
- Janzamin, M., Sedghi, H., and Anandkumar, A. Score function features for discriminative learning: Matrix and tensor framework. abs/1412.2863, 2014. URL <http://arxiv.org/abs/1412.2863>.
- Janzamin, M., Sedghi, H., and Anandkumar, A. Beating the perils of non-convexity: Guaranteed training of neural networks using tensor methods. *arXiv preprint arXiv:1506.08473*, 2015.
- Jin, C., Zhang, Y., Balakrishnan, S., Wainwright, M. J., and Jordan, M. Local maxima in the likelihood of gaussian mixture models: Structural results and algorithmic consequences. *arXiv preprint arXiv:1609.00978*, 2016.
- Jordan, M. I. and Jacobs, R. A. Hierarchical mixtures of experts and the EM algorithm. *Neural computation*, 6(2): 181–214, 1994.
- Jordan, M. I. and Xu, L. Convergence results for the EM approach to mixtures of experts architectures. *Neural Networks*, 8(9):1409–1431, 1995.
- Kruskal, J. B. Three-way arrays: rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics. *Linear algebra and its applications*, 18(2):95–138, 1977.
- Le, P., Dymetman, M., and Renders, J.-M. Lstm-based mixture-of-experts for knowledge-aware dialogues. *arXiv preprint arXiv:1605.01652*, 2016.
- Li, Y. and Yuan, Y. Convergence analysis of two-layer neural networks with relu activation. In *Advances in Neural Information Processing Systems*, pp. 597–607, 2017.
- Liu, Y.-C. and Yeh, I.-C. Using mixture design and neural networks to build stock selection decision support systems. *Neural Computing and Applications*, 28(3): 521–535, 2017. doi: 10.1007/s00521-015-2090-x. URL <https://archive.ics.uci.edu/ml/datasets/Stock+portfolio+performance>.
- Masoudnia, S. and Ebrahimpour, R. Mixture of experts: a literature survey. *Artificial Intelligence Review*, 42(2): 275, 2014.
- McLachlan, G. and Krishnan, T. *The EM algorithm and extensions*, volume 382. John Wiley & Sons, 2007.
- Ng, J. W. and Deisenroth, M. P. Hierarchical mixture-of-experts model for large-scale gaussian process regression. *arXiv preprint arXiv:1412.3078*, 2014.
- Safran, I. and Shamir, O. Spurious local minima are common in two-layer relu neural networks. *arXiv preprint arXiv:1712.08968*, 2017.
- Sedghi, H., Janzamin, M., and Anandkumar, A. Provable tensor methods for learning mixtures of classifiers. *arXiv preprint arXiv:1412.3046*, 2014.

- Sharan, V. and Valiant, G. Orthogonalized ALS: A theoretically principled tensor decomposition algorithm for practical use. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 3095–3104, 06–11 Aug 2017. URL <http://proceedings.mlr.press/v70/sharan17a.html>.
- Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q., Hinton, G., and Dean, J. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017.
- Stein, C. A bound for the error in the normal approximation to the distribution of a sum of dependent random variables. In *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability*, volume 2, pp. 583–602. University of California Press, 1972.
- Sun, X., Peng, X., Ren, F., and Xue, Y. Human-machine conversation based on hybrid neural network. In *Computational Science and Engineering (CSE) and Embedded and Ubiquitous Computing (EUC), 2017 IEEE International Conference on*, volume 1, pp. 260–266. IEEE, 2017.
- Sun, Y., Ioannidis, S., and Montanari, A. Learning mixtures of linear classifiers. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32, pp. 721–729, 2014.
- Theis, L. and Bethge, M. Generative image modeling using spatial lstms. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2, NIPS’15*, pp. 1927–1935, Cambridge, MA, USA, 2015. MIT Press.
- Tresp, V. Mixtures of gaussian processes. NIPS, 2001.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *Advances in Neural Information Processing Systems*, pp. 5998–6008, 2017.
- Wang, X., Yu, F., Wang, R., Ma, Y.-A., Mirhoseini, A., Darrell, T., and Gonzalez, J. E. Deep mixture of experts via shallow embedding. *arXiv preprint arXiv:1806.01531*, 2018.
- Yeh, I.-C. Modeling of strength of high performance concrete using artificial neural networks. *Cement and Concrete Research*, 28(12):1797–1808, 1998. URL <https://archive.ics.uci.edu/ml/datasets/Concrete+Compressive+Strength>.
- Yi, X., Caramanis, C., and Sanghavi, S. Alternating minimization for mixed linear regression. In *International Conference on Machine Learning*, pp. 613–621, 2014.
- Yi, X., Caramanis, C., and Sanghavi, S. Solving a mixture of many random linear equations by tensor decomposition and alternating minimization. *arXiv preprint arXiv:1608.05749*, 2016.
- Yuksel, S. E., Wilson, J. N., and Gader, P. D. Twenty years of mixture of experts. *IEEE Transactions on Neural Networks and Learning Systems*, 23(8):1177–1193, 2012.
- Zhong, K., Jain, P., and Dhillon, I. S. Mixed linear regression with multiple components. pp. 2190–2198. 2016.
- Zhong, K., Song, Z., Jain, P., Bartlett, P. L., and Dhillon, I. S. Recovery guarantees for one-hidden-layer neural networks. *arXiv preprint arXiv:1706.03175*, 2017.

Organization. The appendix is organized as follows:

- Appendix A and Appendix B contain the requisite material for method of moments and the convergence analysis of EM respectively.
- Appendix C details the class of non-linearities for which our results hold.
- Appendix D contains all the proofs of Section 3. Two technical lemmas needed to prove Theorem 2 are relegated to Appendix E and Appendix F.
- Appendix G provides convergence guarantees for Gradient EM.
- Appendix H contains additional experiments for the comparison of joint-EM and our algorithm for the synthetic data.

A. Toolbox for method of moments

In this section, we introduce the key techniques that are useful in parameter estimation of mixture models via the method of moments.

Stein’s identity (Stein’s lemma) is a well-known result in probability and statistics and is widely used in estimation and inference tasks. A refined version of the Stein’s lemma (Stein, 1972) for higher-order moments is the key to parameter estimation in mixture of generalized linear models. We utilize this machinery in proving Theorem 1. We first recall the Stein’s lemma.

Lemma 1 (Stein’s lemma (Stein, 1972)). *Let $\mathbf{x} \sim \mathcal{N}(0, I_d)$ and $g : \mathbb{R}^d \rightarrow \mathbb{R}$ be a function such that both $\mathbb{E}[\nabla_{\mathbf{x}}g(\mathbf{x})]$ and $\mathbb{E}[g(\mathbf{x}) \cdot \mathbf{x}]$ exist and are finite. Then*

$$\mathbb{E}[g(\mathbf{x}) \cdot \mathbf{x}] = \mathbb{E}[\nabla_{\mathbf{x}}g(\mathbf{x})].$$

The following lemma, which can be viewed as an extension of Stein’s lemma for higher-order moments, is the central technique behind parameter estimation in M-GLMs.

Lemma 2 ((Sedghi et al., 2014)). *Let $\mathbf{x} \sim \mathcal{N}(0, I_d)$ and $\mathcal{S}_3(\mathbf{x})$ be as defined in (6) and let $\mathcal{S}_2(\mathbf{x}) \triangleq \mathbf{x} \otimes \mathbf{x} - I_d$. Then for any $g : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfying some regularity conditions, we have*

$$\mathbb{E}[g(\mathbf{x}) \cdot \mathcal{S}_2(\mathbf{x})] = \mathbb{E}[\nabla_{\mathbf{x}}^{(2)}g(\mathbf{x})], \quad \mathbb{E}[g(\mathbf{x}) \cdot \mathcal{S}_3(\mathbf{x})] = \mathbb{E}[\nabla_{\mathbf{x}}^{(3)}g(\mathbf{x})].$$

B. Toolbox for EM convergence analysis

Recall that the domain of our gating parameters is $\Omega = \{\mathbf{w} : \|\mathbf{w}\| \leq 1\}$. Then the population EM for the mixture of experts consists of the following two steps:

- **E-step:** Using the current estimate \mathbf{w}_t to compute the function $Q(\cdot|\mathbf{w}_t)$.
- **M-step:** $\mathbf{w}_{t+1} = \operatorname{argmax}_{\|\mathbf{w}\| \leq 1} Q(\mathbf{w}|\mathbf{w}_t)$.

Thus the EM can be viewed as a deterministic procedure which maps $\mathbf{w}_t \mapsto M(\mathbf{w}_t)$ where

$$M(\mathbf{w}) = \operatorname{argmax}_{\mathbf{w}' \in \Omega} Q(\mathbf{w}'|\mathbf{w}).$$

Our convergence analysis relies on tools from (Balakrishnan et al., 2017) where they provided local convergence results on both the EM and gradient EM algorithms. In particular, they showed that if we initialize EM in a sufficiently small neighborhood around the true parameters, the EM iterates converge geometrically to the true parameters under some strong-concavity and gradient stability conditions. We now formally state the assumptions in (Balakrishnan et al., 2017) under which the convergence guarantees hold. We will show in the next section that these conditions hold *globally* in our setting.

Assumption 1 (Convexity of the domain). Ω is convex.

Assumption 2 (Strong-concavity). $Q(\cdot|\mathbf{w}^*)$ is a λ -strongly concave function over a r -neighborhood of \mathbf{w}^* , i.e. $\mathcal{B}(\mathbf{w}^*, r) \triangleq \{\mathbf{w} \in \Omega : \|\mathbf{w} - \mathbf{w}^*\| \leq r\}$.

Remark 1. An important point to note is that the true parameter \mathbf{w}^* is a fixed point for the EM algorithm, i.e. $M(\mathbf{w}^*) = \mathbf{w}^*$. This is also known as *self-consistency* of the EM algorithm. Hence it is reasonable to expect that in a sufficiently small neighborhood around \mathbf{w}^* there exists a unique maximizer for $Q(\cdot|\mathbf{w}^*)$.

Assumption 3 (First-order stability condition). Assume that

$$\|\nabla Q(M(\mathbf{w})|\mathbf{w}^*) - \nabla Q(M(\mathbf{w})|\mathbf{w})\| \leq \gamma \|\mathbf{w} - \mathbf{w}^*\|, \quad \forall \mathbf{w} \in \mathcal{B}(\mathbf{w}^*, r).$$

Remark 2. Intuitively, the gradient stability condition enforces the gradient maps $\nabla Q(\cdot|\mathbf{w})$ and $\nabla Q(\cdot|\mathbf{w}^*)$ to be close whenever \mathbf{w} lies in a neighborhood of \mathbf{w}^* . This will ensure that the mapped output $M(\mathbf{w})$ stays closer to \mathbf{w}^* .

Theorem 4 (Theorem 1, (Balakrishnan et al., 2017)). *If the above assumptions are met for some radius $r > 0$ and $0 \leq \gamma < \lambda$, then the map $\mathbf{w} \mapsto M(\mathbf{w})$ is contractive over $\mathcal{B}(\mathbf{w}^*, r)$, i.e.*

$$\|M(\mathbf{w}) - \mathbf{w}^*\| \leq \left(\frac{\gamma}{\lambda}\right) \|\mathbf{w} - \mathbf{w}^*\|, \quad \forall \mathbf{w} \in \mathcal{B}(\mathbf{w}^*, r),$$

and consequently, the EM iterates $\{\mathbf{w}_t\}_{t \geq 0}$ converge geometrically to \mathbf{w}^* , i.e.

$$\|\mathbf{w}_t - \mathbf{w}^*\| \leq \left(\frac{\gamma}{\lambda}\right)^t \|\mathbf{w}_0 - \mathbf{w}^*\|,$$

whenever the initialization $\mathbf{w}_0 \in \mathcal{B}(\mathbf{w}^*, r)$.

C. Class of non-linearities

In this section, we characterize the class of non-linearities for which our theoretical results for the recovery of regressors hold. Let $Z \sim \mathcal{N}(0, 1)$ and $Y|Z \sim \mathcal{N}(g(Z), \sigma^2)$, where $g : \mathbb{R} \rightarrow \mathbb{R}$. For $(\alpha, \beta, \gamma) \in \mathbb{R}^3$, define

$$\mathcal{P}_3(y) \triangleq Y^3 + \alpha Y^2 + \beta Y, \quad \mathcal{S}_3(Z) = \mathbb{E}[\mathcal{P}_3(y)|Z] = g(Z)^3 + \alpha g(Z)^2 + g(Z)(\beta + 3\sigma^2) + \alpha\sigma^2,$$

and

$$\mathcal{S}_2(Y) \triangleq Y^2 + \gamma Y, \quad \mathcal{S}_2(Z) = \mathbb{E}[\mathcal{S}_2(Y)|Z] = g(Z)^2 + \gamma g(Z) + \sigma^2.$$

Condition 1. $\mathbb{E}[\mathcal{S}'_3(Z)] = \mathbb{E}[\mathcal{S}''_3(Z)] = 0$ and $\mathbb{E}[\mathcal{S}'''_3(Z)] \neq 0$.

Condition 2. $\mathbb{E}[\mathcal{S}'_2(Z)] = 0$ and $\mathbb{E}[\mathcal{S}''_2(Z)] \neq 0$.

We are now ready to define the (α, β, γ) -valid class of non-linearities.

Definition 1. *We say that the non-linearity g is (α, β, γ) -valid if there exists $(\alpha, \beta, \gamma) \in \mathbb{R}^3$ such that both Condition 1 and Condition 2 are satisfied.*

We have that

$$\begin{aligned} \mathcal{S}'_3(Z) &= 3g(Z)^2 g'(Z) + 2\alpha g(Z) g'(Z) + g'(Z)(\beta + 3\sigma^2) \\ &= 2\alpha g(Z) g'(Z) + \beta g'(Z) + 3g(Z)^2 g'(Z) + 3g'(Z)\sigma^2, \\ \mathcal{S}''_3(Z) &= 2\alpha (g'(Z)^2 + g(Z)g''(Z)) + \beta g''(Z) + 3g''(Z)(g(Z)^2 + \sigma^2) + 6g(Z)g'(Z)^2. \end{aligned}$$

Thus $\mathbb{E}[\mathcal{S}'_3(Z)] = \mathbb{E}[\mathcal{S}''_3(Z)] = 0$ implies that

$$\begin{bmatrix} 2\mathbb{E}(g(Z)g'(Z)) & \mathbb{E}(g'(Z)) \\ 2\mathbb{E}(g'(Z)^2 + g(Z)g''(Z)) & \mathbb{E}(g''(Z)) \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} -3\mathbb{E}(g(Z)^2 g'(Z) + g'(Z)\sigma^2) \\ -3\mathbb{E}(g''(Z)(g(Z)^2 + \sigma^2) + 2g(Z)g'(Z)^2) \end{bmatrix}$$

To ensure Condition 1, we need the pair (α, β) obtained by solving the above linear equation to satisfy $\mathbb{E}[\mathcal{S}'''_3(Z)] \neq 0$. Similarly, $\mathbb{E}[\mathcal{S}'_2(Z)] = 0$ implies that

$$\gamma = \frac{-2\mathbb{E}[g(Z)g'(Z)]}{\mathbb{E}[g'(Z)]}.$$

Thus Condition 2 stipulates that $\mathbb{E}[\mathcal{S}_2''(Z)] \neq 0$ with this choice of γ . It turns out that these conditions hold for a wide class of non-linearities and in particular, when g is either the identity function, or the sigmoid function, or the ReLU. For these three choices of popular non-linearities, the values of the tuple (α, β, γ) are provided below (which are obtained by solving the linear equations mentioned above).

