

# Divergence-Minimization for Latent-Structure Models: Monotone Operators, Contraction Guarantees, and Robust Inference

Lei Li

Department of Statistics, George Mason University

and

Anand N. Vidyashankar

Department of Statistics, George Mason University

November 25, 2025

## Abstract

We develop a divergence-minimization (DM) framework for robust and efficient inference in latent-mixture models. By optimizing a residual-adjusted divergence, the DM approach recovers EM as a special case and yields robust alternatives through different divergence choices. We establish that the sample objective decreases monotonically along the iterates, leading the DM sequence to stationary points under standard conditions, and that at the population level the operator exhibits local contractivity near the minimizer. Additionally, we verify consistency and  $\sqrt{n}$ -asymptotic normality of minimum-divergence estimators and of finitely many DM iterations, showing that under correct specification their limiting covariance matches the Fisher information. Robustness is analyzed via the residual-adjustment function, yielding bounded influence functions and a strictly positive breakdown bound for bounded-RAF divergences, and we contrast this with the non-robust behaviour of KL/EM. Next, we address the challenge of determining the number of mixture components by proposing a penalized divergence criterion combined with repeated sample splitting, which delivers consistent order selection and valid post-selection inference. Empirically, DM instantiations based on Hellinger and negative exponential divergences deliver accurate inference and remain stable under contamination in mixture and image-segmentation tasks. The results clarify connections to MM and proximal-point methods and offer practical defaults, making DM a drop-in alternative to EM for robust latent-structure inference.

Keywords: DM algorithm; divergence-based mixture complexity; sample splitting; computational complexity

# 1 Introduction

Models with latent structure, such as finite mixture models (FMM), are ubiquitous in clustering, image segmentation, and density modeling across fields ranging from astronomy to medical research. A standard approach in many of these settings is the classical expectation–maximization (EM) algorithm and its variants (e.g., stochastic EM) for model fitting and inference. The properties of EM, which are likelihood-based, have been well studied when the model is correctly specified. However, when the model is misspecified or the data contain aberrant observations, likelihood-based methods can perform poorly, and minimum divergence estimators, also referred to as minimum disparity estimators (see [Lindsay \(1994\)](#), [Basu, Sarkar & Vidyashankar \(1997\)](#)), have been proposed as robust and efficient alternatives to likelihood-based inference. For models with latent structure, however, a unified, divergence-agnostic treatment with operator-level guarantees remains underdeveloped.

The primary objective of this paper is to develop a general divergence-minimization (DM) algorithm for latent-structure models and (i) to establish monotone descent and local contractivity properties, (ii) large-sample guarantees (consistency and  $\sqrt{n}$ -asymptotic normality), and (iii) a divergence-based mixture complexity estimator via multiple sample splitting. As a consequence of our operator-theoretic view, existing methods such as EM and proximal-point updates appear as special cases of the DM algorithm. A brief literature review of divergence-based methods for latent structure models appears below; an extended review is deferred to Appendix A.

## 1.1 Background and literature review

Divergence-based methods are widely used due to their ability to unify classical and robust inference. From [Beran \(1977\)](#) through [Lindsay \(1994\)](#) (among others), it is well established

that divergence-based estimators are robust to misspecification and first-order efficient under correct specification. In finite mixtures, [Woodward et al. \(1995\)](#) analyzed a two-component normal mixture and used the minimum Hellinger distance (MHD) to estimate the mixing proportion. [Cutler & Cordero-Braña \(1996\)](#) estimated all parameters in normal mixtures via MHD, referred to as the HMIX algorithm. For discrete mixtures, [Karlis & Xekalaki \(1998\)](#) studied Poisson mixtures using HELMIX, a *Poisson-specific* variant whose update relies on recurrence relations for the Poisson pmf and on the special algebraic form of the of the Hellinger distance; it does not readily generalize either to other divergences or to other count families. In particular, HELMIX does not extend to the Poisson–Gamma (Negative Binomial) or Poisson–lognormal mixtures that we study here, for which no Hellinger-type EM algorithm appears to be available in the literature. Our divergence–minimization (DM) algorithm provides EM-type procedures for these models in a unified way, yielding estimators that are both robust and first-order efficient under the correctly specified model (their limiting covariance equals the Fisher information). Beyond Hellinger distance, the DM framework accommodates negative exponential and related divergences, which are known to improve stability in the presence of outliers and inliers ([Lindsay 1994](#), [Basu, Sarkar & Vidyashankar 1997](#)).

Our finite–step convergence and large–sample results are essentially divergence–agnostic within the DM class: under mild curvature and regularity conditions, the DM iterates are  $\sqrt{n}$ –consistent and first–order efficient under the correctly specified model, with limiting covariance equal to the Fisher information, for any disparity  $G$  in this class. Robustness, however, is driven not by the DM scheme itself but by the residual adjustment factor  $A$  associated with  $G$ : when  $A$  is bounded away from  $-1$  on the relevant range, the resulting DM estimators have a strictly positive breakdown point and bounded gross–error sensitivity to outliers, while uniform boundedness of  $A$  yields additional stability under inlier

contamination, in line with the behaviour of negative exponential and related divergences. In contrast, classical EM, including commonly used weight-flooring schemes, does not yield robustness: its influence function is unbounded and its breakdown point is essentially zero. The remainder of the paper is organized as follows. Section 2 introduces the DM functional and algorithm in a principled way via variational elimination of the latent weights, and clarifies its relationships to EM, MM, and proximal-point methods. Section 3 studies the resulting algorithm, establishing monotone descent, local contraction for the population and sample DM maps, and path properties of the iterates. Section 4 investigates the asymptotic behaviour of the sample-level iterates, showing  $\sqrt{n}$ -consistency, first-order efficiency under the correctly specified model, finite-step CLT and (Godambe-)Wilks properties, and robustness in the form of bounded influence functions and strictly positive breakdown points. Finally, Section 5 addresses unknown mixture complexity: we introduce a generalized divergence information criterion (GDIC), study its robustness, and combine it with a simple sample-splitting scheme (using a majority-vote aggregation across splits) to obtain valid post-selection inference, where the selection split is used to compute GDIC while the estimation split inherits the finite-step CLT and (Godambe-)Wilks guarantees developed earlier. This yields, to our knowledge, the first end-to-end divergence-based procedure that both estimates  $K$  and delivers robust inference in overdispersed latent count mixtures such as Poisson-Gamma and Poisson-lognormal models. Section 6 presents empirical studies and Section 7 case studies; detailed background and historical notes appear in Appendix A. Additional technical details are presented in Appendices B through N.

## 2 DM Algorithm for Models with Latent Structure

We begin this section with a brief description of models with latent structure and minimum divergence estimation.



**Finite mixtures and latent structure.** The random variable  $Y$  with density given by

$$f(y; \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k h_k(y; \boldsymbol{\phi}_k), \quad \boldsymbol{\theta} = (\pi, \boldsymbol{\phi}) \in \Delta_{K-1} \times \Phi, \quad \pi_k > 0, \quad \sum_k \pi_k = 1,$$

is referred to as the finite mixture model (FMM). When

$$h_k(y; \boldsymbol{\phi}_k) = \int_{\mathbb{R}} s_k(y | \lambda) r_k(\lambda | \boldsymbol{\phi}_k) d\rho(\lambda),$$

where  $s_k(\cdot | \lambda)$  is a density for each  $\lambda$  and  $r(\cdot | \boldsymbol{\phi}_k)$  is the density on  $\Lambda$  with parameter  $\boldsymbol{\phi}_k = (\phi_{k1}, \dots, \phi_{kd})$ , we refer to it as Hierarchical Finite Mixture Model (HFMM). In applications, we assume a common conditional kernel, that is,  $s_k \equiv s$ , and allow only the mixing law to vary across components; hence, we write

$$h_k(y; \boldsymbol{\phi}_k) \equiv h(y; \boldsymbol{\phi}_k) = \int_{\Lambda} s(y|\lambda) r(\lambda|\boldsymbol{\phi}_k) d\rho(\lambda),$$

where  $\Lambda \subset \mathbb{R}$ . Turning to the model, let  $Z \in \mathcal{Z} := \{1, \dots, K\}$  denote the random variable representing the unobserved class of  $Y$  and set for  $z = k$ ,  $p(y, z; \boldsymbol{\theta}) = \pi_k h(y; \boldsymbol{\phi}_k)$ . Then, the marginal density of  $Y$  is given by  $f(y; \boldsymbol{\theta}) = \int_{\mathcal{Z}} p(y, z; \boldsymbol{\theta}) d\nu(z)$ , where  $\nu$  is the counting measure on  $(\mathcal{Z}, 2^{\mathcal{Z}})$ .  $\boldsymbol{\pi}$  represents the vector of class-inclusion probabilities and is also referred to as mixing weights. To remove label ambiguity, we order the components by increasing mixing weights and break ties deterministically:  $\pi_1 \leq \pi_2 \leq \dots \leq \pi_K$ , and if  $\pi_j = \pi_{j+1}$  we use a fixed scalar summary  $m(\phi)$  (e.g., a location) with  $m(\phi_j) \leq m(\phi_{j+1})$ . All results are permutation-invariant and hold for any fixed labeling rule; this convention only selects a canonical representative of the label-equivalence class.

