
Learning Mixtures of Experts with EM: A Mirror Descent Perspective

Quentin Fruytier¹ Aryan Mokhtari¹ Sujay Sanghavi¹

Abstract

Classical Mixtures of Experts (MoE) are Machine Learning models that involve partitioning the input space, with a separate “expert” model trained on each partition. Recently, MoE-based model architectures have become popular as a means to reduce training and inference costs. There, the partitioning function and the experts are both learnt jointly via gradient descent-type methods on the log-likelihood. In this paper we study theoretical guarantees of the Expectation Maximization (EM) algorithm for the training of MoE models. We first rigorously analyze EM for MoE where the conditional distribution of the target and latent variable conditioned on the feature variable belongs to an exponential family of distributions and show its equivalence to projected Mirror Descent with unit step size and a Kullback-Leibler Divergence regularizer. This perspective allows us to derive new convergence results and identify conditions for local linear convergence; In the special case of mixture of 2 linear or logistic experts, we additionally provide guarantees for linear convergence based on the signal-to-noise ratio. Experiments on synthetic and (small-scale) real-world data supports that EM outperforms the gradient descent algorithm both in terms of convergence rate and the achieved accuracy.

1. Introduction

Classical Mixtures of Experts (MoE) (Jacobs et al., 1991; Jordan & Jacobs, 1994) are a crucial class of parametric latent variable models that have gained significant popularity in deep learning for their ability to reduce both training and inference costs (Chen et al., 2022). MoE are particu-

larly effective when the feature space can be divided and processed by specialized models, known as experts. Instead of relying on a single, large model to handle all input-output mappings, MoE utilize an ensemble of specialized experts. Each expert is responsible for a specific subset of the input space, allowing the system to efficiently route inputs to the most appropriate expert. This partitioning enables each expert to focus on mapping its designated inputs to outputs using a separate, optimized model. By leveraging multiple specialized experts rather than a monolithic model, MoE can achieve greater scalability and flexibility. This modular approach not only enhances computational efficiency but also allows for improved performance, as each expert can be finely tuned to handle its particular segment of the input space effectively. Real-world applications of MoE such as Sparse MoE span across various domains, including language translation, speech recognition, recommendation systems, and more (Fedus et al., 2022; Ma et al., 2018; Hinton et al., 2012; Liu et al., 2024).

In its most generic form, training an MoE model involves training both (a) the parameters in the individual experts, and (b) the gating function that routes inputs to the appropriate expert. Typically, these are both learnt jointly by first formulating the final loss function as applied to the ensemble output, and then minimizing this joint loss function (via SGD or its variants) over the parameters of the gate and the experts.

In this paper we investigate, primarily from a theoretical perspective, the training of classical MoE as defined by Jacobs et al. (1991) using a classic algorithm: Expectation Maximization (EM). As opposed to SGD-based methods which are agnostic to whether a parameter is in the gate or in an expert, EM first formulates two separate problems – one for the router, and another for the experts – in a specific way. It then solves each problem in isolation and in parallel, and then collates the outputs to arrive at the updated set of gate and expert parameters.

For our theoretical results, we consider the setting of general MoE where the joint distribution of the target and latent variable conditioned on the feature variable belongs to an exponential family of distributions. We then narrow in on two simpler instances of MoE models: Mixture of Linear Experts (where each expert is a linear regressor) and Mix-

¹Electrical and Computer Engineering Department, University of Texas at Austin, Austin, Texas, United States of America. Correspondence to: Quentin Fruytier <qdf76@my.utexas.edu>, Aryan Mokhtari <mokhtari@austin.utexas.edu>, Sujay Sanghavi <sanghavi@mail.utexas.edu>.

ture of Logistic Experts (where each expert is a logistic regressor). The router in each case is a linear softmax.

Main Contributions: The primary finding of this paper is to unveil the correspondence between EM for general MoE to projected Mirror Descent. A similar correspondence was first discovered by [Kunstner et al. \(2021\)](#), but was limited to generative models for which the joint complete data distribution belongs to an exponential family of distributions; this did not include MoE. As such, our contributions can be seen as a generalization of ([Kunstner et al., 2021](#)) to all generative models for which the conditional distribution of the target and latent variables conditioned on the feature variable belongs to an exponential family of distributions. We next state the details of our contributions.

- 1) In Theorem 4.1, we show that when EM is applied to general MoE, the iterates are equivalent to the ones generated by projected Mirror Descent on the conditional likelihood function with unit step-size and KL divergence. Here, the projection is an expectation moment projection over the parameter space. By leveraging this correspondence, in Theorem 4.2, we obtain sufficient conditions for which EM applied to general MoE converges to a stationary point or the true parameters. We further characterize the explicit convergence rate for each of the considered settings.
- 2) Next, in Theorem 5.1, we narrow in on the special cases of mixtures of 2 linear or logistic experts and show EM is equivalent to Mirror Descent with unit step-size and KL divergence, without requiring any extra projection. This perspective allows us to recover classic MD convergence results which we contextualize for our setting in Corollary B.1. Finally, we characterize the sufficient conditions for convergence in terms of the Missing Information Matrix (MIM) in Theorem B.2 and, subsequently, the signal-to-noise ratio (SNR) of the generative model in Theorem B.4.
- 3) Finally, on synthetic and small scale proof of concept real world datasets, we observe that EM outperforms gradient descent both in terms of convergence rate and the achieved performance. As well as supporting our theoretical results, this re-iterates the power of the EM algorithm for fitting MoE previously suggested by ([Jordan & Jacobs, 1994](#); [Jordan & Xu, 1995](#)).

1.1. Related Work

The EM algorithm ([Dempster et al., 1977](#)) is a powerful tool for fitting latent variable models. Previous research on EM has demonstrated that, under mild smoothness assumptions, its parameter iterates converge to a stationary point of the log-likelihood objective ([Dempster et al., 1977](#); [Wu, 1983](#); [Tseng, 2004](#)). Subsequent research on EM introduced new analytical frameworks to provide specialized guarantees,

particularly regarding the convergence of EM iterates to the true model parameters and the rate of this convergence. An adopted framework in the past decade, introduced by [Balakrishnan et al. \(2017\)](#), interprets EM as a variant of the gradient descent algorithm. For latent variable models with a strongly convex EM objective that satisfies a condition known as “first-order stability,” it was shown that EM iterates converge linearly to the true parameters. Subsequent works utilized this framework to show a local linear rate for Mixtures of Gaussians and Mixtures of Linear Regressions ([Balakrishnan et al., 2017](#); [Daskalakis et al., 2017](#); [Dwivedi et al., 2018](#); [Kwon et al., 2019](#); [Kwon & Caramanis, 2020b](#); [Kwon et al., 2021](#); [Kwon & Caramanis, 2020a](#)).

In a recent work, [Kunstner et al. \(2021\)](#) proved that EM – where the complete data distribution belongs in an exponential family – is equivalent to the mirror descent algorithm with unit step-size and Kullback–Leibler (KL) divergence regularizer:

$$\text{KL}[p(\mathbf{x}; \boldsymbol{\theta}) || p(\mathbf{x}, \boldsymbol{\phi})] = \mathbb{E}_{\mathbf{x} | \boldsymbol{\theta}} \left[\log \left(\frac{p(\mathbf{x}; \boldsymbol{\theta})}{p(\mathbf{x}; \boldsymbol{\phi})} \right) \right]. \quad (1)$$

This led to the first non-asymptotic convergence rates for EM, independent of parameterization. While this characterization of EM included Mixtures of Gaussians, it failed to extend to Mixtures of Regression or MoE. Still, the authors analyzed the setting where the distribution of the latent variable conditioned on all other variables is an exponential family distribution for which they showed EM converged sub-linearly to a stationary point. Our work extends these findings to MoE, obtaining sufficient conditions under which EM converges sub-linearly and linearly to the true parameters.

There have also been works attempting to explore and understand how to fit MoE. The foundational paper ([Jordan & Jacobs, 1994](#)) was the first to empirically use EM and EM-variants to fit MoE, yielding encouraging results. Then, the follow-up paper by [Jordan & Xu \(1995\)](#) showed that for MoE and Hierarchical MoE with strongly convex negative log-likelihood objective, EM and EM-like iterations converge linearly to the true parameters where the rate constant depends on the eigenvalues of the Hessian matrix. However, the objective is generally non convex, raising doubts about whether the necessary assumptions for their result hold, even locally. Other works have remarked that the nature of the gating function creates a form of competition between the experts during training that can lead to local minima. Bayesian methods for MoE include variational learning and maximum a posteriori (MAP) estimation. But, [Yuksel et al. \(2012\)](#) noted that these solutions are not trivial due to the softmax gate not admitting a conjugate prior and are prone to getting stuck at local minima. [Makkuva et al. \(2019\)](#) analyzed a variant of the EM algorithm for MoE which consists in 1) first recovering the expert param-

eters using a tensor decomposition method, then 2) whilst freezing the experts, utilizing EM to fit the gating function's parameters only. They proved that their approach recovers the true parameters at a *nearly* linear rate. Then, [Becker et al. \(2020\)](#) experimented with an EM variant algorithm for Gaussian Mixture of Experts called Expectation Information Maximization (EIM) which featured an extra information projection step. They obtained promising empirical results on both synthetic and real world datasets. Our work extends upon previous works by making the direct connection between EM for MoE and projected mirror descent. We also unveil the sufficient conditions for EM to converge sub-linearly to a stationary point or true parameter, and for special cases, characterize these sufficient conditions with respect to the SNR of the generative model.

2. Mixture of Experts

Next, we formally describe the setting under consideration for the Mixture of k Experts. The notation used throughout is summarized in the [Notation Section](#) of the appendix.

Data Generation Model: First, the input or feature vector variable $\mathbf{x} \in \mathbb{R}^d$ is sampled based on a probability density function $p(\mathbf{x})$. Second, given the feature vector \mathbf{x} , a latent variable $z \in [k]$, responsible for routing \mathbf{x} to the appropriate expert is sampled with probability mass function $P(z|\mathbf{x})$. Finally, given the pair (\mathbf{x}, z) , the target y is generated according to the probability distribution $p(y|\mathbf{x}, z)$. Hence, the complete data distribution is $p(\mathbf{x}, y, z) = p(y|\mathbf{x}, z)p(z|\mathbf{x})p(\mathbf{x})$ and the joint input-output probability distribution can be written as

$$p(\mathbf{x}, y) = p(\mathbf{x}) \sum_{z \in [k]} p(y|\mathbf{x}, z)P(z|\mathbf{x}). \quad (2)$$

This is the most general form of data generation under mixture of k -experts, but here we focus on the case where \mathbf{x} is sampled from a unit spherical Gaussian distribution, i.e. $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$, and $P(z|\mathbf{x}) = P(z|\mathbf{x}; \mathbf{w}^*)$ is parameterized by a set of vectors $(\mathbf{w}_1^*, \dots, \mathbf{w}_k^*)$ and can be cast as

$$P(z = i|\mathbf{x}; \mathbf{w}^*) = \frac{e^{\mathbf{x}^\top \mathbf{w}_i^*}}{\sum_{j \in [k]} e^{\mathbf{x}^\top \mathbf{w}_j^*}}, \quad i \in [k]. \quad (3)$$

In other words, the probability mass function of the latent variable is the softmax function between the inner product of \mathbf{x} with the parameter vectors concatenated as $\mathbf{w}^* = (\mathbf{w}_1^*, \dots, \mathbf{w}_k^*) \in \mathbb{R}^{d \times k}$.

Regarding $p(y|\mathbf{x}, z)$, the probability distribution of the target variable y conditioned on the input variable \mathbf{x} and latent routing variable $z = i$ (i.e., expert i), we also assume that it is parameterized by a vector β_i^* and we have $p(y|\mathbf{x}, z = i) = p(y|\mathbf{x}, z = i; \beta_i^*)$. Thus, given the concatenated vector $\beta^* = (\beta_1^*, \dots, \beta_k^*) \in \mathbb{R}^{s \times k}$, we can write $p(y|\mathbf{x}, z) = p(y|\mathbf{x}, z; \beta^*)$.

In this paper, we will consider three different settings. The first, *General MoE*, is the most general setting we consider. It comprises all MoE where the distribution of y, z conditioned on \mathbf{x} belongs to an exponential family of distribution and enjoys a natural re-parameterization $p(y, z|\mathbf{x}, \theta^*) = p(y, z; \theta_x^*)$ with $\theta_x^* = \eta(\mathbf{x}, \theta^*) \in \tilde{\Omega}$. That is to say that the conditional distribution can be written as

$$p(y, z|\mathbf{x}, \theta^*) \propto \exp \{ \langle s(y, z), \theta_x^* \rangle + A(\theta_x^*) \}, \quad (4)$$

where $s(y, z)$, θ_x , and $A(\cdot)$ are, respectively, referred to as the sufficient statistic, natural parameterization, and log partition of the exponential family of distributions. This setting includes the popular cases where $p(y|\mathbf{x}, z)$ is Gaussian or multivariate Bernoulli. In the former case, the exponential family of distribution corresponds to a Gaussian Mixture.

The second setting, *Mixture of Linear Experts*, is when

$$p(y|\mathbf{x}, z = i; \beta_i^*) \propto \exp \left\{ \frac{(y - \mathbf{x}^\top \beta_i^*)^2}{2} \right\} \quad (5)$$

as the density function of the normal distribution. This is equivalent to assume that the target $y = \mathbf{x}^\top \beta_i^* + \epsilon$ where ϵ is additive zero mean Gaussian noise with unit variance.

Finally, *Mixture of Logistic Experts*, is when the density function $p(y|\mathbf{x}, z = i; \beta_i^*)$ can be written as

$$P(y = 1|\mathbf{x}, z = i; \beta_i^*) = \frac{\exp(\mathbf{x}^\top \beta_i^*)}{1 + \exp(\mathbf{x}^\top \beta_i^*)} \quad (6)$$

and $P(y = -1|\mathbf{x}, z = i; \beta_i^*) = 1 - P(y = 1|\mathbf{x}, z = i; \beta_i^*)$.

Maximum Likelihood Loss: Given the assumed data distribution of (\mathbf{x}, y) , our goal is to find the set of feasible parameters $\beta \in \mathbb{R}^{d \times k}$ and $\mathbf{w} \in \mathbb{R}^{d \times k}$ that maximize the log likelihood function. For ease of notation we define $\theta := (\beta, \mathbf{w}) \in \mathbb{R}^{2d \times k}$ as the concatenation of all parameters. Specifically, from (2), the expected log likelihood is

$$\mathbb{E}_{\mathbf{X}} [\log p(\mathbf{x})] + \mathbb{E}_{\mathbf{X}, Y} \left[\log \left(\sum_{z \in [k]} p(y|\mathbf{x}, z)P(z|\mathbf{x}) \right) \right].$$

Given that only the last term of the sum depends on parameters $\theta = (\mathbf{w}, \beta)^\top$ the negative log-likelihood objective function, $\mathcal{L}(\theta)$, that we aim to minimize can be written as

$$-\mathbb{E}_{\mathbf{X}, Y} \left[\log \left(\sum_{z \in [k]} p(y|\mathbf{x}, z; \beta)P(z|\mathbf{x}; \mathbf{w}) \right) \right] \quad (7)$$

Note that as we will discuss in detail, for both mixtures of linear or logistic experts, the above objective function is known to be non convex with respect to θ . In fact, this is generally true for Mixtures of Gaussian, Mixtures of Regressions, and Mixtures of Experts. In the next section we will discuss the use of the EM algorithm for solving this optimization problem.

3. EM for Mixtures of Experts

Next, we present the EM algorithm for MoE. EM takes a structured approach to minimizing the objective $\mathcal{L}(\theta)$ in (7). Each iteration of EM is decomposed into two steps as follows. The first step is called “expectation”: For current parameter estimate θ^t , we compute the expectation of the complete-data log-likelihood with respect to the latent variables, using the current parameter estimates θ^t and denote it by $Q(\theta|\theta^t)$, i.e.,

$$Q(\theta|\theta^t) = -\mathbb{E}_{\mathbf{X}, \mathbf{Y}} [\mathbb{E}_{Z|\mathbf{x}, \mathbf{y}; \theta^t} [\log p(\mathbf{x}, \mathbf{y}, z; \theta)]] . \quad (8)$$

Then, in the second step called “maximization”, we simply minimize the objective $Q(\theta|\theta^t)$ (or maximize $-Q(\theta|\theta^t)$) with respect to $\theta \in \Omega$ and obtain our new parameter as

$$\theta^{t+1} := \operatorname{argmin}_{\theta \in \Omega} Q(\theta|\theta^t). \quad (9)$$

Since $\log p(y, z|\mathbf{x}; \theta) = \log p(y|\mathbf{x}; \beta) + \log p(z|\mathbf{x}; \mathbf{w})$, it follows that the EM objective (8) is linearly separable in the parameters β and \mathbf{w} . Thus, we can rewrite $Q(\theta|\phi)$ as the sum of two functions that depend only on β and \mathbf{w} , respectively. Subsequently, the EM update (9) is obtained as the concatenation $\theta^{t+1} = (\mathbf{w}^{t+1}, \beta^{t+1})^\top$, where

$$\begin{aligned} \mathbf{w}^{t+1} &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} -\mathbb{E}_{\mathbf{X}, \mathbf{Y}} [\mathbb{E}_{Z|\mathbf{x}, \mathbf{y}; \theta^t} [\log p(z|\mathbf{x}; \mathbf{w})]] , \\ \beta^{t+1} &= \operatorname{argmin}_{\beta \in \mathbb{R}^d} -\mathbb{E}_{\mathbf{X}, \mathbf{Y}} [\mathbb{E}_{Z|\mathbf{x}, \mathbf{y}; \theta^t} [\log p(y|z, \mathbf{x}; \beta)]] . \end{aligned}$$

EM has two well understood characteristics: (a) its update always minimize an objective that is an upper bound on the likelihood, and (b) fixed points of the EM update are also stationary points of the likelihood. We now show this below. The original objective, $\mathcal{L}(\theta)$, can be decomposed into the difference between the EM objective and another function that is bounded below by 0. Specifically, for any $\theta, \phi \in \Omega$,

$$\begin{aligned} \mathcal{L}(\theta) &= -\mathbb{E}_{\mathbf{X}, \mathbf{Y}} [\log(p(y|\mathbf{x}; \theta))] \\ &= -\mathbb{E}_{\mathbf{X}, \mathbf{Y}} \mathbb{E}_{Z|\mathbf{x}, \mathbf{y}, \phi} [\log(p(y|\mathbf{x}, z; \theta))] \\ &= -\mathbb{E}_{\mathbf{X}, \mathbf{Y}} \mathbb{E}_{Z|\mathbf{x}, \mathbf{y}, \phi} \left[\log \left(\frac{p(y, z|\mathbf{x}; \theta)}{p(z|\mathbf{x}, \mathbf{y}; \theta)} \right) \right] . \end{aligned}$$

Denoting $H(\theta|\phi) := -\mathbb{E}_{\mathbf{X}, \mathbf{Y}} \mathbb{E}_{Z|\mathbf{x}, \mathbf{y}, \phi} [\log p(z|\mathbf{x}, \mathbf{y}; \theta)]$, it follows that

$$\mathcal{L}(\theta) = Q(\theta|\phi) - H(\theta|\phi). \quad (10)$$

where $H(\theta|\phi)$ is bounded below by 0. Thus, $Q(\theta|\phi)$ acts as an upper bound on the negative log-likelihood.

Next, applying Jensen’s inequality shows that $H(\theta|\phi)$ is minimized at $\theta = \phi$ where $H(\theta|\theta) = 0$. Consequently, it follows that the negative log-likelihood gradient matches that of the surrogate EM objective at $\phi = \theta$, i.e., $\nabla \mathcal{L}(\theta) = \nabla Q(\theta|\theta)$. This suggests that any stationary point of $\mathcal{L}(\theta)$ is also a stationary point of the EM algorithm and vice versa.

3.1. EM for Symmetric Mixture of 2-Experts

So far, we discussed EM for the most general form of MoE. Next, we derive EM for the special case of *Symmetric Mixture of Experts* (SymMoE), the focus of our analysis in Section 5. SymMoE is a simplified version of MoE where (i) the number of experts is restricted to 2, represented as $z \in \{-1, 1\}$, and (ii) the experts are symmetric around the linear separator, i.e., $\tilde{\beta}^* := \beta_1^* = -\beta_{-1}^*$. This symmetric structure simplifies the probability density functions introduced earlier, making the subsequent analysis easier to follow. We explore these simplifications in detail below.

As we are restricted to two experts, the expression for $P(z = 1|\mathbf{x}; \mathbf{w}^*) = 1 - P(z = -1|\mathbf{x}; \mathbf{w}^*)$ is

$$P(z = 1|\mathbf{x}; \mathbf{w}^*) = \frac{e^{\mathbf{x}^\top \mathbf{w}_1^*}}{e^{\mathbf{x}^\top \mathbf{w}_1^*} + e^{\mathbf{x}^\top \mathbf{w}_2^*}} .$$

For ease of notation, we define $\tilde{\mathbf{w}}^* := \mathbf{w}_1^* - \mathbf{w}_2^*$ and reparameterize the probability mass function of z given \mathbf{x} as

$$P(z|\mathbf{x}; \tilde{\mathbf{w}}^*) = \frac{\exp\{\frac{z+1}{2} \mathbf{x}^\top \tilde{\mathbf{w}}^*\}}{1 + e^{\mathbf{x}^\top \tilde{\mathbf{w}}^*}} . \quad (11)$$

Thus, under this simplification, the EM update of the gating parameter \mathbf{w} is now given as

$$\mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \mathbb{E}_{\mathbf{X}, \mathbf{Y}} \mathbb{E}_{Z|\mathbf{x}, \mathbf{y}; \theta^t} \left[\log \left(\frac{1 + e^{\mathbf{x}^\top \mathbf{w}}}{e^{\frac{z+1}{2} \mathbf{x}^\top \mathbf{w}}} \right) \right] .$$

While the above minimization problem is strongly convex, it does not have a closed form solution. In our experiments, we use gradient descent to obtain \mathbf{w}^{t+1} .

For the special case of a *symmetric mixture of linear experts* (SymMoLinE), the expression for $p(y|\mathbf{x}, z; \beta^*)$ given in (5) can be simplified as

$$p(y|\mathbf{x}, z; \beta^*) \propto \exp \left\{ \frac{(y - z \mathbf{x}^\top \tilde{\beta}^*)^2}{2} \right\} . \quad (12)$$

Under this simplification, the EM update of the expert parameter β is now obtained more compactly as

$$\beta^{t+1} = \mathbb{E}_{\mathbf{X}, \mathbf{Y}} [(2p(z = 1|\mathbf{x}, \mathbf{y}; \theta^t) - 1) \mathbf{x} \mathbf{y}] .$$

Similarly, for *symmetric mixtures of logistic experts* (SymMoLogE) with $y \in \{-1, 1\}$, the expression of $p(y|\mathbf{x}, z; \beta^*)$ given in (6) can also be simplified as

$$P(y|\mathbf{x}, z; \beta^*) = \frac{\exp\{\frac{yz+1}{2} \mathbf{x}^\top \beta^*\}}{1 + e^{\mathbf{x}^\top \beta^*}} . \quad (13)$$

Under this simplification, the EM update of the expert parameter β is now given as

$$\beta^{t+1} = \operatorname{argmin}_{\beta \in \mathbb{R}^d} \mathbb{E}_{\mathbf{X}, \mathbf{Y}} \left[\mathbb{E}_{Z|\mathbf{x}, \mathbf{y}; \theta^t} \left[\log \left(\frac{1 + e^{\mathbf{x}^\top \beta}}{e^{\frac{yz+1}{2} \mathbf{x}^\top \beta}} \right) \right] \right] .$$

3.2. EM for Deep and Sparse MoE

So far, we have derived EM for the foundational formulations of MoE, as initially proposed in (Jacobs et al., 1991). While EM is straightforward to derive in these cases, the same does not hold for deep MoE—and especially not for Sparse MoE. A deep MoE consists of $l \geq 2$ MoE blocks, as defined in Section 2, stacked sequentially. In this setup, the input x passes through the first MoE block and then sequentially through all subsequent blocks to produce the output y .