Example 1. If g is the identity mapping, then $\mathcal{P}_3(y) = y^3 - 3y(1 + \sigma^2)$ and $\mathcal{S}_2(y) = y^2$.

Example 2. If g is the sigmoid function, i.e. $g(z) = \frac{1}{1+e^{-z}}$, then α and β can be obtained by solving the following linear equation:

$$\begin{bmatrix} 0.2066 & 0.2066 \\ 0.0624 & -0.0001 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} -0.1755 - 0.6199\sigma^2 \\ -0.0936 \end{bmatrix}$$

The second-order transformation is given by $\mathcal{S}_2(y) = y^2 - y$ (since $\gamma = -1$ when g is sigmoid).

Example 3. If g is the ReLU function, i.e. $g(z) = \max\{0, z\}$, then $\alpha = -3\sqrt{\frac{2}{\pi}}$, $\beta = 3\left(\frac{4}{\pi} - \sigma^2 - 1\right)$ and $\gamma = -2\sqrt{\frac{2}{\pi}}$.

D. Proofs of Section 3

In this section, for the simplicity of the notation we denote the true parameters as \mathbf{w}_i 's and \mathbf{a}_i 's dropping the $*$ sign.

D.1. Proof of Theorem 1 for $k = 2$

Proof. Suppose that g is the linear activation function. For $k = 2$, (1) implies that

$$P_{y|\mathbf{x}} = f(\mathbf{w}^\top \mathbf{x}) \cdot \mathcal{N}(y|\mathbf{a}_1^\top \mathbf{x}, \sigma^2) + (1 - f(\mathbf{w}^\top \mathbf{x})) \cdot \mathcal{N}(y|\mathbf{a}_2^\top \mathbf{x}, \sigma^2), \quad \mathbf{x} \sim \mathcal{N}(0, I_d), \quad (13)$$

where $f(\cdot)$ is the sigmoid function. Using the fact $\mathbb{E}[Z^3] = \mu^3 + 3\mu\sigma^2$ for any Gaussian random variable $Z \sim \mathcal{N}(\mu, \sigma^2)$, we get

$$\mathbb{E}[y^3|\mathbf{x}] = f(\mathbf{w}^\top \mathbf{x})((\mathbf{a}_1^\top \mathbf{x})^3 + 3(\mathbf{a}_1^\top \mathbf{x})\sigma^2) + (1 - f(\mathbf{w}^\top \mathbf{x}))((\mathbf{a}_2^\top \mathbf{x})^3 + 3(\mathbf{a}_2^\top \mathbf{x})\sigma^2).$$

Moreover,

$$\mathbb{E}[y|\mathbf{x}] = f(\mathbf{w}^\top \mathbf{x})(\mathbf{a}_1^\top \mathbf{x}) + (1 - f(\mathbf{w}^\top \mathbf{x}))(\mathbf{a}_2^\top \mathbf{x}).$$

Thus,

$$\mathbb{E}[y^3 - 3y(1 + \sigma^2)|\mathbf{x}] = f(\mathbf{w}^\top \mathbf{x})((\mathbf{a}_1^\top \mathbf{x})^3 - 3(\mathbf{a}_1^\top \mathbf{x})) + (1 - f(\mathbf{w}^\top \mathbf{x}))((\mathbf{a}_2^\top \mathbf{x})^3 - 3(\mathbf{a}_2^\top \mathbf{x})).$$

If we define $\mathcal{P}_3(y) \triangleq y^3 - 3y(1 + \sigma^2)$, in view of Lemma 2 we get that

$$\begin{aligned} \mathcal{T}_3 &= \mathbb{E}[\mathcal{P}_3(y) \cdot \mathcal{S}_3(\mathbf{x})] = \mathbb{E}[(y^3 - 3y(1 + \sigma^2)) \cdot \mathcal{S}_3(\mathbf{x})] \\ &= \mathbb{E}[(f(\mathbf{w}^\top \mathbf{x})((\mathbf{a}_1^\top \mathbf{x})^3 - 3(\mathbf{a}_1^\top \mathbf{x}))) \cdot \mathcal{S}_3(\mathbf{x})] + \mathbb{E}[(1 - f(\mathbf{w}^\top \mathbf{x})((\mathbf{a}_2^\top \mathbf{x})^3 - 3(\mathbf{a}_2^\top \mathbf{x}))) \cdot \mathcal{S}_3(\mathbf{x})] \\ &= \mathbb{E}\left[\nabla_{\mathbf{x}}^{(3)}(f(\mathbf{w}^\top \mathbf{x})((\mathbf{a}_1^\top \mathbf{x})^3 - 3(\mathbf{a}_1^\top \mathbf{x})))\right] + \mathbb{E}\left[\nabla_{\mathbf{x}}^{(3)}(1 - f(\mathbf{w}^\top \mathbf{x})((\mathbf{a}_2^\top \mathbf{x})^3 - 3(\mathbf{a}_2^\top \mathbf{x})))\right]. \end{aligned} \quad (14)$$

Using the chain rule for multi-derivatives, the first term simplifies to

$$\begin{aligned} \mathbb{E}\left[\nabla_{\mathbf{x}}^{(3)}(f(\mathbf{w}^\top \mathbf{x})((\mathbf{a}_1^\top \mathbf{x})^3 - 3(\mathbf{a}_1^\top \mathbf{x})))\right] &= \mathbb{E}[f'''((\mathbf{a}_1^\top \mathbf{x})^3 - 3(\mathbf{a}_1^\top \mathbf{x}))] \cdot \mathbf{w} \otimes \mathbf{w} \otimes \mathbf{w} + \mathbb{E}[f''(3(\mathbf{a}_1^\top \mathbf{x})^2 - 3)] \cdot \\ &\quad (\mathbf{w} \otimes \mathbf{w} \otimes \mathbf{a}_1 + \mathbf{w} \otimes \mathbf{a}_1 \otimes \mathbf{w} + \mathbf{a}_1 \otimes \mathbf{w} \otimes \mathbf{w}) + \\ &\quad \mathbb{E}[f'(6(\mathbf{a}_1^\top \mathbf{x}))] \cdot (\mathbf{a}_1 \otimes \mathbf{a}_1 \otimes \mathbf{w} + \mathbf{a}_1 \otimes \mathbf{w} \otimes \mathbf{a}_1 + \mathbf{w} \otimes \mathbf{a}_1 \otimes \mathbf{a}_1) + 6\mathbb{E}[f] \cdot \mathbf{a}_1 \otimes \mathbf{a}_1 \otimes \mathbf{a}_1. \end{aligned} \quad (15)$$

Since $f(z) = \frac{1}{1+e^{-z}}$, $f'(\cdot)$, $f'''(\cdot)$ are even functions whereas $f''(\cdot)$ is an odd function. Furthermore, both $\mathbf{x} \mapsto (\mathbf{a}_1^\top \mathbf{x})^3 - 3(\mathbf{a}_1^\top \mathbf{x})$ and $\mathbf{x} \mapsto \mathbf{a}_1^\top \mathbf{x}$ are odd functions whereas $\mathbf{x} \mapsto 3(\mathbf{a}_1^\top \mathbf{x})^2 - 3$ is an even function. Since $\mathbf{x} \sim \mathcal{N}(0, I_d)$, $-\mathbf{x} \stackrel{(d)}{=} \mathbf{x}$. Thus all the expectation terms in (15) equal zero except for the last term since $\mathbb{E}[f(\mathbf{w}^\top \mathbf{x})] = \frac{1}{2} > 0$. We have,

$$\mathbb{E}\left[\nabla_{\mathbf{x}}^{(3)}(f(\mathbf{w}^\top \mathbf{x})((\mathbf{a}_1^\top \mathbf{x})^3 - 3(\mathbf{a}_1^\top \mathbf{x})))\right] = 3 \cdot \mathbf{a}_1 \otimes \mathbf{a}_1 \otimes \mathbf{a}_1.$$

Similarly,

$$\mathbb{E} \left[\nabla_{\mathbf{x}}^{(3)} (1 - f(\mathbf{w}^\top \mathbf{x})) ((\mathbf{a}_2^\top \mathbf{x})^3 - 3(\mathbf{a}_2^\top \mathbf{x})) \right] = 3 \cdot \mathbf{a}_2 \otimes \mathbf{a}_2 \otimes \mathbf{a}_2.$$

Together, we have that

$$\mathcal{T}_3 = 3 \cdot \mathbf{a}_1 \otimes \mathbf{a}_1 \otimes \mathbf{a}_1 + 3 \cdot \mathbf{a}_2 \otimes \mathbf{a}_2 \otimes \mathbf{a}_2.$$

Now consider an arbitrary link function g belonging to the class of non-linearities described in Appendix C. Then

$$P_{y|\mathbf{x}} = f(\mathbf{w}^\top \mathbf{x}) \cdot \mathcal{N}(y|g(\mathbf{a}_1^\top \mathbf{x}), \sigma^2) + (1 - f(\mathbf{w}^\top \mathbf{x})) \cdot \mathcal{N}(y|g(\mathbf{a}_2^\top \mathbf{x}), \sigma^2), \quad \mathbf{x} \sim \mathcal{N}(0, I_d),$$

implies that

$$\mathbb{E}[y^3|\mathbf{x}] = f(\mathbf{w}^\top \mathbf{x})(g(\mathbf{a}_1^\top \mathbf{x})^3 + 3g(\mathbf{a}_1^\top \mathbf{x})\sigma^2) + (1 - f(\mathbf{w}^\top \mathbf{x}))(g(\mathbf{a}_2^\top \mathbf{x})^3 + 3g(\mathbf{a}_2^\top \mathbf{x})\sigma^2),$$

and

$$\begin{aligned} \mathbb{E}[y^2|\mathbf{x}] &= f(\mathbf{w}^\top \mathbf{x})(g(\mathbf{a}_1^\top \mathbf{x})^2 + \sigma^2) + (1 - f(\mathbf{w}^\top \mathbf{x}))(g(\mathbf{a}_2^\top \mathbf{x})^2 + \sigma^2), \\ \mathbb{E}[y|\mathbf{x}] &= f(\mathbf{w}^\top \mathbf{x})g(\mathbf{a}_1^\top \mathbf{x}) + (1 - f(\mathbf{w}^\top \mathbf{x}))g(\mathbf{a}_2^\top \mathbf{x}). \end{aligned}$$

If we define $\mathcal{P}_3(y) \triangleq y^3 + \alpha y^2 + \beta y$, we have that

$$\begin{aligned} \mathcal{T}_3 &= \mathbb{E}[\mathcal{P}_3(y) \cdot \mathcal{S}_3(\mathbf{x})] = \mathbb{E}[\mathbb{E}[y^3 + \alpha y^2 + \beta y|\mathbf{x}] \cdot \mathcal{S}_3(\mathbf{x})] \\ &= \mathbb{E} [f(\mathbf{w}^\top \mathbf{x}) (g(\mathbf{a}_1^\top \mathbf{x})^3 + \alpha g(\mathbf{a}_1^\top \mathbf{x})^2 + g(\mathbf{a}_1^\top \mathbf{x})(\beta + 3\sigma^2)) \cdot \mathcal{S}_3(\mathbf{x})] + \\ &\quad \mathbb{E} [(1 - f(\mathbf{w}^\top \mathbf{x})) (g(\mathbf{a}_2^\top \mathbf{x})^3 + \alpha g(\mathbf{a}_2^\top \mathbf{x})^2 + g(\mathbf{a}_2^\top \mathbf{x})(\beta + 3\sigma^2)) \cdot \mathcal{S}_3(\mathbf{x})] \\ &= \mathbb{E} \left[\nabla_{\mathbf{x}}^{(3)} (f(\mathbf{w}^\top \mathbf{x}) (g(\mathbf{a}_1^\top \mathbf{x})^3 + \alpha g(\mathbf{a}_1^\top \mathbf{x})^2 + g(\mathbf{a}_1^\top \mathbf{x})(\beta + 3\sigma^2))) \right] + \\ &\quad \mathbb{E} \left[\nabla_{\mathbf{x}}^{(3)} (f(\mathbf{w}^\top \mathbf{x}) (g(\mathbf{a}_2^\top \mathbf{x})^3 + \alpha g(\mathbf{a}_2^\top \mathbf{x})^2 + g(\mathbf{a}_2^\top \mathbf{x})(\beta + 3\sigma^2))) \right] \\ &\stackrel{(a)}{=} \mathbb{E}[f]\mathbb{E} \left[\nabla_{\mathbf{x}}^{(3)} (g(\mathbf{a}_1^\top \mathbf{x})^3 + \alpha g(\mathbf{a}_1^\top \mathbf{x})^2 + g(\mathbf{a}_1^\top \mathbf{x})(\beta + 3\sigma^2)) \right] \cdot \mathbf{a}_1 \otimes \mathbf{a}_1 \otimes \mathbf{a}_1 + \\ &\quad \mathbb{E}[1 - f]\mathbb{E} \left[\nabla_{\mathbf{x}}^{(3)} (g(\mathbf{a}_2^\top \mathbf{x})^3 + \alpha g(\mathbf{a}_2^\top \mathbf{x})^2 + g(\mathbf{a}_2^\top \mathbf{x})(\beta + 3\sigma^2)) \right] \cdot \mathbf{a}_2 \otimes \mathbf{a}_2 \otimes \mathbf{a}_2 \\ &= c_{g,\sigma} (\mathbb{E}[f] \cdot \mathbf{a}_1 \otimes \mathbf{a}_1 \otimes \mathbf{a}_1 + \mathbb{E}[1 - f] \cdot \mathbf{a}_2 \otimes \mathbf{a}_2 \otimes \mathbf{a}_2), \end{aligned}$$

where (a) follows from the choice of α and β and the fact that $\mathbf{w} \perp \{\mathbf{a}_1, \mathbf{a}_2\}$, and $c_{g,\sigma} \triangleq \mathbb{E} \left[(g(Z)^3 + \alpha g(Z)^2 + g(Z)(\beta + 3\sigma^2))''' \right]$ where $Z \sim \mathcal{N}(0, 1)$. The proof for \mathcal{T}_2 is similar. \square