Several commonly used models are particular cases of the HFMM. First, setting  $s(y|\lambda) = \text{Poisson}(y|\lambda) = e^{-\lambda} \lambda^y / y!$  and  $r_k(\lambda; \boldsymbol{\phi}_k) = r(\lambda; \boldsymbol{\phi}_k) = \text{Gamma}(\lambda|\alpha_k, \beta_k) = \{\beta_k^{\alpha_k} / \Gamma(\alpha_k)\} \lambda^{\alpha_k-1} e^{-\beta_k \lambda}$ , with  $\boldsymbol{\phi}_k = (\alpha_k, \beta_k)$ ,  $\alpha_k > 0$ ,  $\beta_k > 0$ , one obtains a finite mixture of Poisson-Gamma models. The resulting marginal is the same as a mixture of negative binomial models; that is,

$$h(y; \alpha_k, \beta_k) = \frac{\Gamma(y + \alpha_k)}{\Gamma(\alpha_k) y!} \left( \frac{\beta_k}{\beta_k + 1} \right)^{\alpha_k} \left( \frac{1}{\beta_k + 1} \right)^y.$$

Continuing with the Poisson distribution, setting  $s(y \mid \lambda)$  as above and  $r(\lambda; \phi_k) = \text{Lognormal}(\lambda \mid \mu_k, \sigma_k^2) = \{\lambda \sqrt{2\pi\sigma_k^2}\}^{-1} \exp\left[-(\log \lambda - \mu_k)^2 / (2\sigma_k^2)\right]$ , with  $\phi_k = (\mu_k, \sigma_k^2)$ ,  $\sigma_k^2 > 0$ , we obtain

$$h(y; \mu_k, \sigma_k^2) = \int_0^\infty \text{Poisson}(y \mid \lambda) \text{Lognormal}(\lambda \mid \mu_k, \sigma_k^2) d\lambda.$$

The above expression is evaluated numerically as there is no closed-form expression. We will investigate these two HFMMs, along with the standard  $K$ -component Poisson mixture, in the simulation section. In the rest of the manuscript, we use the following ratio convention: *Ratio convention.* For any ratio  $a/b$  used below (e.g.,  $r(y) = g(y)/f(y; \theta')$ ,  $t(z \mid y) = w(z \mid y; \theta)/w(z \mid y; \theta')$ ), we interpret it only where the denominator is positive; on null sets, we define the ratio to be 1.

**Residuals and divergences.** Let  $g$  denote the target density and let the *Pearson residual* be

$$\delta(y; \theta) = \frac{g(y)}{f(y; \theta)} - 1 \quad (\text{so } g = f(1 + \delta)).$$

Given a thrice differentiable convex generator  $G : [-1, \infty) \rightarrow [0, \infty)$  satisfying  $G(0) = G'(0) = 0$  and  $G''(0) = 1$ , its *residual-adjustment function* (RAF) is given by  $A(\delta) := (1 + \delta)G'(\delta) - G(\delta)$ . From the properties of  $G$ ,  $A(\delta)$  is an increasing function on  $[-1, \infty)$  and carries the relevant information about the trade-off between efficiency and robustness. The residual-adjusted divergence

$$D_G(g, f_\theta) = \int G(\delta(y; \theta)) f(y; \theta) dy.$$

The population minimum-divergence estimator (MDE) solves  $\theta^* \in \arg \min_\theta D_G(g, f_\theta)$ ; in practice  $g$  is replaced by an empirical version  $g_n$  yielding  $\hat{\theta}_n^{(G)} \in \arg \min_\theta D_G(g_n, f_\theta)$ . Different generators  $G$  induce different RAFs and hence different weighting of residuals in the estimating equations. Some of the canonical choices are (i) squared Hellinger distance:  $G_{\text{HD}}(t) = 2(\sqrt{1+t} - 1)^2$ ,  $A_{\text{HD}}(\delta) = 2[(1 + \delta)^{1/2} - 1]$  and (ii) the Negative exponential

divergence (NED):  $G_{\text{NED}}(t) = [\exp(-t) - 1 + t]$ ,  $A_{\text{NED}} = 2 - 2e^{-\delta} - \delta e^{-\delta}$ . A Variant of NED (vNED), with  $g$  and  $f_{\theta}$  flipped, equivalently replacing  $\delta = g/f - 1$  by  $\delta^* = f/g - 1$ , yields a softened left-tail for  $\delta < 0$  to avoid over-penalizing “holes” and helps improve segmentation and count-mixture fits. The details are in the appendix.

**Remark (Limitations of KL for robustness).** For the Kullback–Leibler generator  $G_{\text{KL}}(t) = (1+t)\log(1+t) - t$ , the corresponding residual adjustment factor is  $A_{\text{KL}}(\delta) = \delta$ , which is unbounded in the residual  $\delta = g/f_{\theta} - 1$ . As a consequence, large density ratios  $g/f_{\theta}$  can exert arbitrarily large leverage on the estimating equations: the gradient  $\nabla_{\theta} D_G(g, f_{\theta})$  is not uniformly bounded over contaminated distributions, and individual data points can dominate the fit. Common EM heuristics such as *weight-flooring* (enforcing  $\pi_k \geq \pi_{\text{floor}}$  and renormalizing) only constrain the mixing weights and do not cap these per-observation contributions: a single extreme observation can still drive the fitted component parameters and, in model-selection settings, spuriously favor an extra component. In particular, the influence function remains unbounded and the breakdown point is essentially zero. In Sections 4.2–4.3 we make this precise, showing that our general convergence and finite-step CLT results still hold for KL, but the uniform contraction and breakdown lower bounds available for bounded-RAF divergences (such as NED/vNED) do not extend to KL without additional local density-ratio assumptions.

We now turn to a heuristic description of the DM method in terms of divergences, borrowing terminology from the EM literature. As above, let  $Y$  be an observable real-valued random variable and  $Z$  is a latent random variable representing the class membership of  $Y$ . The pair  $(Y, Z)$  is referred to as the complete data. Let  $\mathcal{G}$  denote the densities on  $\mathbb{R}$ , and  $g(\cdot)$  be the true density of  $Y$  and  $f(\cdot; \theta) \in \mathcal{F}_{\theta}$  denote the postulated density. The postulated joint density of  $(Y, Z)$  is denoted by  $p(\cdot, \cdot; \theta)$  and the marginal postulated density is  $f(y; \theta) = \int_{\mathcal{Z}} p(y, z; \theta) d\nu(z)$ . The marginal density is referred to as the *incomplete data*

model. Let  $w(z|y; \boldsymbol{\theta}) := \frac{p(y, z; \boldsymbol{\theta})}{f(y; \boldsymbol{\theta})} = \mathbb{P}_{\boldsymbol{\theta}}(Z = z | Y = y)$  be the conditional density of  $Z$  given  $Y$ , it follows that  $p(y, z; \boldsymbol{\theta}) = f(y; \boldsymbol{\theta})w(z|y; \boldsymbol{\theta})$ . This is referred to as the *responsibility*. Finally, we frequently use the notations  $w(z|y, \boldsymbol{\theta})$  and  $w_z(y; \boldsymbol{\theta})$  interchangeably.