For completeness, and to encourage future work on large-scale applications of EM for MoE, we propose a formulation of EM for deep and sparse MoE in Appendix D.4. Instead of solving a separate latent-variable problem at each layer—as is done in classical MoE—we posit that the latent variable $z \in [k]^l$ should represent the entire sequence of experts selected across the network. This allows us to construct the EM objective in Equation (8).

An EM-like solution can then be derived for Sparse MoE, where the loss is computed solely from the sequences of experts observed through greedy expert selection. This formulation provides a principled approach to training deep and sparse MoE models using EM.

4. Main Result

In the previous section, we derived the EM update for both MoE and SymMoE as the solution to minimizing the EM objective in (8). We further demonstrated that this solution can be decomposed into the concatenation of the respective solutions to two minimization sub-problems. In this section, we will show that this update is exactly equivalent to performing a single step of the projected Mirror Descent (MD) update on $\mathcal{L}(\theta)$ with a unit step size and the KL divergence as a regularizer. To illustrate more clearly that minimizing $Q(\theta|\theta^t)$ in (8) corresponds to a projected MD step on the loss $\mathcal{L}(\theta)$ at the point θ^t , we first provide a brief overview of the core concept behind the MD update.

In most gradient-based methods, the next iteration is obtained by minimizing an upper bound of the objective function. For example, in Gradient Descent (GD), the next iterate is found by minimizing the first-order Taylor expansion of the objective at θ^t , with a squared norm regularizer. Specifically, for minimizing $\mathcal{L}(\cdot)$, the GD update with step size η at θ^t is equivalent to minimizing the following function:

$$\mathcal{L}(\theta^t) + \langle \nabla \mathcal{L}(\theta^t), \theta - \theta^t \rangle + \frac{1}{2\eta} \|\theta - \theta^t\|_2^2.$$

This function indeed serves as an upper bound for $\mathcal{L}(\cdot)$ if $\eta \leq 1/L$, where L is the Lipschitz constant of $\nabla \mathcal{L}(\theta)$.

Mirror Descent solves a similar sub-problem where instead of a squared norm regularizer, we employ the Bregman Di-

vergence regularizer: The Bregman divergence induced by a differentiable, convex function $h : \mathbb{R}^d \rightarrow \mathbb{R}$, measures the difference between $h(\theta)$ and its first-order approximation at θ^t , i.e.,

$$D_h(\theta, \theta^t) = h(\theta) - h(\theta^t) - \langle \nabla h(\theta^t), \theta - \theta^t \rangle. \quad (14)$$

Thus, the iterations of MD are derived by minimizing the following expression:

$$\mathcal{L}(\theta^t) + \langle \nabla \mathcal{L}(\theta^t), \theta - \theta^t \rangle + \frac{1}{\eta} D_h(\theta^t, \theta). \quad (15)$$

Finally, in projected mirror descent, the update is completed by projecting the solution obtained from minimizing (15) onto a subspace. As mentioned, this scheme is reasonable if the function approximation using the Bregman divergence serves as an upper bound for the function $\mathcal{L}(\theta)$. This can be ensured when the step size η is sufficiently small, and the condition of relative smoothness is satisfied (Lu et al., 2018).

In the upcoming theorem, we formally establish that for general MoE, minimizing the EM objective function defined in equation (9) is exactly equivalent to minimizing the subproblem associated with a single step of MD, as defined in equation (15), followed by a projection step. This result demonstrates that the EM update for General MoE is essentially performing an MD step in a specific natural re-parameterization space, then projecting the resulting solution onto Ω . The proof of this result is involved and is deferred to Appendix A.1.

Theorem 4.1. *For General MoE, there exists a natural reparameterization $\theta_x \in \{\eta(\cdot, \theta) : \theta \in \Omega\}$ with*

$$\mathcal{L}(\theta) = \mathbb{E}_X [L(\theta_x)] \quad (16)$$

and a mirror map $A(\theta_x)$ such that the EM update in (9) simplifies and is equivalent to the expectation moment projection,

$$\operatorname{argmin}_{\theta \in \Omega} \mathbb{E}_X \left[KL \left[p \left(y, z \mid \tilde{\theta}_x^{t+1} \right) \parallel p(y, z \mid \eta(x, \theta)) \right] \right], \quad (17)$$

where for each x , $\tilde{\theta}_x^{t+1}$ is obtained from the following MD step,

$$\operatorname{argmin}_{\psi \in \tilde{\Omega}} \langle \nabla L(\theta_x^t), \psi - \theta_x^t \rangle + D_A(\psi, \theta_x^t), \quad (18)$$

with $L(\theta_x)$ being 1-smooth relative to $A(\theta_x)$. Further, $\forall \psi_1, \psi_2 \in \tilde{\Omega}$, the divergence function $D_A(\psi_1, \psi_2)$ is equal to the KL divergence on $p(y, z, \psi)$:

$$D_A(\psi_1, \psi_2) = KL[p(y, z \mid \psi_2) \parallel p(y, z \mid \psi_1)]. \quad (19)$$

Proof sketch. We utilize the assumed property that the conditional distribution $p(y, z \mid x, \theta)$ belongs to an exponential

family of distributions to decompose the EM surrogate as

$$Q(\theta|\theta^t) - Q(\theta^t|\theta^t) = \mathbb{E}_X [\langle \nabla L(\theta_x^t), \theta_x - \theta_x^t \rangle + D_A(\theta_x, \theta_x^t)]. \quad (20)$$

The above derivation follows similar steps to Kunstner et al. Theorem 1 and is provided in Appendix A for completeness. We note that because $\theta_x = \eta(x, \theta)$ is not necessarily linear in x , we cannot conclude that minimizing (1) with respect to θ results in a direct MD step. Instead, recall the point-wise (in x) MD iterate, $\tilde{\theta}_x^{t+1} := \arg\min_{\theta_x} \langle \nabla L(\theta_x^t), \theta_x - \theta_x^t \rangle + D_A(\theta_x, \theta_x^t)$. Differentiating and setting equal to 0, it holds that

$$\nabla A(\tilde{\theta}_x^{t+1}) = s(x; \theta^t). \quad (21)$$

Using (2) and the decomposing $\nabla L(\theta_x^t) = \nabla A(\theta_x^t) - s(x; \theta^t)$, we see that (1) can be further simplified to

$$Q(\theta|\theta^t) - Q(\theta^t|\theta^t) = \mathbb{E}_X [-\langle \nabla A(\tilde{\theta}_x^{t+1}), \theta_x - \theta_x^t \rangle + A(\theta_x) - A(\theta_x^t)].$$

Because $A(\tilde{\theta}_x^{t+1})$ and $A(\theta_x^t)$ only depend on θ^t , minimizing the above with respect to θ is equivalent to minimizing the following with respect to θ

$$\mathbb{E}_X [D_A(\theta_x, \tilde{\theta}_x^{t+1})]. \quad (22)$$

Finally, substituting the Bregman Divergence induced by A by the KL divergence yields the claim (this derivation is also included in Appendix A.1 for completeness). \square

We remark that our proof for this result is not merely limited to the setting of MoE, but is satisfied for any mixture for which the distribution $p(y, z|x, \theta)$ of the target and latent variable conditioned on the feature variable belongs to an exponential family of distribution.

4.1. Convergence Guarantees for General MoE

Before moving into the convergence results, we specify the necessary assumptions previously outlined in (Kunstner et al., 2021) for the iterations of the EM Algorithm to be well-defined for our class of MoE models. See Appendix A.2 for a more in-depth discussion of the implications of these assumptions.

A_1 . The conditional distribution $p(y, z|\theta_x)$ is a steep, minimal exponential family of distribution and $\eta(x, \cdot)$ is a continuously differentiable function.

A_2 . The optimal objective function value is bounded below, i.e., $\mathcal{L}(\theta^*) > -\infty$, on the constraint set Ω .

A_3 . The following sub-level sets $\Omega_\theta := \{\phi \in \Omega : Q(\phi|\theta) \leq Q(\theta|\theta)\}$ are compact.

Next, we briefly introduce key definitions that will be used later. We say θ^1 is initialized in a locally average-convex region of $\mathcal{L}(\theta)$ with respect to the random variable X , if there exists a convex set $\Theta \subseteq \Omega$ containing θ^1, θ^* such that for all $\phi, \theta \in \Theta$,

$$\mathbb{E}_X [L(\phi_x)] \geq \mathbb{E}_X [L(\theta_x) + \langle \nabla L(\theta_x), \phi_x - \theta_x \rangle] \quad (23)$$

where $\theta_x := \eta(x, \theta)$. Furthermore, Θ is called α -average-convex relative to h if

$$\mathbb{E}_X [L(\phi_x)] \geq \mathbb{E}_X [L(\theta_x)] + \mathbb{E}_X [\langle \nabla L(\theta_x), \phi_x - \theta_x \rangle + \alpha D_h(\phi_x, \theta_x)] \quad (24)$$

Now, thanks to the previously shown correspondence between EM for general MoE and projected MD, we are able to present novel convergence properties of the EM Algorithm when applied to General MoE. The theorem that follows provides sufficient conditions and explicit rates for the convergence of EM iterations to a stationary point or true parameters. The proof, adapted from (Lu et al., 2018), is a bit more involved due to the nature of the extra projection step. It is included in Appendix A.2.

Theorem 4.2 (Convergence of EM). *Assuming $A_1 - A_3$. For general MoE with re-parameterization given by $\theta_x := \eta(x, \theta)$, strictly convex mirror map $A(\theta_x)$, and if for all $\theta_x^{t+1}, \theta_x^{t+1}, \phi_x \in \{\theta_x^t, \theta_x^*\}$,*

$$\mathbb{E}_X [D_A(\phi_x, \tilde{\theta}_x^{t+1})] \geq \mathbb{E}_X [D_A(\theta_x^{t+1}, \tilde{\theta}_x^{t+1}) + D_A(\phi_x, \theta_x^{t+1})], \quad (25)$$

then, the EM iterates $\{\theta^t\}_{t \in [T]}$ satisfy:

1) **Stationnarity**. For no additional conditions,

$$\min_{t \in [T]} \mathbb{E}_X [D_A(\theta_x^t, \theta_x^{t+1})] \leq \frac{\mathcal{L}(\theta^1) - \mathcal{L}(\theta^*)}{T}; \quad (26)$$

2) **Sub-linear Rate to θ^*** . If θ^1 is initialized in Θ , a locally average-convex region of $\mathcal{L}(\theta)$ containing θ^* , then

$$\mathcal{L}(\theta^T) - \mathcal{L}(\theta^*) \leq \frac{\mathbb{E}_X [D_A(\theta_x^*, \theta_x^1)]}{T} \quad (27)$$

3) **Linear Rate to θ^*** . If θ^1 is initialized in $\Theta \subseteq \Omega$, a locally average-convex region of $\mathcal{L}(\theta)$ relative to $A(\theta)$ that contains θ^* , then

$$\mathcal{L}(\theta^T) - \mathcal{L}(\theta^*) \leq (1 - \alpha)^T \mathbb{E}_X [D_A(\theta_x^*, \theta_x^1)] \quad (28)$$

The above condition on the initialization to belong in a locally average-convex region is satisfied trivially if $L(\theta_x)$ is almost surely convex relative to $A(\theta_x)$. Such an assumption is stronger, but more in line with standard sufficient conditions for optimality of MD.

As noted in the literature, EM’s convergence is sensitive to initialization. If θ^1 is initialized within a locally average-convex region of $\mathcal{L}(\theta)$, the EM iterates for the MoE problem will converge sub-linearly to the true parameter. However, if θ^1 is in a region where $\mathcal{L}(\theta)$ is strongly average-convex relative to A , the iterates will converge linearly. This last assumption is different from that of prior work, which typically require θ^1 to be initialized in a locally strongly convex region.

4.2. Discussion of Main Result

The results show that the EM update (9) for general MoE is equivalent to projected mirror descent with a unit step-size and KL divergence regularizer on the complete data distribution. We offer the following additional remarks.

First, if we have oracle access to the EM updates for w and β , EM requires no hyper-parameters, unlike GD, which is sensitive to the step size. This can be especially advantageous for cases where the β -update has a closed-form solution (as is the case for linear experts), making EM’s benefits over GD more evident. Additionally, while GD regularizes progress based on the Euclidean distance between iterates, EM adjusts progress based on the divergence between probability distributions across iterations. This is often more suitable for latent variable models, where small Euclidean changes may cause large shifts in the mixture distribution, and vice versa.

Second, whereas previous analysis of EM for various settings hinged on various types of analyses ranging from verifying obscure conditions to – less reproducible at scale – direct proofs, the connection to MD that we unveil greatly unifies the process of analysis and provides more intuition as to the inner workings of EM. In particular, as we will discuss in Section 5, our framework for analysis allows to easily provide intuitive conditions for linear convergence to the true parameters that are based on the MIM and subsequently, the SNR of the generative model.

Third, while Jordan & Xu (1995) also demonstrated that the EM algorithm for MoE converges linearly to the true parameters, the sufficient conditions they provided are more restrictive. Specifically, their analysis requires the Hessian to be negative definite, and the convergence rate depends explicitly on its eigenvalues. These conditions are similar in nature to those typically required for GD-type methods. In contrast, our sufficient conditions for optimality align with those of MD, which more accurately captures the convergence behavior of EM, as established by the equivalence shown in Theorem 4.1.

Finally, large-scale applications often favor a mini-batch training paradigm, as it tends to yield better performance for a given computational cost. Large-scale implementation

of EM can directly benefit from this paradigm for solving each iteration’s convex optimization subproblem (i.e., Equation (9)) where a GD-style method is typically used. Scaling laws for the mini-batch paradigm suggest that reducing the batch size should be accompanied by a proportional reduction in the learning rate (Shuai et al., 2024; Malladi et al., 2022; Goyal et al., 2017). However, since EM is equivalent to MD with a fixed learning rate of 1, this kind of modular tuning is not directly applicable to EM. This does not imply that there is no optimal batch size for EM. Rather, extending theoretical guarantees to stochastic and mini-batch settings can be approached through the framework of stochastic mirror descent (SMD) and mini-batch MD (MBMD), both of which have been studied in the context of composite optimization (Duchi et al., 2010). We highlight this as an open direction for future research on scalable implementations of EM for MoE.

5. Special case of SymMoLinE and SymMoLogE

In this section, we narrow in on two special cases: the Symmetric Mixture of Linear Experts (SymMoLinE) and the Symmetric Mixture of Logistic Experts (SymMoLogE) models. We first show that minimizing the EM objective function defined in (9) is exactly equivalent to minimizing the subproblem associated with a single step of MD as defined in equation (15), with a step size of $\eta = 1$. This time, we fully specify the mirror map and show it is strictly convex. Unlike the previous result, this correspondence between EM and MD does not feature the extra projection step that was present for general MoE. This allows us to more easily characterize the sufficient condition of EM for linear convergence to the true parameters, relating it to the MIM and SNR of the generative model (see Appendix B.5 and B.6). The proof is provided in Appendix B.2.

Theorem 5.1. *For SymMoLinE and SymMoLogE, there is a mirror map $A(\theta)$ such that the EM update in (9) is equivalent to*

$$\operatorname{argmin}_{\theta \in \Omega} \langle \nabla \mathcal{L}(\theta^t), \theta - \theta^t \rangle + D_A(\theta, \theta^t), \quad (29)$$

where, $\forall \phi, \theta \in \Omega$, the divergence function $D_A(\theta, \theta^t)$ is equal to the KL divergence on the complete data:

$$D_A(\phi, \theta) = KL[p(\mathbf{x}, y, z; \theta) \| p(\mathbf{x}, y, z; \phi)]. \quad (30)$$

In particular, in the case of SymMoLinE,

$$A(\theta) = \mathbb{E}_{\mathbf{X}} \left[\frac{(\mathbf{x}^\top \beta)^2}{2} + \log(1 + e^{\mathbf{x}^\top w}) \right], \quad (31)$$

while in the case of SymMoLogE,

$$A(\theta) = \mathbb{E}_{\mathbf{X}} \left[\log \left((1 + e^{\mathbf{x}^\top \beta}) (1 + e^{\mathbf{x}^\top w}) \right) \right]. \quad (32)$$

Finally, in both cases, the map $A(\theta)$ is strictly convex in θ and $\mathcal{L}(\theta)$ is 1-smooth relative to $A(\theta)$.

As shown in the proof of the result, it is evident from (41) and (45) that $p(x, y, z; \theta)$ does not belong to an exponential family of distributions for either SymMoLinE or SymMoLogE. Therefore, this result does not simply follow as a corollary of (Kunstner et al., 2021, Proposition 1), but stands as an independent finding, introducing another class of latent variable models where EM is equivalent to MD. A couple of follow-up remarks are in order.

First, as the loss is 1-smooth relative to A , this validates the choice of $\eta = 1$ for the Mirror Descent update, and subsequent convergence results from the MD literature. Specifically, in Corollary B.1 of Appendix B.4, we contextualize convergence results from (Lu et al., 2018; Kunstner et al., 2021) for SymMoLinE and SymMoLogE that feature 1) a guaranteed sub-linear rate of convergence to a stationary point at no additional assumption, 2) a sub-linear rate of convergence to the true parameter if initialized within a convex region of the loss-function that includes the true parameters, and 3) a linear rate of convergence to the true parameter if initialized within a region of the loss function that contains the true parameters and is strongly-convex relative to the mirror map. Then, in Theorem B.2 of Appendix B.5, we further characterize the assumptions required for linear convergence by relating the relative strong convexity of the objective to the eigenvalues of the MIM. Lastly, in Theorem B.4 of Appendix B.6, we characterize the existence of the local region of convergence as a function of the SNR of the generative model for the cases of SymMoLinE and SymMoLogE. We then conclude in Appendix B.7 with a discussion of the implications of the results.

6. Experiments

In this section, we empirically validate our theoretical results by comparing the performance of EM with Gradient EM (Algorithm 2), and Gradient Descent (GD, Algorithm 3). Recall that EM for MoE obtains its next parameter iterate as the concatenation to the solutions of two minimization problems. Instead, Gradient EM obtains its next parameter iterate as the concatenation of a single gradient update on the respective sub-minimization problems of EM. This differs from GD that obtains its next parameter iterate as the gradient update on the negative log-likelihood objective θ . We evaluate these methods on both a synthetic dataset and the real-world Fashion MNIST dataset, consistently reporting significant improvements for EM and Gradient EM over GD. We also provide mini-batch CIFAR-10 experiments with more than 2-experts in Appendix C.1. Note that our aim is not to achieve state-of-the-art accuracy, but to reiterate that EM can be more suitable than GD for fitting specific models.

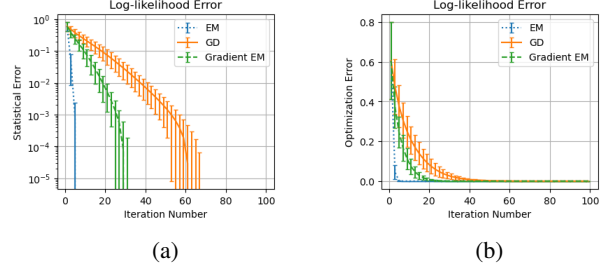


Figure 1. Convergence of objective errors $\mathcal{L}(\theta^t) - \mathcal{L}(\theta^*)$ and $\mathcal{L}(\theta^t) - \mathcal{L}(\theta^T)$ in Fig 1a and Fig 1b, respectively, averaged over 50 instances when fitting a SymMoLinE.

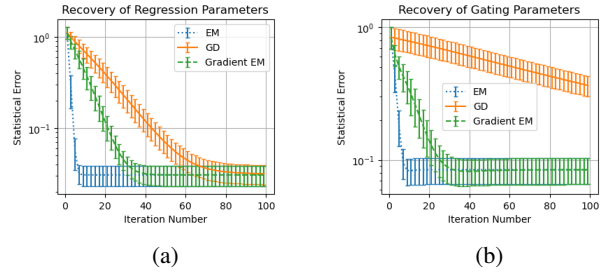


Figure 2. This figure shows the progress made towards the true parameters, $\frac{\|\beta^t - \beta^*\|_2}{\|\beta^*\|_2}$ and $\frac{\|w^t - w^*\|_2}{\|w^*\|_2}$ in figures 2a and 2b respectively, averaged over 50 instances when fitting a SymMoLinE

Synthetic Dataset. We created the synthetic dataset so as to simulate a population setting of SymMoLinE. We sampled 10^3 data points from an SymMoLinE with known additive unit Gaussian noise (i.e. $\mathcal{N}(0, 1)$) and true parameters $\beta^*, w^* \in \mathbb{R}^{10}$ that satisfy $\|\beta^*\|_2 = \|w^*\|_2 = 4$. Subsequently, we run full-batch EM, Gradient EM, and GD for 50 iterations and report the results on the training set averaged over 50 instances. Each time, re-sampling the true parameters, initial parameters, and whole dataset. The initial parameters, are randomly initialize within a neighborhood of the true parameters, and are consistent across all benchmarks.

Figure 1 shows the objective function progress. EM requires fewer iterations to fit the mixture compared to both Gradient EM and GD, with Gradient EM also outperforming GD in fitting time. Figure 2 illustrates the progress toward recovering the true SymMoLinE parameters. Once again, EM requires significantly fewer iterations to fit the mixture compared to both Gradient EM and GD, with Gradient EM also taking considerably less time than GD.

Overall, we observe that all three algorithms exhibit a linear convergence rate, both in optimizing the objective function and fitting the true parameters. This aligns with our theoretical results for MoE and is consistent with findings for Mixtures of Gaussians and Mixtures of Linear Regression

in high SNR scenarios. To validate our results further, we perform a paired t-test (Ross & Willson, 2018). For EM and Gradient EM compared to GD, we obtain a T-statistic ≥ 22 indicating that the difference in final accuracy is statistically significant (p-value ~ 0.000).

Validation Experiment on Fashion MNIST. For the small scale proof of concept experiment on Fashion MNIST (Xiao et al., 2017), we alter the dataset so as to simulate a mixture of 2 Logistic Experts. To do so, we perform an inversion transformation on the images at random with probability $\frac{1}{2}$. Effectively, the transformation inverts the images from a white article of clothing on a black background to a black article of clothing on a white background. As shown in Table 1, the single expert on the original Fashion MNIST dataset reaches an accuracy of 83.2% on the test set. Meanwhile, the single expert cannot achieve better than an accuracy of 10.2% on the altered dataset. This suggests a 2-component MoLogE is appropriate for fitting the altered dataset, so long as the ground truth partitioning is linear in image space.

The 2-component MoLogE to be trained consists of one Linear gating layer of dimension $2 \times 28 \times 28$, and 2 logistic experts of dimension $10 \times 28 \times 28$ each. We randomly initialize each linear layer to be unit norm and execute the algorithms on the same datasets and with the same initializations. For Gradient EM, the only additional code needed over GD is to define the EM Loss function appropriately, and then perform a Gradient Step on the Gating parameters and the Expert parameters separately as describe in Algorithm 2. For EM, for each iteration, we perform several gradient steps in an inner loop to approximately recover the solutions to the sub-problems described in (9). We report our findings for the full-batch iteration of the respective algorithms in Table 2 and Figure 3.