D.2. Proof of Theorem 1 for general k

Proof. The proof for general k closely follows that of $k = 2$, described in Appendix D.1. For the general k , we first prove the theorem when g is the identity function, i.e.

$$P_{y|\mathbf{x}} = \sum_{i \in [k]} P_{i|\mathbf{x}} P_{y|\mathbf{x},i} = \sum_{i \in [k]} \frac{e^{\mathbf{w}_i^\top \mathbf{x}}}{\sum_{i \in [k]} e^{\mathbf{w}_i^\top \mathbf{x}}} \cdot \mathcal{N}(y|\mathbf{a}_i^\top \mathbf{x}, \sigma^2), \quad \mathbf{x} \sim \mathcal{N}(0, I_d).$$

Denoting $P_{i|\mathbf{x}}$ by $p_i(\mathbf{x})$, we have that

$$\begin{aligned} \mathbb{E}[y^3|\mathbf{x}] &= \sum_{i \in [k]} p_i(\mathbf{x}) ((\mathbf{a}_i^\top \mathbf{x})^3 + 3(\mathbf{a}_i^\top \mathbf{x})\sigma^2), \\ \mathbb{E}[y|\mathbf{x}] &= \sum_{i \in [k]} p_i(\mathbf{x}) (\mathbf{a}_i^\top \mathbf{x}). \end{aligned}$$

Hence

$$\mathbb{E}[y^3 - 3y(1 + \sigma^2)|\mathbf{x}] = \sum_{i \in [k]} p_i(\mathbf{x}) ((\mathbf{a}_i^\top \mathbf{x})^3 - 3(\mathbf{a}_i^\top \mathbf{x}))$$

If we let $\mathcal{P}_3(y) \triangleq y^3 - 3y(1 + \sigma^2)$, we get

$$\mathbb{E}[\mathcal{P}_3(y) \cdot \mathcal{S}_3(\mathbf{x})] = \sum_{i \in [k]} \mathbb{E} \left[\nabla_{\mathbf{x}}^{(3)} (p_i(\mathbf{x}) ((\mathbf{a}_i^\top \mathbf{x})^3 - 3(\mathbf{a}_i^\top \mathbf{x}))) \right]$$

Since $\mathbf{x} \sim \mathcal{N}(0, I_d)$ and $\mathbf{a}_i \perp \text{span}\{\mathbf{w}_1, \dots, \mathbf{w}_{k-1}\}$, we have that $\mathbf{a}_i^\top \mathbf{x} \perp (\mathbf{w}_1^\top \mathbf{x}, \dots, \mathbf{w}_{k-1}^\top \mathbf{x})$. Moreover, $\mathbb{E}[(\mathbf{a}_i^\top \mathbf{x})^3 - 3(\mathbf{a}_i^\top \mathbf{x})] = \mathbb{E}[(\mathbf{a}_i^\top \mathbf{x})^2 - 1] = \mathbb{E}[\mathbf{a}_i^\top \mathbf{x}] = 0$ for each i . Using the chain-rule for multi-derivatives, the above equation thus simplifies to

$$\mathbb{E}[\mathcal{P}_3(y) \cdot \mathcal{S}_3(\mathbf{x})] = \sum_{i \in [k]} \mathbb{E}[p_i(\mathbf{x})] \cdot \mathbb{E} \left[\nabla_{\mathbf{x}}^{(3)} ((\mathbf{a}_i^\top \mathbf{x})^3 - 3(\mathbf{a}_i^\top \mathbf{x})) \right] = \sum_{i \in [k]} 6\mathbb{E}[p_i(\mathbf{x})] \cdot \mathbf{a}_i \otimes \mathbf{a}_i \times \mathbf{a}_i.$$

For a generic $g : \mathbb{R} \rightarrow \mathbb{R}$ which is (α, β, γ) -valid, let $\mathcal{P}_3(y) = y^3 + \alpha y^2 + \beta y$. Then it is easy to see that the same proof goes through except for a change in the coefficients of rank-1 terms, i.e.

$$\mathbb{E}[\mathcal{P}_3(y) \cdot \mathcal{S}_3(\mathbf{x})] = \sum_{i \in [k]} \alpha_i \mathbb{E}[p_i(\mathbf{x})] \cdot \mathbf{a}_i \otimes \mathbf{a}_i \otimes \mathbf{a}_i,$$

where $\alpha_i \triangleq \mathbb{E} \left[(g(Z)^3 + \alpha g(Z)^2 + g(Z)(\beta + 3\sigma^2))''' \right]$ where $Z \sim \mathcal{N}(0, 1)$ and $'''$ denotes the third-derivative with respect to Z . Note that Condition 2 together with the fact that $\mathbb{E}[p_i(\mathbf{x})] > 0$ ensures that $\alpha_i \neq 0$ and thus the coefficients of the rank-1 terms are non-zero. The proof for \mathcal{T}_2 is similar. \square

D.3. Proof of Theorem 2

The following two lemmas are central to the proof of Theorem 2. Let $\mathbf{A}^\top = [\mathbf{a}_1 | \dots | \mathbf{a}_k] \in \mathbb{R}^{d \times k}$ denote the matrix of regressor parameters whereas $\mathbf{W}^\top = [\mathbf{w}_1 | \dots | \mathbf{w}_{k-1}] \in \mathbb{R}^{d \times (k-1)}$ denote the matrix of gating parameters. With a slight change of notation, when $\mathbf{A} = \mathbf{A}^*$, we denote the EM operator $M(\mathbf{W})$ as either $M(\mathbf{W}, \mathbf{A}^*)$ or $M(\mathbf{w})$, introduced in Section 3. For the general case, we simply denote it by $M(\mathbf{W}, \mathbf{A})$. In the following lemmas, we use the norm $\|\mathbf{A}\| = \max_{i \in [k]} \|\mathbf{A}_i^\top\|_2$ where $\mathbf{A} \in \mathbb{R}^{k \times d}$ is a matrix of regressors, similarly for any matrix of classifiers $\mathbf{W} \in \mathbb{R}^{(k-1) \times d}$.

Lemma 3 (Contraction of the EM operator). *Under the assumptions of Theorem 2, we have that*

$$\|M(\mathbf{W}, \mathbf{A}^*) - \mathbf{W}^*\| \leq \kappa_\sigma \|\mathbf{W} - \mathbf{W}^*\|.$$

Moreover, $\mathbf{W} = \mathbf{W}^*$ is a fixed point for $M(\mathbf{W}, \mathbf{A}^*)$.

Lemma 4 (Robustness of the EM operator). *Let the matrix of regressors \mathbf{A} be such that $\max_{i \in [k]} \|\mathbf{A}_i^\top - (\mathbf{A}_i^*)^\top\|_2 = \sigma^2 \varepsilon_1$. Then for any $\mathbf{W} \in \Omega$, we have that*

$$\|M(\mathbf{W}, \mathbf{A}) - M(\mathbf{W}, \mathbf{A}^*)\| \leq \kappa \varepsilon_1,$$

where κ is a constant depending on g, k and σ . In particular, $\kappa \leq (k-1) \frac{\sqrt{6(2+\sigma^2)}}{2}$ for $g = \text{linear, sigmoid and ReLU}$.

We are now ready to prove Theorem 2.

Proof. We first note that the EM iterates $\{\mathbf{W}_t\}_{t \geq 1}$ evolve according to

$$\mathbf{W}_t = M(\mathbf{W}_{t-1}, \mathbf{A}), \quad t \geq 1$$

Thus

$$\begin{aligned} \|\mathbf{W}_t - \mathbf{W}^*\| &= \|M(\mathbf{W}_{t-1}, \mathbf{A}) - \mathbf{W}^*\| = \|M(\mathbf{W}_{t-1}, \mathbf{A}) - M(\mathbf{W}_{t-1}, \mathbf{A}^*)\| \\ &\leq \|M(\mathbf{W}_{t-1}, \mathbf{A}) - M(\mathbf{W}_{t-1}, \mathbf{A}^*)\| + \|M(\mathbf{W}_{t-1}, \mathbf{A}^*) - \mathbf{W}^*\| \\ &\leq \kappa \varepsilon_1 + \kappa_\sigma \|\mathbf{W}_{t-1} - \mathbf{W}^*\|, \end{aligned}$$

where the last inequality follows from Lemma 3 and Lemma 4. Recursively using the above inequality, we obtain that

$$\|\mathbf{W}_t - \mathbf{W}^*\| \leq (\kappa_\sigma)^t \|\mathbf{W}_0 - \mathbf{W}^*\| + \kappa_\sigma \varepsilon_1 (1 + \kappa_\sigma + \dots + \kappa_\sigma^{t-1}) \leq (\kappa_\sigma)^t \|\mathbf{W}_0 - \mathbf{W}^*\| + \frac{\kappa_\sigma \varepsilon_1}{1 - \kappa_\sigma}.$$

□

D.4. Proof of Theorem 3

Proof. We are given that $(\mathbf{a}_1, \mathbf{a}_2) = (\mathbf{a}_1^*, \mathbf{a}_2^*)$. Denoting \mathbf{w}^* with \mathbf{w} , from (13), we have that

$$\mathbb{E}[y|\mathbf{x}] = f(\mathbf{w}^\top \mathbf{x}) \cdot \mathbf{a}_1^\top \mathbf{x} + (1 - f(\mathbf{w}^\top \mathbf{x})) \cdot \mathbf{a}_2^\top \mathbf{x}, \quad (16)$$

$$= \mathbf{a}_2^\top \mathbf{x} + f(\mathbf{w}^\top \mathbf{x}) \cdot (\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}. \quad (17)$$

Thus,

$$\frac{\mathbb{E}[y|\mathbf{x}] - \mathbf{a}_2^\top \mathbf{x}}{(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}} = f(\mathbf{w}^\top \mathbf{x}).$$

Notice that in the above equation we have $(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}$ in the denominator. But this equals zero with zero probability whenever \mathbf{x} is generated from a continuous distribution; in our case \mathbf{x} is Gaussian. Thus we may write

$$\begin{aligned} \mathbb{E} \left[\left(\frac{y - \mathbf{a}_2^\top \mathbf{x}}{(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}} \right) \cdot \mathbf{x} \right] &\stackrel{\text{X}}{=} \mathbb{E} \left[\left(\frac{\mathbb{E}[y|\mathbf{x}] - \mathbf{a}_2^\top \mathbf{x}}{(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}} \right) \cdot \mathbf{x} \right] = \mathbb{E} [f(\mathbf{w}^\top \mathbf{x}) \cdot \mathbf{x}] \\ &= \mathbb{E} [f'(\mathbf{w}^\top \mathbf{x})] \cdot \mathbf{w} \\ &= \mathbb{E}_{Z \sim \mathcal{N}(0,1)} f'(\|\mathbf{w}\| Z) \cdot \mathbf{w} \\ &\propto \mathbf{w}. \end{aligned}$$

However, it turns out that the above chain of equalities does not hold. Surprisingly, the first equality, which essentially is the law of iterated expectations, is not valid in this case as $\frac{y - \mathbf{a}_2^\top \mathbf{x}}{(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}}$ is not integrable. To see this, notice that the model in (13) can also be written as

$$y \stackrel{(d)}{=} Z(\mathbf{a}_1^\top \mathbf{x}) + (1 - Z)(\mathbf{a}_2^\top \mathbf{x}) + \sigma N, \quad Z \sim \text{Bern}(f(\mathbf{w}^\top \mathbf{x})), N \sim \mathcal{N}(0, 1).$$

Thus,

$$\text{Ratio} \triangleq \frac{y - \mathbf{a}_2^\top \mathbf{x}}{(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}} \stackrel{(d)}{=} Z + \frac{\sigma N}{(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}}.$$

Since Z is independent of N and $\frac{N}{(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}}$ is a Cauchy random variable, it follows that the random variable Ratio is not integrable. To deal with the non-integrability of Ratio, we look at its conditional cdf, given by

$$\mathbb{P}[\text{Ratio} \leq z|\mathbf{x}] = f(\mathbf{w}^\top \mathbf{x}) \Phi \left((z - 1) \frac{|\Delta_x|}{\sigma} \right) + (1 - f(\mathbf{w}^\top \mathbf{x})) \Phi \left(z \frac{|\Delta_x|}{\sigma} \right), \quad \Delta_x = (\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x},$$

where $\Phi(\cdot)$ is the standard Gaussian cdf. Substituting $z = 0.5$ and using the fact that $\Phi(z) + \Phi(-z) = 1$, we obtain

$$\begin{aligned} \mathbb{P}[\text{Ratio} \leq 0.5|\mathbf{x}] &= f(\mathbf{w}^\top \mathbf{x}) \Phi \left(-\frac{|\Delta_x|}{2\sigma} \right) + (1 - f(\mathbf{w}^\top \mathbf{x})) \Phi \left(\frac{|\Delta_x|}{2\sigma} \right) \\ &= \Phi \left(\frac{|(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}|}{2\sigma} \right) + f(\mathbf{w}^\top \mathbf{x}) \left(1 - 2\Phi \left(\frac{|(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}|}{2\sigma} \right) \right). \end{aligned}$$

Since $\Phi\left(\frac{|(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}|}{2\sigma}\right)$ is a symmetric function in \mathbf{x} its first moment with \mathbf{x} equals zero. Furthermore, if we assume that \mathbf{w} is orthogonal to \mathbf{a}_1 and \mathbf{a}_2 , we have

$$\begin{aligned}
 \mathbb{E}[\mathbb{1}\{\text{Ratio} \leq 0.5\} \cdot \mathbf{x}] &= \mathbb{E}[\mathbb{P}[\text{Ratio} \leq 0.5 | \mathbf{x}] \cdot \mathbf{x}] \\
 &= \mathbb{E}\left[f(\mathbf{w}^\top \mathbf{x}) \left(1 - 2\Phi\left(\frac{|(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}|}{2\sigma}\right)\right) \cdot \mathbf{x}\right] \\
 &= \mathbb{E}[f'(\mathbf{w}^\top \mathbf{x})] \cdot \mathbb{E}\left(1 - 2\Phi\left(\frac{|(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}|}{2\sigma}\right)\right) \cdot \mathbf{w} + \\
 &\quad \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})] \cdot \underbrace{\mathbb{E}\left[\nabla_{\mathbf{x}}\left(1 - 2\Phi\left(\frac{|(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}|}{2\sigma}\right)\right)\right]}_{=0, \text{ since derivative of a even function is odd}} \\
 &= \mathbb{E}[f'(\mathbf{w}^\top \mathbf{x})] \cdot \mathbb{E}\left(1 - 2\Phi\left(\frac{|(\mathbf{a}_1 - \mathbf{a}_2)^\top \mathbf{x}|}{2\sigma}\right)\right) \cdot \mathbf{w} \\
 &\propto \mathbf{w}.
 \end{aligned}$$

Thus, if $\|\mathbf{w}\| = 1$, we have that

$$\frac{\mathbb{E}[\mathbb{1}\{\text{Ratio} \leq 0.5\} \cdot \mathbf{x}]}{\|\mathbb{E}[\mathbb{1}\{\text{Ratio} \leq 0.5\} \cdot \mathbf{x}]\|} = \mathbf{w}.$$

In the finite sample regime, $\mathbb{E}[\mathbb{1}\{\text{Ratio} \leq 0.5\} \cdot \mathbf{x}]$ can be estimated from samples using the empirical moments and its normalized version will be an estimate of \mathbf{w} . \square

E. Proof of Lemma 4

We need the following lemma which establishes the stability of the minimizers for strongly convex functions under Lipschitz perturbations.