The cross-entropy loss for the joint distribution of  $(Y, Z)$ , referred to as the *EM Q-functional*, is  $\mathbf{E}[-\log p(Y, Z; \boldsymbol{\theta})]$ , and can be expressed as  $\mathbf{E}_Y[\mathbf{E}_{Z|Y; \boldsymbol{\theta}}[-\log p(Y, Z; \boldsymbol{\theta}')] ]$ . The conditional expectation, conditioned on  $Y = y$ , is referred to as the EM surrogate. To obtain a robust variant of the EM surrogate, we replace the cross-entropy loss with a residual adjusted divergence on complete data. Specifically, we introduce the complete-data residual (parameterized by an auxiliary conditional  $\tilde{q}$ )

$$\delta(y, z; \boldsymbol{\theta}', \tilde{q}) = \frac{g(y) \tilde{q}(z | y)}{p(y, z; \boldsymbol{\theta}')} - 1, \quad Q_G(\boldsymbol{\theta}', \tilde{q} | \boldsymbol{\theta}) := \iint G(\delta(y, z; \boldsymbol{\theta}', \tilde{q})) p(y, z; \boldsymbol{\theta}') d\nu(z) dy.$$

We begin with a simple lemma describing the variational elimination of the auxiliary conditional from the complete-data residual.

**Lemma 1** (Variational elimination of the auxiliary conditional). *For any  $\boldsymbol{\theta}'$ ,*

$$\min_{\tilde{q}} \widetilde{Q}_G(\boldsymbol{\theta}', \tilde{q} | \boldsymbol{\theta}) = D_G(g, f_{\boldsymbol{\theta}'} ) \quad \text{attained at} \quad \tilde{q}(\cdot | y) = w(\cdot | y; \boldsymbol{\theta}') := \frac{p(y, \cdot; \boldsymbol{\theta}')}{f(y; \boldsymbol{\theta}')}.$$

*Proof.* Let  $\boldsymbol{\theta}'$  be fixed and set  $w(\cdot | y; \boldsymbol{\theta}')$  as above. Then, by Jensen's inequality, for each  $y$ ,

$$\int G\left(\frac{g(y)}{f(y; \boldsymbol{\theta}')} \cdot \frac{\tilde{q}(z | y)}{w(z | y; \boldsymbol{\theta}')} - 1\right) w(z | y; \boldsymbol{\theta}') d\nu(z) \geq G\left(\frac{g(y)}{f(y; \boldsymbol{\theta}')} - 1\right),$$

and  $\int \frac{\tilde{q}}{w} w d\nu = \int \tilde{q} d\nu = 1$ . Next, integrating over  $y$ , we obtain  $\min_{\tilde{q}} \widetilde{Q}_G(\boldsymbol{\theta}', \tilde{q} | \boldsymbol{\theta}) = D_G(g, f_{\boldsymbol{\theta}'} )$ , attained at  $\tilde{q} = w(\cdot | y; \boldsymbol{\theta}')$ .  $\square$

We notice that fixing the auxiliary conditional at the current responsibilities yields a valid majorizer. Hence, we set

$$Q_G(\boldsymbol{\theta}' | \boldsymbol{\theta}) := Q_G(\boldsymbol{\theta}', w(\cdot | \cdot; \boldsymbol{\theta}) | \boldsymbol{\theta}) \geq D_G(g, f_{\boldsymbol{\theta}'}) \quad \text{with equality at } \boldsymbol{\theta}' = \boldsymbol{\theta}.$$

Thus, if  $\theta_{m+1} \in \arg \min_{\theta'} Q_G(\theta' | \theta_m)$  then  $D_G(g, f_{\theta_{m+1}}) \leq D_G(g, f_{\theta_m})$ . In other words,  $Q_G(\theta' | \theta)$  inherits the properties of EM surrogate. Next, setting  $t(z | y) = \frac{w(z|y;\theta)}{w(z|y;\theta')}$  and recalling  $\delta(y; \theta') + 1 = \frac{g(y)}{f(y;\theta')}$ , observe that

$$Q_G(\theta' | \theta) = \int_Y f(y; \theta') \left\{ \int_Z G((1 + \delta(y; \theta')) t(z | y) - 1) w(z | y; \theta') d\nu(z) \right\} dy. \quad (2.1)$$

By subtracting  $G(-1 + \frac{g(y)}{f(y;\theta')})$  from the inner integral, it follows that the difference  $Q_G(\theta' | \theta) - D_G(g, \theta')$  represents average latent-conditional gap. Applying Jensen's inequality w.r.t.  $w(\cdot | y; \theta')$  (noting  $\int t(\cdot | y) w(\cdot | y; \theta') d\nu = 1$ ) it follows that  $H_G(\theta' | \theta) \geq 0$ , and  $H_G(\theta | \theta) = 0$ . This yields the separation-majorization identity, namely,  $Q_G(\theta' | \theta) = \Psi_G(\theta') + H_G(\theta' | \theta) \geq \Psi_G(\theta') = D_G(g_n, f_{\theta'})$  with equality iff  $w(\cdot | y; \theta) = w(\cdot | y; \theta')$  for  $f(\cdot; \theta')$ -a.e.  $y$ .

**The DM Algorithm:** The expression (2.1) can alternatively be expressed as follows:

$$\begin{aligned} Q_G(\theta' | \theta) &= \mathbf{E}_Y \left[ \mathbf{E}_{Z|Y, \theta'} \left[ G \left( -1 + \frac{g(Y)w(Z|Y; \theta)}{f(Y; \theta')w(Z|Y; \theta')} \right) \right] \right] \\ &= \mathbf{E}_Y \left[ \mathbf{E}_{Z|Y, \theta'} \left[ G \left( -1 + \frac{g(Y)w(Z|Y; \theta)}{p(Y, Z; \theta')} \right) \right] \right]. \end{aligned} \quad (2.2)$$

This suggests a natural approach to obtain the MDE by iterating the RHS of (2.2). We will show, in the following lemma, that if  $\theta'$  is MDE, then  $Q^{(G)}(\theta' | \theta') = D^{(G)}(\theta') \leq Q^{(G)}(\theta' | \theta)$  for  $\theta \in \Theta$ . Using Supplement Lemma 9.4 (variational elimination), choosing  $\tilde{q}(\cdot | y) = w(\cdot | y; \theta)$  yields the following majorization.

**Lemma 2.** *For any  $\theta, \theta' \in \Theta$ ,  $D_G(\theta') \leq Q_G(\theta' | \theta)$  for all  $\theta \in \Theta$  and the equality holds if and only if  $\theta' = \theta$ . Hence, for any  $\theta$ , the map  $Q_G(\cdot | \theta)$  majorizes  $D^{(G)}(\cdot)$ .*

Based on Lemma 2, it follows that  $Q_G(\theta' | \theta)$  is a majorizing function for the MM algorithm (see Hunter & Lange (2000b)). Also, note that the divergence in  $Q_G(\theta' | \theta)$  includes additional information regarding the unobserved conditional density  $w(z|y; \theta)$ , which enables the calculation of the posterior distribution, as in the EM algorithm. This is based on knowing

the true density  $g(\cdot)$  and hence  $D_G(\boldsymbol{\theta})$  and  $Q_G(\boldsymbol{\theta}'|\boldsymbol{\theta})$  will be referred to as the *population level* DM objective functions.

When  $g(\cdot)$  is unknown it can be replaced by  $g_n(\cdot)$  (estimated from observed data), leading to the *sample-level* objective functions. One choice of  $g_n(\cdot)$  is the kernel density estimator

$$g_n(y) = \frac{1}{nc_n} \sum_{i=1}^n \mathcal{K}\left(\frac{y - Y_i}{c_n}\right), \quad (2.3)$$

where  $\mathcal{K}(\cdot)$  is a probability density and  $c_n$  is the bandwidth. Other choices of  $g_n(\cdot)$  include wavelet density estimator (see [Doukhan & León \(1990\)](#)) and local polynomial regression estimator (see [Fan & Gijbels \(1996\)](#)). We now introduce notations to describe the sample-level objective function and surrogate. Let  $\delta_n(y; \boldsymbol{\theta}) = g_n(y)/f(y; \boldsymbol{\theta}) - 1$  denote the sample Pearson's residual and

$$D_G(g_n, f_{\boldsymbol{\theta}}) = \int G(\delta_n(y; \boldsymbol{\theta})) f(y; \boldsymbol{\theta}) dy \quad \text{and} \quad A(\delta_n) = (1 + \delta_n)G'(\delta_n) - G(\delta_n).$$

denote the sample-level objective function and the RAF. Given responsibilities  $w_k(y; \boldsymbol{\theta})$ , we set

$$\delta_{k,n}(y; \boldsymbol{\theta}, \boldsymbol{\theta}') = \frac{g_n(y) w_k(y; \boldsymbol{\theta})}{\pi'_k h(y; \phi'_k)} - 1, \quad \text{and} \quad Q_G(\boldsymbol{\theta}' | \boldsymbol{\theta}; g_n) = \sum_{k=1}^K \int G(\delta_{k,n}(y; \boldsymbol{\theta}, \boldsymbol{\theta}')) \pi'_k h(y; \phi'_k) dy.$$

Thus, the DM algorithm can be broadly divided into two steps. Below, we fix  $\mu \in \{g_n, g\}$ . Then the steps of the DM algorithm are:

1. **D-step.** Determine  $Q_G(\boldsymbol{\theta}'|\boldsymbol{\theta}; \mu)$ .
2. **M-step.** Choose  $\boldsymbol{\theta}_{m+1} \in \boldsymbol{\Theta}$  to minimize  $Q_G(\boldsymbol{\theta}'|\boldsymbol{\theta}_m; \mu)$  over  $\boldsymbol{\theta}' \in \boldsymbol{\Theta}$ .