In Table 2, we report the respective final test accuracy and cross-entropy loss values after 100 iterations of EM, Gradient EM and GD for fitting a 2-component MoLogE on the altered Fashion MNIST dataset, averaged over 25 instances. We see that EM boasts a much improved final test accuracy that nearly recovers the single expert accuracy on the original unaltered Fashion MNIST dataset of 79.2%. Meanwhile, Gradient EM also registers an improvement over GD. In Figure 3, we report the progress made on the accuracy and objective function for the test set over the 100 iterations, averaged over 25 instances. As was observed in our synthetic experiment, EM takes considerably less iterations to fit the mixture than both Gradient EM and GD, where the former also takes considerably less time to fit the mixture than GD. To validate our results further, we perform a paired t-test. For EM and Gradient EM compared to GD, we obtain a T-statistic ≥ 17 indicating that the difference in final accuracy is statistically significant (p-value ~ 0.000).

Table 1. Performance for single Logistic Expert

	Accuracy	Random Invert
Single Expert	83.2%	No
Single Expert	10.2%	Yes

Table 2. Performance for 2-Component MoLogE

	Accuracy	Cross Entropy
EM	78.5%	0.827
Gradient EM	66.0%	1.29
Gradient Descent	62.4%	1.30

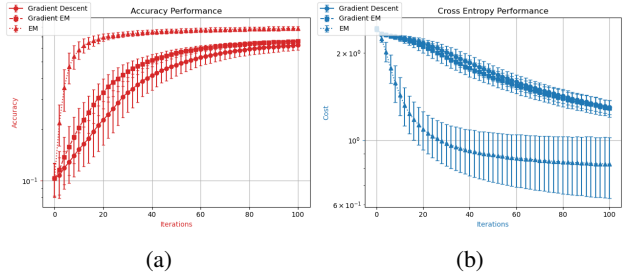


Figure 3. Test accuracy and objective function, $\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\hat{y}_i = y_i}$ and $\mathcal{L}(\theta^t)$ in 3a and 3b, respectively, averaged over 25 instances for a 2-component MoLogE train on Random Invert FMNIST.

7. Conclusion

In this paper, we theoretically addressed the problem of Mixtures of Experts (MoE) with the use of the EM algorithm. We first showed that the EM update for MoE could be interpreted as a projected Mirror Descent step on the log-likelihood with a unit step size and a KL divergence regularizer, extending the result of (Kunstner et al., 2021) beyond complete data distribution belonging to an exponential family. Building on this, we characterized different convergence rates for EM in this setting under various assumptions about the log-likelihood function and specified when these assumptions held. Lastly, we empirically observed that EM can outperform gradient descent in both convergence rate and final performance.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

Acknowledgements

This research was supported by NSF Encore Tripods (2217069), NSF’s AI Institute IFML (2019844) and NSF Grant 2007668.

References

- Balakrishnan, S., Wainwright, M. J., and Yu, B. Statistical guarantees for the EM algorithm: From population to sample-based analysis. 2017.
- Becker, P., Arenz, O., and Neumann, G. Expected information maximization: Using the i-projection for mixture density estimation. *arXiv preprint arXiv:2001.08682*, 2020.
- Chen, Z., Deng, Y., Wu, Y., Gu, Q., and Li, Y. Towards understanding mixture of experts in deep learning. *arXiv preprint arXiv:2208.02813*, 2022.
- Daskalakis, C., Tzamos, C., and Zampetakis, M. Ten steps of EM suffice for mixtures of two Gaussians. In *Conference on Learning Theory*, pp. 704–710. PMLR, 2017.
- Dempster, A. P., Laird, N. M., and Rubin, D. B. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the royal statistical society: series B (methodological)*, 39(1):1–22, 1977.
- Duchi, J. C., Shalev-Shwartz, S., Singer, Y., and Tewari, A. Composite objective mirror descent. In *Colt*, volume 10, pp. 14–26. Citeseer, 2010.
- Dwivedi, R., Khamaru, K., Wainwright, M. J., Jordan, M. I., et al. Theoretical guarantees for em under misspecified gaussian mixture models. *Advances in Neural Information Processing Systems*, 31, 2018.
- Fedus, W., Zoph, B., and Shazeer, N. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022.
- Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- Hinton, G., Deng, L., Yu, D., Dahl, G. E., Mohamed, A.-r., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath,
- T. N., et al. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal processing magazine*, 29(6):82–97, 2012.
- Jacobs, R. A., Jordan, M. I., Nowlan, S. J., and Hinton, G. E. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991.
- 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.
- Kunstner, F., Kumar, R., and Schmidt, M. Homeomorphic-invariance of em: Non-asymptotic convergence in kl divergence for exponential families via mirror descent. In *International Conference on Artificial Intelligence and Statistics*, pp. 3295–3303. PMLR, 2021.
- Kwon, J. and Caramanis, C. EM converges for a mixture of many linear regressions. In *International Conference on Artificial Intelligence and Statistics*, pp. 1727–1736. PMLR, 2020a.
- Kwon, J. and Caramanis, C. The EM algorithm gives sample-optimality for learning mixtures of well-separated gaussians. In *Conference on Learning Theory*, pp. 2425–2487. PMLR, 2020b.
- Kwon, J., Qian, W., Caramanis, C., Chen, Y., and Davis, D. Global convergence of the EM algorithm for mixtures of two component linear regression. In *Conference on Learning Theory*, pp. 2055–2110. PMLR, 2019.
- Kwon, J., Ho, N., and Caramanis, C. On the minimax optimality of the EM algorithm for learning two-component mixed linear regression. In *International Conference on Artificial Intelligence and Statistics*, pp. 1405–1413. PMLR, 2021.
- Liu, A., Feng, B., Xue, B., Wang, B., Wu, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C., et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- Lu, H., Freund, R. M., and Nesterov, Y. Relatively smooth convex optimization by first-order methods, and applications. *SIAM Journal on Optimization*, 28(1):333–354, 2018.
- Ma, J., Zhao, Z., Yi, X., Chen, J., Hong, L., and Chi, E. H. Modeling task relationships in multi-task learning with multi-gate mixture-of-experts. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 1930–1939, 2018.

- Makkuva, A., Viswanath, P., Kannan, S., and Oh, S. Breaking the gridlock in mixture-of-experts: Consistent and efficient algorithms. In *International Conference on Machine Learning*, pp. 4304–4313. PMLR, 2019.
- Malladi, S., Lyu, K., Panigrahi, A., and Arora, S. On the sdes and scaling rules for adaptive gradient algorithms. *Advances in Neural Information Processing Systems*, 35: 7697–7711, 2022.
- Orchard, T. and Woodbury, M. A. A missing information principle: theory and applications. In *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Theory of Statistics*, volume 6, pp. 697–716. University of California Press, 1972.
- Ross, A. and Willson, V. L. *Basic and advanced statistical tests: Writing results sections and creating tables and figures*. Springer, 2018.
- Shuai, X., Wang, Y., Wu, Y., Jiang, X., and Ren, X. Scaling law for language models training considering batch size. *arXiv preprint arXiv:2412.01505*, 2024.
- Tseng, P. An analysis of the EM algorithm and entropy-like proximal point methods. *Mathematics of Operations Research*, 29(1):27–44, 2004.
- Wu, C. J. On the convergence properties of the EM algorithm. *The Annals of statistics*, pp. 95–103, 1983.
- Xiao, H., Rasul, K., and Vollgraf, R. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017. URL <https://arxiv.org/abs/1708.07747>.
- 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.

Appendix

Appendix Contents

- A) EM, Projected Mirror Descent, and General MoE
 - A.1) EM is Projected Mirror Descent for General MoE.
 - A.2) Convergence Results for EM applied to General MoE.
- B) EM, Mirror Descent, and SymMoLogE and SymMoLinE.
 - B.1) EM is Mirror Descent for SymMoLogE and SymMoLinE
 - B.2) Proof of Theorem 5.1 for SymMoLinE
 - B.3) Proof of Theorem 5.1 for SymMoLogE
 - B.4) Convergence Guarantees of EM for SymMoLogE and SymMoLinE
 - B.5) Satisfiability of Conditions from Corollary B.1
 - B.6) Correspondence Between the MIM and SNR for SymMoLinE and SymMoLogE
 - B.7) Existence of Locally Convex Region
- C) Additional Experiments
 - C.1) Experiment on Grayscale CIFAR-10
- D) Algorithms
 - D.1) EM for MoE
 - D.2) Gradient EM for MoE
 - D.3) Gradient Descent for MoE
 - D.4) EM for Deep and Sparse MoE

Notations

We summarize here the notations used throughout the paper. For clarity, we distinguish between different types of mathematical objects (e.g., vectors, random variables, distributions) and follow standard conventions where possible.

The Kullback-Leibler (KL) divergence of a distribution p from a distribution q is denoted by $\text{KL}[q \| p] := \int q(x) \log(q(x)/p(x)) dx$. We use lowercase letters such as p to denote continuous probability density functions and uppercase letters such as P to denote discrete probability mass functions. The Euclidean (or ℓ_2) norm of a vector is denoted by $\|\cdot\|_2$. We use the compact notation $[k] := \{1, 2, \dots, k\}$.

We denote vectors using bold lowercase letters (e.g., \mathbf{x}), and random variables using uppercase letters (e.g., X). Bold uppercase letters (e.g., \mathbf{X}) are used to represent either vector-valued random variables or matrices; the distinction between the two is clear from context. For a matrix \mathbf{M} , we denote its i^{th} eigenvalue by λ_i , and the corresponding eigenvector by \mathbf{v}_i . The minimum and maximum eigenvalues of \mathbf{M} are denoted by λ_{\min} and λ_{\max} , respectively. We use \mathbf{I}_d to denote the $d \times d$ identity matrix, and \mathbf{e}_i to denote the i^{th} standard basis (unit) vector in \mathbb{R}^d .

Expectations are written as $\mathbb{E}_{\mathbf{X}}[f(\mathbf{X})] = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x}$, where the distribution of \mathbf{X} is implicitly defined by the context. When needed, we make the dependence on parameters explicit by writing $\mathbf{X}; \boldsymbol{\theta}$, where $\boldsymbol{\theta}$ denotes the parameters of the distribution of \mathbf{X} . The notation $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ refers to the multivariate normal distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$.

A. EM, Projected Mirror Descent, and General MoE.

In this Section, we provide all the complete proofs and discussions relating to results from Section 4.

A.1. EM is Projected Mirror Descent for General MoE.

In this section, we provide the full and detailed proof of the main result, Theorem 4.1. For ease of comprehension and in the hope that this will provide useful insights into other types of non-exponential family mixtures for which EM is also connected to MD, we prove our result following the same general ideas as that of [Kunstner et al. \(2021, Proposition 1\)](#).

For ease of reading, we re-state the theorem below:

Theorem 4.1: For General MoE, there exists a natural re-parameterization $\theta_x \in \{\eta(\cdot, \theta) : \theta \in \Omega\}$ with

$$\mathcal{L}(\theta) = \mathbb{E}_X [L(\theta_x)]$$

and a mirror map $A(\theta_x)$ such that the EM update in (9) simplifies and is equivalent to the expectation moment projection,

$$\operatorname{argmin}_{\theta \in \Omega} \mathbb{E}_X \left[\text{KL} \left[p(y, z | \tilde{\theta}_x^{t+1}) \parallel p(y, z | \eta(x, \theta)) \right] \right],$$

where for each x , $\tilde{\theta}_x^{t+1}$ is obtained from the following MD step,

$$\operatorname{argmin}_{\psi \in \tilde{\Omega}} \langle \nabla L(\theta_x^t), \psi - \theta_x^t \rangle + D_A(\psi, \theta_x^t),$$

with $L(\theta_x)$ being 1-smooth relative to $A(\theta_x)$. Further, $\forall \psi_1, \psi_2 \in \tilde{\Omega}$, the divergence function $D_A(\psi_1, \psi_2)$ is equal to the KL divergence on $p(y, z, |\psi)$:

$$D_A(\psi_1, \psi_2) = \text{KL}[p(y, z | \psi_2) \parallel p(y, z | \psi_1)].$$

Proof. The EM is centered around iterative minimization of the surrogate upper-bound $Q(\phi | \theta)$. For conditionally exponential family of distribution, We can decompose it in terms of the sufficient statistic and log-partition:

$$\begin{aligned} Q(\theta | \theta^t) &= -\mathbb{E}_{Y,X} \left[\sum_z \ln(p(y, x, z; \theta)) P(z | y, x; \theta^t) \right] \\ &= -\mathbb{E}_{Y,X} \left[\sum_z (\ln(p(y, z | x; \theta)) + \ln(p(x))) P(z | y, x; \theta^t) \right] \\ &= -\mathbb{E}_{Y,X} \left[\sum_z (\langle S(y, z), \theta_x \rangle - A(\theta_x^t)) + \ln(p(x)) \right] P(z | y, x; \theta^t) \\ &= \mathbb{E}_X \left[-\underbrace{\langle \mathbb{E}_{Y|x, \theta^*} \mathbb{E}_{Z|y, x, \theta^t} [S(y, z)], \theta_x \rangle}_{s(x; \theta^t)} + A(\theta_x) - \ln(p(x)) \right] \\ &= \mathbb{E}_X [-\langle s(x; \theta^t), \theta_x \rangle + A(\theta_x) - \ln(p(x))]. \end{aligned}$$

Therefore, the above also implies

$$\mathbb{E}_X [\nabla L(\theta_x^t)] = \mathbb{E}_X [\nabla Q(\theta_x^t | \theta_x^t)] \quad (33)$$

$$= \mathbb{E}_X [\nabla A(\theta_x) - s(x; \theta^t)] \quad (34)$$

Simple algebra then shows that

$$\begin{aligned} Q(\theta | \theta^t) - Q(\theta^t | \theta^t) &= \mathbb{E}_X [-\langle s(x; \theta^t), \theta_x \rangle + A(\theta_x) + \langle s(x; \theta^t), \theta_x^t \rangle - A(\theta_x^t)] \\ &= \mathbb{E}_X [-\langle s(x; \theta^t), \theta_x - \theta_x^t \rangle + A(\theta_x) - A(\theta_x^t)] \\ &\stackrel{i)}{=} \mathbb{E}_X [-\langle s(x; \theta^t) - \nabla A(\theta_x^t), \theta_x - \theta_x^t \rangle + D_A(\theta_x, \theta_x^t)] \\ &= \mathbb{E}_X [\langle \nabla L(\theta_x^t), \theta_x - \theta_x^t \rangle + D_A(\theta_x, \theta_x^t)] \end{aligned}$$

where i) adds and subtracts $\langle \nabla A(\theta_x^t), \theta_x - \theta_x^t \rangle$. This is especially important as EM minimizes $Q(\theta|\theta^t) - Q(\theta^t|\theta^t)$ in each iteration. Now recall that $\mathcal{L}(\theta) = Q(\theta|\theta^t) - H(\theta|\theta^t)$ and $H(\theta^t|\theta^t) - H(\theta|\theta^t) \leq 0$, it follows that

$$\begin{aligned} \mathcal{L}(\theta) - \mathcal{L}(\theta^t) &= Q(\theta|\theta^t) - Q(\theta^t|\theta^t) - H(\theta|\theta^t) - H(\theta^t|\theta^t) \\ &\leq Q(\theta|\theta^t) - Q(\theta^t|\theta^t) \\ &= \mathbb{E}_X [\langle \nabla L(\theta_x^t), \theta_x - \theta_x^t \rangle + D_A(\theta_x, \theta_x^t)] \end{aligned}$$

Thus far, we have shown that the EM iteration under the considered setting is equivalent to minimizing the following upper bound on $\mathcal{L}(\theta)$, w.r.t. θ :

$$\mathcal{L}(\theta^t) + \mathbb{E}_X [\langle \nabla L(\theta_x^t), \theta_x - \theta_x^t \rangle + D_A(\theta_x, \theta_x^t)] \quad (35)$$

Now, recall $\tilde{\theta}_x^{t+1} := \operatorname{argmin}_{\theta_x} \langle \nabla L(\theta_x^t), \theta_x - \theta_x^t \rangle + D_A(\theta_x, \theta_x^t)$ which is the outcome of a mirror descent step. Differentiating and setting equal to 0, it holds that

$$\nabla A(\tilde{\theta}_x^{t+1}) = s(x; \theta_x^t). \quad (36)$$

Using this and decomposing $\nabla L(\theta_x)$ above, we see that (35) is equal to

$$\begin{aligned} &= \mathcal{L}(\theta^t) + \mathbb{E}_X [\langle \nabla A(\theta_x^t) - s(x; \theta_x^t), \theta_x - \theta_x^t \rangle + D_A(\theta_x, \theta_x^t)] \\ &= \mathcal{L}(\theta^t) + \mathbb{E}_X [\langle \nabla A(\theta_x^t) - \nabla A(\tilde{\theta}_x^{t+1}), \theta_x - \theta_x^t \rangle + D_A(\theta_x, \theta_x^t)] \\ &= \mathcal{L}(\theta^t) + \mathbb{E}_X [-\langle \nabla A(\tilde{\theta}_x^{t+1}), \theta_x - \theta_x^t \rangle + A(\theta_x) - A(\theta_x^t)]. \end{aligned}$$

Thus, minimizing (35) with respect to θ_x is equivalent to minimizing

$$\mathbb{E}_X [D_A(\theta_x, \tilde{\theta}_x^{t+1})]. \quad (37)$$

Substituting the Bregman Divergence induced by A by the KL divergence yields the claim.

It remains to verify that $D_A(\theta_x, \tilde{\theta}_x^{t+1}) = \text{KL} [p(y, z|\tilde{\theta}_x^{t+1})||p(y, z|\theta_x)]$ and that the function $L(\theta_x)$ is 1-smooth relative to $A(\theta_x)$. This follows directly from previous work by [Kunstner et al. \(2021\)](#) since $L(\theta_x)$ is the expected negative log-likelihood of $y|x$ where $A(\theta_x)$ is the log-partition of the exponential distribution $p(y, z|\theta_x)$. For completeness, the derivation is as follows. For $Y, Z|x$ belonging to an exponential family of distribution, the KL divergence can be decomposed directly using (36) as follows to obtain the Bregman Divergence:

$$\begin{aligned} \text{KL} [p(y, z|\tilde{\theta}_x^{t+1})||p(y, z|\theta_x)] &:= \mathbb{E}_{Y, Z|x; \tilde{\theta}_x^{t+1}} \left[\log \left(\frac{p(y, z; \tilde{\theta}_x^{t+1})}{p(y, z; \theta_x)} \right) \right] \\ &= \mathbb{E}_{Y, Z|x; \tilde{\theta}_x^{t+1}} [\langle S(y, z), \tilde{\theta}_x^{t+1} - \theta_x \rangle + A(\theta_x) - A(\tilde{\theta}_x^{t+1})] \\ &= \langle \nabla A(\tilde{\theta}_x^{t+1}), \tilde{\theta}_x^{t+1} - \theta_x \rangle + A(\theta_x) - A(\tilde{\theta}_x^{t+1}) \\ &= D_A(\theta_x, \tilde{\theta}_x^{t+1}). \end{aligned}$$

□

A.2. Convergence Results for EM applied to General MoE.

In this section, we provide the full proof of Theorem 4.2. Before we begin, we recall and discuss the regularity assumptions previously made. Recall assumptions A_1 , A_2 , and A_3 :

- A_1 The conditional distribution $p(y, z|\theta_x)$ is a steep, minimal exponential family of distribution and $\eta(x, \cdot)$ is a continuously differentiable function,
- A_2 The optimal objective function values is bounded below, i.e., $\mathcal{L}(\theta^*) > -\infty$, on the constraint set Ω ,
- A_3 The following sub-level sets $\Omega_\theta := \{\phi \in \Omega : Q(\phi|\theta) \leq Q(\theta|\theta)\}$ are compact. .

Assumption A_1 serves to ensure that $\mathcal{L}(\theta) = \mathbb{E}_X [L(\theta_x)]$ is differentiable, that the EM surrogate is also differentiable and has a solution. It further serves to guarantee the mirror map A is smooth, ensuring that projecting into the dual space is well-defined. We note that if the re-parametrization function is continuously differentiable in θ , it will hold that A_1 is satisfied for the popular case that $p(y, z|x)$ is a Gaussian mixture (this includes MoE with Gaussian experts). Next, A_2 and A_3 are classical optimization assumptions that serve to guarantee the solution of the M-step is unique and exists within the constraint set Ω , thereby ensuring the EM iterations are well defined.

Further, we make the additional remark that the projection, (17), in Theorem 4.1 can be seen to be equivalent to the following projection over the space of functions on x .

$$\eta(x, \theta^{t+1}) = \theta_x^{t+1} = \underset{\phi_x \in \{\eta(\cdot, \theta) : \theta \in \Omega\}}{\operatorname{argmin}} \mathbb{E}_X [D_A(\phi_x, \tilde{\theta}_x^{t+1})] \quad (38)$$

We can see that the set $\{\eta(\cdot, \theta) : \theta \in \Omega\}$ is not necessarily guaranteed to be convex. Such a result would require the re-parametrization function $\eta(x, \theta)$ to be affine. In situations where this set is not convex, we cannot take our weak generalized Pythagorean identity (25) for granted, and thus we include it as an extra assumption to be satisfied for these convergence results. Still, convexity is not necessary for our generalized Pythagorean inequality to hold. For instance, ensuring $\tilde{\theta}_x^{t+1}$ is in the relative interior of $\{\eta(\cdot, \theta) : \theta \in \Omega\}$, or simply satisfying the inequality in expectation will suffice, but may be difficult to show.

For ease of reading, we re-state the result below:

Theorem 4.2: For general MoE with re-parameterization given by $\theta_x := \eta(x, \theta)$, strictly convex mirror map $A(\theta_x)$, and if for all $\tilde{\theta}_x^{t+1}, \theta_x^{t+1}, \phi_x \in \{\theta_x, \theta_x^*\}$,

$$\mathbb{E}_X [D_A(\phi_x, \tilde{\theta}_x^{t+1})] \geq \mathbb{E}_X [D_A(\theta_x^{t+1}, \tilde{\theta}_x^{t+1}) + D_A(\phi_x, \theta_x^{t+1})],$$

then, the EM iterates $\{\theta^t\}_{t \in [T]}$ satisfy:

- 1) **Stationnarity.** For no additional conditions,

$$\min_{t \in [T]} \mathbb{E}_X [D_A(\theta_x^t, \theta_x^{t+1})] \leq \frac{\mathcal{L}(\theta^1) - \mathcal{L}(\theta^*)}{T};$$

- 2) **Sub-linear Rate to θ^* .** If θ^1 is initialized in Θ , a locally average-convex region of $\mathcal{L}(\theta)$ containing θ^* , then

$$\mathcal{L}(\theta^T) - \mathcal{L}(\theta^*) \leq \frac{\mathbb{E}_X [D_A(\theta_x^*, \theta_x^1)]}{T}$$

- 3) **Linear Rate to θ^* .** If θ^1 is initialized in $\Theta \subseteq \Omega$, a locally average-strongly convex region of $\mathcal{L}(\theta)$ relative to $A(\theta)$ that contains θ^* , then

$$\mathcal{L}(\theta^T) - \mathcal{L}(\theta^*) \leq (1 - \alpha)^T \mathbb{E}_X [D_A(\theta_x^*, \theta_x^1)]$$

Proof. The proof is divided into three parts that correspond to each of the three sub-results of the corollary.