Lemma 5. *Suppose $\Omega \subseteq \mathbb{R}^d$ is a closed convex subset, $f : \Omega \rightarrow \mathbb{R}$ is a λ -strongly convex function for some $\lambda > 0$ and B is an L -Lipschitz continuous function on Ω . Let $\mathbf{w}_f = \operatorname{argmin}_{\mathbf{w} \in \Omega} f(\mathbf{w})$ and $\mathbf{w}_{f+B} = \operatorname{argmin}_{\mathbf{w} \in \Omega} f(\mathbf{w}) + B(\mathbf{w})$. Then*

$$\|\mathbf{w}_f - \mathbf{w}_{f+B}\| \leq \frac{L}{\lambda}.$$

Proof. Let $\mathbf{w}' \in \Omega$ be such that $\|\mathbf{w}' - \mathbf{w}_f\| > \frac{L}{\lambda}$. Let $\mathbf{w}_\alpha = \alpha \mathbf{w}_f + (1 - \alpha) \mathbf{w}'$ for $0 < \alpha < 1$. From the fact that \mathbf{w}_f is the minimizer of f on Ω and that f is strongly convex, we have that

$$f(\mathbf{w}') \geq f(\mathbf{w}_f) + \frac{\lambda \|\mathbf{w}' - \mathbf{w}_f\|^2}{2}.$$

Furthermore, the strong-convexity of f implies that

$$\begin{aligned}
 f(\mathbf{w}_\alpha) &\leq \alpha f(\mathbf{w}_f) + (1 - \alpha) f(\mathbf{w}') - \frac{\alpha(1 - \alpha)\lambda}{2} \|\mathbf{w}' - \mathbf{w}_f\|^2 \\
 &= f(\mathbf{w}') + \alpha(f(\mathbf{w}_f) - f(\mathbf{w}')) - \frac{\alpha(1 - \alpha)\lambda}{2} \|\mathbf{w}' - \mathbf{w}_f\|^2 \\
 &\leq f(\mathbf{w}') - \alpha \frac{\lambda \|\mathbf{w}' - \mathbf{w}_f\|^2}{2} - \frac{\alpha(1 - \alpha)\lambda}{2} \|\mathbf{w}' - \mathbf{w}_f\|^2 \\
 &= f(\mathbf{w}') - \lambda \alpha \left(1 - \frac{\alpha}{2}\right) \|\mathbf{w}' - \mathbf{w}_f\|^2
 \end{aligned} \tag{18}$$

Since B is L -Lipschitz, we have

$$B(\mathbf{w}_\alpha) \leq B(\mathbf{w}') + L\alpha \|\mathbf{w}' - \mathbf{w}_f\|. \tag{19}$$

Adding (18) and (19), we get

$$\begin{aligned} f(\mathbf{w}_\alpha) + B(\mathbf{w}_\alpha) &\leq f(\mathbf{w}') + B(\mathbf{w}') + L\alpha \|\mathbf{w}' - \mathbf{w}_f\| - \lambda\alpha \left(1 - \frac{\alpha}{2}\right) \|\mathbf{w}' - \mathbf{w}_f\|^2 \\ &= f(\mathbf{w}') + B(\mathbf{w}') + \alpha\lambda \|\mathbf{w}' - \mathbf{w}_f\| \left(\frac{L}{\lambda} - \left(1 - \frac{\alpha}{2}\right) \|\mathbf{w}' - \mathbf{w}_f\|\right) \end{aligned}$$

By the assumption that $\|\mathbf{w}' - \mathbf{w}_f\| > \frac{L}{\lambda}$, the term $\frac{L}{\lambda} - \left(1 - \frac{\alpha}{2}\right) \|\mathbf{w}' - \mathbf{w}_f\|$ will be negative for sufficiently small α . This in turn implies that $f(\mathbf{w}_\alpha) + B(\mathbf{w}_\alpha) < f(\mathbf{w}') + B(\mathbf{w}')$ for such α . Consequently \mathbf{w}' is not a minimizer of $f + B$ for any \mathbf{w}' such that $\|\mathbf{w}' - \mathbf{w}_f\| > \frac{L}{\lambda}$. The conclusion follows. \square

We are now ready to prove Lemma 4. Fix any $\mathbf{W} \in \Omega$ and let $\mathbf{A} = \begin{bmatrix} \mathbf{a}_1^\top \\ \dots \\ \mathbf{a}_k^\top \end{bmatrix} \in \mathbb{R}^{k \times d}$ be such that $\max_{i \in [k]} \|\mathbf{a}_i - \mathbf{a}_i^*\|_2 = \sigma^2 \varepsilon_1$ for some $\varepsilon_1 > 0$. Let

$$\mathbf{W}' = M(\mathbf{W}, \mathbf{A}), \quad (\mathbf{W}')^* = M(\mathbf{W}, \mathbf{A}^*),$$

where,

$$M(\mathbf{W}, \mathbf{A}) = \arg \max_{\mathbf{W}' \in \Omega} Q(\mathbf{W}' | \mathbf{W}, \mathbf{A}),$$

and,

$$Q(\mathbf{W}' | \mathbf{W}, \mathbf{A}) = \mathbb{E} \left[\sum_{i \in [k-1]} p^{(i)}(\mathbf{W}, \mathbf{A}) ((\mathbf{W}'_i)^\top \mathbf{x}) - \log \left(1 + \sum_{i \in [k-1]} e^{(\mathbf{W}'_i)^\top \mathbf{x}} \right) \right].$$

Here $p^{(i)}(\mathbf{A}, \mathbf{W}) \triangleq \frac{p_i(\mathbf{x}) N_i}{\sum_{i \in [k]} p_i(\mathbf{x}) N_i}$ denotes the posterior probability of choosing the i^{th} expert, where

$$p_i(\mathbf{x}) = \frac{e^{\mathbf{w}_i^\top \mathbf{x}}}{1 + \sum_{k \in [k-1]} e^{\mathbf{w}_k^\top \mathbf{x}}}, \quad N_i \triangleq \mathcal{N}(y | g(\mathbf{a}_i^\top \mathbf{x}), \sigma^2), \quad N_i^* = \mathcal{N}(y | g((\mathbf{a}_i^*)^\top \mathbf{x}), \sigma^2).$$

Since both $Q(\cdot | \mathbf{W}, \mathbf{A})$ and $Q(\cdot | \mathbf{W}, \mathbf{A}^*)$ are strongly concave functions over Ω with some strong-concavity parameter λ , Lemma 5 implies that

$$\|M(\mathbf{W}, \mathbf{A}) - M(\mathbf{W}, \mathbf{A}^*)\| \leq \frac{L}{\lambda},$$

where L is the Lipschitz-constant for the function $l(\cdot) \triangleq Q(\cdot | \mathbf{W}, \mathbf{A}) - Q(\cdot | \mathbf{W}, \mathbf{A}^*)$. We have that

$$l(\mathbf{W}') = \sum_{i \in [k-1]} \mathbb{E}[(p^{(i)}(\mathbf{W}, \mathbf{A}) - p^{(i)}(\mathbf{W}, \mathbf{A}^*)) (\mathbf{W}'_i)^\top \mathbf{x}]$$

Without loss of generality let $i = 1$. Since $l(\cdot)$ is linear in \mathbf{W}' , it suffices to show for each i that

$$\left\| \mathbb{E}[(p^{(1)}(\mathbf{W}, \mathbf{A}) - p^{(1)}(\mathbf{W}, \mathbf{A}^*)) \mathbf{x}] \right\| \leq L,$$

We show that $L = \kappa \varepsilon_1$, or equivalently,

$$\left\| \mathbb{E}[(p^{(1)}(\mathbf{W}, \mathbf{A}) - p^{(1)}(\mathbf{W}, \mathbf{A}^*)) \mathbf{x}] \right\| \leq \kappa \varepsilon_1,$$

Let

$$\mathbf{A}_t = \mathbf{A}^* + t\Delta, \quad \Delta = \mathbf{A} - \mathbf{A}^* \in \mathbb{R}^{k \times d}.$$

By hypothesis, we have that $\|\Delta_i\|_2 \leq \sigma^2 \varepsilon_1$ for all $i \in [k]$. Thus in order to show that

$$\left\| \mathbb{E}[(p^{(1)}(\mathbf{A}, \mathbf{W}) - p^{(1)}(\mathbf{A}^*, \mathbf{W}))\mathbf{x}] \right\|_2 \leq \kappa \varepsilon_1,$$

it suffices to show that

$$\langle \mathbb{E}[(p^{(1)}(\mathbf{A}, \mathbf{W}) - p^{(1)}(\mathbf{A}^*, \mathbf{W}))\mathbf{x}], \tilde{\Delta} \rangle \leq \kappa \|\Delta/\sigma^2\|_2 \|\tilde{\Delta}\|_2, \quad \text{for all } \tilde{\Delta} \in \mathbb{R}^d.$$

Or equivalently,

$$\mathbb{E}[(p^{(1)}(\mathbf{A}, \mathbf{W}) - p^{(1)}(\mathbf{A}^*, \mathbf{W}))\langle \mathbf{x}, \tilde{\Delta} \rangle] \leq \kappa \|\Delta/\sigma^2\|_2 \|\tilde{\Delta}\|_2.$$

We can rewrite the difference of the posteriors as

$$p^{(1)}(\mathbf{A}, \mathbf{W}) - p^{(1)}(\mathbf{A}^*, \mathbf{W}) = \int_0^1 \frac{d}{dt} p^{(1)}(\mathbf{A}^* + t\Delta, \mathbf{W}) dt = \sum_{i \in [k]} \int_0^1 \langle \nabla_{\mathbf{a}_i} p^{(1)}(\mathbf{A}_t, \mathbf{W}), \Delta_i \rangle dt. \quad (20)$$

Since $N_i = \mathcal{N}(y|g(\mathbf{a}_i^\top \mathbf{x}), \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-g(\mathbf{a}_i^\top \mathbf{x}))^2/2\sigma^2}$, we have that

$$\nabla_{\mathbf{a}_i} N_i = N_i \left(\frac{y - g(\mathbf{a}_i^\top \mathbf{x})}{\sigma^2} \right) g'(\mathbf{a}_i^\top \mathbf{x}).$$

Thus,

$$\begin{aligned} \nabla_{\mathbf{a}_i} p^{(1)}(\mathbf{A}_t, \mathbf{W}) &= \nabla_{\mathbf{a}_i} \left(\frac{p_1(\mathbf{x})N_1}{\sum_{i \in [k]} p_i(\mathbf{x})N_i} \right) \\ &= \begin{cases} \frac{(\sum_{i \neq 1} p_i(\mathbf{x})N_i)p_1(\mathbf{x})N_1}{(\sum_i p_i(\mathbf{x})N_i)^2} \left(\frac{y - g(\mathbf{a}_1^\top \mathbf{x})}{\sigma^2} \right) g'(\mathbf{a}_1^\top \mathbf{x}), & \text{if } i = 1 \\ \frac{-p_i(\mathbf{x})p_1(\mathbf{x})N_iN_1}{(\sum_i p_i(\mathbf{x})N_i)^2} \left(\frac{y - g(\mathbf{a}_i^\top \mathbf{x})}{\sigma^2} \right) g'(\mathbf{a}_i^\top \mathbf{x}), & \text{if } i \neq 1 \end{cases} \end{aligned}$$

Hence,

$$\begin{aligned} \mathbb{E}[(p^{(1)}(\mathbf{A}, \mathbf{W}) - p^{(1)}(\mathbf{A}^*, \mathbf{W}))\langle \mathbf{x}, \tilde{\Delta} \rangle] &= \sum_{i \in [k]} \int_0^1 \mathbb{E}[\langle \nabla_{\mathbf{a}_i} p^{(1)}(\mathbf{A}_t, \mathbf{W}), \Delta_i \rangle \langle \mathbf{x}, \tilde{\Delta} \rangle] dt \quad (21) \\ &= \int_0^1 \mathbb{E} \left[\frac{(\sum_{i \neq 1} p_i(\mathbf{x})N_i)p_1(\mathbf{x})N_1}{(\sum_i p_i(\mathbf{x})N_i)^2} \left(\frac{y - g(\mathbf{a}_1^\top \mathbf{x})}{\sigma^2} \right) g'(\mathbf{a}_1^\top \mathbf{x}) \langle \mathbf{x}, \Delta_1 \rangle \langle \mathbf{x}, \tilde{\Delta} \rangle \right] dt \quad (22) \\ &\quad + \sum_{i \neq 1} \int_0^1 \mathbb{E} \left[\frac{-p_i(\mathbf{x})p_1(\mathbf{x})N_iN_1}{(\sum_i p_i(\mathbf{x})N_i)^2} \left(\frac{y - g(\mathbf{a}_i^\top \mathbf{x})}{\sigma^2} \right) g'(\mathbf{a}_i^\top \mathbf{x}) \langle \mathbf{x}, \Delta_i \rangle \langle \mathbf{x}, \tilde{\Delta} \rangle \right] dt, \quad (23) \end{aligned}$$

where we denoted $(\mathbf{a}_i)_t$ by \mathbf{a}_i in the integrals above (with a slight abuse of notation) for the sake of notational simplicity. For any $i \neq 1$, we have that

$$\begin{aligned} &\left| \frac{-p_i(\mathbf{x})p_1(\mathbf{x})N_iN_1}{(\sum_i p_i(\mathbf{x})N_i)^2} \left(\frac{y - g(\mathbf{a}_i^\top \mathbf{x})}{\sigma^2} \right) g'(\mathbf{a}_i^\top \mathbf{x}) \langle \mathbf{x}, \Delta_i \rangle \langle \mathbf{x}, \tilde{\Delta} \rangle \right| \\ &\leq \frac{p_i(\mathbf{x})p_1(\mathbf{x})N_iN_1}{(p_1(\mathbf{x})N_1 + p_i(\mathbf{x})N_i)^2} |(y - g(\mathbf{a}_i^\top \mathbf{x}))g'(\mathbf{a}_i^\top \mathbf{x}) \langle \mathbf{x}, \Delta_i/\sigma^2 \rangle \langle \mathbf{x}, \tilde{\Delta} \rangle| \end{aligned}$$