As a *standing convention* for the rest of the paper, for every fixed  $n$ , we denote by  $D_{G,n}(\boldsymbol{\theta}) := D_G(g_n, f_{\boldsymbol{\theta}})$  and  $Q_{G,n}(\boldsymbol{\theta}' | \boldsymbol{\theta}) := Q_G(\boldsymbol{\theta}' | \boldsymbol{\theta}; g_n)$ . Also, for a fixed  $G$ , we suppress the dependence on  $G$  and simply refer to it as  $D_n(\boldsymbol{\theta})$  and  $Q_n(\boldsymbol{\theta}' | \boldsymbol{\theta})$ .

**The generalized DM Algorithm:** It is possible that in the M-step, the minimizer  $\boldsymbol{\theta}_{m+1}$  is not unique. Let  $\boldsymbol{\Theta}_m$  denote the set of all minimizers at step  $m$ ; that is,  $\boldsymbol{\theta}_{m+1} \in \boldsymbol{\Theta}_{m+1}$ .

Sometimes it may be difficult to perform M-step numerically; in this case, we can define a generalized DM algorithm (referred to as the G-DM algorithm) as follows: Let  $M : \boldsymbol{\theta}_m \rightarrow \boldsymbol{\Theta}_{m+1}$  be a point to set map: then the G-DM algorithm is an iterative scheme such that

$$Q^{(G)}(\boldsymbol{\theta}'|\boldsymbol{\theta}_m) \leq Q^{(G)}(\boldsymbol{\theta}_m|\boldsymbol{\theta}_m) \quad \text{for all } \boldsymbol{\theta}' \in M(\boldsymbol{\theta}_m).$$

We notice here that for any G-DM sequence  $\{\boldsymbol{\theta}_m\}$ ,  $D_G(\boldsymbol{\theta}_{m+1}) \leq D_G(\boldsymbol{\theta}_m)$  and DM algorithm is a special case of G-DM algorithm. We emphasize here that by choosing different divergences, we obtain many existing algorithms, including the EM, HMIX, and HELMIX algorithms. These examples are given in the Supplementary analysis [9.4](#).

The key ingredient for the proposed population DM algorithm is the conditional Pearson-type ratio between  $g(y)w(z|y; \boldsymbol{\theta})$  and  $f(y; \boldsymbol{\theta}')w(z|y; \boldsymbol{\theta}')$ , which we denote by

$$\tau(y, z; g, \boldsymbol{\theta}, \boldsymbol{\theta}') = \frac{g(y)w(z|y; \boldsymbol{\theta})}{f(y; \boldsymbol{\theta}')w(z|y; \boldsymbol{\theta}')}.$$

When  $z = k$ , the denominator of  $\tau(y, k; g, \boldsymbol{\theta}, \boldsymbol{\theta}')$  reduces to  $\pi'_k h(y; \boldsymbol{\phi}'_k)$ . We observe here that, unlike the Pearson residual,  $\tau$  does not subtract 1; it plays the role of a conditional analogue instead. When there is no scope for confusion, for notational simplicity, we suppress  $g$ ,  $\boldsymbol{\theta}$  and  $\boldsymbol{\theta}'$  in  $\tau_k(y, g, \boldsymbol{\theta}, \boldsymbol{\theta}')$  and denote it as  $\tau_k(y)$  and when  $g$  is replaced by  $g_n$  we refer to it as  $\tau_{k,n}(y)$ .

## 2.1 The DM-MIX Algorithm for FMM

A general description of the DM algorithm for HFMM (i.e., the DM-MIX algorithm) is given below. As noted above,  $\tau_k(y) := \tau_k(y, g; \boldsymbol{\theta}, \boldsymbol{\theta}') = \frac{g(y)w_k(y; \boldsymbol{\theta})}{\pi'_k h(y; \boldsymbol{\theta}')}$ . Also, set

$$Q^{(k)}(\boldsymbol{\theta}'|\boldsymbol{\theta}) = \int_{\mathbb{R}} G(\tau_k(y)) \pi'_k h(y; \boldsymbol{\phi}'_k) dy.$$

When  $g$  is replaced by  $g_n$ , we denote the above quantities by  $\tau_{k,n}(\cdot, \cdot; \cdot)$  and  $Q^{(k)}(\cdot|\cdot)$  by  $Q_n^{(k)}(\cdot|\cdot)$ . Finally, define  $B(u) = G(u) - uG'(u)$ . We now turn to the algorithm. For a given starting point  $\boldsymbol{\theta}^{(0)}$  and a density estimate  $g_n(y)$ ,

---

**Algorithm 1** The DM-Mix Algorithm

---

Set  $m = 0$ .

**repeat**

1. Compute

$$Q_n(\boldsymbol{\theta}^{(m)}|\boldsymbol{\theta}^{(m)}) = \sum_{k=1}^K \int_{\mathcal{Y}} G\left(-1 + \frac{g_n(y)w_k(y;\boldsymbol{\theta}^{(m)})}{\pi_k^{(m)}h(y;\boldsymbol{\phi}_k^{(m)})}\right) \pi_k^{(m)} h(y|\boldsymbol{\phi}_k^{(m)}) dy, \quad \text{where } w_k(y;\boldsymbol{\theta}) = \frac{\pi_k h(y;\boldsymbol{\phi}_k)}{\sum_{l=1}^K \pi_l h(y;\boldsymbol{\phi}_l)}.$$

2. Update  $\boldsymbol{\phi}_k^{(m+1)}$ :

$$\boldsymbol{\phi}_k^{(m+1)} = \underset{\boldsymbol{\phi}'_k \in \Theta}{\operatorname{argmin}} \int_{\mathcal{Y}} G\left(-1 + \frac{g_n(y)w_k(y;\boldsymbol{\theta}^{(m)})}{\pi_k^{(m)}h(y;\boldsymbol{\phi}'_k)}\right) \pi_k^{(m)} h(y|\boldsymbol{\phi}'_k) dy.$$

3. Update the mixing probabilities using the Lagrangian multiplier step:

$$\pi_k^{(m+1)} = \frac{\Phi_k(\boldsymbol{\pi}^{(m)}, \boldsymbol{\phi}^{(m+1)})}{\sum_{\ell=1}^K \Phi_{\ell}(\boldsymbol{\pi}^{(m)}, \boldsymbol{\phi}^{(m+1)})}, \quad k = 1, \dots, K,$$

where

$$\begin{aligned} \Phi_k(\boldsymbol{\pi}, \boldsymbol{\phi}) &:= -\pi_k \frac{\partial Q_n^{(k)}(\pi_k, \boldsymbol{\phi}_k | \boldsymbol{\theta})}{\partial \pi_k}, \\ \frac{\partial Q_n^{(k)}}{\partial \pi_k}(\pi_k, \boldsymbol{\phi}_k | \boldsymbol{\theta}) &= \int_{\mathcal{Y}} h(y; \boldsymbol{\phi}_k) \left\{ B(\tau_{k,n}(y)) \right\} dy, \\ \Phi_k(\boldsymbol{\pi}, \boldsymbol{\phi}) &:= \pi_k \int_{\mathcal{Y}} h(y; \boldsymbol{\phi}_k) \left\{ B(\tau_{k,n}(y)) \right\} dy. \end{aligned}$$

4. Update  $w_k(y; \boldsymbol{\theta}^{(m+1)})$ , and compute  $Q_n(\boldsymbol{\theta}^{(m+1)}|\boldsymbol{\theta}^{(m+1)})$  and the difference  $\epsilon_{m+1} = |Q_n(\boldsymbol{\theta}^{(m+1)}|\boldsymbol{\theta}^{(m+1)}) - Q_n(\boldsymbol{\theta}^{(m)}|\boldsymbol{\theta}^{(m)})|$ .