To aid the reader's comprehension, we re-state the cosine law for Bregman divergence (also known as 3-point lemma).

Lemma A.1 (cosine law for Bregman divergence). *Assume the mapping A is proper and convex. Then, for all $a, b, c \in \tilde{\Omega}$, it holds that*

$$D_A(a, b) = D_A(a, c) + D_A(c, b) - \langle \nabla A(b) - \nabla A(c), a - c \rangle. \quad (39)$$

Part 1): Stationarity.

We begin by utilizing the result from Theorem 4.1 that the conditional log-likelihood $L(\theta_x)$ is 1-smooth relative to the mirror map $A(\theta_x)$. For all $\theta_x \in \tilde{\Omega}$,

$$L(\theta_x^{t+1}) \leq L(\theta_x^t) + \langle \nabla L(\theta_x^t), \theta_x^{t+1} - \theta_x^t \rangle + D_A(\theta_x^{t+1}, \theta_x^t).$$

Taking the expectation on both sides with respect to the random feature variable x yields:

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^t) + \mathbb{E}_X [\langle \nabla L(\theta_x^t), \theta_x^{t+1} - \theta_x^t \rangle + D_A(\theta_x^{t+1}, \theta_x^t)].$$

Combining (33) and (36) then plugging into the above, we obtain

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^t) + \mathbb{E}_X \left[\langle \nabla A(\theta_x^t) - \nabla A(\tilde{\theta}_x^{t+1}), \theta_x^{t+1} - \theta_x^t \rangle + D_A(\theta_x^{t+1}, \theta_x^t) \right].$$

Decomposing $D_A(\theta_x^{t+1}, \theta_x^t)$ above and canceling appropriate terms yields

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^t) + \underbrace{\mathbb{E}_X \left[-\langle \nabla A(\tilde{\theta}_x^{t+1}), \theta_x^{t+1} + A(\theta_x^{t+1}) - A(\theta_x^t) \rangle \right]}_{i)}$$

It can then be checked that i) is equal to $\mathbb{E}_X \left[D_A(\theta_x^{t+1}, \tilde{\theta}_x^{t+1}) - D_A(\theta_x^t, \tilde{\theta}_x^{t+1}) \right]$. Therefore, substituting the above and utilizing the inequality (25), it follows that

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^t) - \mathbb{E}_X \left[D_A(\theta_x^t, \theta_x^{t+1}) \right].$$

Re-arranging the terms and averaging over T -iterations yields the claim:

$$\begin{aligned} \mathbb{E}_X \left[D_A(\theta_x^t, \theta_x^{t+1}) \right] &\leq \mathcal{L}(\theta^{t+1}) - \mathcal{L}(\theta^t) \\ \implies \min_{t \leq T} \mathbb{E}_X \left[D_A(\theta_x^t, \theta_x^{t+1}) \right] &\leq \frac{1}{T} \sum_{t=1}^T \mathcal{L}(\theta^{t+1}) - \mathcal{L}(\theta^t) = \frac{\mathcal{L}(\theta^1) - \mathcal{L}(\theta^T)}{T} \end{aligned}$$

Part 2): Sub-linear rate to θ^* .

We begin by utilizing the result from Theorem 4.1 that the conditional log-likelihood $L(\theta_x)$ is 1-smooth relative to the mirror map $A(\theta_x)$. For all $\theta_x \in \tilde{\Omega}$, then apply expectation with respect to x on both sides yielding:

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^t) + \mathbb{E}_X \left[\langle \nabla L(\theta_x^t), \theta_x^{t+1} - \theta_x^t \rangle + D_A(\theta_x^{t+1}, \theta_x^t) \right].$$

We then add and subtract $\mathbb{E}_X \left[\langle \nabla L(\theta_x^t), \theta_x^* \rangle \right]$ to obtain

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^t) + \mathbb{E}_X \left[\langle \nabla L(\theta_x^t), \theta_x^{t+1} - \theta_x^* + \theta_x^* - \theta_x^t \rangle + D_A(\theta_x^{t+1}, \theta_x^t) \right].$$

We then use the average local convexity assumption, (23), and obtain

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^*) + \mathbb{E}_X \left[\langle \nabla L(\theta_x^t), \theta_x^{t+1} - \theta_x^* \rangle + D_A(\theta_x^{t+1}, \theta_x^t) \right].$$

Combining (33) and (36) then plugging into the above, we obtain

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^*) + \mathbb{E}_X \left[\langle \nabla A(\theta_x^t) - \nabla A(\tilde{\theta}_x^{t+1}), \theta_x^{t+1} - \theta_x^* \rangle + D_A(\theta_x^{t+1}, \theta_x^t) \right]. \quad (40)$$

We now decompose $D_A(\theta_x^{t+1}, \theta_x^t)$ using Lemma A.1 with $a = \tilde{\theta}_x^{t+1}$, $b = \theta_x^t$, and $c = \theta_x^{t+1}$ and obtain

$$\begin{aligned} &\langle \nabla A(\theta_x^t) - \nabla A(\tilde{\theta}_x^{t+1}), \theta_x^{t+1} - \theta_x^* \rangle + D_A(\theta_x^{t+1}, \theta_x^t) \\ &= \langle \nabla A(\theta_x^t) - \nabla A(\tilde{\theta}_x^{t+1}), \theta_x^{t+1} - \theta_x^* \rangle + D_A(\tilde{\theta}_x^{t+1}, \theta_x^t) - D_A(\tilde{\theta}_x^{t+1}, \theta_x^{t+1}) + \langle \nabla A(\theta_x^t) - \nabla A(\theta_x^{t+1}), \tilde{\theta}_x^{t+1} - \theta_x^{t+1} \rangle \\ &= A(\tilde{\theta}_x^{t+1}) - A(\theta_x^t) + \langle \nabla A(\theta_x^t), \theta_x^t - \theta_x^* \rangle + \langle -\nabla A(\tilde{\theta}_x^{t+1}), \theta_x^{t+1} - \theta_x^* \rangle + A(\theta_x^{t+1}) - A(\tilde{\theta}_x^{t+1}) \end{aligned}$$

We now add and subtract $\langle \nabla A(\theta_x^t), \tilde{\theta}_x^{t+1} \rangle$ and $\langle \nabla A(\tilde{\theta}_x^{t+1}), \tilde{\theta}_x^{t+1} \rangle$ and group terms to obtain

$$D_A(\theta_x^{t+1}, \tilde{\theta}_x^{t+1}) + \underbrace{D_A(\tilde{\theta}_x^{t+1}, \theta_x^t) + \langle \nabla A(\theta_x^t) - \nabla A(\tilde{\theta}_x^{t+1}), \tilde{\theta}_x^{t+1} - \theta_x^* \rangle}_{ii)}$$

We now apply Lemma A.1 again to ii) with $a = \theta_x^*$, $b = \theta_x^t$, $c = \tilde{\theta}_x^{t+1}$ and obtain the sub-result:

$$\langle \nabla A(\theta_x^t) - \nabla A(\tilde{\theta}_x^{t+1}), \theta_x^{t+1} - \theta_x^* \rangle + D_A(\theta_x^{t+1}, \theta_x^t) = D_A(\theta_x^*, \theta_x^t) - D_A(\theta_x^*, \tilde{\theta}_x^{t+1}) + D_A(\theta_x^{t+1}, \tilde{\theta}_x^{t+1})$$

Plugging the above equality into (40), we obtain

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^*) + \mathbb{E}_X \left[D_A(\theta_x^*, \theta_x^t) - D_A(\theta_x^*, \tilde{\theta}_x^{t+1}) + D_A(\theta_x^{t+1}, \tilde{\theta}_x^{t+1}) \right].$$

Then, using 25, we obtain

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^*) + \mathbb{E}_X \left[D_A(\theta_x^*, \theta_x^t) - D_A(\theta_x^*, \theta_x^{t+1}) \right].$$

Re-arranging, then averaging over T iterations yields the claim:

$$\begin{aligned} T(\mathcal{L}(\theta^T) - \mathcal{L}(\theta^*)) &\leq \sum_{t=1}^T (\mathcal{L}(\theta^t) - \mathcal{L}(\theta^*)) \leq \sum_{t=1}^T \mathbb{E}_X [D_A(\theta_x^*, \theta_x^t) - D_A(\theta_x^*, \theta_x^{t+1})] \\ \implies \mathcal{L}(\theta^T) - \mathcal{L}(\theta^*) &\leq \frac{\mathbb{E}_X [D_A(\theta_x^*, \theta_x^1) - D_A(\theta_x^*, \theta_x^T)]}{T} \leq \frac{\mathbb{E}_X [D_A(\theta_x^*, \theta_x^1)]}{T} \end{aligned}$$

Part 3): Linear rate to θ^* .

We begin by utilizing the result from Theorem 4.1 that the conditional log-likelihood $L(\theta_x)$ is 1-smooth relative to the mirror map $A(\theta_x)$. For all $\theta_x \in \hat{\Omega}$, then apply expectation with respect to x on both sides yielding:

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^t) + \mathbb{E}_X [\langle \nabla L(\theta_x^t), \theta_x^{t+1} - \theta_x^t \rangle + D_A(\theta_x^{t+1}, \theta_x^t)].$$

We then add and subtract $\mathbb{E}_X [\langle \nabla L(\theta_x^t), \theta_x^* \rangle]$ to obtain

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^t) + \mathbb{E}_X [\langle \nabla L(\theta_x^t), \theta_x^{t+1} - \theta_x^* + \theta_x^* - \theta_x^t \rangle + D_A(\theta_x^{t+1}, \theta_x^t)].$$

We then use the local α -strongly average-convexity assumption, (23), and obtain

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^*) + \mathbb{E}_X [\langle \nabla L(\theta_x^t), \theta_x^{t+1} - \theta_x^* \rangle + D_A(\theta_x^{t+1}, \theta_x^t) - \alpha D_A(\theta_x^*, \theta_x^t)].$$

Then, following the same steps as for the sub-linear case, we Combining (33) and (36), utilize the cosine law for Bregman Divergence $D_A(\cdot, \cdot)$ twice, then apply (25) to obtain:

$$\begin{aligned} \mathcal{L}(\theta^{t+1}) &\leq \mathcal{L}(\theta^*) + \mathbb{E}_X [D_A(\theta_x^*, \theta_x^t) - D_A(\theta_x^*, \theta_x^{t+1}) - \alpha D_A(\theta_x^*, \theta_x^t)] \\ &\leq \mathcal{L}(\theta^*) + \mathbb{E}_X [(1 - \alpha) D_A(\theta_x^*, \theta_x^t)]. \end{aligned}$$

Unraveling the recurrence over T iterations, yields the result:

$$\mathcal{L}(\theta^T) - \mathcal{L}(\theta^*) \leq (1 - \alpha)^T \mathbb{E}_X [D_A(\theta_x^*, \theta_x^1)]$$

This completes the proof. □

B. EM, Mirror Descent, and SymMoLogE and SymMoLinE.

In this section, we provide all results, proofs, and discussion pertaining to EM for symmetric mixtures of logistic or linear experts.

B.1. EM is Mirror Descent for SymMoLogE and SymMoLinE

In this section, we provide the full and detailed proof of Theorem 5.1. For ease of comprehension and in the hope that this will provide useful insights into other types of non-exponential family mixtures for which EM is also connected to MD, we prove our result following the same general ideas as that of (Kunstner et al., 2021, Proposition 1). We split the proof into two parts (SymMoLinE and SymMoLogE) which can be found in Appendix B.2 and B.3.

For ease of reading, we re-state the theorem below:

Theorem 5.1: For SymMoLinE and SymMoLogE, there is a mirror map $A(\theta)$ such that the EM update in (9) simplifies and is equivalent to

$$\underset{\theta \in \Omega}{\operatorname{argmin}} \langle \nabla \mathcal{L}(\theta), \theta - \theta^t \rangle + D_A(\theta, \theta^t),$$

where $\forall \phi, \theta \in \Omega$ the divergence function $D_A(\theta, \theta^t)$ is equal to the KL divergence on the complete data:

$$D_A(\phi, \theta) = \text{KL}[p(\mathbf{x}, y, z; \theta) \| p(\mathbf{x}, y, z; \phi)].$$

In particular, in the case of SymMoLinE,

$$A(\theta) = \mathbb{E}_{\mathbf{X}} \left[\frac{(\mathbf{x}^\top \beta)^2}{2} + \log(1 + e^{\mathbf{x}^\top \mathbf{w}}) \right],$$

while in the case of SymMoLogE,

$$A(\theta) = \mathbb{E}_{\mathbf{X}} \left[\log \left((1 + e^{\mathbf{x}^\top \beta}) (1 + e^{\mathbf{x}^\top \mathbf{w}}) \right) \right].$$

Finally, in both cases, the map $A(\theta)$ is strictly convex in θ and $\mathcal{L}(\theta)$ is 1-smooth relative to $A(\theta)$.

B.2. Proof of Theorem 5.1 for SymMoLinE

Proof. Recall that we consider a 2 component SymMoLinE (see Section 3.1) where $z \in \{-1, 1\}$ is the latent unobserved variable, and

- 1) $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$,
- 2) $P(z|\mathbf{x}; \mathbf{w}) = \frac{\exp\{\frac{z+1}{2} \mathbf{x}^\top \mathbf{w}\}}{1 + e^{\mathbf{x}^\top \mathbf{w}}}$,
- 3) $p(y|\mathbf{x}, z; \beta) = \frac{\exp\left\{-\frac{(y - z\mathbf{x}^\top \beta)^2}{2}\right\}}{\sqrt{2\pi}}$.

We begin by deriving a near exponential form of the complete data probability density function $p(\mathbf{x}, z, y; \theta)$:

$$\begin{aligned} p(\mathbf{x}, y, z; \theta) &= p(y|\mathbf{x}, z; \theta) P(z|\mathbf{x}; \theta) p(\mathbf{x}) \\ &= \exp\{\log p(y|\mathbf{x}, z; \beta) + \log P(z|\mathbf{x}; \mathbf{w}) + \log p(\mathbf{x})\} \\ &= \exp\left\{\frac{-(y - z\mathbf{x}^\top \beta)^2}{2} - \frac{1}{2} \log(2\pi) + \left(\frac{z+1}{2}\right) \mathbf{x}^\top \mathbf{w} - \log(1 + e^{\mathbf{x}^\top \mathbf{w}}) + \log p(\mathbf{x})\right\} \\ &= \exp\left\{\frac{-y^2}{2} + yz\mathbf{x}^\top \beta - \frac{z^2(\mathbf{x}^\top \beta)^2}{2} + \left(\frac{z+1}{2}\right) \mathbf{x}^\top \mathbf{w} - \log(1 + e^{\mathbf{x}^\top \mathbf{w}}) + \log p(\mathbf{x}) - \frac{1}{2} \log(2\pi)\right\} \\ &= \exp\left\{\left\langle \begin{bmatrix} yz\mathbf{x} \\ \frac{z\mathbf{x}}{2} \end{bmatrix}, \begin{bmatrix} \beta \\ \mathbf{w} \end{bmatrix} \right\rangle + \frac{\mathbf{x}^\top \mathbf{w}}{2} - \frac{(\mathbf{x}^\top \beta)^2}{2} - \log(1 + e^{\mathbf{x}^\top \mathbf{w}}) + \log p(\mathbf{x}) - \frac{y^2}{2} - \frac{1}{2} \log(2\pi)\right\}. \end{aligned}$$

Thus we have recovered the decomposition,

$$p(\mathbf{x}, y, z; \boldsymbol{\theta}) = \exp \left\{ \left\langle \underbrace{\begin{bmatrix} \frac{z\mathbf{x}}{2} \\ yz\mathbf{x} \end{bmatrix}}_{S(\mathbf{x}, y, z)}, \begin{bmatrix} \mathbf{w} \\ \boldsymbol{\beta} \end{bmatrix} \right\rangle + a(\mathbf{x}, y, \boldsymbol{\theta}) \right\}, \quad (41)$$

where in $a(\mathbf{x}, y, \boldsymbol{\theta})$, the feature variable \mathbf{x} cannot be linearly separated from the parameter $\boldsymbol{\theta}$:

$$a(\mathbf{x}, y, \boldsymbol{\theta}) = \frac{\mathbf{x}^\top \mathbf{w}}{2} - \frac{(\mathbf{x}^\top \boldsymbol{\beta})^2}{2} - \log(1 + e^{\mathbf{x}^\top \mathbf{w}}) + \log p(\mathbf{x}) - \frac{y^2}{2} - \frac{1}{2} \log(2\pi).$$

At this point, we pause and discuss the implications of the obtained form. First we recall that for a random variable \mathbf{U} to belong to an exponential family, it must satisfy

$$p(\mathbf{u}; \boldsymbol{\theta}) = h(\mathbf{u}) \exp \{ \langle s(\mathbf{u}), \boldsymbol{\theta} \rangle - A(\boldsymbol{\theta}) \}$$

for some $h(\cdot)$, $s(\cdot)$, $\boldsymbol{\theta}$, $A(\cdot)$ that are called the normalization function, sufficient statistics, natural parameters, and log-partition function respectively. To clarify, we note that 1) $A(\boldsymbol{\theta})$ must be a function of the parameters only and cannot depend on \mathbf{u} and 2) $h(\mathbf{u})$ must be a function of the variable \mathbf{u} only and cannot depend on the parameters. In other words, it must be that $h(\mathbf{u})$ and $A(\boldsymbol{\theta})$ are linearly separated inside the exp. With the above in mind, we see that SymMoLinE is not an exponential family. Further, we remark that the above formulation showing $S(\mathbf{x}, y, z)$ is linear with $(\mathbf{w}, \boldsymbol{\beta})^\top$ does not extend beyond the **symmetric** setting of the Mixture of Linear Experts; for $k \geq 3$, this relationship becomes non-linear. This turns out to be problematic for showing EM is equivalent to MD for $k \geq 3$. Lastly, note that taking the expectation of $a(\mathbf{x}, y, \boldsymbol{\theta})$ over $(\mathbf{x}, y) \sim p(\mathbf{x}, y; \boldsymbol{\theta}^*)$ yields,

$$\begin{aligned} \mathbb{E}_{\mathbf{X}, Y}[a(\mathbf{x}, y, \boldsymbol{\theta})] &= -\mathbb{E}_{\mathbf{X}} \left[\frac{(\mathbf{x}^\top \boldsymbol{\beta})^2}{2} + \log(1 + e^{\mathbf{x}^\top \mathbf{w}}) \right] - C \\ &= -A(\boldsymbol{\theta}) - C, \end{aligned}$$

where the above follows from $\mathbb{E}_{\mathbf{X}}[\frac{\mathbf{x}^\top \mathbf{w}}{2}] = 0$ and $C := -\mathbb{E}_{\mathbf{X}, Y}[\log p(\mathbf{x}) - \frac{y^2}{2} - \frac{1}{2} \log(2\pi)]$ is not a function of the parameter $\boldsymbol{\theta}$. With the obtained form (41), we now continue with the proof.

Part a): Show EM is MD, i.e., $\operatorname{argmin}_{\boldsymbol{\theta} \in \Omega} Q(\boldsymbol{\theta} | \boldsymbol{\theta}^t) = \operatorname{argmin}_{\boldsymbol{\theta} \in \Omega} \langle \nabla \mathcal{L}(\boldsymbol{\theta}), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle + D_A(\boldsymbol{\theta}, \boldsymbol{\theta}^t)$.

Taking the appropriate expectation, the EM objective Q can be written as

$$\begin{aligned} Q(\boldsymbol{\theta} | \boldsymbol{\theta}^t) &= -\mathbb{E}_{\mathbf{X}, Y} [\mathbb{E}_{Z | \mathbf{x}, y; \boldsymbol{\theta}^t} [\log p(\mathbf{x}, y, z; \boldsymbol{\theta})]] \\ &= -\mathbb{E}_{\mathbf{X}, Y} [\mathbb{E}_{Z | \mathbf{x}, y; \boldsymbol{\theta}^t} [\langle S(\mathbf{x}, y, z), \boldsymbol{\theta} \rangle + a(\mathbf{x}, y, \boldsymbol{\theta})]] \\ &= -\mathbb{E}_{\mathbf{X}, Y} [a(\mathbf{x}, y, \boldsymbol{\theta})] - \mathbb{E}_{\mathbf{X}, Y} [\mathbb{E}_{Z | \mathbf{x}, y; \boldsymbol{\theta}^t} [\langle S(\mathbf{x}, y, z), \boldsymbol{\theta} \rangle]] \\ &= A(\boldsymbol{\theta}) - \langle s(\boldsymbol{\theta}^t), \boldsymbol{\theta} \rangle + C \end{aligned}$$

where $s(\boldsymbol{\theta}^t) := \mathbb{E}_{\mathbf{X}, Y} \mathbb{E}_{Z | \mathbf{x}, y; \boldsymbol{\theta}^t} [S(\mathbf{x}, y, z)]$. As a consequence, it is also true that

$$\nabla Q(\boldsymbol{\theta}^t | \boldsymbol{\theta}^t) = \nabla A(\boldsymbol{\theta}^t) - s(\boldsymbol{\theta}^t). \quad (42)$$

Continuing, we use the above to simplify the expression for $Q(\boldsymbol{\theta} | \boldsymbol{\theta}^t) - Q(\boldsymbol{\theta}^t | \boldsymbol{\theta}^t)$ that will subsequently give us the MD loss:

$$\begin{aligned} Q(\boldsymbol{\theta} | \boldsymbol{\theta}^t) - Q(\boldsymbol{\theta}^t | \boldsymbol{\theta}^t) &= A(\boldsymbol{\theta}) - \langle s(\boldsymbol{\theta}^t), \boldsymbol{\theta} \rangle - A(\boldsymbol{\theta}^t) + \langle s(\boldsymbol{\theta}^t), \boldsymbol{\theta}^t \rangle \\ &= -\langle s(\boldsymbol{\theta}^t), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle + \langle \nabla A(\boldsymbol{\theta}^t), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle - \langle \nabla A(\boldsymbol{\theta}^t), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle + A(\boldsymbol{\theta}) - A(\boldsymbol{\theta}^t) \\ &\stackrel{i)}{=} \langle \nabla Q(\boldsymbol{\theta}^t | \boldsymbol{\theta}^t), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle + D_A(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \\ &\stackrel{ii)}{=} \langle \nabla \mathcal{L}(\boldsymbol{\theta}^t), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle + D_A(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \end{aligned}$$

where we first adding and subtracting $\langle \nabla A(\theta^t), \theta - \theta^t \rangle$ then *i*) follows from (42) and *ii*) follows from $\nabla \mathcal{L}(\theta) = \nabla Q(\theta|\theta)$ (see Section 3 for the derivation).