For g = linear, sigmoid and ReLU, we have that $|g'(\cdot)| \leq 1$. Moreover, $\frac{p_i(\mathbf{x})p_1(\mathbf{x})N_iN_1}{(p_1(\mathbf{x})N_1 + p_i(\mathbf{x})N_i)^2} \leq 1/4$. Thus we have

$$\frac{p_i(\mathbf{x})p_1(\mathbf{x})N_iN_1}{(p_1(\mathbf{x})N_1 + p_i(\mathbf{x})N_i)^2} |(y - g(\mathbf{a}_i^\top \mathbf{x}))g'(\mathbf{a}_i^\top \mathbf{x}) \langle \mathbf{x}, \Delta_i/\sigma^2 \rangle \langle \mathbf{x}, \tilde{\Delta} \rangle| \leq \frac{1}{4} |y - g(\mathbf{a}_i^\top \mathbf{x})| |\langle \mathbf{x}, \Delta_i/\sigma^2 \rangle \langle \mathbf{x}, \tilde{\Delta} \rangle|.$$

We thus get

$$\mathbb{E} \left[\frac{-p_i(\mathbf{x})p_1(\mathbf{x})N_iN_1}{(\sum_i p_i(\mathbf{x})N_i)^2} \left(\frac{y - g(\mathbf{a}_i^\top \mathbf{x})}{\sigma^2} \right) g'(\mathbf{a}_i^\top \mathbf{x}) \langle \mathbf{x}, \Delta_i \rangle \langle \mathbf{x}, \tilde{\Delta} \rangle \right] \leq \frac{1}{4} \mathbb{E} [|y - g(\mathbf{a}_i^\top \mathbf{x})| \langle \mathbf{x}, \Delta_i / \sigma^2 \rangle \langle \mathbf{x}, \tilde{\Delta} \rangle] \quad (24)$$

$$\leq \frac{1}{4} \sqrt{\mathbb{E}[(y - g(\mathbf{a}_i^\top \mathbf{x}))^2] \mathbb{E}[\langle \mathbf{x}, \Delta_i / \sigma^2 \rangle^2 \langle \mathbf{x}, \tilde{\Delta} \rangle^2]} \quad (25)$$

$$\leq \frac{\sqrt{3}}{4} \sqrt{\mathbb{E}[(y - g(\mathbf{a}_i^\top \mathbf{x}))^2]} \|\Delta_i / \sigma^2\|_2 \|\tilde{\Delta}\|_2 \quad (26)$$

Now it remains to bound $\sqrt{\mathbb{E}[(y - g(\mathbf{a}_i^\top \mathbf{x}))^2]}$. Since $\|\mathbf{a}_i\|_2 \leq 1$, one can show that $\mathbb{E}[g(\mathbf{a}_i^\top \mathbf{x})^2] \leq 1$ for the given choice of non-linearities for g . Also, we have that

$$\mathbb{E}[y^2] = \mathbb{E}[\mathbb{E}[y^2 | \mathbf{x}]] = \mathbb{E} \left[\sum_{i \in [k]} p_i^*(\mathbf{x}) g(\langle \mathbf{a}_i^*, \mathbf{x} \rangle)^2 + \sigma^2 \right] = \mathbb{E} \left[\sum_{i \in [k]} p_i^*(\mathbf{x}) \mathbb{E}[g(\langle \mathbf{a}_i^*, \mathbf{x} \rangle)^2] + \sigma^2 \right] \leq 1 + \sigma^2,$$

where we used the following facts: (i) $\langle \mathbf{a}_i^*, \mathbf{x} \rangle$ is independent of the random variable $p_i^*(\mathbf{x})$ for each $i \in [k]$, (ii) $\langle \mathbf{a}_i^*, \mathbf{x} \rangle \stackrel{(d)}{=} \langle \mathbf{a}_1^*, \mathbf{x} \rangle$ and (iii) $\mathbb{E}[g(\langle \mathbf{a}_1^*, \mathbf{x} \rangle)^2] \leq 1$. Since $\mathbb{E}[(y - g(\mathbf{a}_i^\top \mathbf{x}))^2] \leq 2\mathbb{E}[y^2] + \mathbb{E}[g(\mathbf{a}_i^\top \mathbf{x})^2]$, after substituting these bounds in (26), we get

$$\mathbb{E} \left[\frac{-p_i(\mathbf{x})p_1(\mathbf{x})N_iN_1}{(\sum_i p_i(\mathbf{x})N_i)^2} \left(\frac{y - g(\mathbf{a}_i^\top \mathbf{x})}{\sigma^2} \right) g'(\mathbf{a}_i^\top \mathbf{x}) \langle \mathbf{x}, \Delta_i \rangle \langle \mathbf{x}, \tilde{\Delta} \rangle \right] \leq \frac{\sqrt{6(2 + \sigma^2)}}{4} \|\Delta_i / \sigma^2\|_2 \|\tilde{\Delta}\|_2.$$

Similarly,

$$\mathbb{E} \left[\frac{p_i(\mathbf{x})N_i p_1(\mathbf{x})N_1}{(\sum_i p_i(\mathbf{x})N_i)^2} \left(\frac{y - g(\mathbf{a}_1^\top \mathbf{x})}{\sigma^2} \right) g'(\mathbf{a}_1^\top \mathbf{x}) \langle \mathbf{x}, \Delta_1 \rangle \langle \mathbf{x}, \tilde{\Delta} \rangle \right] \leq \frac{\sqrt{6(2 + \sigma^2)}}{4} \|\Delta_1 / \sigma^2\|_2 \|\tilde{\Delta}\|_2.$$

Substituting the above two inequalities in (23), we obtain that

$$\mathbb{E}[(p^{(1)}(\mathbf{A}, \mathbf{W}) - p^{(1)}(\mathbf{A}^*, \mathbf{W})) \langle \mathbf{x}, \tilde{\Delta} \rangle] \leq 2(k-1) \frac{\sqrt{6(2 + \sigma^2)}}{4} \|\Delta_1 / \sigma^2\|_2 \|\tilde{\Delta}\|_2.$$

Defining $\kappa \triangleq (k-1) \frac{\sqrt{6(2 + \sigma^2)}}{2}$ and using the fact that $\|\Delta / \sigma^2\|_2 \leq \varepsilon_1$, we thus obtain

$$\left\| \mathbb{E}[(p^{(1)}(\mathbf{A}, \mathbf{W}) - p^{(1)}(\mathbf{A}^*, \mathbf{W})) \mathbf{x}] \right\|_2 \leq \kappa \varepsilon_1.$$

F. Proof of Lemma 3

F.1. Proof for $k = 2$

Proof. We first prove the lemma for $k = 2$. We show that the assumptions in Appendix B hold *globally* in our setting yielding a geometric convergence. Here we simply denote $M(\mathbf{W}, \mathbf{A}^*)$ as $M(\mathbf{w})$ dropping the explicit dependence on \mathbf{A}^* . Recall that

$$Q(\mathbf{w} | \mathbf{w}_t) = \mathbb{E}_{p_{\mathbf{w}^*}(\mathbf{x}, y)} \left[p_1(\mathbf{x}, y, \mathbf{w}_t) \cdot (\mathbf{w}^\top \mathbf{x}) - \log(1 + e^{\mathbf{w}^\top \mathbf{x}}) \right],$$

where

$$p_1(\mathbf{x}, y, \mathbf{w}_t) = \frac{f(\mathbf{w}_t^\top \mathbf{x}) \mathcal{N}(y | g(\mathbf{a}_1^\top \mathbf{x}), \sigma^2)}{f(\mathbf{w}^\top \mathbf{x}) \mathcal{N}(y | g(\mathbf{a}_1^\top \mathbf{x}), \sigma^2) + (1 - f(\mathbf{w}^\top \mathbf{x})) \mathcal{N}(y | g(\mathbf{a}_2^\top \mathbf{x}), \sigma^2)}. \quad (27)$$

For simplicity we drop the subscript in the above expectation with respect to the distribution $p_{\mathbf{w}^*}(\mathbf{x}, y)$. Now we verify each of the assumptions.

- Convexity of Ω easily follows from its definition.

- We have that

$$Q(\mathbf{w}|\mathbf{w}^*) = \mathbb{E} \left[p_1(\mathbf{x}, y, \mathbf{w}^*) \cdot (\mathbf{w}^\top \mathbf{x}) - \log(1 + e^{\mathbf{w}^\top \mathbf{x}}) \right].$$

Note that the strong-concavity of $Q(\cdot|\mathbf{w}^*)$ is equivalent to the strong-convexity of $-Q(\cdot|\mathbf{w}^*)$. Denoting the sigmoid function by f , we have that for all $\mathbf{w} \in \Omega$,

$$\begin{aligned} -\nabla^2 Q(\mathbf{w}|\mathbf{w}^*) &= \mathbb{E} [f'(\mathbf{w}^\top \mathbf{x}) \cdot \mathbf{x} \mathbf{x}^\top], \\ &\stackrel{\text{(Stein's lemma)}}{=} \mathbb{E} [f'''(\mathbf{w}^\top \mathbf{x}) \cdot \mathbf{w} \mathbf{w}^\top + \mathbb{E}[f'(\mathbf{w}^\top \mathbf{x})] \cdot I] \\ &= \mathbb{E}[f'''(\|\mathbf{w}\| Z)] \cdot \mathbf{w} \mathbf{w}^\top + \mathbb{E}[f'(\|\mathbf{w}\| Z)] \cdot I, \quad Z \sim \mathcal{N}(0, 1) \\ &\stackrel{(a)}{\succcurlyeq} \inf_{0 \leq \alpha \leq 1} \min \{ \mathbb{E}[f'(\alpha Z)], \mathbb{E}[f'(\alpha Z)] + \alpha^2 \mathbb{E}[f'''(\alpha Z)] \} \cdot I \\ &= \underbrace{0.14}_{\lambda} \cdot I \end{aligned} \tag{28}$$

where (a) follows from finding the two possible eigenvalues of the positive-definite matrix in the previous step and considering the minimum among them to ensure strong-convexity. Here the value of λ is found numerically to be approximately around 0.1442.

- For any $\mathbf{w}, \mathbf{w}_t \in \Omega$,

$$\nabla Q(\mathbf{w}|\mathbf{w}_t) = \mathbb{E} [p_1(\mathbf{x}, y, \mathbf{w}_t) \cdot \mathbf{x} - f(\mathbf{w}^\top \mathbf{x}) \cdot \mathbf{x}].$$

Thus,

$$\|\nabla Q(M(\mathbf{w})|\mathbf{w}^*) - \nabla Q(M(\mathbf{w})|\mathbf{w})\| = \|\mathbb{E} [(p_1(\mathbf{x}, y, \mathbf{w}_t) - p_1(\mathbf{x}, y, \mathbf{w}^*)) \cdot \mathbf{x}]\| \stackrel{(a)}{\leq} \gamma_\sigma \|\mathbf{w} - \mathbf{w}^*\|,$$

where we want to prove in (a) that γ_σ is smaller than 0.14 for all $\mathbf{w} \in \Omega$. Intuitively, this means that the posterior probability in (27) is smooth with respect to the parameter \mathbf{w} . We will now show that this can be achieved in the high-SNR regime when σ is sufficiently small. This will ensure that $\kappa_\sigma \triangleq \frac{\gamma_\sigma}{\lambda} < 1$. In particular, the value of γ_σ is dimension-independent and depends only on the choice of the non-linearity g .