5. Set  $m = m+1$ .

**until**  $\epsilon_m < \text{threshold}$ .

---

The EM, HMIX, and HELMIX algorithms are special cases of this framework for particular choices of  $G(\cdot)$  and probability distributions (see Supplementary analysis 9.4), wherein we also illustrate that by reformulating the DM objective function algorithm can be interpreted as an MM algorithm, a proximal point algorithm, and a coordinate descent algorithm.

The DM algorithm applied to HFMM has a special structure for updating the mixing probability  $\boldsymbol{\pi}$ . To see this, we partition the parameter vector  $\boldsymbol{\theta}' = (\boldsymbol{\pi}', \boldsymbol{\phi}')$  into two sub-vectors  $\boldsymbol{\pi}' = (\pi'_1, \dots, \pi'_K)$  and  $\boldsymbol{\phi}' = (\boldsymbol{\phi}'_1, \dots, \boldsymbol{\phi}'_K)$ , and we can update  $\boldsymbol{\phi}'$  and  $\boldsymbol{\pi}'$  iteratively.



Specifically, starting with updating parameters  $\phi'_k$ , we update  $\pi'_k$  using the Lagrangian multiplier and then iterate the steps. This process leads to the estimate of  $\pi_k$  as

$$\pi'_k = \frac{\pi'_k \frac{\partial Q_n^{(k)}(\boldsymbol{\theta}' | \boldsymbol{\theta})}{\partial \pi'_k}}{\sum_{\ell=1}^K \pi'_\ell \frac{\partial Q_n^{(\ell)}(\boldsymbol{\theta}' | \boldsymbol{\theta})}{\partial \pi'_\ell}}. \quad (2.4)$$

Details of the derivation are provided in the Supplementary analysis 9.4. The update of the mixing probability in (2.4) naturally arises from leveraging the majorizing divergence function  $Q(\cdot|\cdot)$ .

### 3 Convergence Guarantees of the DM Algorithm

Figure 1 summarizes the DM map and the descent property of  $D_G$ .

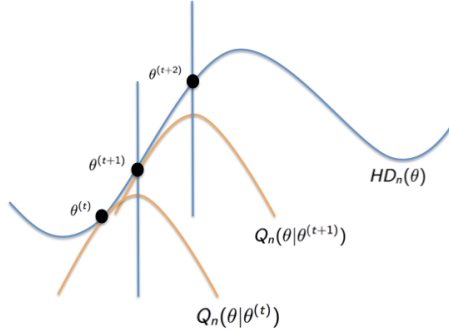


Figure 1: The DM Algorithm Illustration

Let  $M_n(\boldsymbol{\theta}) : \boldsymbol{\Theta} \rightrightarrows \boldsymbol{\Theta}$  and  $M(\boldsymbol{\theta}) : \boldsymbol{\Theta} \rightrightarrows \boldsymbol{\Theta}$  be two set-valued maps defined as follows:

$$M_n(\boldsymbol{\theta}) := \arg \min_{\boldsymbol{\theta}' \in \boldsymbol{\Theta}} Q_n(\boldsymbol{\theta}' | \boldsymbol{\theta}) \quad \text{and} \quad M(\boldsymbol{\theta}) := \arg \min_{\boldsymbol{\theta}' \in \boldsymbol{\Theta}} Q(\boldsymbol{\theta}' | \boldsymbol{\theta}).$$

**Definition 1** (DM sequences). *A sequence  $\{\boldsymbol{\theta}_m\}$  is a population DM sequence if  $\boldsymbol{\theta}_{m+1} \in M(\boldsymbol{\theta}_m)$ ; a sequence  $\{\boldsymbol{\theta}_{m,n}\}$  is a sample DM sequence if  $\boldsymbol{\theta}_{m+1,n} \in M_n(\boldsymbol{\theta}_{m,n})$ .*

In this section, we establish the monotone descent property and stationarity of the limit

points for the sample-level iterates. Using these, we establish local contraction of the population map  $M(\cdot)$  together with the finite-sample perturbation bounds.

**DM surrogate decomposition (recall).** With  $G$  fixed, let  $D_n(\boldsymbol{\theta}) := D_G(g_n, f_{\boldsymbol{\theta}})$  and  $Q_n(\boldsymbol{\theta}' | \boldsymbol{\theta}) := Q_G(\boldsymbol{\theta}' | \boldsymbol{\theta}; g_n)$ . Define the remainder

$$H_G(\boldsymbol{\theta}' | \boldsymbol{\theta}) := Q_n(\boldsymbol{\theta}' | \boldsymbol{\theta}) - D_n(\boldsymbol{\theta}).$$

Then  $Q_n(\boldsymbol{\theta}' | \boldsymbol{\theta}) = D_n(\boldsymbol{\theta}') + H_G(\boldsymbol{\theta}' | \boldsymbol{\theta})$  with  $H_G(\boldsymbol{\theta}' | \boldsymbol{\theta}) \geq 0$  and  $H_G(\boldsymbol{\theta} | \boldsymbol{\theta}) = 0$ . (See Supplementary material Section D for a derivation.)

**Lemma 3** (Monotone descent for DM updates). *Let  $\boldsymbol{\theta}^+ \in M_n(\boldsymbol{\theta})$  be a DM update (set-valued selection). Then*

$$D_n(\boldsymbol{\theta}^+) \leq Q_n(\boldsymbol{\theta}^+ | \boldsymbol{\theta}) \leq Q_n(\boldsymbol{\theta} | \boldsymbol{\theta}) = D_n(\boldsymbol{\theta}),$$

hence  $D_n$  is nonincreasing along DM iterates.

*Proof.* By the decomposition above,  $Q_n(\boldsymbol{\theta}^+ | \boldsymbol{\theta}) \geq D_n(\boldsymbol{\theta}^+)$  and  $Q_n(\boldsymbol{\theta} | \boldsymbol{\theta}) = D_n(\boldsymbol{\theta})$ . Since  $\boldsymbol{\theta}^+ \in \arg \min_{\boldsymbol{\theta}'} Q_n(\boldsymbol{\theta}' | \boldsymbol{\theta})$ ,  $Q_n(\boldsymbol{\theta}^+ | \boldsymbol{\theta}) \leq Q_n(\boldsymbol{\theta} | \boldsymbol{\theta})$ .  $\square$

*Stationarity* will be established below via our self-consistency and first-order conditions.

### 3.1 Global Properties of the DM Sequences

Our first proposition concerns the convergence properties of the sample-level DM sequence.

**Proposition 1.** *Assume (B1)-(B3) and (D1) holds and let  $\{\boldsymbol{\theta}_{m,n}\}$  be any sample-level DM sequence. Then  $D_n(\boldsymbol{\theta}_{m+1,n}) \leq D_n(\boldsymbol{\theta}_{m,n})$ , with equality if and only if  $\boldsymbol{\theta}_{m,n}$  is a fixed point of  $M_n$ . Moreover, every limit point  $\bar{\boldsymbol{\theta}}$  of  $\{\boldsymbol{\theta}_{m,n}\}$  is stationary for  $D_n$ , i.e.  $\nabla D_n(\bar{\boldsymbol{\theta}}) = 0$ .*