Finally, the first part of our result follows trivially as

$$\operatorname{argmin}_{\theta \in \Omega} Q(\theta|\theta^t) = \operatorname{argmin}_{\theta \in \Omega} Q(\theta|\theta^t) - Q(\theta^t|\theta^t).$$

Part b): Show $D_A(\phi, \theta) = \text{KL}[p(\mathbf{x}, y, z; \theta) \| p(\mathbf{x}, y, z; \phi)]$

This result follows simply from decomposing $\text{KL}[p(\mathbf{x}, y, z; \theta) \| p(\mathbf{x}, y, z; \phi)]$ as follows:

$$\begin{aligned} \text{KL}[p(\mathbf{x}, y, z; \theta) \| p(\mathbf{x}, y, z; \phi)] &= \mathbb{E}_{\mathbf{X}, Y, Z | \theta} \left[\log \frac{p(\mathbf{x}, y, z; \theta)}{p(\mathbf{x}, y, z; \phi)} \right] \\ &\stackrel{(41)}{=} \langle s(\theta), \theta - \phi \rangle - A(\theta) + A(\phi) \pm \mathbb{E}_{\mathbf{X}, Y | \theta} \left[\log p(\mathbf{x}) - \frac{y^2}{2} - \frac{1}{2} \log(2\pi) \right] \\ &= A(\phi) - A(\theta) - \langle s(\theta), \phi - \theta \rangle \\ &\stackrel{i)}{=} A(\phi) - A(\theta) - \langle \nabla A(\theta), \phi - \theta \rangle. \end{aligned}$$

where *i*) follows from the fact that $\phi = \theta$ minimizes $-\mathbb{E}_{\mathbf{X}, Z, Y | \theta} [\log p(\mathbf{x}, y, z; \phi)]$. To see this, we use Jensen's inequality:

$$\begin{aligned} 0 &\leq -\mathbb{E}_{\mathbf{X}, Z, Y | \theta} [\log p(\mathbf{x}, y, z; \theta)] \\ &\stackrel{\text{Jensen's}}{\leq} -\log \mathbb{E}_{\mathbf{X}, Z, Y | \theta} [p(\mathbf{x}, y, z; \theta)] = -\log \int_{\mathbf{x}, y, z} p(\mathbf{x}, y, z; \theta)^2 d\mathbf{x} dz dy \\ &\stackrel{\text{Jensen's}}{\leq} -\log \left(\int_{\mathbf{x}, y, z} p(\mathbf{x}, y, z; \theta) d\mathbf{x} dz dy \right)^2 = -\log(1) = 0. \end{aligned}$$

Finally, taking the derivative with respect to ϕ and setting equal to 0 completes the proof:

$$\begin{aligned} 0 &= \frac{\partial}{\partial \theta} \mathbb{E}_{\mathbf{X}, Z, Y | \theta} [\log p(\mathbf{x}, y, z; \phi)] |_{\phi=\theta} \\ &= \mathbb{E}_{\mathbf{X}, Z, Y | \theta} \left[S(\mathbf{x}, y, z) + \frac{\partial}{\partial \phi} a(\mathbf{x}, y, \phi) |_{\phi=\theta} \right] \\ &= s(\theta) - \nabla A(\theta). \end{aligned}$$

Part c): Show $\mathcal{L}(\theta)$ is 1-smooth relative to $A(\theta)$.

The function $\mathcal{L}(\theta)$ is said to be 1-smooth relative to $A(\theta)$ if for all ϕ, θ , it holds that

$$\mathcal{L}(\theta) \leq \mathcal{L}(\phi) + \langle \nabla \mathcal{L}(\phi), \theta - \phi \rangle + D_A(\theta, \phi).$$

Recall the following from Section 3. The objective function $\mathcal{L}(\theta)$ is related to the EM objective $Q(\phi|\theta)$ by (10),

$$\mathcal{L}(\theta) = Q(\theta|\phi) - H(\theta|\phi),$$

where $H(\phi|\theta) \geq 0$ and $H(\theta|\theta) = 0$ for all $\phi, \theta \in \Omega$. Consequently, it then holds that for all $\phi, \theta \in \Omega$,

$$\mathcal{L}(\theta) = Q(\theta|\theta) \tag{43}$$

$$\mathcal{L}(\theta) \leq Q(\theta|\phi). \tag{44}$$

Recall also from part a) that $Q(\theta|\phi) - Q(\phi|\phi) = \langle \nabla \mathcal{L}(\phi), \theta - \phi \rangle + D_A(\theta, \phi)$. Then, the claim follows naturally from the above as follows:

$$\begin{aligned} \mathcal{L}(\theta) &\stackrel{(44)}{\leq} Q(\theta|\phi) \\ &\stackrel{a)}{=} Q(\phi|\phi) + \langle \nabla \mathcal{L}(\phi), \theta - \phi \rangle + D_A(\theta, \phi) \\ &\stackrel{(43)}{=} \mathcal{L}(\phi) + \langle \nabla \mathcal{L}(\phi), \theta - \phi \rangle + D_A(\theta, \phi) \end{aligned}$$

It follows that $\mathcal{L}(\boldsymbol{\theta})$ is 1-smooth relative to $A(\boldsymbol{\theta})$.

Part d): $A(\boldsymbol{\theta})$ and the MD objective is convex with respect to $\boldsymbol{\theta}$.

Here, we will show that the mirror descent objective is strongly convex in $\boldsymbol{\theta}$. It is important for this to hold so that the iterations of MD are well-defined; the minimizer of a strongly convex objective exists and is unique.

Note that the mirror descent objective, $\langle \nabla \mathcal{L}(\boldsymbol{\theta}), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle + D_A(\boldsymbol{\theta}, \boldsymbol{\theta}^t)$, is strongly convex in $\boldsymbol{\theta}$ if $A(\boldsymbol{\theta})$ is strongly convex in $\boldsymbol{\theta}$. Therefore, since $A(\boldsymbol{\theta})$ given in (31) is twice continuously differentiable, it is strongly convex with respect to $\boldsymbol{\theta}$ if and only if $\nabla^2 A(\boldsymbol{\theta}) \succeq r \mathbf{I}_{2d}$, for some $r > 0$. We begin:

$$\begin{aligned} \nabla^2 A(\boldsymbol{\theta}) &= \frac{\partial^2}{\partial \boldsymbol{\theta}^2} \mathbb{E}_{\mathbf{X}} \left[\frac{(\mathbf{x}^\top \boldsymbol{\beta})^2}{2} + \log(1 + e^{\mathbf{x}^\top \mathbf{w}}) \right] \\ &= \mathbb{E}_{\mathbf{X}} \left[\frac{\partial^2}{\partial \boldsymbol{\theta}^2} \left(\frac{(\mathbf{x}^\top \boldsymbol{\beta})^2}{2} + \log(1 + e^{\mathbf{x}^\top \mathbf{w}}) \right) \right] \\ &= \begin{pmatrix} \mathbb{E}_{\mathbf{X}} \left[\mathbf{x} \mathbf{x}^\top \frac{e^{\mathbf{x}^\top \mathbf{w}}}{(1 + e^{\mathbf{x}^\top \mathbf{w}})^2} \right] & \mathbf{0} \\ \mathbf{0} & \mathbb{E}_{\mathbf{X}} [\mathbf{x} \mathbf{x}^\top] \end{pmatrix} \\ &= \begin{pmatrix} \mathbb{E}_{\mathbf{X}} \left[\mathbf{x} \mathbf{x}^\top \frac{e^{\mathbf{x}^\top \mathbf{w}}}{(1 + e^{\mathbf{x}^\top \mathbf{w}})^2} \right] & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_d \end{pmatrix} \end{aligned}$$

where the last line follows from the assumption that \mathbf{x} is sampled from a unit spherical Gaussian distribution: $\mathbb{E}_{\mathbf{X}} [\mathbf{x} \mathbf{x}^\top] = \mathbf{I}_d$ for $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$.

From the above, we see that $A(\boldsymbol{\theta})$ is strictly convex, and it is strongly convex if $\mathbb{E}_{\mathbf{X}} \left[\mathbf{x} \mathbf{x}^\top \frac{e^{\mathbf{x}^\top \mathbf{w}}}{(1 + e^{\mathbf{x}^\top \mathbf{w}})^2} \right] \succeq r \mathbf{I}_d$ for some $r > 0$. This follows from Lemma B.5 where we show its eigenvalues are bounded below by $\min \left\{ \Omega \left(\frac{1}{\|\mathbf{w}\|_2} \right), \Omega \left(\frac{1}{\|\mathbf{w}\|_2^3} \right) \right\}$. Thus, it holds that

$$\nabla^2 A(\boldsymbol{\theta}) \succeq \min \left\{ \Omega \left(\frac{1}{\|\mathbf{w}\|_2} \right), \Omega \left(\frac{1}{\|\mathbf{w}\|_2^3} \right), 1 \right\} \mathbf{I}_{2d}.$$

Restricting the feasible set Ω to be all $\boldsymbol{\theta} \in \mathbb{R}^{2d}$ with $\|\boldsymbol{\theta}\|_2 \leq N$ for some $N \in [0, \infty)$, it holds that $A(\boldsymbol{\theta})$ is strongly convex with respect to $\boldsymbol{\theta}$ on Ω .

With part d) proven, this concludes the proof of Theorem 5.1 for SymMoLinE. We now prove the same for SymMoLogE, referring to this section where necessary.

□

B.3. Proof of Theorem 5.1 for SymMoLogE

Proof. Recall that we consider a 2 component SymMoLogE (see Section 3.1) where $z \in \{-1, 1\}$ is the latent unobserved variable, and

$$1) \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d),$$

$$2) P(z|\mathbf{x}; \mathbf{w}) = \frac{\exp\{\frac{z+1}{2} \mathbf{x}^\top \mathbf{w}\}}{1 + e^{\mathbf{x}^\top \mathbf{w}}}.$$

$$3) P(y|\mathbf{x}, z; \boldsymbol{\beta}) = \frac{\exp\{(\frac{yz+1}{2}) \mathbf{x}^\top \boldsymbol{\beta}\}}{1 + e^{\mathbf{x}^\top \boldsymbol{\beta}}}.$$

We begin by deriving a near exponential form of the complete data probability density function $p(\mathbf{x}, z, y; \boldsymbol{\theta})$:

$$\begin{aligned}
 p(\mathbf{x}, z, y; \boldsymbol{\theta}) &= P(y|\mathbf{x}, z; \boldsymbol{\theta})P(z|\mathbf{x}; \boldsymbol{\theta})p(\mathbf{x}) \\
 &= \exp\{\log P(y|\mathbf{x}, z; \boldsymbol{\beta}) + \log P(z|\mathbf{x}; \mathbf{w}) + \log p(\mathbf{x})\} \\
 &= \exp\left\{\log\left(\left(\frac{\exp\{\mathbf{x}^\top \boldsymbol{\beta}\}}{1 + e^{\mathbf{x}^\top \boldsymbol{\beta}}}\right)^{\frac{yz+1}{2}} \left(\frac{1}{1 + e^{\mathbf{x}^\top \boldsymbol{\beta}}}\right)^{1 - \frac{yz+1}{2}}\right) + \left(\frac{z+1}{2}\right) \mathbf{x}^\top \mathbf{w} - \log(1 + e^{\mathbf{x}^\top \mathbf{w}}) + \log p(\mathbf{x})\right\} \\
 &= \exp\left\{\frac{yz+1}{2} \log\left(\frac{\exp\{\mathbf{x}^\top \boldsymbol{\beta}\}}{1 + e^{\mathbf{x}^\top \boldsymbol{\beta}}}\right) + \left(1 - \frac{yz+1}{2}\right) \log\left(\frac{1}{1 + e^{\mathbf{x}^\top \boldsymbol{\beta}}}\right) + \left(\frac{z+1}{2}\right) \mathbf{x}^\top \mathbf{w} - \log(1 + e^{\mathbf{x}^\top \mathbf{w}}) + \log p(\mathbf{x})\right\} \\
 &= \exp\left\{\left(\frac{yz+1}{2}\right) \mathbf{x}^\top \boldsymbol{\beta} - \log(1 + e^{\mathbf{x}^\top \boldsymbol{\beta}}) + \left(\frac{z+1}{2}\right) \mathbf{x}^\top \mathbf{w} - \log(1 + e^{\mathbf{x}^\top \mathbf{w}}) + \log p(\mathbf{x})\right\} \\
 &= \exp\left\{\left\langle \left[\frac{yz\mathbf{x}}{2}, \frac{z\mathbf{x}}{2}\right], \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{w} \end{bmatrix} \right\rangle + \frac{\mathbf{x}^\top (\mathbf{w} + \boldsymbol{\beta})}{2} - \log\left[\left(1 + e^{\mathbf{x}^\top \boldsymbol{\beta}}\right) \left(1 + e^{\mathbf{x}^\top \mathbf{w}}\right)\right] + \log p(\mathbf{x})\right\}.
 \end{aligned}$$

Thus we have recovered the decomposition,

$$p(\mathbf{x}, y, z; \boldsymbol{\theta}) = \exp\left\{\left\langle \left[\frac{yz\mathbf{x}}{2}, \frac{z\mathbf{x}}{2}\right], \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{w} \end{bmatrix} \right\rangle + a(\mathbf{x}, y, \boldsymbol{\theta})\right\}, \quad (45)$$

where in $a(\mathbf{x}, y, \boldsymbol{\theta})$, \mathbf{x} cannot be linearly separated from the parameter $\boldsymbol{\theta}$:

$$a(\mathbf{x}, y, \boldsymbol{\theta}) = \frac{\mathbf{x}^\top (\mathbf{w} + \boldsymbol{\beta})}{2} - \log\left[\left(1 + e^{\mathbf{x}^\top \boldsymbol{\beta}}\right) \left(1 + e^{\mathbf{x}^\top \mathbf{w}}\right)\right] + \log p(\mathbf{x}).$$

Similar to SymMoLinE, we can see that $p(\mathbf{x}, y, z; \boldsymbol{\theta})$ does not belong to an exponential family of distribution. Also, note that taking the expectation of $a(\mathbf{x}, y, \boldsymbol{\theta})$ over $(\mathbf{x}, y) \sim p(\mathbf{x}, y; \boldsymbol{\theta}^*)$ yields,

$$\begin{aligned}
 \mathbb{E}_{\mathbf{X}, Y}[a(\mathbf{x}, y, \boldsymbol{\theta})] &= -\mathbb{E}_{\mathbf{X}}\left[\log\left[\left(1 + e^{\mathbf{x}^\top \boldsymbol{\beta}}\right) \left(1 + e^{\mathbf{x}^\top \mathbf{w}}\right)\right]\right] - C \\
 &= -A(\boldsymbol{\theta}) - C,
 \end{aligned}$$

where the above follows from $\mathbb{E}_{\mathbf{X}}\left[\frac{\mathbf{x}^\top (\mathbf{w} + \boldsymbol{\beta})}{2}\right] = 0$ and $C := -\mathbb{E}_{\mathbf{X}, Y|\boldsymbol{\theta}^*}[\log p(\mathbf{x})]$ is not a function of the parameter $\boldsymbol{\theta}$. We now continue with the proof.

From here on, the proofs of part a), part b) and part c) follow identically from that of SymMoLinE, so we will refer to Appendix B.2 for those proofs. We will now show part d).

Part d): $A(\boldsymbol{\theta})$ and the MD objective is convex with respect to $\boldsymbol{\theta}$.

Here, we will show that the mirror descent objective is strongly convex in $\boldsymbol{\theta}$. It is important for this to hold so that the iterations of MD are well-defined; the minimizer of a strongly convex objective exists and is unique.

Note that the mirror descent objective, $\langle \nabla \mathcal{L}(\boldsymbol{\theta}), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle + D_A(\boldsymbol{\theta}, \boldsymbol{\theta}^t)$, is strongly convex in $\boldsymbol{\theta}$ if $A(\boldsymbol{\theta})$ is strongly convex in $\boldsymbol{\theta}$. Therefore, since $A(\boldsymbol{\theta})$ given in (31) is twice continuously differentiable, it is strongly convex with respect to $\boldsymbol{\theta}$ if and only if $\nabla^2 A(\boldsymbol{\theta}) \succeq r \mathbf{I}_{2d}$, for some $r > 0$. We begin:

$$\begin{aligned}
 \nabla^2 A(\boldsymbol{\theta}) &= \frac{\partial^2}{\partial \boldsymbol{\theta}^2} \mathbb{E}_{\mathbf{X}} \left[\log \left(\left(1 + e^{\mathbf{x}^\top \boldsymbol{\beta}}\right) \left(1 + e^{\mathbf{x}^\top \mathbf{w}}\right) \right) \right] \\
 &= \mathbb{E}_{\mathbf{X}} \left[\frac{\partial^2}{\partial \boldsymbol{\theta}^2} \left(\log \left(1 + e^{\mathbf{x}^\top \boldsymbol{\beta}}\right) + \log \left(1 + e^{\mathbf{x}^\top \mathbf{w}}\right) \right) \right] \\
 &= \begin{pmatrix} \mathbb{E}_{\mathbf{X}} \left[\mathbf{x} \mathbf{x}^\top \frac{e^{\mathbf{x}^\top \mathbf{w}}}{(1 + e^{\mathbf{x}^\top \mathbf{w}})^2} \right] & \mathbf{0} \\ \mathbf{0} & \mathbb{E}_{\mathbf{X}} \left[\mathbf{x} \mathbf{x}^\top \frac{e^{\mathbf{x}^\top \boldsymbol{\beta}}}{(1 + e^{\mathbf{x}^\top \boldsymbol{\beta}})^2} \right] \end{pmatrix} \\
 &= \begin{pmatrix} \mathbb{E}_{\mathbf{X}} \left[\mathbf{x} \mathbf{x}^\top \frac{e^{\mathbf{x}^\top \mathbf{w}}}{(1 + e^{\mathbf{x}^\top \mathbf{w}})^2} \right] & \mathbf{0} \\ \mathbf{0} & \mathbb{E}_{\mathbf{X}} \left[\mathbf{x} \mathbf{x}^\top \frac{e^{\mathbf{x}^\top \boldsymbol{\beta}}}{(1 + e^{\mathbf{x}^\top \boldsymbol{\beta}})^2} \right] \end{pmatrix}.
 \end{aligned}$$

From the above, we see that $A(\theta)$ is strictly convex, and it is strongly convex if $\mathbb{E}_{\mathbf{X}} \left[\mathbf{x} \mathbf{x}^\top \frac{e^{\mathbf{x}^\top \mathbf{w}}}{(1+e^{\mathbf{x}^\top \mathbf{w}})^2} \right] \succeq r \mathbf{I}_d$ and $\mathbb{E}_{\mathbf{X}} \left[\mathbf{x} \mathbf{x}^\top \frac{e^{\mathbf{x}^\top \beta}}{(1+e^{\mathbf{x}^\top \beta})^2} \right] \succeq r \mathbf{I}_d$ for some $r > 0$. This follows from Lemma B.5 where we show their respective eigenvalues are bounded below by, $\min \left\{ \Omega \left(\frac{1}{\|\mathbf{w}\|_2} \right), \Omega \left(\frac{1}{\|\mathbf{w}\|_2^3} \right) \right\}$ and $\min \left\{ \Omega \left(\frac{1}{\|\beta\|_2} \right), \Omega \left(\frac{1}{\|\beta\|_2^3} \right) \right\}$. Thus, it holds that

$$\nabla^2 A(\theta) \succeq \min \left\{ \Omega \left(\frac{1}{\|\mathbf{w}\|_2} \right), \Omega \left(\frac{1}{\|\mathbf{w}\|_2^3} \right), \Omega \left(\frac{1}{\|\beta\|_2} \right), \Omega \left(\frac{1}{\|\beta\|_2^3} \right) \right\} \mathbf{I}_{2d}.$$

Restricting Ω to be all θ with $\|\theta\|_2 \leq N$ for some $N \in [0, \infty)$, it holds that $A(\theta)$ is strongly convex with respect to θ on Ω .

With part d) proven, this concludes the proof of Theorem 5.1 for SymMoLinE. We now prove the same for SymMoLogE, referring to this section where necessary. □

B.4. Convergence Guarantees of EM for SymMoLogE and SymMoLinE

In this section, we provide the proofs of Corollary B.1. Building on prior work (Lu et al., 2018), we contextualize convergence properties of MD for SymMoLinE and SymMoLogE. Before presenting the result, we briefly review key concepts. We say θ^1 is initialized in a locally convex region of $\mathcal{L}(\theta)$ if there exists a convex set $\Theta \subseteq \Omega$ containing θ^1, θ^* such that for all $\phi, \theta \in \Theta$,

$$\mathcal{L}(\phi) \geq \mathcal{L}(\theta) + \langle \nabla \mathcal{L}(\theta), \phi - \theta \rangle. \quad (46)$$

Furthermore, Θ is called α -strongly convex relative to h if

$$\mathcal{L}(\phi) \geq \mathcal{L}(\theta) + \langle \nabla \mathcal{L}(\theta), \phi - \theta \rangle + \alpha D_h(\phi, \theta). \quad (47)$$

The corollary that follows provides conditions for convergence of EM to (1) a stationary point in the KL divergence, (2) the true parameters at a sub-linear rate, and (3) the true parameters at a linear rate. We further note that the proof is adapted from (Kunstner et al., 2021, Proposition 2, Corollary 1, and Corollary 3) and (Lu et al., 2018, Theorem 3.1), we provide it here for completeness.

Corollary B.1 (Convergence of EM). *For SymMoLinE, SymMoLogE with mirror map $A(\theta)$ given as (31) (32) respectively, and denoting $D_A(\theta^t, \theta^{t+1}) := \text{KL}[p(\mathbf{x}, y, z; \theta^{t+1}) \| p(\mathbf{x}, y, z; \theta^t)]$, the EM iterates $\{\theta^t\}_{t \in [T]}$ satisfy:*

1) **Stationarity.** *For no additional conditions,*

$$\min_{t \in [T]} D_A(\theta^{t+1}, \theta^t) \leq \frac{\mathcal{L}(\theta^1) - \mathcal{L}(\theta^*)}{T}; \quad (48)$$

2) **Sub-linear Rate to θ^* .** *If θ^1 is initialized in Θ , a locally convex region of $\mathcal{L}(\theta)$ containing θ^* , then*

$$\mathcal{L}(\theta^T) - \mathcal{L}(\theta^*) \leq \frac{D_A(\theta^*, \theta^1)}{T} \quad (49)$$

3) **Linear Rate to θ^* .** *If θ^1 is initialized in $\Theta \subseteq \Omega$, a locally strongly convex region of $\mathcal{L}(\theta)$ relative to $A(\theta)$ that contains θ^* , then*

$$\mathcal{L}(\theta^T) - \mathcal{L}(\theta^*) \leq (1 - \alpha)^T (\mathcal{L}(\theta^1) - \mathcal{L}(\theta^*)). \quad (50)$$

Proof. The proof is divided into three parts that correspond to each of the three sub-results of the corollary.

Part 1): Stationarity.

Given Theorem 5.1, this proof follows from identical arguments to that of (Kunstner et al., 2021, Proposition 2). We write it below for completeness.

Recall from Theorem 5.1 that θ^{t+1} is obtained as the minimizer of the convex objective, (15):

$$\langle \nabla \mathcal{L}(\theta^t), \theta - \theta^t \rangle + D_A(\theta, \theta^t).$$

As such, differentiating and setting equal to 0, it holds that θ^{t+1} satisfies

$$\nabla \mathcal{L}(\theta^t) = \nabla A(\theta^t) - \nabla A(\theta^{t+1}) \quad (51)$$

Further, by the above together with relative smoothness, it holds that

$$\begin{aligned} \mathcal{L}(\theta^{t+1}) &\leq \mathcal{L}(\theta^t) + \langle \nabla \mathcal{L}(\theta^t), \theta^{t+1} - \theta^t \rangle + D_A(\theta^{t+1}, \theta^t) \\ &= \mathcal{L}(\theta^t) + \langle \nabla A(\theta^t) - \nabla A(\theta^{t+1}), \theta^{t+1} - \theta^t \rangle + D_A(\theta^{t+1}, \theta^t) \\ &= \mathcal{L}(\theta^t) - \langle \nabla A(\theta^{t+1}), \theta^{t+1} - \theta^t \rangle + A(\theta^{t+1}) - A(\theta^t) \\ &= \mathcal{L}(\theta^t) - D_A(\theta^t, \theta^{t+1}). \end{aligned}$$

Thus it we have shown that

$$D_A(\theta^t, \theta^{t+1}) \leq \mathcal{L}(\theta^t) - \mathcal{L}(\theta^{t+1}). \quad (52)$$

The claim then follows from taking the mean over T iterations:

$$\min_{t \leq T} D_A(\theta^t, \theta^{t+1}) \leq \frac{1}{T} \sum_{t=1}^T D_A(\theta^t, \theta^{t+1}) \leq \frac{1}{T} \sum_{t=1}^T \mathcal{L}(\theta^t) - \mathcal{L}(\theta^{t+1}) = \frac{\mathcal{L}(\theta^1) - \mathcal{L}(\theta^T)}{T} \leq \frac{\mathcal{L}(\theta^1) - \mathcal{L}(\theta^*)}{T}.$$

Part 2): Sub-linear Rate to θ^* .