To prove that

$$\|\mathbb{E} [(p_1(\mathbf{x}, y, \mathbf{w}) - p_1(\mathbf{x}, y, \mathbf{w}^*)) \cdot \mathbf{x}]\| \leq \gamma \|\mathbf{w} - \mathbf{w}^*\| = \gamma \|\Delta\|,$$

it suffices to show

$$\langle \mathbb{E} [(p_1(\mathbf{x}, y, \mathbf{w}) - p_1(\mathbf{x}, y, \mathbf{w}^*)) \cdot \mathbf{x}], \tilde{\Delta} \rangle \leq \gamma \|\Delta\| \|\tilde{\Delta}\|, \quad \forall \tilde{\Delta} \in \mathbb{R}^d.$$

Or equivalently,

$$\mathbb{E} \left[(p_1(\mathbf{x}, y, \mathbf{w}) - p_1(\mathbf{x}, y, \mathbf{w}^*)) \langle \mathbf{x}, \tilde{\Delta} \rangle \right] \leq \gamma \|\Delta\| \|\tilde{\Delta}\|.$$

Let $\Delta \triangleq \mathbf{w} - \mathbf{w}^*$ and $f(u) \triangleq p_1(\mathbf{x}, y, \mathbf{w}_u)$ where $\mathbf{w}_u = \mathbf{w}^* + u\Delta, u \in [0, 1]$. Thus $f(1) = p_1(\mathbf{x}, y, \mathbf{w})$ and $f(0) = p_1(\mathbf{x}, y, \mathbf{w}^*)$. So we get

$$p_1(\mathbf{x}, y, \mathbf{w}) - p_1(\mathbf{x}, y, \mathbf{w}^*) = f(1) - f(0) = \int_0^1 f'(u) du = \int_0^1 \langle \nabla p_1(\mathbf{x}, y, \mathbf{w}_u), \Delta \rangle du,$$

where the gradient is evaluated with respect to \mathbf{w}_u . Differentiating (27) with respect to \mathbf{w} , we get that

$$\begin{aligned} \nabla_{\mathbf{w}} p_1(\mathbf{x}, y, \mathbf{w}) &= \frac{f(\mathbf{w}^\top \mathbf{x})(1 - f(\mathbf{w}^\top \mathbf{x}))\mathcal{N}(y|g(\mathbf{a}_1^\top \mathbf{x}), \sigma^2)\mathcal{N}(y|g(\mathbf{a}_2^\top \mathbf{x}), \sigma^2)}{(f(\mathbf{w}^\top \mathbf{x})\mathcal{N}(y|g(\mathbf{a}_1^\top \mathbf{x}), \sigma^2) + (1 - f(\mathbf{w}^\top \mathbf{x}))\mathcal{N}(y|g(\mathbf{a}_2^\top \mathbf{x}), \sigma^2))^2} \cdot \mathbf{x} \\ &\triangleq R(\mathbf{x}, y, \mathbf{w}, \sigma) \cdot \mathbf{x}. \end{aligned}$$

Thus,

$$\begin{aligned}
 \mathbb{E} \left[(p_1(\mathbf{x}, y, \mathbf{w}) - p_1(\mathbf{x}, y, \mathbf{w}^*)) \langle \mathbf{x}, \tilde{\Delta} \rangle \right] &= \mathbb{E} \left[\left(\int_0^1 R(\mathbf{x}, y, \mathbf{w}_u, \sigma) \langle \mathbf{x}, \Delta \rangle du \right) \langle \mathbf{x}, \tilde{\Delta} \rangle \right] \\
 &= \int_0^1 \mathbb{E} \left[R(\mathbf{x}, y, \mathbf{w}_u, \sigma) \langle \mathbf{x}, \Delta \rangle \langle \mathbf{x}, \tilde{\Delta} \rangle \right] du \\
 &\leq \left(\int_0^1 \sqrt{\mathbb{E}[R(\mathbf{x}, y, \mathbf{w}_u, \sigma)^2]} du \right) \sqrt{\mathbb{E} \left[\langle \mathbf{x}, \Delta \rangle^2 \langle \mathbf{x}, \tilde{\Delta} \rangle^2 \right]} \\
 &\leq \underbrace{\sqrt{3} \left(\int_0^1 \sqrt{\mathbb{E}[R(\mathbf{x}, y, \mathbf{w}_u, \sigma)^2]} du \right)}_{\gamma_\sigma} \|\Delta\| \|\tilde{\Delta}\| \\
 &= \gamma_\sigma \|\Delta\| \|\tilde{\Delta}\|,
 \end{aligned}$$

where the last inequality follows from Lemma 5 of (Balakrishnan et al., 2017). Our goal is to now prove that $\gamma_\sigma \rightarrow 0$ as $\sigma \rightarrow 0$. First observe that

$$\begin{aligned}
 R(\mathbf{x}, y, \mathbf{w}, \sigma) &= \frac{f(\mathbf{w}^\top \mathbf{x})(1 - f(\mathbf{w}^\top \mathbf{x}))e^{-(y-g(\mathbf{a}_1^\top \mathbf{x}))/2\sigma^2} e^{-(y-g(\mathbf{a}_1^\top \mathbf{x}))/2\sigma^2}}{(f(\mathbf{w}^\top \mathbf{x})e^{-(y-g(\mathbf{a}_1^\top \mathbf{x}))/2\sigma^2} + (1 - f(\mathbf{w}^\top \mathbf{x}))e^{-(y-g(\mathbf{a}_2^\top \mathbf{x}))/2\sigma^2})^2} \leq \frac{1}{4} \left(\text{since } \frac{ab}{(a+b)^2} \leq 1/4 \right) \\
 &= \frac{f(1-f)e^{\frac{(y-g(\mathbf{a}_1^\top \mathbf{x}))^2 - (y-g(\mathbf{a}_2^\top \mathbf{x}))^2}{2\sigma^2}}}{\left(f + (1-f)e^{\frac{(y-g(\mathbf{a}_1^\top \mathbf{x}))^2 - (y-g(\mathbf{a}_2^\top \mathbf{x}))^2}{2\sigma^2}} \right)^2} \rightarrow 0 \text{ as } \sigma \rightarrow 0,
 \end{aligned}$$

where the key observation is that irrespective of the sign of $(y - g(\mathbf{a}_1^\top \mathbf{x}))^2 - (y - g(\mathbf{a}_2^\top \mathbf{x}))^2$, the ratio still goes to zero and hence by dominated convergence theorem $\mathbb{E}[R(\mathbf{x}, y, \mathbf{w}_u, \sigma)^2] \rightarrow 0$ for each $u \in [0, 1]$. Now we show that this convergence is uniform in u and thus $\gamma_\sigma \rightarrow 0$. For simplicity, define

$$\Delta_1 \triangleq (y - g(\mathbf{a}_1^\top \mathbf{x}))^2, \quad \Delta_2 \triangleq (y - g(\mathbf{a}_2^\top \mathbf{x}))^2 \text{ and } \sigma = \frac{1}{n}. \quad (29)$$

Thus,

$$R(\mathbf{x}, y, \mathbf{w}_u, \sigma) = \frac{f(1-f)e^{\frac{n^2}{2}(\Delta_1 - \Delta_2)}}{\left(f + (1-f)e^{\frac{n^2}{2}(\Delta_1 - \Delta_2)} \right)^2} \quad (30)$$

$$\leq \frac{f(1-f)e^{\frac{n^2}{2}(\Delta_1 - \Delta_2)}}{\left((1-f)e^{\frac{n^2}{2}(\Delta_1 - \Delta_2)} \right)^2} = \frac{f}{1-f} e^{-\frac{n^2}{2}(\Delta_1 - \Delta_2)}. \quad (31)$$

Similarly,

$$R(\mathbf{x}, y, \mathbf{w}_u, \sigma) \leq \frac{1-f}{f} e^{-\frac{n^2}{2}(\Delta_2 - \Delta_1)}. \quad (32)$$

Thus, we get

$$R(\mathbf{x}, y, \mathbf{w}_u, \sigma) \leq \max \left(\frac{1-f}{f}, \frac{f}{1-f} \right) e^{-\frac{n^2}{2}(|\Delta_1 - \Delta_2|)}. \quad (33)$$

Hence

$$\frac{\gamma_\sigma}{\sqrt{3}} = \int_0^1 \sqrt{\mathbb{E}[\text{Ratio}(\mathbf{x}, y, \mathbf{w}_u, \sigma)^2]} du \quad (34)$$

$$\leq \int_0^1 \sqrt{\mathbb{E} \left[\max \left(\frac{1-f}{f}, \frac{f}{1-f} \right)^2 e^{-n^2|\Delta_1 - \Delta_2|} \right]} du \quad (35)$$

$$\leq \int_0^1 \sqrt{\mathbb{E} \left[\left(\frac{1-f}{f} \right)^2 e^{-n^2|\Delta_1 - \Delta_2|} + \left(\frac{f}{1-f} \right)^2 e^{-n^2|\Delta_1 - \Delta_2|} \right]} du \quad (36)$$

$$= \int_0^1 \sqrt{2\mathbb{E} [e^{2\mathbf{w}_u^\top \mathbf{x}} e^{-n^2|\Delta_1 - \Delta_2|}]} du \quad (37)$$

$$\leq \int_0^1 \sqrt{2\sqrt{\mathbb{E}[e^{4\mathbf{w}_u^\top \mathbf{x}}]} \mathbb{E}[e^{-2n^2|\Delta_1 - \Delta_2|}]} du \quad (38)$$

$$\stackrel{(a)}{\leq} \sqrt{2e^4} \sqrt{\mathbb{E}[e^{-2n^2|\Delta_1 - \Delta_2|}]}, \quad (39)$$

where (a) follows from the fact $\|\mathbf{w}_u\| \leq 1$ and $\mathbb{E}[e^{4\mathbf{w}_u^\top \mathbf{x}}] = e^{8\|\mathbf{w}_u\|^2} \leq e^8$, for each $u \in [0, 1]$. Now we analyze the convergence rate of the last term $\mathbb{E}[e^{-2n^2|\Delta_1 - \Delta_2|}]$ for the case of linear regression, i.e. $g(z) = z$. Notice that for the two-mixtures, we have

$$y \stackrel{(d)}{=} Z(\mathbf{a}_1^\top \mathbf{x}) + (1-Z)\mathbf{a}_2^\top \mathbf{x} + \sigma N = Z(\mathbf{a}_1^\top \mathbf{x}) + (1-Z)\mathbf{a}_2^\top \mathbf{x} + \frac{N}{n}, \quad Z|\mathbf{x} \sim \text{Bern}(f(\mathbf{w}_*^\top \mathbf{x})). \quad (40)$$

Thus,

$$\Delta_1 - \Delta_2 \stackrel{(d)}{=} (y - \mathbf{a}_1^\top \mathbf{x})^2 - (y - \mathbf{a}_2^\top \mathbf{x})^2 \quad (41)$$

$$= (\mathbf{a}_1^\top \mathbf{x} - \mathbf{a}_2^\top \mathbf{x})^2 (1-2Z) + \frac{2N}{n} (\mathbf{a}_2^\top \mathbf{x} - \mathbf{a}_1^\top \mathbf{x}) \quad (42)$$

$$= \langle \mathbf{x}, \mathbf{v} \rangle^2 (1-2Z) + \frac{2N}{n} \langle \mathbf{x}, \mathbf{v} \rangle, \quad \mathbf{v} = \mathbf{a}_1 - \mathbf{a}_2. \quad (43)$$

Since Z can equal either 0 or 1, we have

$$\gamma_\sigma \leq \sqrt{3}\sqrt{2e^4} \left(\mathbb{E}[e^{-2n^2|\langle \mathbf{x}, \mathbf{v} \rangle^2(1-2Z) + \frac{2N}{n}\langle \mathbf{x}, \mathbf{v} \rangle}|] \right)^{1/4} \quad (44)$$

$$\leq \sqrt{6e^4} \left(\mathbb{E} \left[\max \left(e^{-2n^2|\langle \mathbf{x}, \mathbf{v} \rangle^2 + \frac{2N}{n}\langle \mathbf{x}, \mathbf{v} \rangle}, e^{-2n^2|-\langle \mathbf{x}, \mathbf{v} \rangle^2 + \frac{2N}{n}\langle \mathbf{x}, \mathbf{v} \rangle} \right) \right] \right)^{1/4} \quad (45)$$

$$\leq \sqrt{6\sqrt{2}e^4} \left(\mathbb{E} \left[e^{-2n^2|\langle \mathbf{x}, \mathbf{v} \rangle^2 + \frac{2N}{n}\langle \mathbf{x}, \mathbf{v} \rangle} \right] \right)^{1/4} \quad (46)$$

$$= \sqrt{6\sqrt{2}e^4} \left(\mathbb{E} \left[e^{-2n^2|Z^2 + \frac{2ZN}{n}|} \right] \right)^{1/4}, \quad Z \sim \mathcal{N}(0, \|\mathbf{a}_1 - \mathbf{a}_2\|), N \sim \mathcal{N}(0, 1). \quad (47)$$

$$= O \left(\sqrt{6\sqrt{2}e^4} \left(\mathbb{E}[e^{-2n^2 Z^2}] \right)^{1/4} \right) \quad (48)$$

$$= \sqrt{6\sqrt{2}e^4} \left(\sqrt{\frac{1}{4n^2 \|\mathbf{a}_1 - \mathbf{a}_2\|^2 + 1}} \right)^{1/4} \quad (49)$$

$$= O \left(\frac{1}{(n \|\mathbf{a}_1 - \mathbf{a}_2\|)^{1/4}} \right) \quad (50)$$

$$= O \left(\left(\frac{\sigma}{\|\mathbf{a}_1 - \mathbf{a}_2\|} \right)^{1/4} \right). \quad (51)$$

□

F.2. Proof for general k

Proof. The proof strategy for general k is similar. First let $\varepsilon_1 = 0$. Our task is to show that the assumptions of Appendix B hold globally in our setting. The domain Ω is clearly convex since

$$\Omega = \{\mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_{k-1}) : \|\mathbf{w}_i\| \leq 1, \forall i \in [k-1]\}.$$

Now we verify Assumption 2. The function $Q(\cdot|\mathbf{w}_t)$ is given by

$$Q(\mathbf{w}|\mathbf{w}_t) = \mathbb{E} \left[\sum_{i \in [k-1]} p_{\mathbf{w}_t}^{(i)}(\mathbf{w}_i^\top \mathbf{x}) - \log \left(1 + \sum_{i \in [k-1]} e^{\mathbf{w}_i^\top \mathbf{x}} \right) \right],$$

where $p_{\mathbf{w}_t}^{(i)} \triangleq \mathbb{P}[z = i | \mathbf{x}, y, \mathbf{w}_t]$ corresponds to the posterior probability for the i^{th} expert, given by

$$p_{\mathbf{w}_t}^{(i)} = \frac{p_{i,t}(\mathbf{x}) \mathcal{N}(y | g(\mathbf{a}_i^\top \mathbf{x}), \sigma^2)}{\sum_{j \in [k]} p_{j,t}(\mathbf{x}) \mathcal{N}(y | g(\mathbf{a}_j^\top \mathbf{x}), \sigma^2)}, \quad p_{i,t}(\mathbf{x}) = \frac{e^{(\mathbf{w}_t)_i^\top \mathbf{x}}}{1 + \sum_{j \in [k-1]} e^{(\mathbf{w}_t)_j^\top \mathbf{x}}}.$$

Throughout we follow the convention that $\mathbf{w}_k = 0$. Thus the gradient of Q with respect to the i^{th} gating parameter \mathbf{w}_i is given by

$$\nabla_{\mathbf{w}_i} Q(\mathbf{w}|\mathbf{w}_t) = \mathbb{E} \left[\left(p_{\mathbf{w}_t}^{(i)} - \frac{e^{\mathbf{w}_i^\top \mathbf{x}}}{1 + \sum_{j \in [k-1]} e^{\mathbf{w}_j^\top \mathbf{x}}} \right) \cdot \mathbf{x} \right], \quad i \in [k-1].$$