*Proof (majorization-minimization).* By construction  $D_n(\boldsymbol{\theta}) \leq Q_n(\boldsymbol{\theta} | \boldsymbol{\theta}')$  for all  $\boldsymbol{\theta}, \boldsymbol{\theta}'$ , with equality at  $\boldsymbol{\theta} = \boldsymbol{\theta}'$ . Hence, noticing that  $\boldsymbol{\theta}_{m+1,n}$  is a minimizer, it follows that  $D_n(\boldsymbol{\theta}_{m+1,n}) \leq Q_n(\boldsymbol{\theta}_{m+1,n} | \boldsymbol{\theta}_{m,n}) \leq Q_n(\boldsymbol{\theta}_{m,n} | \boldsymbol{\theta}_{m,n}) = D_n(\boldsymbol{\theta}_{m,n})$ . Next, if equality holds, then  $\boldsymbol{\theta}_{m,n}$

minimizes  $Q_n(\cdot \mid \boldsymbol{\theta}_{m,n})$  and is thus a fixed point. Now let  $\bar{\boldsymbol{\theta}}$  be any limit point of  $\{\boldsymbol{\theta}_{m,n}\}$ . By the closed-graph property (D1), any such limit point satisfies  $\bar{\boldsymbol{\theta}} \in M_n(\bar{\boldsymbol{\theta}})$ , so  $\bar{\boldsymbol{\theta}}$  minimizes  $Q_n(\cdot \mid \bar{\boldsymbol{\theta}})$ . On the other hand, for every  $\boldsymbol{\theta}$ , using the tangency identity  $Q_n(\boldsymbol{\theta} \mid \boldsymbol{\theta}) = D_n(\boldsymbol{\theta})$ , and differentiating this with respect to the first argument at  $\boldsymbol{\theta} = \bar{\boldsymbol{\theta}}$  gives

$$\nabla_{\boldsymbol{\theta}'} Q_n(\boldsymbol{\theta}' \mid \bar{\boldsymbol{\theta}}) \Big|_{\boldsymbol{\theta}' = \bar{\boldsymbol{\theta}}} = \nabla D_n(\bar{\boldsymbol{\theta}}).$$

Since  $\bar{\boldsymbol{\theta}}$  minimizes  $Q_n(\cdot \mid \bar{\boldsymbol{\theta}})$ , the left-hand side is zero, and hence  $\nabla D_n(\bar{\boldsymbol{\theta}}) = 0$ .  $\square$

**Remark 1** (When “stationary” coincides with “fixed point”). *If, in addition, for each  $\boldsymbol{\theta}$  the function  $\boldsymbol{\theta}' \mapsto Q_n(\boldsymbol{\theta}' \mid \boldsymbol{\theta})$  is strongly convex near  $\boldsymbol{\theta}$  and tangent to  $D_n$  at  $\boldsymbol{\theta}$  in the sense that  $\nabla_{\boldsymbol{\theta}'} Q_n(\boldsymbol{\theta}' \mid \boldsymbol{\theta})|_{\boldsymbol{\theta}' = \boldsymbol{\theta}} = \nabla D_n(\boldsymbol{\theta})$ , then  $\bar{\boldsymbol{\theta}}$  is stationary for  $D_n$  if and only if it minimizes  $Q_n(\cdot \mid \bar{\boldsymbol{\theta}})$ . Under the uniqueness of that minimizer, this is equivalent to  $\bar{\boldsymbol{\theta}}$  being a fixed point of the DM update map.*

The next proposition concerns the self-consistency of the DM updates, namely, any minimizer of the observed-data objective also minimizes the complete-data surrogate at the same parameter value. That is, if  $\boldsymbol{\theta}_n^*$  minimizes  $D_n(\boldsymbol{\theta}')$  over  $\boldsymbol{\theta}' \in \boldsymbol{\Theta}$ , then  $\boldsymbol{\theta}_n^*$  also minimizes  $Q_n(\cdot \mid \boldsymbol{\theta}_n^*)$ .

**Proposition 2.** *Let  $\boldsymbol{\theta}_n^* \in \arg \min_{\boldsymbol{\theta}' \in \boldsymbol{\Theta}} D_n(\boldsymbol{\theta}')$  and assume that  $Q_n(\cdot \mid \boldsymbol{\theta})$  attains a minimum for each fixed  $\boldsymbol{\theta}$ , then  $\boldsymbol{\theta}_n^* \in \arg \min_{\boldsymbol{\theta}' \in \boldsymbol{\Theta}} Q_n(\boldsymbol{\theta}' \mid \boldsymbol{\theta}_n^*)$ .*

*Proof.* For any  $\boldsymbol{\theta}'$ ,  $Q_n(\boldsymbol{\theta}' \mid \boldsymbol{\theta}_n^*) \geq D_n(\boldsymbol{\theta}') \geq D_n(\boldsymbol{\theta}_n^*) = Q_n(\boldsymbol{\theta}_n^* \mid \boldsymbol{\theta}_n^*)$ , so  $Q_n(\cdot \mid \boldsymbol{\theta}_n^*)$  is minimized at  $\boldsymbol{\theta}_n^*$ .  $\square$

*An immediate consequence of the above proposition is the following corollary.*

**Corollary 1** (Fixed point under uniqueness). *If, in addition,  $Q_n(\cdot \mid \boldsymbol{\theta}_n^*)$  has a unique minimizer, then the update map satisfies  $M_n(\boldsymbol{\theta}_n^*) = \boldsymbol{\theta}_n^*$ .*

We now turn to the cyclical behavior of the iterates of the algorithm.

## Cyclical behavior

For fixed  $n$ , a finite set of distinct points  $\{\boldsymbol{\theta}_{1,n}^*, \dots, \boldsymbol{\theta}_{t,n}^*\} \subset \Theta$  is a *cycle of length*  $t \geq 2$  for  $M_n$  if  $M_n(\boldsymbol{\theta}_{i,n}^*) = \boldsymbol{\theta}_{i+1,n}^*$  for  $i = 1, \dots, t-1$  and  $M_n(\boldsymbol{\theta}_{t,n}^*) = \boldsymbol{\theta}_{1,n}^*$ . We use assumptions (B1)-(B4) and the isolated-stationary-points condition (B4) throughout this subsection.

**Lemma 4** (Finiteness at a level). *Assume (B1)-(B4). Fix any value  $D_n^* \in \mathbb{R}$ . Then the set of stationary points of  $D_n$  with objective value  $D_n^*$  is at most finite.*

*Proof sketch.* If infinitely many stationary points shared the same value, compactness (B2) would yield a convergent subsequence to another stationary point, contradicting isolation (B4). More detailed proof is in the appendix.  $\square$

**Proposition 3** (Limit-set structure: convergence or finite cycle). *Let  $\{\boldsymbol{\theta}_{m,n}\}$  be any sample-level DM sequence with  $\boldsymbol{\theta}_{m+1,n} \in M_n(\boldsymbol{\theta}_{m,n})$ , and assume (B1)-(B4), and (D1)'. Then  $D_n(\boldsymbol{\theta}_{m,n}) \downarrow D_n^*$ , and every limit point of  $\{\boldsymbol{\theta}_{m,n}\}$  is a stationary point of  $D_n$  with value  $D_n^*$ . Moreover, the update  $M_n$  permutes the (finite) set of limit points; hence either  $\boldsymbol{\theta}_{m,n} \rightarrow \boldsymbol{\theta}_n^*$  (a stationary point with  $M_n(\boldsymbol{\theta}_n^*) = \boldsymbol{\theta}_n^*$ ), or the limit points form a finite cycle.*

*Proof sketch.* Monotone majorization gives  $D_n(\boldsymbol{\theta}_{m+1,n}) \leq D_n(\boldsymbol{\theta}_{m,n})$  with a finite limit  $D_n^*$ ; any limit point is stationary by Prop. 1. By Lemma 4, the limit set is finite. Continuity (D1) implies  $M_n$  maps the set to itself, so  $M_n$  acts as a permutation; this yields either a singleton (convergence) or a finite cycle.  $\square$

**Proposition 4** (No cycles under uniqueness). *Assume (B1)-(B4), and (D2) (for every  $\boldsymbol{\theta}$ ,  $Q_n(\cdot | \boldsymbol{\theta})$  has a unique minimizer). Then any sample-level DM sequence  $\{\boldsymbol{\theta}_{m,n}\}$  converges to a stationary point  $\boldsymbol{\theta}_n^*$  with  $M_n(\boldsymbol{\theta}_n^*) = \boldsymbol{\theta}_n^*$ . Moreover, if  $\boldsymbol{\theta}_{m,n} \neq \boldsymbol{\theta}_n^*$  for all  $m$ , then  $D_n(\boldsymbol{\theta}_{m+1,n}) < D_n(\boldsymbol{\theta}_{m,n})$  and  $D_n(\boldsymbol{\theta}_{m,n}) \searrow D_n(\boldsymbol{\theta}_n^*)$ .*

*Proof sketch.* Uniqueness excludes nontrivial cycles and enforces single-valuedness/continuity of  $M_n$  near stationary points; strict inequality follows from majorization unless at a

minimizer of  $Q_n(\cdot \mid \boldsymbol{\theta}_{m,n})$ , which coincides with stationarity.  $\square$

## 3.2 Contraction Properties

We analyze the DM operator at both the population and sample levels. Throughout this subsection, assume the true model holds,  $g(y) = f(y; \boldsymbol{\theta}^*)$ , and define  $q(\cdot) := Q(\cdot \mid \boldsymbol{\theta}^*)$ . We will use the strong convexity of  $q$  near  $\boldsymbol{\theta}^*$  and a first-order stability (FOS) condition that controls the sensitivity of the gradient map  $\nabla Q(\cdot \mid \boldsymbol{\theta})$  with respect to its second argument.