Given Theorem 5.1, this proof follows from identical arguments to that of [Kunstner et al. \(2021, Corollary 1\)](#) and [Lu et al. \(2018, Theorem 3.1\)](#). We write it below for completeness.

Here, we assume that $\mathcal{L}(\theta)$ is convex on the set Θ . In part 1), we used (51) to show,

$$\langle \nabla \mathcal{L}(\theta^t), \theta^{t+1} - \theta^t \rangle + D_A(\theta^{t+1}, \theta^t) = -D_A(\theta^t, \theta^{t+1}),$$

where the right hand side is non-positive since the Bregman divergence is non-negative if the inducing function A is convex – which it is. Now, starting from relative smoothness, we see that

$$\begin{aligned} \mathcal{L}(\theta^{t+1}) &\leq \mathcal{L}(\theta^t) + \langle \nabla \mathcal{L}(\theta^t), \theta^{t+1} - \theta^t \rangle + D_A(\theta^{t+1}, \theta^t) \\ &= \mathcal{L}(\theta^t) + \langle \nabla \mathcal{L}(\theta^t), \theta^{t+1} - \theta^* + \theta^* - \theta^t \rangle + D_A(\theta^{t+1}, \theta^t) \\ &= \mathcal{L}(\theta^t) + \langle \nabla \mathcal{L}(\theta^t), \theta^{t+1} - \theta^* \rangle + \langle \nabla \mathcal{L}(\theta^t), \theta^* - \theta^t \rangle + D_A(\theta^{t+1}, \theta^t) \\ &\stackrel{i)}{\leq} \mathcal{L}(\theta^*) + \langle \nabla \mathcal{L}(\theta^t), \theta^{t+1} - \theta^* \rangle + D_A(\theta^{t+1}, \theta^t) \\ &\stackrel{(51)}{=} \mathcal{L}(\theta^*) + \langle \nabla A(\theta^t) - \nabla A(\theta^{t+1}), \theta^{t+1} - \theta^* \rangle + D_A(\theta^{t+1}, \theta^t) \end{aligned}$$

where i) follows from convexity of $\mathcal{L}(\theta)$ on the set Θ . Subsequently, we apply the 3-point lemma, $D_A(\theta^*, \theta^t) = D_A(\theta^*, \theta^{t+1}) + \langle \theta^* - \theta^{t+1}, \nabla A(\theta^{t+1}) - A(\theta^t) \rangle + D_A(\theta^{t+1}, \theta^t)$, and obtain,

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^*) + D_A(\theta^*, \theta^t) - D_A(\theta^*, \theta^{t+1}). \quad (53)$$

Finally, the result follows from summing the left and right hand side over T iterations:

$$T(\mathcal{L}(\theta^T) - \mathcal{L}(\theta^*)) \leq \sum_{t=1}^T \mathcal{L}(\theta^{t+1}) - \mathcal{L}(\theta^*) \leq \sum_{t=1}^T D_A(\theta^*, \theta^t) - D_A(\theta^*, \theta^{t+1}) \leq D_A(\theta^*, \theta^1)$$

Part 3): Linear Rate to θ^* .

Given Theorem 5.1, this proof follows from identical arguments to that of [Kunstner et al. \(2021, Corollary 3\)](#) and [Lu et al. \(2018, Theorem 3.1\)](#). We write it below for completeness.

In addition to convexity, we now assume that $\mathcal{L}(\theta)$ is α -strongly convex relative to $A(\theta)$ on the set Θ . Specifically, we have that for any $\phi, \theta \in \Theta$,

$$\mathcal{L}(\theta) \geq \mathcal{L}(\phi) + \langle \nabla \mathcal{L}(\phi), \theta - \phi \rangle + \alpha D_A(\theta, \phi). \quad (54)$$

Using the three point lemma again, we have

$$\begin{aligned}
 D_A(\theta^*, \theta^{t+1}) &= D_A(\theta^*, \theta^t) + \langle \theta^* - \theta^t, \nabla A(\theta^t) - \nabla A(\theta^{t+1}) \rangle + D_A(\theta^t, \theta^{t+1}) \\
 &\stackrel{(51)}{=} D_A(\theta^*, \theta^t) + \langle \nabla \mathcal{L}(\theta^t), \theta^* - \theta^t \rangle + D_A(\theta^t, \theta^{t+1}) \\
 &\stackrel{(54)}{\leq} D_A(\theta^*, \theta^t) + \mathcal{L}(\theta^*) - \mathcal{L}(\theta^t) - \alpha D_A(\theta^*, \theta^t) + D_A(\theta^t, \theta^{t+1}) \\
 &= (1 - \alpha) D_A(\theta^*, \theta^t) + \mathcal{L}(\theta^*) - \mathcal{L}(\theta^t) + D_A(\theta^t, \theta^{t+1}) \\
 &\stackrel{(52)}{\leq} (1 - \alpha) D_A(\theta^*, \theta^t) + \mathcal{L}(\theta^*) - \mathcal{L}(\theta^t) + \mathcal{L}(\theta^t) - \mathcal{L}(\theta^{t+1}) \\
 &\leq (1 - \alpha) D_A(\theta^*, \theta^t) \\
 &\leq (1 - \alpha)^T D_A(\theta^*, \theta^1).
 \end{aligned}$$

Finally, from (53), we see that

$$\begin{aligned}
 \mathcal{L}(\theta^{t+1}) - \mathcal{L}(\theta^*) &\leq D_A(\theta^*, \theta^t) - D_A(\theta^*, \theta^{t+1}) \\
 &\leq (1 - \alpha)^T D_A(\theta^*, \theta^1) - D_A(\theta^*, \theta^{t+1}) \\
 &\leq (1 - \alpha)^T D_A(\theta^*, \theta^1).
 \end{aligned}$$

□

B.5. Satisfiability of Conditions from Corollary B.1

The above result raises an important question: when does a locally convex or relatively strongly convex region of $\mathcal{L}(\theta)$ containing θ^* exist? Interestingly, this is closely tied to the Signal-to-Noise Ratio (SNR). Before exploring this connection, we first introduce the concept of the Missing Information Matrix (MIM) introduced in (Orchard & Woodbury, 1972).

The MIM relates the level of information the pair (x, y) holds about the latent expert label z given parameters θ , and it is formally defined as

$$M(\theta) = I_{x,z,y|\theta}^{-1} I_{z|x,y,\theta}. \quad (55)$$

Here, $I_{x,z,y|\theta}$ is the Fisher information matrix of the complete data distribution and $I_{z|x,y,\theta}$ is the Fisher information matrix of the conditional distribution of the latent unobserved variable given the observed ones, denoted by

$$\begin{aligned}
 I_{x,z,y|\theta} &:= -\mathbb{E}_{\mathbf{X},Y,Z|\theta} \left[\frac{\partial^2}{\partial \theta^2} \log p(\mathbf{x}, z, y; \theta) \right] \\
 &= \nabla^2 Q(\phi|\theta)|_{\phi=\theta} = \nabla^2 A(\theta) \\
 I_{z|x,y,\theta} &:= -\mathbb{E}_{\mathbf{X},Y} \mathbb{E}_{Z|x,y,\theta} \left[\frac{\partial^2}{\partial \theta^2} \log P(z|x, y; \theta) \right] \\
 &= \nabla^2 H(\phi|\theta)|_{\phi=\theta}.
 \end{aligned}$$

Thus, it also holds that the MIM is a function of $A(\theta)$ and $H(\phi|\theta)$, i.e.,

$$M(\theta) = \nabla^2 A(\theta)^{-1} \nabla^2 H(\phi|\theta)|_{\phi=\theta}. \quad (56)$$

Due to the pairwise independence of \mathbf{X} and its rotational invariance, there exists an orthonormal matrix R such that $\Delta := R I_{x,z,y|\theta}^{-1} R^\top$ is positive semi-definite and diagonal and $J := R I_{z|x,y,\theta} R^\top$ is symmetric positive semi-definite (see Lemma B.3). As such, the MIM in (55) is also symmetric and positive semi-definite: $M(\theta) = R^\top \Delta R R^\top J R = R^\top \Delta J R$. Note that the MIM quantifies the difficulty of estimating parameters when only x, y are observed. To understand its significance, consider the following: large eigenvalues of M indicate that x, y contain little information about the true value of the latent variable z , making estimation more difficult. Conversely, small eigenvalues suggest that x, y provide enough information to effectively constrain the possible values of z . Thus, the MIM can be seen as analogous to the Signal-to-Noise Ratio (SNR).

In the theorem below, we show how the eigenvalues of $M(\theta)$ are related to the satisfiability of the conditions for Corollary B.1 regarding the relative strong convexity of $\mathcal{L}(\theta)$ with respect to the mirror map $A(\theta)$.

Theorem B.2. For SymMoLinE and SymMoLogE and their respective mirror mappings (31) and (32), the objective $\mathcal{L}(\theta)$ is α -strongly convex relative to the mirror map $A(\theta)$ on the convex set Θ if and only if

$$\lambda_{\max}(\mathbf{M}(\theta)) \leq (1 - \alpha) \text{ for all } \theta \in \Theta. \quad (57)$$

Proof. Recall that $\mathcal{L}(\theta)$ is strongly convex relative to $A(\theta)$ on Θ if for all $\theta, \phi \in \Theta$, it holds that

$$\mathcal{L}(\phi) \geq \mathcal{L}(\theta) + \langle \nabla \mathcal{L}(\theta), \phi - \theta \rangle + \alpha D_A(\phi, \theta).$$

For $\mathcal{L}(\theta)$ and $A(\theta)$ twice continuously differentiable, it was shown by (Lu et al., 2018) that this is equivalent to the following bound on the Hessian:

$$\nabla^2 \mathcal{L}(\theta) \succeq \alpha \nabla^2 A(\theta).$$

Now, using $\mathcal{L}(\phi) = Q(\phi|\theta) - H(\phi|\theta)$, we see that

$$\begin{aligned} \nabla^2 \mathcal{L}(\theta) &= \nabla^2 (Q(\phi|\theta) - H(\phi|\theta))|_{\phi=\theta} \\ &= \nabla^2 Q(\phi|\theta)|_{\phi=\theta} - \nabla^2 H(\phi|\theta)|_{\phi=\theta} \\ &= \nabla^2 A(\theta) - \nabla^2 H(\phi|\theta)|_{\phi=\theta} \end{aligned}$$

Therefore, since $\nabla^2 A(\theta)$ is symmetric positive definite (proven in Appendix B.1), our condition simplifies to

$$\begin{aligned} (1 - \alpha) \nabla^2 A(\theta) &\succeq \nabla^2 H(\phi|\theta)|_{\phi=\theta} \\ \iff (1 - \alpha) \mathbf{I}_{2d} &\succeq \nabla^2 A(\theta)^{-1} \nabla^2 H(\phi|\theta)|_{\phi=\theta} = \mathbf{M}(\theta). \end{aligned}$$

Finally, for \mathbf{x} from a unit spherical Gaussian distribution, we know that $\mathbf{M}(\theta)$ is symmetric positive-definite (see Lemma B.3). As a result, the above inequality is equivalent to the following bound on the eigenvalues of the MIM:

$$1 - \alpha \geq \lambda_{\max}(\mathbf{M}(\theta)).$$

□

We now provide the simple Lemma that the MIM is a symmetric matrix for SymMoLinE and SymMoLogE.

Lemma B.3 ($\mathbf{M}(\theta)$ is symmetric). For SymMoLinE and SymMoLogE, the MIM is a symmetric matrix, i.e. $\mathbf{M}(\theta) = \mathbf{M}(\theta)^\top$.

Proof. Recall the assumption that \mathbf{x} is sampled from a unit spherical Gaussian distribution: $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. As such, for any orthonormal $d \times d$ matrix \mathbf{R} , we know that $\mathbf{R}\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$; this is called rotational invariance of the Gaussian distribution. Thus, consider orthonormal $\mathbf{R}_u \in \mathbb{R}^{d \times d}$, such that $\mathbf{R}_u \mathbf{u} = \mathbf{e}_1 \|\mathbf{u}\|_2$ where $\mathbf{e}_j \in \mathbb{R}^d$ is the $\mathbf{0}$ vector with a 1 at index j . Now, we can observe that for any $\mathbf{w} \in \mathbb{R}^d$, $\mathbb{E}_{\mathbf{x}} \left[\mathbf{x} \mathbf{x}^\top \frac{e^{\mathbf{x}^\top \mathbf{u}}}{(1 + e^{\mathbf{x}^\top \mathbf{u}})^2} \right]$, is diagonalizable by an orthonormal matrix \mathbf{R} :

$$\begin{aligned} \mathbf{R}_u \left(\mathbb{E}_{\mathbf{x}} \left[\mathbf{x} \mathbf{x}^\top \frac{e^{\mathbf{x}^\top \mathbf{u}}}{(1 + e^{\mathbf{x}^\top \mathbf{u}})^2} \right] \right) \mathbf{R}_u^\top &= \mathbb{E}_{\mathbf{x}} \left[\mathbf{R}_u \mathbf{x} \mathbf{x}^\top \mathbf{R}_u^\top \frac{e^{\mathbf{x}^\top \mathbf{R}_u^\top \mathbf{R}_u \mathbf{u}}}{(1 + e^{\mathbf{x}^\top \mathbf{R}_u^\top \mathbf{R}_u \mathbf{u}})^2} \right] \\ &= \mathbb{E}_{\mathbf{R}_u \mathbf{x}} \left[\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \frac{e^{\tilde{\mathbf{x}}^\top \mathbf{e}_1 \|\mathbf{u}\|_2}}{(1 + e^{\tilde{\mathbf{x}}^\top \mathbf{e}_1 \|\mathbf{u}\|_2})^2} \right] \\ &= \begin{pmatrix} \mathbb{E}_{\mathbf{R}_u \mathbf{x}} \left[\tilde{x}_1^2 \frac{e^{\tilde{x}_1 \|\mathbf{u}\|_2}}{(1 + e^{\tilde{x}_1 \|\mathbf{u}\|_2})^2} \right] & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{E}_{\mathbf{R}_u \mathbf{x}} \left[\tilde{x}_d^2 \frac{e^{\tilde{x}_d \|\mathbf{u}\|_2}}{(1 + e^{\tilde{x}_d \|\mathbf{u}\|_2})^2} \right] \end{pmatrix} \end{aligned}$$

Where in the above, 1) $\mathbf{R}^\top \mathbf{R} = \mathbf{I}_d$ since $\mathbf{R}^\top = \mathbf{R}^{-1}$ for orthonormal matrices and 2) non diagonal elements evaluate to 0 because for all $i \neq j$, the 0 mean random variables \tilde{X}_j is independent from \tilde{X}_i .

Finally, we put the above together to show $M(\theta)$ is symmetric. We define the block diagonal orthonormal matrix R_M as

$$R_M := \begin{pmatrix} R_w & 0 \\ 0 & R_\beta \end{pmatrix}.$$

From the above, it follows that both for SymMoLinE and SymMoLogE, $R_M \nabla^2 A(\theta)^{-1} R_M^\top$ is a diagonal matrix. We can now use this change of basis matrix to show the MIM is a symmetric matrix:

$$\begin{aligned} M(\theta) &= \nabla^2 A(\theta)^{-1} \nabla^2 H(\theta|\theta) \\ &= R_M^\top \left(R_M \nabla^2 A(\theta)^{-1} R_M^\top \right) (R_M \nabla^2 H(\theta|\theta) R_M^\top) R_M \\ &= R_M^\top (R_M \nabla^2 H(\theta|\theta) R_M^\top) \left(R_M \nabla^2 A(\theta)^{-1} R_M^\top \right) R_M \\ &= \nabla^2 H(\theta|\theta)^\top (\nabla^2 A(\theta)^{-1})^\top \\ &= M(\theta)^\top. \end{aligned}$$

□

B.6. Correspondence Between the MIM and SNR for SymMoLinE and SymMoLogE

The above result states if $M(\theta) \prec I_{2d}$ for all $\theta \in \Omega$ and $\theta^* \in \Omega$, then the EM updates will converge linearly to θ^* . This offers a unified framework for analyzing EM for MoE, linking the rate of convergence to a classical statistical metric. To determine whether EM achieves linear or sub-linear convergence, we need to understand the behavior of $M(\theta)$ and $M(\theta^*)$, which indicates the existence and size of the local region where EM enjoys such convergence.

Theorem B.4. For SymMoLinE, the eigenvalues of $I_{x,z,y|\theta}^{-1}$ belong to the set

$$\lambda \left(I_{x,z,y|\theta}^{-1} \right) = \{ \Theta(\|w\|_2^3), \Theta(\|w\|_2), 1 \} \quad (58)$$

and $I_{z|x,y,\theta}$ is given as the expectation over (X, Y) of a function that is decreasing as a function of $\|\theta\|_2$:

$$I_{z|x,y,\theta} = \mathbb{E}_{X,Y} \left[\frac{\exp \left\langle \begin{bmatrix} x \\ 2yx \end{bmatrix}, \theta \right\rangle \begin{bmatrix} x \\ 2yx \end{bmatrix} \begin{bmatrix} x \\ 2yx \end{bmatrix}^\top}{\left(1 + \exp \left\langle \begin{bmatrix} x \\ 2yx \end{bmatrix}, \theta \right\rangle \right)^2} \right]. \quad (59)$$

Similarly, For SymMoLogE, the eigenvalues of $I_{x,z,y|\theta}^{-1}$ belong to the set

$$\lambda \left(I_{x,z,y|\theta}^{-1} \right) = \{ \Theta(\|w\|_2^3), \Theta(\|w\|_2), \Theta(\|\beta\|_2^3), \Theta(\|\beta\|_2) \} \quad (60)$$

and $I_{z|x,y,\theta}$ is given as the expectation over (X, Y) of a function that is decreasing as a function of $\|\theta\|_2$:

$$I_{z|x,y,\theta} = \mathbb{E}_{X,Y} \left[\frac{\exp \left\langle \begin{bmatrix} x \\ yx \end{bmatrix}, \theta \right\rangle \begin{bmatrix} x \\ yx \end{bmatrix} \begin{bmatrix} x \\ yx \end{bmatrix}^\top}{\left(1 + \exp \left\langle \begin{bmatrix} x \\ yx \end{bmatrix}, \theta \right\rangle \right)^2} \right]. \quad (61)$$

Proof. We divide the proof into two parts. In the first part, we consider the SymMoLinE setting and, in the second part, we consider the SymMoLogE setting.

Part a): SymMoLinE.

For $I_{x,z,y|\theta}$, recall that it has the following form:

$$I_{x,y,z|\theta} = \begin{pmatrix} \mathbb{E}_X \left[xx^\top \frac{e^{x^\top w}}{(1 + e^{x^\top w})^2} \right] & 0 \\ 0 & I_d \end{pmatrix}.$$

By Lemma B.5, we see that $\mathbf{I}_{\mathbf{x},z,y|\boldsymbol{\theta}}$ can be diagonalized into the following form:

$$\mathbf{I}_{\mathbf{x},y,z|\boldsymbol{\theta}} = \mathbf{R}_M \begin{pmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \dots & 0 & 0 \\ 0 & 0 & \lambda_2 & 0 \\ 0 & 0 & 0 & \mathbf{I}_d \end{pmatrix} \mathbf{R}_M^\top$$

where \mathbf{R}_M is an orthonormal rotation matrix and $\lambda_1 = \Theta\left(\frac{1}{\|\mathbf{w}\|_2^3}\right)$ and $\lambda_2 = \Theta\left(\frac{1}{\|\mathbf{w}\|_2}\right)$. Therefore, $\mathbf{I}_{\mathbf{x},z,y|\boldsymbol{\theta}}$ has three eigenvalues given as $\left\{\Theta\left(\frac{1}{\|\mathbf{w}\|_2^3}\right), \Theta\left(\frac{1}{\|\mathbf{w}\|_2}\right), 1\right\}$. It follows that $\mathbf{I}_{\mathbf{x},z,y|\boldsymbol{\theta}}^{-1}$ has the form

$$\mathbf{I}_{\mathbf{x},y,z|\boldsymbol{\theta}}^{-1} = \mathbf{R}_M^\top \begin{pmatrix} 1/\lambda_1 & 0 & 0 & 0 \\ 0 & \dots & 0 & 0 \\ 0 & 0 & 1/\lambda_2 & 0 \\ 0 & 0 & 0 & \mathbf{I}_d \end{pmatrix} \mathbf{R}_M.$$

Therefore, $\mathbf{I}_{\mathbf{x},y,z|\boldsymbol{\theta}}^{-1}$ has three eigenvalues given as $\left\{\Theta\left(\|\mathbf{w}\|_2^3\right), \Theta\left(\|\mathbf{w}\|_2\right), 1\right\}$.