Thus the $(i, j)^{\text{th}}$ block of the negative Hessian $-\nabla_{\mathbf{w}}^{(2)} Q(\mathbf{w}|\mathbf{w}^*) \in \mathbb{R}^{d(k-1) \times d(k-1)}$ is given by

$$-\nabla_{\mathbf{w}_i, \mathbf{w}_j} Q(\mathbf{w}|\mathbf{w}^*) = \begin{cases} \mathbb{E}[p_i(\mathbf{x})(1 - p_i(\mathbf{x})) \cdot \mathbf{x} \mathbf{x}^\top], & j = i \\ \mathbb{E}[-p_i(\mathbf{x})p_j(\mathbf{x}) \cdot \mathbf{x} \mathbf{x}^\top], & j \neq i \end{cases} \quad (52)$$

where $p_i(\mathbf{x}) = \frac{e^{\mathbf{w}_i^\top \mathbf{x}}}{1 + \sum_{j \in [k-1]} e^{\mathbf{w}_j^\top \mathbf{x}}}$. It is clear from (52) that $-\nabla_{\mathbf{w}}^{(2)} Q(\mathbf{w}|\mathbf{w}^*)$ is positive semi-definite. Since we are interested in the strong convexity of $-Q(\cdot|\mathbf{w}^*)$ which is equivalent to positive definiteness of the negative Hessian, it suffices to show that

$$\lambda \triangleq \inf_{\mathbf{w} \in \Omega} \lambda_{\min} \left(-\nabla_{\mathbf{w}}^{(2)} Q(\mathbf{w}|\mathbf{w}^*) \right) > 0.$$

Since the Hessian is continuous with respect to \mathbf{w} and consequently the minimum eigenvalue of it, there exists a $\mathbf{w}' \in \Omega$ such that

$$\lambda = \lambda_{\min} \left(-\nabla_{\mathbf{w}'}^{(2)} Q(\mathbf{w}'|\mathbf{w}^*) \right) = \inf_{\|\mathbf{a}\|=1} \mathbf{a}^\top \left(-\nabla_{\mathbf{w}'}^{(2)} Q(\mathbf{w}'|\mathbf{w}^*) \right) \mathbf{a},$$

where $\mathbf{a} = (\mathbf{a}_1^\top, \dots, \mathbf{a}_{k-1}^\top)^\top \in \mathbb{R}^{d(k-1)}$. In view of (52), the above equation can be further simplified to

$$\lambda = \inf_{\|\mathbf{a}\|=1} \mathbb{E}[\mathbf{a}_x^\top M_x \mathbf{a}_x], \quad (53)$$

where $\mathbf{a}_x = (\mathbf{a}_1^\top \mathbf{x}, \dots, \mathbf{a}_{k-1}^\top \mathbf{x})^\top \in \mathbb{R}^{k-1}$ and M_x is given by

$$M_x(i, j) = \begin{cases} p_i(\mathbf{x})(1 - p_i(\mathbf{x})), & i = j \\ -p_i(\mathbf{x})p_j(\mathbf{x}), & i \neq j \end{cases}$$

Let the infimum in (53) is attained by \mathbf{a}^* , i.e. $\lambda = \mathbb{E}[(\mathbf{a}_x^*)^\top M_x \mathbf{a}_x^*]$. For each \mathbf{x} , M_x is strictly diagonally dominant since $|M_x(i, i)| = p_i(\mathbf{x})(1 - p_i(\mathbf{x})) = p_i(\mathbf{x}) \left(\sum_{j \neq i, j \in [k]} p_j(\mathbf{x}) \right) > p_i(\mathbf{x}) \left(\sum_{j \neq i, j \in [k-1]} p_j(\mathbf{x}) \right) = \sum_{j \neq i} M(i, j)$. Thus M_x is positive-definite and $(\mathbf{a}_x^*)^\top M_x \mathbf{a}_x^* > 0$ whenever $\mathbf{a}_x^* \neq 0$. Since \mathbf{x} follows a continuous distribution it follows that $\mathbf{a}_x^* \neq 0$ with probability 1 and thus $\lambda = \mathbb{E}[(\mathbf{a}_x^*)^\top M_x \mathbf{a}_x^*] > 0$.

Now it remains to show that Assumption 3 too holds, i.e.

$$\|\nabla Q(M(\mathbf{w})|\mathbf{w}^*) - \nabla Q(M(\mathbf{w})|\mathbf{w})\| \leq \gamma \|\mathbf{w} - \mathbf{w}^*\|.$$

Note that $\mathbf{w} = (\mathbf{w}_1^\top, \dots, \mathbf{w}_{k-1}^\top)^\top \in \mathbb{R}^{d(k-1)}$. We will show that

$$\|(\nabla Q(M(\mathbf{w})|\mathbf{w}^*))_i - (\nabla Q(M(\mathbf{w})|\mathbf{w}))_i\| \leq \gamma_\sigma \|\mathbf{w} - \mathbf{w}^*\|, \quad i \in [k-1],$$

where $(\nabla Q(M(\mathbf{w})|\mathbf{w}))_i \in \mathbb{R}^d$ refers to the i^{th} block of the gradient and $\gamma_\sigma \rightarrow 0$. Observe that

$$(\nabla Q(M(\mathbf{w})|\mathbf{w}^*))_i - (\nabla Q(M(\mathbf{w})|\mathbf{w}))_i = \mathbb{E} \left[(p_{\mathbf{w}}^{(i)} - p_{\mathbf{w}^*}^{(i)}) \cdot \mathbf{x} \right]$$

Let $\Delta = \mathbf{w} - \mathbf{w}^*$ and correspondingly $\Delta = (\Delta_1^\top, \dots, \Delta_{k-1}^\top)^\top$ where $\Delta_i = \mathbf{w}_i - \mathbf{w}_i^*$. Thus it suffices to show that

$$\left\| \mathbb{E}[(p_{\mathbf{w}}^{(i)} - p_{\mathbf{w}^*}^{(i)}) \cdot \mathbf{x}] \right\| \leq \gamma_\sigma \|\Delta\|.$$

Or equivalently,

$$\mathbb{E}[(p_{\mathbf{w}}^{(i)} - p_{\mathbf{w}^*}^{(i)}) \langle \mathbf{x}, \tilde{\Delta} \rangle] \leq \gamma_\sigma \|\Delta\| \|\tilde{\Delta}\|, \quad \forall \tilde{\Delta} \in \mathbb{R}^d.$$

We consider the case $i = 1$. The proof for the other cases is similar. Recall that

$$p_{\mathbf{w}}^{(1)} = \frac{p_1(\mathbf{x}) \mathcal{N}(y|g(\mathbf{a}_1^\top \mathbf{x}), \sigma^2)}{\sum_{j \in [k]} p_j(\mathbf{x}) \mathcal{N}(y|g(\mathbf{a}_j^\top \mathbf{x}), \sigma^2)}, \quad p_i(\mathbf{x}) = \frac{e^{\mathbf{w}_i^\top \mathbf{x}}}{1 + \sum_{j \in [k-1]} e^{\mathbf{w}_j^\top \mathbf{x}}}, \quad i \in [k-1].$$

For simplicity we define $N_i = \mathcal{N}(y|g(\mathbf{a}_i^\top \mathbf{x}), \sigma^2)$. It is straightforward to verify that

$$\nabla_{\mathbf{w}_j} p_i(\mathbf{x}) = \begin{cases} p_i(\mathbf{x})(1 - p_i(\mathbf{x})) \cdot \mathbf{x}, & j = i \\ -p_i(\mathbf{x}) p_j(\mathbf{x}) \cdot \mathbf{x}, & j \neq i \end{cases}$$

Thus

$$\begin{aligned} \nabla_{\mathbf{w}_1} (p_{\mathbf{w}}^{(1)}) &= \nabla_{\mathbf{w}_1} \left(\frac{p_1(\mathbf{x}) N_1}{\sum_{i=1}^N p_i(\mathbf{x}) N_i} \right) \\ &= \frac{\left(\sum_{i=1}^N p_i(\mathbf{x}) N_i \right) p_1(\mathbf{x})(1 - p_1(\mathbf{x})) N_1 - p_1(\mathbf{x}) N_1 \left(-\sum_{j \neq 1} p_j(\mathbf{x}) p_1(\mathbf{x}) N_j + p_1(\mathbf{x})(1 - p_1(\mathbf{x})) N_1 \right)}{\left(\sum_{i=1}^N p_i(\mathbf{x}) N_i \right)^2} \cdot \mathbf{x} \\ &= \frac{p_1(\mathbf{x}) N_1 \left(\sum_{j \geq 2} p_j(\mathbf{x}) N_j \right)}{\left(\sum_{i=1}^N p_i(\mathbf{x}) N_i \right)^2} \cdot \mathbf{x} \\ &\triangleq R_1(\mathbf{x}, y, \mathbf{w}, \sigma) \cdot \mathbf{x} \end{aligned}$$

Similarly,

$$\begin{aligned} \nabla_{\mathbf{w}_i} (p_{\mathbf{w}}^{(1)}) &= \frac{p_1(\mathbf{x}) p_i(\mathbf{x}) N_1 N_i}{\left(\sum_{i=1}^N p_i(\mathbf{x}) N_i \right)^2} \cdot \mathbf{x}, \quad i \neq 1, \\ &\triangleq R_i(\mathbf{x}, y, \mathbf{w}, \sigma) \cdot \mathbf{x}. \end{aligned}$$

Let $\mathbf{w}_u \triangleq \mathbf{w}^* + u\Delta$, $u \in [0, 1]$ and $f(u) \triangleq p_{\mathbf{w}_u}^{(1)}$. Thus

$$\begin{aligned} p_{\mathbf{w}}^{(1)} - p_{\mathbf{w}^*}^{(1)} &= f(1) - f(0) = \int_0^1 f'(u) du \\ &= \int_0^1 \left(\sum_{i \in [k-1]} \langle \nabla_{\mathbf{w}_i} (p_{\mathbf{w}_u}^{(1)}), \Delta_i \rangle \right) du \\ &= \sum_{i \in [k-1]} \int_0^1 R_i(\mathbf{x}, y, \mathbf{w}, \sigma) \langle \mathbf{x}, \Delta_i \rangle du. \end{aligned}$$

So we get

$$\begin{aligned}
 \mathbb{E}[(p_{\mathbf{w}}^{(1)} - p_{\mathbf{w}^*}^{(1)})\langle \mathbf{x}, \tilde{\Delta} \rangle] &= \sum_{i \in [k-1]} \int_0^1 \mathbb{E}[R_i(\mathbf{x}, y, \mathbf{w}_u, \sigma)\langle \mathbf{x}, \Delta_i \rangle\langle \mathbf{x}, \tilde{\Delta} \rangle] du \\
 &\leq \sum_{i \in [k-1]} \int_0^1 \sqrt{\mathbb{E}[R_i(\mathbf{x}, y, \mathbf{w}_u, \sigma)^2] \mathbb{E}[\langle \mathbf{x}, \Delta_i \rangle^2 \langle \mathbf{x}, \tilde{\Delta} \rangle^2]} du \\
 &\leq \sum_{i \in [k-1]} \int_0^1 \sqrt{\mathbb{E}[R_i(\mathbf{x}, y, \mathbf{w}_u, \sigma)^2]} \left(\sqrt{3} \|\Delta_i\| \|\tilde{\Delta}\| \right) du \\
 &\leq \sum_{i \in [k-1]} \int_0^1 \sqrt{\mathbb{E}[R_i(\mathbf{x}, y, \mathbf{w}_u, \sigma)^2]} \left(\sqrt{3} \|\Delta\| \|\tilde{\Delta}\| \right) du \\
 &= \underbrace{\left(\sum_{i \in [k-1]} \int_0^1 \sqrt{\mathbb{E}[R_i(\mathbf{x}, y, \mathbf{w}_u, \sigma)^2]} du \right)}_{\gamma_\sigma^{(1)}} \left(\sqrt{3} \|\Delta\| \|\tilde{\Delta}\| \right)
 \end{aligned}$$

Now our goal is to show that $\mathbb{E}[R_i(\mathbf{x}, y, \mathbf{w}_u, \sigma)^2] \rightarrow 0$ as $\sigma \rightarrow 0$. For $i = 1$, we have

$$R_1(\mathbf{x}, y, \mathbf{w}_u, \sigma)^2 = \left(\frac{\sum_{j \geq 2} p_1(\mathbf{x}) p_j(\mathbf{x}) N_1 N_j}{\left(\sum_{i=1}^N p_i(\mathbf{x}) N_i \right)^2} \right)^2 \leq k \sum_{j \geq 2} \left(\frac{p_1(\mathbf{x}) p_j(\mathbf{x}) N_1 N_j}{\left(\sum_{i=1}^N p_i(\mathbf{x}) N_i \right)^2} \right)^2 \leq k \sum_{j \geq 2} \left(\frac{p_1(\mathbf{x}) p_j(\mathbf{x}) N_1 N_j}{(p_1(\mathbf{x}) N_1 + p_j(\mathbf{x}) N_j)^2} \right)^2$$

Similarly,

$$R_i(\mathbf{x}, y, \mathbf{w}_u, \sigma)^2 \leq \left(\frac{p_1(\mathbf{x}) p_i(\mathbf{x}) N_1 N_i}{(p_1(\mathbf{x}) N_1 + p_i(\mathbf{x}) N_i)^2} \right)^2, \quad \forall i \neq 1, i \in [k-1].$$

For $\mathbf{w} = \mathbf{w}_u$ and $i \neq 1$, we have that

$$\begin{aligned}
 \frac{p_1(\mathbf{x}) p_i(\mathbf{x}) N_1 N_i}{(p_1(\mathbf{x}) N_1 + p_i(\mathbf{x}) N_i)^2} &= \frac{e^{\mathbf{w}_1^\top \mathbf{x}} e^{\mathbf{w}_i^\top \mathbf{x}} e^{-\frac{(y-g(\mathbf{a}_1^\top \mathbf{x}))^2}{2\sigma^2}} e^{-\frac{(y-g(\mathbf{a}_i^\top \mathbf{x}))^2}{2\sigma^2}}}{\left(e^{\mathbf{w}_1^\top \mathbf{x}} e^{-\frac{(y-g(\mathbf{a}_1^\top \mathbf{x}))^2}{2\sigma^2}} + e^{\mathbf{w}_i^\top \mathbf{x}} e^{-\frac{(y-g(\mathbf{a}_i^\top \mathbf{x}))^2}{2\sigma^2}} \right)^2} \leq \frac{1}{4} \\
 &= \frac{e^{\mathbf{w}_1^\top \mathbf{x}} e^{\mathbf{w}_i^\top \mathbf{x}} e^{\frac{(y-g(\mathbf{a}_1^\top \mathbf{x}))^2 - (y-g(\mathbf{a}_i^\top \mathbf{x}))^2}{2\sigma^2}}}{\left(e^{\mathbf{w}_1^\top \mathbf{x}} + e^{\mathbf{w}_i^\top \mathbf{x}} e^{\frac{(y-g(\mathbf{a}_1^\top \mathbf{x}))^2 - (y-g(\mathbf{a}_i^\top \mathbf{x}))^2}{2\sigma^2}} \right)^2} \\
 &\xrightarrow{\sigma \rightarrow 0} 0.
 \end{aligned}$$