### Geometric Convergence of the DM Sequences

In this section, we focus on the guarantees for the population-level DM algorithm. For  $q : \Theta \rightarrow \mathbb{R}$ ,  $q$  is  $\lambda$ -strongly convex on  $B_2(r; \boldsymbol{\theta}^*)$  if  $q(\boldsymbol{\theta}_1) - q(\boldsymbol{\theta}_2) - \langle \nabla q(\boldsymbol{\theta}_2), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 \rangle \geq \frac{\lambda}{2} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2^2$ . By self-consistency,  $\boldsymbol{\theta}^* = M(\boldsymbol{\theta}^*)$ , and the first-order optimality conditions are, for all  $\boldsymbol{\theta}$ ,  $\langle \nabla_{\boldsymbol{\theta}'} Q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^*), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle \geq 0$  and  $\langle \nabla_{\boldsymbol{\theta}'} Q(M(\boldsymbol{\theta}) \mid \boldsymbol{\theta}), \boldsymbol{\theta} - M(\boldsymbol{\theta}) \rangle \geq 0$ .

**Definition 2.** (*First-order Stability (FOS)*) We say that  $\{Q(\cdot \mid \boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta\}$  satisfies FOS( $\gamma$ ) on  $B_2(r'; \boldsymbol{\theta}^*)$  if

$$\|\nabla Q(M(\boldsymbol{\theta}) \mid \boldsymbol{\theta}^*) - \nabla Q(M(\boldsymbol{\theta}) \mid \boldsymbol{\theta})\|_2 \leq \gamma \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2 \quad \text{for all } \boldsymbol{\theta} \in \mathbb{B}_2(r'; \boldsymbol{\theta}^*). \quad (3.5)$$

**Theorem 1.** For some radius  $r' > 0$  and pair  $(\gamma, \lambda)$  such that  $0 < \gamma < \lambda$ , suppose that the function  $Q(\cdot \mid \boldsymbol{\theta}^*)$  is  $\lambda$ -strongly convex and that the FOS( $\gamma$ ) condition (3.5) holds on the ball  $\mathbb{B}_2(r'; \boldsymbol{\theta}^*)$ . Then the population DM operator  $M$  is contractive over  $\mathbb{B}_2(r'; \boldsymbol{\theta}^*)$ ; in particular, the following inequality holds:

$$\|M(\boldsymbol{\theta}) - \boldsymbol{\theta}^*\|_2 \leq \frac{\gamma}{\lambda} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2 \quad \text{for all } \boldsymbol{\theta} \in \mathbb{B}_2(r'; \boldsymbol{\theta}^*).$$

As an immediate consequence, under the conditions of the theorem, for any initial point  $\boldsymbol{\theta}_0 \in \mathbb{B}_2(r'; \boldsymbol{\theta}^*)$ , the population DM sequence  $\{\boldsymbol{\theta}_m\}$  exhibits geometric convergence; that is,

$$\|\boldsymbol{\theta}_m - \boldsymbol{\theta}^*\|_2 \leq \left(\frac{\gamma}{\lambda}\right)^m \|\boldsymbol{\theta}^0 - \boldsymbol{\theta}^*\|_2 \quad \text{for all } m = 1, 2, \dots$$

The  $\text{FOS}(\gamma)$  condition is standard in contraction analyses of EM-type algorithms (cf. [Balakrishnan et al. \(2017\)](#)) and, in our mixture setting, follows from Lipschitz regularity of the responsibilities and component scores; see the Component Lipschitz regularity assumption and Corollary 2 for an explicit bound on  $\gamma_K$ . We emphasize that  $\text{FOS}(\gamma)$  is only required locally in a neighbourhood of  $\theta^*$ , and is not needed for the global monotone descent or stationarity results in Section 3.1; it is used solely to obtain geometric rates for the population map and its noisy sample analogue. We now turn to theoretical results on the sample-level DM algorithm. To this end, let

$$\mathcal{M}_{\text{unif}}(n, \rho) := \sup_{\theta \in B_2(r; \theta^*)} \inf_{\eta \in M_n(\theta)} \|\eta - M(\theta)\|_2, \quad \mathbb{P}[\mathcal{M}_{\text{unif}}(n, \rho) \leq \varepsilon] \geq 1 - \rho. \quad (3.6)$$

Our next result concerns the rate of convergence of the sample-level DM sequence to the population level minimizer of  $D(\theta)$ .

**Theorem 2** (Noisy contraction). *Assume that the conditions of Theorem 1 with contraction factor  $\kappa \in (0, 1)$  on  $B_2(r; \theta^*)$  hold. If  $\mathcal{M}_{\text{unif}}(n, \rho) \leq (1 - \kappa)r$ , then*

$$\|\theta_{m,n} - \theta^*\|_2 \leq \kappa^m \|\theta_{0,n} - \theta^*\|_2 + \frac{\mathcal{M}_{\text{unif}}(n, \rho)}{1 - \kappa} \quad \text{with prob.} \geq 1 - \rho.$$

We next provide an explicit form of the bound in Theorem 2. We need one additional assumption and a related notation.

**Assumption (Component Lipschitz regularity).** There exists  $L_{\text{comp}} > 0$  such that for all  $\theta, \theta' \in B_2(r'; \theta^*)$ ,

$$\sup_y \max_k \left\{ |w_k(y; \theta) - w_k(y; \theta')| + \|\nabla_\phi \log h(y; \phi_k(\theta)) - \nabla_\phi \log h(y; \phi_k(\theta'))\|_2 \right\} \leq L_{\text{comp}} \|\theta - \theta'\|_2.$$

Also, set  $J(\eta, \theta) := D_2(\nabla_1 Q)(\eta \mid \theta)$  where  $\nabla_1 Q(\eta \mid \theta)$  is the gradient of  $Q(\cdot \mid \theta)$  wrt the first argument and  $D_2$  is the Fréchet derivative wrt the second argument. Set

$$C_{\text{fos}} := \sup_{\theta \in B_2(r'; \theta^*)} \sup_{\eta \in M(\theta)} \|J(\eta \mid \theta)\|_{\text{op}} = \sup_{\theta \in B_2(r'; \theta^*)} \sup_{\eta \in M(\theta)} \sqrt{\lambda_{\max}(J^T J)},$$

where  $\lambda_{\max}(J^T J)$  is the maximal eigen-value of the matrix  $J^T J$ . Also, set  $\gamma_K = \inf\{\gamma : \text{FOS}(\gamma) \text{ holds for the } K\text{-component family on } \mathbb{B}_2(r'; \boldsymbol{\theta}^*)\}$ . Let  $\pi_{\min}$  denote the minimal mixing weight on the neighborhood. That is,  $\pi_{\min} = \inf_{\boldsymbol{\theta} \in \mathbb{B}(r'; \boldsymbol{\theta}^*)} \min_{1 \leq k \leq K} \pi_k(\boldsymbol{\theta})$ . Let  $p(K)$  denote the model complexity; that is  $p(K) = (K - 1)(\text{mixing weights}) + Kd_\phi(\text{component parameters})$ .