Now, for $\mathbf{I}_{z|\mathbf{x},y,\boldsymbol{\theta}}$, we first derive a more compact form for the conditional distribution of the latent variable, $p(z|\mathbf{x}, y; \boldsymbol{\theta})$. From simple Bayes rule and algebraic manipulation, we see that

$$\begin{aligned} p(z|\mathbf{x}, y; \boldsymbol{\theta}) &= \frac{p(y|\mathbf{x}, z; \boldsymbol{\theta})p(z|\mathbf{x}; \boldsymbol{\theta})p(\mathbf{x})}{p(\mathbf{x}, y; \boldsymbol{\theta})} \\ &= \frac{\frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(y-z\mathbf{x}^\top\boldsymbol{\beta})^2}{2}\right\} \frac{\exp\left\{\frac{z+1}{2}\mathbf{x}^\top\mathbf{w}\right\}}{1+e^{\mathbf{x}^\top\mathbf{w}}}}{\frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(y-\mathbf{x}^\top\boldsymbol{\beta})^2}{2}\right\} \frac{\exp\{\mathbf{x}^\top\mathbf{w}\}}{1+e^{\mathbf{x}^\top\mathbf{w}}} + \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(y+\mathbf{x}^\top\boldsymbol{\beta})^2}{2}\right\} \frac{1}{1+e^{\mathbf{x}^\top\mathbf{w}}}} \\ &= \frac{\exp\left\{-\frac{(y-z\mathbf{x}^\top\boldsymbol{\beta})^2}{2} + \frac{z+1}{2}\mathbf{x}^\top\mathbf{w}\right\}}{\exp\left\{-\frac{(y-\mathbf{x}^\top\boldsymbol{\beta})^2}{2} + \mathbf{x}^\top\mathbf{w}\right\} + \exp\left\{-\frac{(y+\mathbf{x}^\top\boldsymbol{\beta})^2}{2}\right\}} \\ &= \frac{\exp\left\{zy\mathbf{x}^\top\boldsymbol{\beta} + \frac{z+1}{2}\mathbf{x}^\top\mathbf{w}\right\}}{\exp\{y\mathbf{x}^\top\boldsymbol{\beta} + \mathbf{x}^\top\mathbf{w}\} + \exp\{-y\mathbf{x}^\top\boldsymbol{\beta}\}} \\ &= \frac{\exp\left\{\frac{z+1}{2}(2y\mathbf{x}^\top\boldsymbol{\beta} + \mathbf{x}^\top\mathbf{w})\right\}}{\exp\{2y\mathbf{x}^\top\boldsymbol{\beta} + \mathbf{x}^\top\mathbf{w}\} + 1} \\ &= \frac{\exp\left\{\frac{z+1}{2}\left\langle \begin{bmatrix} \mathbf{x} \\ 2y\mathbf{x} \end{bmatrix}, \begin{bmatrix} \mathbf{w} \\ \boldsymbol{\beta} \end{bmatrix} \right\rangle\right\}}{\exp\left\{\left\langle \begin{bmatrix} \mathbf{x} \\ 2y\mathbf{x} \end{bmatrix}, \begin{bmatrix} \mathbf{w} \\ \boldsymbol{\beta} \end{bmatrix} \right\rangle\right\} + 1} \end{aligned}$$

Now that we have this simplified form, we are able to derive (59) for $\mathbf{I}_{z|\mathbf{x},y,\boldsymbol{\theta}}$:

$$\begin{aligned} \mathbf{I}_{z|\mathbf{x},y,\boldsymbol{\theta}} &= -\mathbb{E}_{\mathbf{X},Y} \mathbb{E}_{Z|\mathbf{x},y,\boldsymbol{\theta}} \left[\frac{\partial^2}{\partial \boldsymbol{\theta}^2} \log p(z|\mathbf{x}, y; \boldsymbol{\theta}) \right] \\ &= -\mathbb{E}_{\mathbf{X},Y} \mathbb{E}_{Z|\mathbf{x},y,\boldsymbol{\theta}} \left[\frac{\partial^2}{\partial \boldsymbol{\theta}^2} \left(\frac{z+1}{2} \left\langle \begin{bmatrix} \mathbf{x} \\ 2y\mathbf{x} \end{bmatrix}, \boldsymbol{\theta} \right\rangle - \log \left(1 + \exp \left\langle \begin{bmatrix} \mathbf{x} \\ 2y\mathbf{x} \end{bmatrix}, \boldsymbol{\theta} \right\rangle \right) \right) \right] \\ &= \mathbb{E}_{\mathbf{X},Y} \mathbb{E}_{Z|\mathbf{x},y,\boldsymbol{\theta}} \left[\frac{\partial^2}{\partial \boldsymbol{\theta}^2} \log \left(1 + \exp \left\langle \begin{bmatrix} \mathbf{x} \\ 2y\mathbf{x} \end{bmatrix}, \boldsymbol{\theta} \right\rangle \right) \right] \\ &= \mathbb{E}_{\mathbf{X},Y} \left[\frac{\partial^2}{\partial \boldsymbol{\theta}^2} \log \left(1 + \exp \left\langle \begin{bmatrix} \mathbf{x} \\ 2y\mathbf{x} \end{bmatrix}, \boldsymbol{\theta} \right\rangle \right) \right] \\ &= \mathbb{E}_{\mathbf{X},Y} \left[\frac{\exp \left\langle \begin{bmatrix} \mathbf{x} \\ 2y\mathbf{x} \end{bmatrix}, \boldsymbol{\theta} \right\rangle \begin{bmatrix} \mathbf{x} \\ 2y\mathbf{x} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ 2y\mathbf{x} \end{bmatrix}^\top}{\left(1 + \exp \left\langle \begin{bmatrix} \mathbf{x} \\ 2y\mathbf{x} \end{bmatrix}, \boldsymbol{\theta} \right\rangle \right)^2} \right]. \end{aligned}$$

This expression depends only on the random variable (X, Y) and the parameter iterate θ .

Part b): SymMoLogE.

Recall that for SymMoLogE, $I_{x,z,y|\theta}$ has the following form:

$$I_{x,z,y|\theta} = \begin{pmatrix} \mathbb{E}_X \left[\mathbf{x} \mathbf{x}^\top \frac{e^{\mathbf{x}^\top \mathbf{w}}}{(1+e^{\mathbf{x}^\top \mathbf{w}})^2} \right] & \mathbf{0} \\ \mathbf{0} & \mathbb{E}_X \left[\mathbf{x} \mathbf{x}^\top \frac{e^{\mathbf{x}^\top \boldsymbol{\beta}}}{(1+e^{\mathbf{x}^\top \boldsymbol{\beta}})^2} \right] \end{pmatrix}.$$

By Lemma B.5, we see that $I_{x,z,y|\theta}$ can be diagonalized into the following form:

$$I_{x,z,y|\theta} = R_M \begin{pmatrix} \lambda_1 & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \lambda_2 & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \lambda_3 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \lambda_4 \end{pmatrix} R_M^\top$$

where R_M is an orthonormal rotation matrix and $\lambda_1 = \Theta\left(\frac{1}{\|\mathbf{w}\|_2^3}\right)$, $\lambda_2 = \Theta\left(\frac{1}{\|\mathbf{w}\|_2}\right)$, $\lambda_3 = \Theta\left(\frac{1}{\|\boldsymbol{\beta}\|_2^3}\right)$, and $\lambda_4 = \Theta\left(\frac{1}{\|\boldsymbol{\beta}\|_2}\right)$. Therefore, $I_{x,z,y|\theta}$ has four eigenvalues given as $\left\{ \Theta\left(\frac{1}{\|\mathbf{w}\|_2^3}\right), \Theta\left(\frac{1}{\|\mathbf{w}\|_2}\right), \Theta\left(\frac{1}{\|\boldsymbol{\beta}\|_2^3}\right), \Theta\left(\frac{1}{\|\boldsymbol{\beta}\|_2}\right) \right\}$. It follows that $I_{x,z,y|\theta}^{-1}$ has the form

$$I_{x,z,y|\theta}^{-1} = R_M \begin{pmatrix} 1/\lambda_1 & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & 1/\lambda_2 & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & 1/\lambda_3 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & 1/\lambda_4 \end{pmatrix} R_M^\top.$$

Therefore, $I_{x,z,y|\theta}^{-1}$ has four eigenvalues given as $\left\{ \Theta\left(\|\mathbf{w}\|_2^3\right), \Theta\left(\|\mathbf{w}\|_2\right), \Theta\left(\|\boldsymbol{\beta}\|_2^3\right), \Theta\left(\|\boldsymbol{\beta}\|_2\right) \right\}$.

Now, for $I_{z|x,y,\theta}$, we first derive a more compact form for the conditional distribution of the latent variable, $P(z|x, y; \theta)$. From simple Bayes rule and algebraic manipulation, we see that

$$\begin{aligned} P(z|x, y; \theta) &= \frac{P(y|x, z; \theta)P(z|x; \theta)p(x)}{p(x, y; \theta)} \\ &= \frac{\frac{\exp\{\frac{yz+1}{2}\mathbf{x}^\top \boldsymbol{\beta}\}}{1+e^{\mathbf{x}^\top \boldsymbol{\beta}}} \frac{\exp\{\frac{z+1}{2}\mathbf{x}^\top \mathbf{w}\}}{1+e^{\mathbf{x}^\top \mathbf{w}}}}{\frac{\exp\{\frac{y+1}{2}\mathbf{x}^\top \boldsymbol{\beta}\}}{1+e^{\mathbf{x}^\top \boldsymbol{\beta}}} \frac{\exp\{\mathbf{x}^\top \mathbf{w}\}}{1+e^{\mathbf{x}^\top \mathbf{w}}} + \frac{\exp\{\frac{-y+1}{2}\mathbf{x}^\top \boldsymbol{\beta}\}}{1+e^{\mathbf{x}^\top \boldsymbol{\beta}}} \frac{1}{1+e^{\mathbf{x}^\top \mathbf{w}}}} \\ &= \frac{\exp\left\{\frac{yz+1}{2}\mathbf{x}^\top \boldsymbol{\beta} + \frac{z+1}{2}\mathbf{x}^\top \mathbf{w}\right\}}{\exp\left\{\frac{y+1}{2}\mathbf{x}^\top \boldsymbol{\beta} + \mathbf{x}^\top \mathbf{w}\right\} + \exp\left\{\frac{-y+1}{2}\mathbf{x}^\top \boldsymbol{\beta}\right\}} \\ &= \frac{\exp\left\{\frac{z+1}{2}(y\mathbf{x}^\top \boldsymbol{\beta} + \mathbf{x}^\top \mathbf{w})\right\}}{\exp\{y\mathbf{x}^\top \boldsymbol{\beta} + \mathbf{x}^\top \mathbf{w}\} + 1} \\ &= \frac{\exp\left\{\frac{z+1}{2}\left\langle \begin{bmatrix} \mathbf{x} \\ y\mathbf{x} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{w} \end{bmatrix} \right\rangle\right\}}{\exp\left\{\left\langle \begin{bmatrix} \mathbf{x} \\ y\mathbf{x} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{w} \end{bmatrix} \right\rangle\right\} + 1} \end{aligned}$$

Now that we have this simplified form, we are able to derive (61) for $I_{z|x,y,\theta}$:

$$\begin{aligned}
 I_{z|x,y,\theta} &= -\mathbb{E}_{\mathbf{X},Y}\mathbb{E}_{Z|x,y,\theta}\left[\frac{\partial^2}{\partial\theta^2}\log P(z|\mathbf{x},y;\theta)\right] \\
 &= -\mathbb{E}_{\mathbf{X},Y}\mathbb{E}_{Z|x,y,\theta}\left[\frac{\partial^2}{\partial\theta^2}\left(\frac{z+1}{2}\left\langle\begin{bmatrix}\mathbf{x} \\ y\mathbf{x}\end{bmatrix},\boldsymbol{\theta}\right\rangle - \log\left(1 + \exp\left\langle\begin{bmatrix}\mathbf{x} \\ y\mathbf{x}\end{bmatrix},\boldsymbol{\theta}\right\rangle\right)\right)\right] \\
 &= \mathbb{E}_{\mathbf{X},Y}\mathbb{E}_{Z|x,y,\theta}\left[\frac{\partial^2}{\partial\theta^2}\log\left(1 + \exp\left\langle\begin{bmatrix}\mathbf{x} \\ y\mathbf{x}\end{bmatrix},\boldsymbol{\theta}\right\rangle\right)\right] \\
 &= \mathbb{E}_{\mathbf{X},Y}\left[\frac{\partial^2}{\partial\theta^2}\log\left(1 + \exp\left\langle\begin{bmatrix}\mathbf{x} \\ y\mathbf{x}\end{bmatrix},\boldsymbol{\theta}\right\rangle\right)\right] \\
 &= \mathbb{E}_{\mathbf{X},Y}\left[\frac{\exp\left\langle\begin{bmatrix}\mathbf{x} \\ y\mathbf{x}\end{bmatrix},\boldsymbol{\theta}\right\rangle\begin{bmatrix}\mathbf{x} \\ y\mathbf{x}\end{bmatrix}\begin{bmatrix}\mathbf{x} \\ y\mathbf{x}\end{bmatrix}^\top}{\left(1 + \exp\left\langle\begin{bmatrix}\mathbf{x} \\ y\mathbf{x}\end{bmatrix},\boldsymbol{\theta}\right\rangle\right)^2}\right]
 \end{aligned}$$

This expression depends only on the random variable (\mathbf{X}, Y) and the parameter iterate $\boldsymbol{\theta}$. \square

Lemma B.5. For $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and $\mathbf{u} \in \mathbb{R}^d$ and $\|\mathbf{u}\|_2 \geq \sqrt{2}$, the symmetric positive definite matrix

$$\mathbb{E}_{\mathbf{X}}\left[\mathbf{x}\mathbf{x}^\top \frac{e^{\mathbf{x}^\top \mathbf{u}}}{(1 + e^{\mathbf{x}^\top \mathbf{u}})^2}\right] \quad (62)$$

is diagonalizable by an orthonormal matrix $\mathbf{R}_{\mathbf{u}} \in \mathbb{R}^{d \times d}$ and has two eigenvalues, $\lambda_1, \lambda_2 \geq 0$, that satisfy

$$\lambda_1 = \Theta\left(\frac{1}{\|\mathbf{u}\|_2^3}\right) \quad (63)$$

$$\lambda_2 = \Theta\left(\frac{1}{\|\mathbf{u}\|_2}\right). \quad (64)$$

Proof. Recall that Gaussian random variables are rotationally invariant. Specifically, for orthonormal matrix $\mathbf{R} \in \mathbb{R}^{d \times d}$ and $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$, it follows that $\mathbf{R}\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. Moreover, $\mathbf{R}^\top \mathbf{R} = \mathbf{R}\mathbf{R}^\top = \mathbf{I}_d$. Using this notion, we will 1) diagonalize (62), then 2) evaluate the eigenvalues of (62) as the diagonal elements.

Consider the orthonormal rotation matrix $\mathbf{R}_{\mathbf{u}} \in \mathbb{R}^{d \times d}$ that is such that $\mathbf{R}_{\mathbf{u}}\mathbf{u} = e_1\|\mathbf{u}\|_2$ where e_j is the j^{th} canonical vector of \mathbb{R}^d . Using this change of basis matrix, we can now obtain the diagonal matrix,

$$\begin{aligned}
 \mathbf{R}_{\mathbf{u}}\left(\mathbb{E}_{\mathbf{X}}\left[\mathbf{x}\mathbf{x}^\top \frac{e^{\mathbf{x}^\top \mathbf{u}}}{(1 + e^{\mathbf{x}^\top \mathbf{u}})^2}\right]\right)\mathbf{R}_{\mathbf{u}}^\top &= \mathbb{E}_{\mathbf{X}}\left[\mathbf{R}_{\mathbf{u}}\mathbf{x}\mathbf{x}^\top \mathbf{R}_{\mathbf{u}}^\top \frac{e^{\mathbf{x}^\top \mathbf{R}_{\mathbf{u}}^\top \mathbf{R}_{\mathbf{u}}\mathbf{u}}}{(1 + e^{\mathbf{x}^\top \mathbf{R}_{\mathbf{u}}^\top \mathbf{R}_{\mathbf{u}}\mathbf{u}})^2}\right] \\
 &= \mathbb{E}_{\mathbf{R}_{\mathbf{u}}\mathbf{x}}\left[\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top \frac{e^{\tilde{\mathbf{x}}^\top e_1\|\mathbf{u}\|_2}}{(1 + e^{\tilde{\mathbf{x}}^\top e_1\|\mathbf{u}\|_2})^2}\right] \\
 &= \begin{pmatrix} \mathbb{E}_{\mathbf{R}_{\mathbf{u}}\mathbf{x}}\left[\frac{\tilde{x}_1^2 e^{\tilde{x}_1\|\mathbf{u}\|_2}}{(1 + e^{\tilde{x}_1\|\mathbf{u}\|_2})^2}\right] & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{E}_{\mathbf{R}_{\mathbf{u}}\mathbf{x}}\left[\frac{e^{\tilde{x}_1\|\mathbf{u}\|_2}}{(1 + e^{\tilde{x}_1\|\mathbf{u}\|_2})^2}\right] \end{pmatrix}.
 \end{aligned}$$

It has only two eigenvalues given in closed form as

$$\begin{aligned}
 \lambda_1 &= \mathbb{E}_{\tilde{X}_1}\left[\frac{\tilde{x}_1^2 e^{\tilde{x}_1\|\mathbf{u}\|_2}}{(1 + e^{\tilde{x}_1\|\mathbf{u}\|_2})^2}\right] = \int_{-\infty}^{\infty} \frac{\tilde{x}_1^2 e^{\tilde{x}_1\|\mathbf{u}\|_2}}{(1 + e^{\tilde{x}_1\|\mathbf{u}\|_2})^2} p(\tilde{x}_1) d\tilde{x}_1 \\
 \lambda_2 &= \mathbb{E}_{\tilde{X}_1}\left[\frac{e^{\tilde{x}_1\|\mathbf{u}\|_2}}{(1 + e^{\tilde{x}_1\|\mathbf{u}\|_2})^2}\right] = \int_{-\infty}^{\infty} \frac{e^{\tilde{x}_1\|\mathbf{u}\|_2}}{(1 + e^{\tilde{x}_1\|\mathbf{u}\|_2})^2} p(\tilde{x}_1) d\tilde{x}_1.
 \end{aligned}$$

The rest of the proof is spent evaluating tight lower and upper bounds on λ_1, λ_2 in terms of $\|\mathbf{u}\|_2$.

Part a): Bounds for λ_1 .

For $\tilde{x}_1 \sim \mathcal{N}(0, 1)$ the probability density function is bounded above by 1: $p(\tilde{x}_1) \leq 1$. Then we can upper bound λ_1 as follows:

$$\begin{aligned} \lambda_1 &= \int_{-\infty}^{\infty} \frac{\tilde{x}_1^2 e^{\tilde{x}_1 \|\mathbf{u}\|_2}}{(1 + e^{\tilde{x}_1 \|\mathbf{u}\|_2})^2} p(\tilde{x}_1) d\tilde{x}_1 \\ &\leq \int_{-\infty}^{\infty} \frac{\tilde{x}_1^2 e^{\tilde{x}_1 \|\mathbf{u}\|_2}}{(1 + e^{\tilde{x}_1 \|\mathbf{u}\|_2})^2} d\tilde{x}_1 \\ &= 2 \int_0^{\infty} \tilde{x}_1^2 e^{-\tilde{x}_1 \|\mathbf{u}\|_2} d\tilde{x}_1 \\ &= \frac{4}{\|\mathbf{u}\|_2^3} \\ &= \mathcal{O}\left(\frac{1}{\|\mathbf{u}\|_2^3}\right) \end{aligned}$$

For the lower bounds, we will use the fact that $e^{-(\|\mathbf{u}\|_2 \tilde{x} + \tilde{x}^2/2)} \geq e^{-\|\mathbf{u}\|_2^2 \tilde{x}^2/2}$ for $x \in \left[\frac{4}{\|\mathbf{u}\|_2}, \infty\right]$ and $\|\mathbf{u}\|_2 \geq \sqrt{2}$. Then, we can lower bound λ_1 as follows:

$$\begin{aligned} \lambda_1 &= \int_{-\infty}^{\infty} \frac{\tilde{x}_1^2 e^{\tilde{x}_1 \|\mathbf{u}\|_2}}{(1 + e^{\tilde{x}_1 \|\mathbf{u}\|_2})^2} p(\tilde{x}_1) d\tilde{x}_1 \\ &\geq 2 \int_0^{\infty} \frac{\tilde{x}_1^2 e^{\tilde{x}_1 \|\mathbf{u}\|_2}}{(2e^{\tilde{x}_1 \|\mathbf{u}\|_2})^2} \left(\frac{e^{-\frac{\tilde{x}_1^2}{2}}}{\sqrt{2\pi}}\right) d\tilde{x}_1 \\ &= \frac{1}{2\sqrt{2\pi}} \int_0^{\infty} \tilde{x}_1^2 e^{-\tilde{x}_1 \|\mathbf{u}\|_2 - \frac{\tilde{x}_1^2}{2}} d\tilde{x}_1 \\ &\geq \frac{1}{2\sqrt{2\pi}} \int_{\frac{4}{\|\mathbf{u}\|_2}}^{\infty} \tilde{x}_1^2 e^{-\|\mathbf{u}\|_2^2 \tilde{x}_1^2/2} d\tilde{x}_1 \\ &\geq \frac{1}{4\sqrt{2\pi} \|\mathbf{u}\|_2^3} \left(\sqrt{\pi} \operatorname{erf}(\tilde{x}_1 \|\mathbf{u}\|_2) - 2\|\mathbf{u}\|_2 \tilde{x}_1 e^{-\|\mathbf{u}\|_2^2 \tilde{x}_1^2/2} \right)_{\frac{4}{\|\mathbf{u}\|_2}}^{\infty} \\ &\geq \Omega\left(\frac{1}{\|\mathbf{u}\|_2^3}\right). \end{aligned}$$

Therefore, it holds that $\lambda_1 = \Theta\left(\frac{1}{\|\mathbf{u}\|_2^3}\right)$.

Part b): Bounds for λ_2 .

For $\tilde{x}_1 \sim \mathcal{N}(0, 1)$ the probability density function is bounded above by 1: $p(\tilde{x}_1) \leq 1$. Then we can upper bound λ_2 as follows:

$$\begin{aligned} \lambda_2 &= \int_{-\infty}^{\infty} \frac{e^{\tilde{x}_1 \|\mathbf{u}\|_2}}{(1 + e^{\tilde{x}_1 \|\mathbf{u}\|_2})^2} p(\tilde{x}_1) d\tilde{x}_1 \\ &\leq \int_{-\infty}^{\infty} \frac{e^{\tilde{x}_1 \|\mathbf{u}\|_2}}{(1 + e^{\tilde{x}_1 \|\mathbf{u}\|_2})^2} d\tilde{x}_1 \\ &= 2 \int_0^{\infty} e^{-\tilde{x}_1 \|\mathbf{u}\|_2} d\tilde{x}_1 \\ &= \frac{1}{\|\mathbf{u}\|_2} \\ &= \mathcal{O}\left(\frac{1}{\|\mathbf{u}\|_2}\right) \end{aligned}$$

For the lower bounds, we will use the fact that $e^{-(\|\mathbf{u}\|_2 \tilde{\mathbf{x}}_1 + \tilde{\mathbf{x}}_1^2/2)} \geq e^{-\|\mathbf{u}\|_2^2 \tilde{\mathbf{x}}_1^2/2}$ for $x \in \left[\frac{4}{\|\mathbf{u}\|_2}, \infty\right]$ and $\|\mathbf{u}\|_2 \geq \sqrt{2}$. Then, we can lower bound λ_2 as follows:

$$\begin{aligned}
 \lambda_2 &= \int_{-\infty}^{\infty} \frac{e^{\tilde{\mathbf{x}}_1 \|\mathbf{u}\|_2}}{(1 + e^{\tilde{\mathbf{x}}_1 \|\mathbf{u}\|_2})^2} p(\tilde{\mathbf{x}}_1) d\tilde{\mathbf{x}}_1 \\
 &\geq 2 \int_0^{\infty} \frac{e^{\tilde{\mathbf{x}}_1 \|\mathbf{u}\|_2}}{(2e^{\tilde{\mathbf{x}}_1 \|\mathbf{u}\|_2})^2} \left(\frac{e^{-\frac{\tilde{\mathbf{x}}_1^2}{2}}}{\sqrt{2\pi}} \right) d\tilde{\mathbf{x}}_1 \\
 &= \frac{1}{2\sqrt{2\pi}} \int_0^{\infty} e^{-\tilde{\mathbf{x}}_1 \|\mathbf{u}\|_2 - \frac{\tilde{\mathbf{x}}_1^2}{2}} d\tilde{\mathbf{x}}_1 \\
 &\geq \frac{1}{2\sqrt{2\pi}} \int_{\frac{4}{\|\mathbf{u}\|_2}}^{\infty} e^{-\|\mathbf{u}\|_2^2 \tilde{\mathbf{x}}_1^2/2} d\tilde{\mathbf{x}}_1 \\
 &\geq \frac{1}{4\sqrt{2\pi} \|\mathbf{u}\|_2} (\sqrt{\pi} \operatorname{erf}(\tilde{\mathbf{x}}_1 \|\mathbf{u}\|_2))_{\frac{4}{\|\mathbf{u}\|_2}}^{\infty} \\
 &\geq \Omega\left(\frac{1}{\|\mathbf{u}\|_2}\right).
 \end{aligned}$$

Therefore, it holds that $\lambda_2 = \Theta\left(\frac{1}{\|\mathbf{u}\|_2}\right)$. □

B.7. Existence of Locally Convex Region

In this section, we further discuss the consequences of Corollary B.1, Theorem B.2, and Theorem B.4.