Thus, by Dominated Convergence Theorem, $\mathbb{E}[R_i(\mathbf{x}, y, \mathbf{w}_u, \sigma)^2] \rightarrow 0$ for each $u \in [0, 1]$. To show that $\int_0^1 \mathbb{E}[R_i(\mathbf{x}, y, \mathbf{w}_u, \sigma)^2] du \rightarrow 0$, we can now follow the same analysis as in the proof of Theorem 2 from (29) on-wards (replacing \mathbf{w} there with $\mathbf{w}_1 - \mathbf{w}_i$) which ensures that $\gamma_\sigma^{(1)}$ in our case converges to zero. Similarly for other $i \in [k-1]$, we get that $\gamma_\sigma^{(i)} \rightarrow 0$. Taking $\gamma_\sigma = \gamma_\sigma^{(1)} + \dots + \gamma_\sigma^{(k-1)}$ and $\kappa_\sigma = \frac{\gamma_\sigma}{\lambda}$ completes the proof. \square

G. Gradient EM algorithm

In this section, we provide the convergence guarantees for the gradient EM algorithm. For simplicity, we prove the results for $k = 2$ and $(\mathbf{a}_1, \mathbf{a}_2) = (\mathbf{a}_1^*, \mathbf{a}_2^*)$. Thus we want to learn the gating parameter \mathbf{w}^* in this setting. The results for the general case follow essentially the same proof as that of Theorem 2. In particular, our Theorem 5 can be viewed as a generalization of Lemma 3. Together with Lemma 4, extension to general k is straightforward.

Note that in the M-step of the EM algorithm, instead of maximizing $Q(\cdot|\mathbf{w}_t)$, we can choose an iterate so that it increases the Q value instead of fully maximizing it, i.e. $Q(\mathbf{w}_{t+1}|\mathbf{w}_t) \geq Q(\mathbf{w}_t|\mathbf{w}_t)$. Such a procedure is termed as *generalized EM*. *Gradient EM* is an example of generalized EM in which we take an ascent step in the direction of the gradient of $Q(\cdot|\mathbf{w}_t)$ to produce the next iterate, i.e.

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha \nabla Q(\mathbf{w}_t|\mathbf{w}_t),$$

where $\alpha > 0$ is a suitably chosen step size and the gradient is with respect to the first argument. To account for the constrained optimization, we can include a projection step. Mathematically,

$$\mathbf{w}_{t+1} = G(\mathbf{w}_t), \quad G(\mathbf{w}) = \Pi_{\Omega}(\mathbf{w} + \alpha \nabla Q(\mathbf{w}|\mathbf{w})),$$

where Π_{Ω} refers to the projection operator. Our next result establishes that the iterates of the gradient EM algorithm too converge geometrically for an appropriately chosen step size α .

Theorem 5. *Suppose that the domain $\Omega = \{\mathbf{w} \in \mathbb{R}^d : \|\mathbf{w}\|_2 \leq 1\}$ and $(\mathbf{a}_1, \mathbf{a}_2) = (\mathbf{a}_1^*, \mathbf{a}_2^*)$. Then there exist constants $\alpha_0 > 0$ and $\sigma_0 > 0$ such that for any step size $0 < \alpha \leq \alpha_0$ and noise variance $\sigma < \sigma_0$, the gradient EM updates on the gating parameter $\{\mathbf{w}\}_{t \geq 0}$ converge geometrically to the true parameter \mathbf{w}^* , i.e.*

$$\|\mathbf{w}_t - \mathbf{w}^*\| \leq (\rho_{\sigma})^t \|\mathbf{w}_0 - \mathbf{w}^*\|,$$

where ρ_{σ} is a dimension-independent constant depending on g and σ .

Remark 3. The condition $\sigma < \sigma_0$ ensures that the Lipschitz constant ρ_{σ} for the map G is strictly less than 1. The constant α_0 depends only on two universal constants which are nothing but the strong-concavity and the smoothness parameters for the function $Q(\cdot|\mathbf{w}^*)$.

Proof. In addition to the assumptions of Appendix B, if we can ensure that the map $-Q(\cdot|\mathbf{w}^*)$ is μ -smooth, then the proof follows from Theorem 3 of (Balakrishnan et al., 2017) if we choose $\alpha_0 = \frac{2}{\mu + \lambda}$ where λ is the strong-convexity parameter of $-Q(\cdot|\mathbf{w}^*)$. The strong-convexity is already established in Appendix D.3. To find the smoothness parameter, note that

$$\begin{aligned} -\nabla^2 Q(\mathbf{w}|\mathbf{w}^*) &= \mathbb{E} [f'(\mathbf{w}^{\top} \mathbf{x}) \cdot \mathbf{x} \mathbf{x}^{\top}], \\ &= \mathbb{E} [f'''(\mathbf{w}^{\top} \mathbf{x}) \cdot \mathbf{w} \mathbf{w}^{\top}] + \mathbb{E} [f'(\mathbf{w}^{\top} \mathbf{x})] \cdot I \\ &= \mathbb{E} [f'''(\|\mathbf{w}\| Z)] \cdot \mathbf{w} \mathbf{w}^{\top} + \mathbb{E} [f'(\|\mathbf{w}\| Z)] \cdot I, \quad Z \sim \mathcal{N}(0, 1) \\ &\leq \sup_{0 \leq \alpha \leq 1} \min \{ \mathbb{E} [f'(\alpha Z)], \mathbb{E} [f'(\alpha Z)] + \alpha^2 \mathbb{E} [f'''(\alpha Z)] \} \cdot I \\ &= \underbrace{0.25}_{\mu} \cdot I. \end{aligned}$$

The contraction parameter is then given by

$$\rho_{\sigma} = 1 - \frac{2\lambda + 2\gamma_{\sigma}}{\mu + \lambda}.$$

Since $\gamma_{\sigma} \xrightarrow{\sigma \rightarrow 0} 0$, $\rho_{\sigma} < 1$ whenever $\sigma < \sigma_0$ for a constant σ_0 . □

H. Additional experiments

H.1. Synthetic data

In Figure 4, we varied the number of samples our data set and fixed the other set of parameters to $k = 3, d = 5, \sigma = 0.5$.

In Figure 5 we repeated our experiments for the choice of $n = 10000, d = 5, k = 3$ for two different popular choices of non-linearities: sigmoid and ReLU. The same conclusion as in the linear setting holds in this case too with our algorithm outperforming the EM consistently.

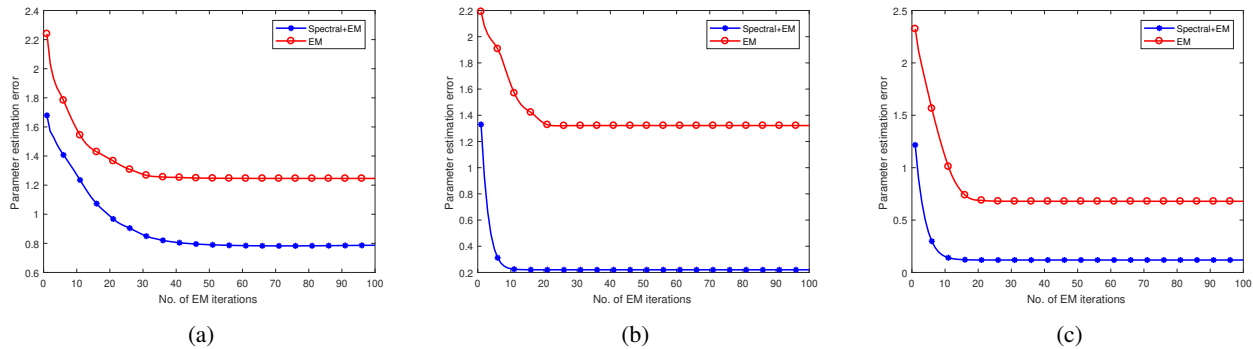


Figure 4: Plot of parameter estimation error with varying number of samples(n): (a) $n = 1000$ (b) $n = 5000$. (c) $n = 10000$.

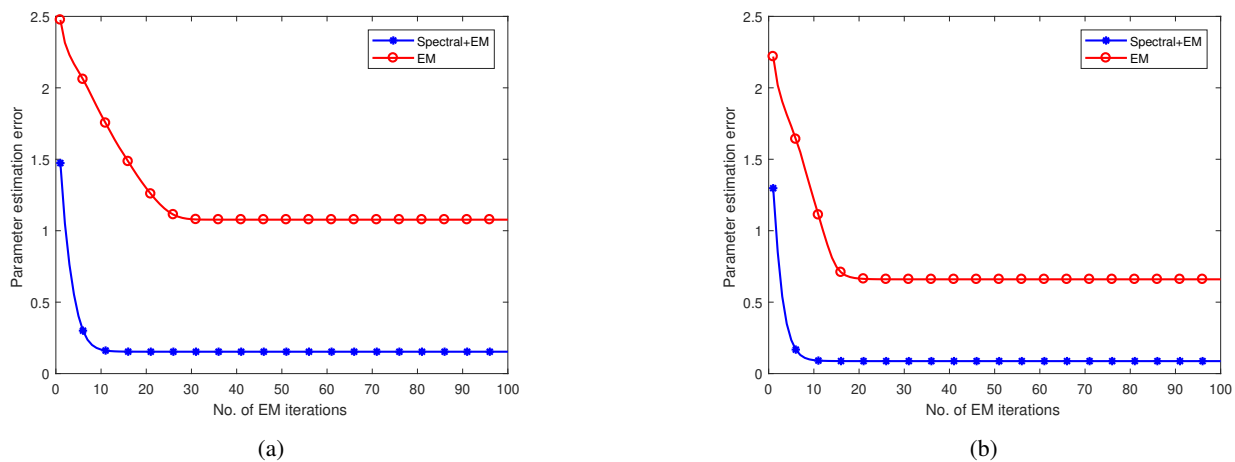


Figure 5: Parameter estimation error for the sigmoid and ReLU nonlinearities respectively.

H.2. Real data

For real data experiments, we choose the 3 standard regression data sets from the UCI Machine Learning Repository: Concrete Compressive Strength Data Set, Stock portfolio performance Data Set, and Airfoil Self-Noise Data Set (Yeh, 1998; Liu & Yeh, 2017; Brooks et al., 1989). In all the three tasks, the goal is to predict the outcome or the response y for each input x , which typically contains some task specific attributes. For example, in the concrete compressive strength, the task is to predict the compressive strength of the concrete given its various attributes such as the component of cement, water, age, etc. For this data, the input $x \in \mathbb{R}^8$ corresponds to 8 different attributes of the concrete and the output $y \in \mathbb{R}$ corresponds to its concrete strength. Similarly, for the stock portfolio data set the input $x \in \mathbb{R}^6$ contains the weights of several stock-picking concepts such as weight of the Large S/P concept, weight of the Small systematic Risk concept, etc., and the output y is the corresponding excess return. The airfoil data set is obtained from a series of aerodynamic and acoustic tests of two and three-dimensional airfoil blade sections and the goal is predict the scaled sound pressure level (in dB) given the frequency, angle of attack, etc., For all the tasks, we pre-processed the data by whitening the input and scaling the output to lie in $(-1, 1)$. We randomly allotted 75% of the data samples for training and the rest for testing. Our evaluation metric is the prediction error on the test set $(x_i, y_i)_{i=1}^n$ defined as

$$\mathcal{E} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2,$$

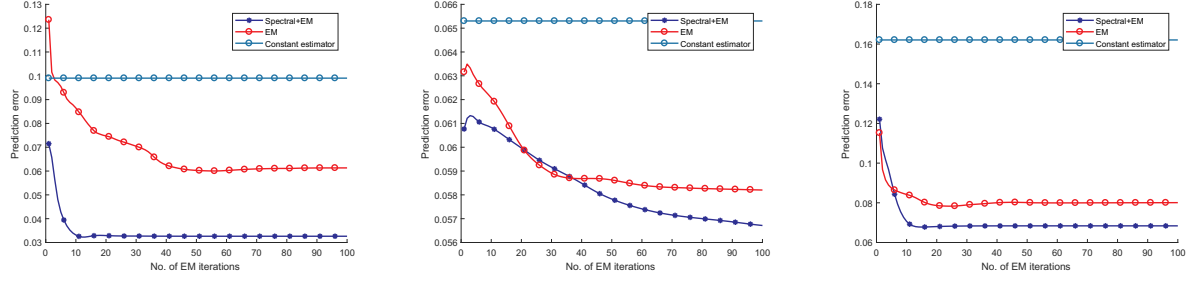


Figure 6: Prediction error for the concrete, stock portfolio and the airfoil data sets respectively.

where \hat{y}_i corresponds to the predicted output response using the learned parameters. In other words,

$$\hat{y} = \sum_{i \in [k]} \frac{e^{\hat{w}_i^\top \mathbf{x}}}{\sum_{j \in [k]} e^{\hat{w}_j^\top \mathbf{x}}} \cdot g(\hat{a}_i^\top \mathbf{x}).$$

We ran the joint-EM algorithm (with 10 different trails) on these tasks with various choices for $k \in \{2, \dots, 10\}$, $\sigma \in \{0.1, 0.4, 0.8, 1\}$, $g \in \{\text{linear}, \text{sigmoid}, \text{ReLU}\}$ and found the best hyper-parameters to be $(k = 3, \sigma = 0.1$ and $g = \text{linear})$, $(k = 3, \sigma = 0.4, g = \text{sigmoid})$ and $(k = 3, \sigma = 0.1, g = \text{linear})$ for the three datasets respectively. For this choice of best hyper-parameters found for joint-EM, we ran our algorithm. Figure 6 highlights the predictive performance of our algorithm as compared to that of the EM. We also plotted the variance of the test data for reference and to gauge the performance of our algorithm. In all the settings our algorithm is able to obtain a better set of parameters resulting in smaller prediction error.