To begin, we will discuss the sufficient condition for Gradient Descent to converge linearly to the true parameters θ^* and compare it to that of EM. For Gradient Descent, it is well understood that if θ^1 is initialized in a convex set Θ that contains θ^* and where $\mathcal{L}(\theta)$ is strongly convex, i.e.

$$\nabla^2 \mathcal{L}(\theta) \succeq \alpha \mathbf{I}_{2d} \quad \text{for all } \theta \in \Theta, \quad (65)$$

then the parameter iterates converge linearly to θ^* . However, as we have shown for SymMoLinE and SymMoLogE, the sufficient condition for EM to converge linearly to θ^* is slightly different. Instead, we require Θ to satisfy that $\mathcal{L}(\theta)$ is strongly convex relative to $A(\theta)$, i.e.,

$$\nabla^2 \mathcal{L}(\theta) \succeq \alpha \nabla^2 A(\theta) \quad \text{for all } \theta \in \Theta. \quad (66)$$

Interestingly, if it holds that $A(\theta)$ is 1-smooth, we see that (66) is weaker than (65), i.e.,

$$\nabla^2 A(\theta) \preceq \mathbf{I}_{2d} \text{ and (65) holds} \implies (66). \quad (67)$$

But more interestingly, as long as $A(\theta)$ is μ -smooth for some $\mu > 0$, it will hold that any set Θ satisfying (65) for some $\alpha > 0$ will also satisfy (66) with $\tilde{\alpha} = \frac{\alpha}{\mu} > 0$. We note that the converse holds for $A(\theta)$ strongly convex. In summary, it then holds that, for SymMoLinE and SymMoLogE, EM's sufficient conditions for a linear rate is strictly weaker than that of Gradient Descent when the mirror map $A(\theta)$ is convex, but not strongly convex.

Next, we will further discuss the implications of Theorem B.4. In the theorem, we obtain clear lower and upper bounds for the eigenvalues of $\mathbf{I}_{\mathbf{x}, \mathbf{y}, \mathbf{z} | \theta}$. However, it is not clear how to do the same for $\mathbf{I}_{\mathbf{z} | \mathbf{x}, \mathbf{y}, \theta}$ as the trick to rotate the axis with an orthonormal matrix R to simplify the expression will not work here because the distribution of the vector $(\mathbf{x}, \mathbf{y}\mathbf{x})^\top$ is not invariant to rotation. Still, there are some things that we can say for special cases. For this discussion, we will constrain ourselves to SymMoLogE. However, the same lines of reasoning also apply to SymMoLinE. First we recall from Theorem B.4 that

$$\mathbf{I}_{\mathbf{z} | \mathbf{x}, \mathbf{y}, \theta} = \mathbb{E}_{\mathbf{X}, \mathbf{Y}} \left[\frac{\exp \left\langle \begin{bmatrix} \mathbf{x} \\ \mathbf{y}\mathbf{x} \end{bmatrix}, \theta \right\rangle \begin{bmatrix} \mathbf{x} \\ \mathbf{y}\mathbf{x} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y}\mathbf{x} \end{bmatrix}^\top}{\left(1 + \exp \left\langle \begin{bmatrix} \mathbf{x} \\ \mathbf{y}\mathbf{x} \end{bmatrix}, \theta \right\rangle \right)^2} \right].$$

From here, one easy way to approach bounding the above is to 1) recall that any outer product of the form $\mathbf{u}\mathbf{u}^\top$ has a single eigenvalue given as $\|\mathbf{u}\|_2^2$, and 2) the inner product between two vectors is equal to the product of their norms and the cosine of the angle between them, i.e. $\mathbf{s}^\top \mathbf{u} = \|\mathbf{s}\|_2 \|\mathbf{u}\|_2 \cos(\phi_{\mathbf{s}, \mathbf{u}})$. Thus, denoting ϕ to be the angle between $\boldsymbol{\theta}$ and the vector $(\mathbf{x}, y\mathbf{x})^\top$, we obtain

$$\begin{aligned} \mathbf{I}_{z|\mathbf{x}, y, \boldsymbol{\theta}} &\preceq \mathbb{E}_{\mathbf{X}, Y, \Phi} \left[\frac{e^{\sqrt{(1+y^2)}\|\mathbf{x}\|_2^2 \|\boldsymbol{\theta}\|_2 \cos(\phi)} (1+y^2) \|\mathbf{x}\|_2^2}{\left(1 + e^{\sqrt{(1+y^2)}\|\mathbf{x}\|_2^2 \|\boldsymbol{\theta}\|_2 \cos(\phi)}\right)^2} \right] \mathbf{I}_{2d} \\ &\preceq \mathbb{E}_{\mathbf{X}, Y, \Phi} \left[e^{-|\sqrt{(1+y^2)}\|\mathbf{x}\|_2^2 \|\boldsymbol{\theta}\|_2 \cos(\phi)} (1+y^2) \|\mathbf{x}\|_2^2 \right] \mathbf{I}_{2d}. \end{aligned}$$

Denoting $s := \left| \sqrt{(1+y^2)}\|\mathbf{x}\|_2^2 \right|$, we can write the above expectation as

$$8 \int_0^{\pi/2} \int_0^\infty s^2 e^{-s\|\boldsymbol{\theta}\|_2 \cos(\phi)} p(s, \phi) ds d\phi.$$

Subsequently, the idea is bound $p(s, \phi) = p(s)p(\phi|s)$ in a way that makes integration easy.

The case where $x, w, \beta \in \mathbb{R}$ is fairly easy. Under this scenario, $\mathbf{x} \sim \mathcal{N}(0, 1)$ and $\mathbf{I}_{\mathbf{x}, y, z|\boldsymbol{\theta}}$ can be upper-bounded as

$$\begin{aligned} \mathbf{I}_{z|\mathbf{x}, y, \boldsymbol{\theta}} &= \mathbb{E}_{X, Y} \left[\begin{pmatrix} x^2 & yx^2 \\ yx^2 & y^2 x^2 \end{pmatrix} \frac{e^{x(w+y\beta)}}{(1 + e^{x(w+y\beta)})^2} \right] \\ &\leq \mathbb{E}_{X, Y} \left[x^2 (1+y^2) e^{-|x(w+y\beta)|} \right] \mathbf{I}_2. \end{aligned}$$

Now, recall that $y \in \{-1, 1\}$, $P(y|x) \leq 1$, $p(x) = \frac{e^{-x^2/2}}{\sqrt{2\pi}} \leq \frac{2e^{-|x|}}{\sqrt{2\pi}}$, and see that

$$\begin{aligned} \mathbf{I}_{z|\mathbf{x}, y, \boldsymbol{\theta}} &= \mathbb{E}_X \left[2x^2 \left(e^{-|x(w-\beta)|} P(y=1|x) + e^{-|x(w+\beta)|} P(y=-1|x) \right) \right] \\ &\leq \mathbb{E}_X \left[2x^2 \left(e^{-|x(w-\beta)|} + e^{-|x(w+\beta)|} \right) \right] \\ &= 4 \int_0^\infty x^2 \left(e^{-|x(w-\beta)|} + e^{-|x(w+\beta)|} \right) p(x) dx \\ &\leq 4 \int_0^\infty x^2 \left(e^{-x(|w-\beta|+1)} + e^{-x(|w+\beta|+1)} \right) dx \\ &\leq 8 \left(\frac{1}{(1+|w-\beta|)^3} + \frac{1}{(1+|w+\beta|)^3} \right) \\ &\leq \mathcal{O} \left(\frac{1}{(1+|w-\beta|)^3} + \frac{1}{(1+|w+\beta|)^3} \right). \end{aligned}$$

Subsequently, together with the fact that $\mathbf{I}_{\mathbf{x}, y, z, \boldsymbol{\theta}}^{-1} \leq \max \{ \mathcal{O}(\|\mathbf{w}\|_2^3), (\|\boldsymbol{\beta}\|_2^3) \}$, it holds that the eigenvalues of the MIM are upper-bounded by

$$\max \left\{ \mathcal{O} \left(\left(\frac{w}{1+|w-\beta|} \right)^3 + \left(\frac{w}{1+|w+\beta|} \right)^3 \right), \mathcal{O} \left(\left(\frac{\beta}{1+|w-\beta|} \right)^3 + \left(\frac{\beta}{1+|w+\beta|} \right)^3 \right) \right\}.$$

This special case is closely related to the case where $\boldsymbol{\beta}$ is parallel to \mathbf{w} ; a similar approach will work.

C. Additional Experiments

C.1. Experiment on Grayscale CIFAR-10

In this section, we consider an additional mini-batch experiment on grayscale CIFAR-10 with a 5-component MoE. The MoE to be trained consists of individual experts that each consist of a single hidden layer MLP with hidden dimension 100 and ReLU activation. The gating function also consists of a single hidden layer MLP with hidden dimension 100 and ReLU activation. We randomly initialize each linear layer to have rows that are unit-norm and execute the algorithms on the same datasets and with the same initializations. For Gradient EM, the only additional code needed over GD is to define the EM Loss function appropriately, and then perform a Gradient Step on the Gating parameters and the Expert parameters separately as describe in Algorithm 2. For EM, for each iteration, we perform several gradient steps in an inner loop to approximately recover the solutions to the sub-problems described in (9). We report our findings for the mini-batch iteration of the respective algorithms in Figure 4.

We report the respective final test accuracy and cross-entropy loss values after 50 epochs of EM, Gradient EM and GD for fitting a 5-component MoE on the CIFAR-10. We see that EM boasts a much improved final test accuracy of 41.6%. Meanwhile, Gradient EM also registers an improvement over GD at 37.4%. Finally, GD obtained a final accuracy of 34.5%. While these accuracies are themselves not good, the challenge was to find an optimal partitioning of the data and utilize very weak experts. In Figure 4, we report the progress made on the accuracy for the test set over the 50 epochs, averaged over 25 instances. As was observed in our synthetic experiment and Fashion MNIST experiments, EM takes considerably less iterations to fit the mixture than both Gradient EM and GD, where the former also takes considerably less time to fit the mixture than GD. To validate our results further, we perform a paired t-test. For EM and Gradient EM compared to GD, we obtain a T-statistic ≥ 22 indicating that the difference in final accuracy is statistically significant (p-value ~ 0.000).

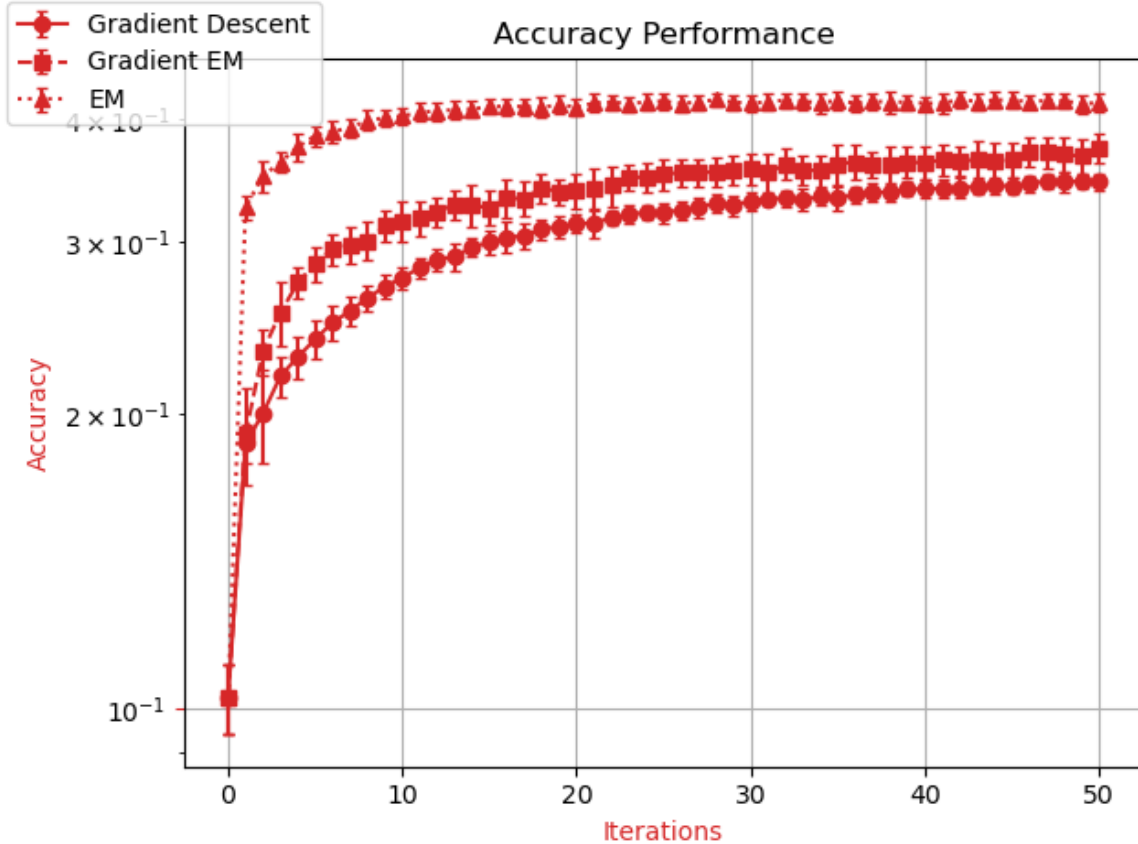


Figure 4. Convergence in predicted label Accuracy for different optimization methods.

D. Algorithms

In this section, we provide explicit formulations of EM, Gradient EM and Gradient Descent for the context of MoE optimization and EM for deep and sparse MoE.

D.1. EM for MoE

Expectation-Maximization (EM): EM takes a structured approach to minimizing the objective $\mathcal{L}(\theta)$ in (7). Each iteration of EM is decomposed into two steps as follows. The first step is called “expectation”: For current parameter estimate θ^t , we compute the expectation of the complete-data log-likelihood with respect to the latent variables, using the current parameter estimates θ^t and denote it by $Q(\theta|\theta^t)$, i.e.,

$$Q(\theta|\theta^t) = -\mathbb{E}_{\mathbf{X}, Y} [\mathbb{E}_{Z|\mathbf{x}, y; \theta^t} [\log p(\mathbf{x}, y, z; \theta)]] . \quad (68)$$

Then, in the second step called “maximization”, we simply minimize the objective $Q(\theta|\theta^t)$ (or maximize $-Q(\theta|\theta^t)$) with respect to $\theta \in \Omega$ and obtain our new parameter as

$$\theta^{t+1} := \underset{\theta \in \Omega}{\operatorname{argmin}} Q(\theta|\theta^t). \quad (69)$$

For MoE described in Section 2, $\log p(y, z|\mathbf{x}; \theta) = \log p(y|z, \mathbf{x}; \beta) + \log p(z|\mathbf{x}; \mathbf{w})$. It follows that the EM objective (8) is linearly separable in the parameters β and \mathbf{w} . Thus, we can rewrite $Q(\theta|\phi)$ as the sum of two functions that depend only on β and \mathbf{w} , respectively. Subsequently, the EM update (9) is obtained as the concatenation $\theta^{t+1} = (\mathbf{w}^{t+1}, \beta^{t+1})^\top$, where

$$\begin{aligned} \mathbf{w}^{t+1} &= \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} -\mathbb{E}_{\mathbf{X}, Y} [\mathbb{E}_{Z|\mathbf{x}, y; \theta^t} [\log p(z|\mathbf{x}; \mathbf{w})]] , \\ \beta^{t+1} &= \underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} -\mathbb{E}_{\mathbf{X}, Y} [\mathbb{E}_{Z|\mathbf{x}, y; \theta^t} [\log p(y|z, \mathbf{x}; \beta)]] . \end{aligned}$$

Algorithm 1 EM for MoE

Input: Initial $\theta^1 \in \Omega$, data: $(\mathbf{X}, Y) \sim p(\mathbf{x}, y; \theta^*)$
for $t = 1$ to T **do**
 θ -Update: Obtain θ^{t+1} as
 $\theta^{t+1} := \arg \min_{\theta \in \Omega} Q(\theta | \theta^t)$
end for
Output: $\theta^T = (\mathbf{w}^T, \beta^T)$

D.2. Gradient EM for MoE

Gradient EM. Whereas EM performs the global minimization of the EM objective given in (68), Gradient EM obtains its next parameter iterate as the concatenation of two gradient updates on the sub objectives,

$$-\mathbb{E}_{\mathbf{X}, Y} [\mathbb{E}_{Z|\mathbf{x}, y; \theta^t} [\log P(z|\mathbf{x}; \mathbf{w})]] \quad (70)$$

$$-\mathbb{E}_{\mathbf{X}, Y} [\mathbb{E}_{Z|\mathbf{x}, y; \theta^t} [\log p(y|z, \mathbf{x}; \beta)]] \quad (71)$$

where the EM objective is given as the summation of (70) and (71).

Algorithm 2 Gradient EM for MOE

Input: Initial $\theta^1 \in \Omega$, data: $(\mathbf{X}, Y) \sim p(\mathbf{x}, y; \theta^*)$, step-size: $\gamma_1, \gamma_2 \in (0, \infty)$.

for $t = 1, \dots, T$: **do**

β -Update: Obtain β^{t+1} as

$$\beta^{t+1} = \beta^t + \gamma_1 \mathbb{E}_{\mathbf{X}, Y} \mathbb{E}_{Z|\mathbf{x}, y; \theta^t} \left[\frac{\partial}{\partial \beta} \log p(y|z, \mathbf{x}; \beta) \right].$$

w -Update: Obtain w^{t+1} as

$$w^{t+1} = w^t + \gamma_2 \mathbb{E}_{\mathbf{X}, Y} \mathbb{E}_{Z|\mathbf{x}, y; \theta^t} \left[\frac{\partial}{\partial w} \log P(z|\mathbf{x}; w) \right].$$

end for

Output: $\theta^T = (w^T, \beta^T)$

D.3. Gradient Descent for MoE

Gradient Descent. Gradient descent is given as the global minimizer of the first order approximation of $\mathcal{L}(\theta)$ at θ plus a quadratic regularizer, i.e.,

$$\mathcal{L}(\theta^t) + \langle \nabla \mathcal{L}(\theta^t), \theta - \theta^t \rangle + \frac{1}{2\eta} \|\theta - \theta^t\|_2^2.$$

Differentiating, and solving for equality at 0 yields the well known gradient update.

Algorithm 3 Gradient Descent for MoE

Input: Initial $\theta^1 \in \Omega$, data: $(\mathbf{X}, Y) \sim p(\mathbf{x}, y; \theta^*)$, step-size: $\gamma \in \mathbb{R}^+$

for $t = 1$ to T **do**

θ -Update:

$$\theta^{t+1} := \theta^t - \gamma \nabla \mathcal{L}(\theta^t)$$

end for

Output: $\theta^T = (w^T, \beta^T)$

D.4. EM for Deep and Sparse MoE

We begin by formalizing the concepts of Deep and Sparse Mixtures of Experts (MoE), extending the classical MoE framework.

Deep MoE: A *deep MoE* is a composition of $l \geq 2$ MoE blocks, denoted as $\text{MoE}_1, \text{MoE}_2, \dots, \text{MoE}_l$, stacked sequentially. Each block MoE_i consists of a gating function $g_i(\mathbf{x}; w_i)$ and a set of k experts $\{f_{i,j}(\mathbf{x}; \beta_{i,j})\}_{j=1}^k$. The input \mathbf{x} is processed through each MoE block in sequence, producing intermediate representations $\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_l$, where:

$$\begin{aligned} \mathbf{h}_1 &= \sum_{j=1}^k g_1(\mathbf{x}; w_1)_j f_{1,j}(\mathbf{x}; \beta_{1,j}), \\ \mathbf{h}_i &= \sum_{j=1}^k g_i(\mathbf{h}_{i-1}; w_i)_j f_{i,j}(\mathbf{h}_{i-1}; \beta_{i,j}) \quad \text{for } i = 2, \dots, l. \end{aligned}$$

The final output is $\mathbf{y} = \mathbf{h}_l$.

Sparse MoE: The *Sparse MoE* is a variant of the MoE that is popular in deep MoE applications where, at each MoE block, only a small subset of experts (typically one or a few) are activated per input both during training and inference. This is achieved via a deterministic or stochastic selection mechanism (e.g., top- k gating), resulting in a sparse latent variable $\mathbf{z} = (z_1, z_2, \dots, z_l) \in [k]^l$ that encodes the sequence of selected experts across the l blocks.

EM for Deep and Sparse MoE: We now describe an EM-like algorithm for training deep and sparse MoE models. The key idea is to treat the expert selection sequence $z = (z_1, \dots, z_l)$ as a latent variable and optimize the expected complete-data log-likelihood. We let $\theta = (\mathbf{w}_1, \dots, \mathbf{w}_l, \beta_{1,1}, \dots, \beta_{l,k})$ denote all model parameters.

In the classical *non-sparse MoE setting* where we do not choose a subset of experts to go through at each layer, the latent variables can only be resolved at the last layer. The EM surrogate is then given as follows:

$$Q(\theta|\theta^t) = -\mathbb{E}_{X,Y} [\mathbb{E}_{Z_l|\mathbf{x},y;\theta^t} [\log p(y|\mathbf{x}, z_l; \theta) P(z_l|\mathbf{x}; \theta)]] . \quad (72)$$

where the given probability functions is given at the last layer as

$$\begin{aligned} p(y|\mathbf{x}, z_1, \dots, z_l; \theta) &= p(y|\mathbf{h}_{l-1}) \\ P(z_1, \dots, z_l|\mathbf{x}; \theta) &= p(z_l|\mathbf{h}_{l-1}; \theta). \end{aligned}$$

In the *sparse MoE setting* where p experts are chosen at each layer, we re-define the latent variables Z_i to be defined over the set of all possible ordered combinations of p experts that could have been chosen out of the k available experts. The latent space embeddings is now given to be $\hat{h}_i := \sum_{j \in z_i} g_i(\mathbf{h}_{i-1}; \mathbf{w}_i)_j f_{i,j}(\mathbf{h}_{i-1}; \beta_{i,j})$. The EM surrogate is given by

$$Q(\theta|\theta^t) = -\mathbb{E}_{X,Y} [\mathbb{E}_{Z_1, \dots, Z_l|\mathbf{x},y;\theta^t} [\log p(y|\mathbf{x}, z_1, \dots, z_l; \theta) P(z_1, \dots, z_l|\mathbf{x}; \theta)]] . \quad (73)$$

where the given probability functions are decomposed per layer as

$$\begin{aligned} p(y|\mathbf{x}, z_1, \dots, z_l; \theta) &= p(y|\hat{\mathbf{h}}_{l-1}) \\ P(z_1, \dots, z_l|\mathbf{x}; \theta) &= P(z_1|\mathbf{x}; \theta) \cdots p(z_l|z_1, \dots, z_{l-1}; \theta). \end{aligned}$$

In practice, constructing the above EM surrogate this would require to make a combinatorial number of forward passes through the model to evaluate the surrogate loss function. To remediate this, we will use the strategy to only evaluate the loss on the greedily chosen sequence through the model as is the case for GD-type solutions for Sparse MoE.