**Corollary 2** (Explicit  $K$ -scaling). *If in addition  $\gamma_K \leq (C_{\text{fos}}/\pi_{\min})KL_{\text{comp}}$  so that  $\kappa_K = \gamma_K/\lambda < 1$ , and*

$$M_{\text{unif}}(n, \rho) \leq C_{\text{op}} A_{\max} \sqrt{\frac{p(K) + \log(1/\rho)}{n}},$$

*then, with probability at least  $(1 - \rho)$ ,*

$$\|\boldsymbol{\theta}_{m,n} - \boldsymbol{\theta}^*\|_2 \leq \kappa_K^m \|\boldsymbol{\theta}_{0,n} - \boldsymbol{\theta}^*\|_2 + \frac{C_{\text{op}} A_{\max}}{1 - \kappa_K} \sqrt{\frac{p(K) + \log(1/\rho)}{n}}.$$

**Remark 2** (How  $d$  enters  $M_{\text{unif}}(n, r)$  when  $M_n$  and  $M$  are set-valued). *Notice that  $\Psi(\boldsymbol{\theta}; g) = \nabla_{\boldsymbol{\theta}} D_G(g, f_{\boldsymbol{\theta}}) = -\int A\left(\frac{g(y)}{f_{\boldsymbol{\theta}}(y)} - 1\right) s_{\boldsymbol{\theta}}(y) f_{\boldsymbol{\theta}}(y) dy$ ,  $s_{\boldsymbol{\theta}} = \nabla_{\boldsymbol{\theta}} \log f_{\boldsymbol{\theta}}$ . For any signed perturbation  $h$  with  $\int h = 0$ ,*

$$\partial_g \Psi(\boldsymbol{\theta}; g)[h] = -\int A'\left(\frac{g(y)}{f_{\boldsymbol{\theta}}(y)} - 1\right) s_{\boldsymbol{\theta}}(y) h(y) dy.$$

*Thus,  $\Psi(\boldsymbol{\theta}; g_n) - \Psi(\boldsymbol{\theta}; g) = \partial_g \Psi(\boldsymbol{\theta}; g)[g_n - g] + r_n(\boldsymbol{\theta})$ , and  $\sup_{\boldsymbol{\theta} \in \mathbb{B}_2(r; \boldsymbol{\theta}^*)} \|r_n(\boldsymbol{\theta})\|_2 = o_p(\|g_n - g\|_{\mathcal{H}})$ , where the remainder bound follows from local Lipschitz continuity of  $A'$ . Under the calibration  $A'(0) = 1$  and at the model ( $g = f_{\boldsymbol{\theta}^*}$ ), this simplifies to  $\partial_g \Psi(\boldsymbol{\theta}^*; g)[h] = -\int s_{\boldsymbol{\theta}^*}(y) h(y) dy$ . Hence, the leading plug-in effect depends only on the score class  $\{s_{\boldsymbol{\theta}} : \boldsymbol{\theta} \in \mathbb{B}_2(r; \boldsymbol{\theta}^*)\}$ , and is bounded by its dual-norm envelope:*

$$\sup_{\boldsymbol{\theta} \in \mathbb{B}_2(r; \boldsymbol{\theta}^*)} \|\Psi(\boldsymbol{\theta}; g_n) - \Psi(\boldsymbol{\theta}; g)\|_2 \leq A'_{\max} \text{Env}(K) \|g_n - g\|_{\mathcal{H}} + o_p(\|g_n - g\|_{\mathcal{H}}).$$

*Hence, for an envelope/entropy constant  $\text{Env}(K)$  of  $\{s_{\xi} : \xi \in M(\theta), \theta \in \mathbb{B}_2(r; \boldsymbol{\theta}^*)\}$  and the seminorm*

$$\|h\|_{\mathcal{H}} := \sup_{\theta \in \mathbb{B}_2(r; \boldsymbol{\theta}^*)} \sup_{\xi \in M(\theta)} \int \|s_{\xi}(y)\|_2 |h(y)| dy,$$

one obtains the high-probability bound

$$M_{\text{unif}}(n, r) \lesssim C_{\text{fos}} A'_{\text{max}} \text{Env}(K) \|g_n - g\|_{\mathcal{H}} \quad \text{with prob.} \geq 1 - \rho, \quad (3.7)$$

where  $C_{\text{fos}}$  is the local FOS modulus and  $A'_{\text{max}} := \sup_{\delta \geq -1} |A'(\delta)|$  is the RAF envelope (both as used in Corollary 2). Strong convexity then converts score-level perturbations to argmin-level deviations, yielding the noisy-contraction bound with  $M_{\text{unif}}(n, r)$  in place of  $\|M_n - M\|$ . All dependence on the data dimension  $d$  enters only through the plug-in rate  $\|g_n - g\|_{\mathcal{H}}$ . For the discrete (finite or countable support)  $\|g_n - g\|_{\mathcal{H}} = O_p(n^{-1/2})$ ; thus  $M_{\text{unif}}(n, r) \lesssim C_{\text{fos}} A'_{\text{max}} \text{Env}(K) n^{-1/2}$  and hence there is no dependence on  $d$ . In the continuous  $d$ -variate (kernel plug-in) case, for  $\beta$ -Hölder  $g$  and a product kernel,  $\|g_n - g\|_{\mathcal{H}} = O_p(h^\beta + \sqrt{1/(nh^d)})$ , optimized at  $n^{-\beta/(2\beta+d)}$ —the standard nonparametric rate. Consequently, the model-side constants— $C_{\text{fos}}$ ,  $A_{\text{max}}$ , local curvature/strong convexity, and any mixture-specific envelopes such as  $\text{Env}(K) \lesssim C_1 \sqrt{K}$  or  $C_2 K/\pi_{\min}$ —govern the explicit  $p(K)$ -dependence, while  $d$  affects only  $\|g_n - g\|_{\mathcal{H}}$  via known density-estimation rates.

## 4 Asymptotic Results

Throughout Sections 4.1–4.4 we take the number of mixture components  $K$  as fixed and known. Section 4.5 treats the case of unknown  $K$  via split–select–estimate and dimension matching. For each sample size  $n$ , let  $m_n \in \mathbb{N}$  denote the number of DM iterations we run on  $D_n$ . Accordingly,  $\boldsymbol{\theta}_{m_n, n}$  denotes the iterate after  $m_n$  updates at sample size  $n$ .

### 4.1 Properties of Truncated Iterates

**Definition 3.** A population (resp. sample) level DM algorithm sequence  $\{\boldsymbol{\theta}_m\}$  (resp.  $\{\boldsymbol{\theta}_{m,n}\}$ ) is called a population (resp. sample) level optimal sequence if  $\boldsymbol{\theta}^* = \lim_{m \rightarrow \infty} \boldsymbol{\theta}_m \in \underset{\boldsymbol{\theta}' \in \Theta}{\text{argmin}} D(\boldsymbol{\theta}')$  (resp.  $\boldsymbol{\theta}_n^* = \lim_{m \rightarrow \infty} \boldsymbol{\theta}_{m,n} \in \underset{\boldsymbol{\theta}' \in \Theta}{\text{argmin}} D_n(\boldsymbol{\theta}')$ ).



The following theorem is concerned with the consistency and asymptotic normality of the finitely iterated sample-level optimal DM sequence  $\{\boldsymbol{\theta}_{m,n}\}$ .

**Theorem 3.** *Assume the number of components  $K$  is fixed and the model is correctly specified  $g = f_{\boldsymbol{\theta}^*}$ .*

1. **Consistency.** *Under (C0)–(C3), any sample optimal DM sequence  $\{\boldsymbol{\theta}_{m,n}\}$  satisfies*

$$\lim_{n \rightarrow \infty} \lim_{m \rightarrow \infty} \boldsymbol{\theta}_{m,n} = \boldsymbol{\theta}^* \text{ almost surely. If } m_n \rightarrow \infty \text{ then } \boldsymbol{\theta}_{m_n,n} \xrightarrow{p} \boldsymbol{\theta}^*.$$

2.  **$\sqrt{n}$ -normality (truncated iterates).** *Assume Theorem 1 holds:  $M$  is  $\kappa$ -contractive on  $B_2(r; \boldsymbol{\theta}^*)$  with  $\kappa \in (0, 1)$ . If  $\boldsymbol{\theta}_{0,n} \in B_2(r; \boldsymbol{\theta}^*)$ ,*

$$m_n \geq \left\lceil \frac{(\frac{1}{2} + \delta) \log n + c_0}{|\log \kappa|} \right\rceil \quad (c_0 > 0, \delta > 0 \text{ fixed}),$$

*then under the conditions (F1), (C1)–(C2), (K1)–(K2), (M1)–(M8)*

$$\sqrt{n}(\boldsymbol{\theta}_{m_n,n} - \boldsymbol{\theta}^*) \Rightarrow \mathcal{N}(0, I(\boldsymbol{\theta}^*)^{-1}).$$

In particular, for some  $\delta > 0$   $m_n = \lceil ((\frac{1}{2} + \delta) \log n) / |\log \kappa| \rceil + O(1)$ , so only  $O(\log n)$  iterations are required.

**Corollary 3** (Finite-step Godambe CLT). *Let  $\hat{\boldsymbol{\theta}}_n$  solve  $\Psi_n(\boldsymbol{\theta}) := \nabla_{\boldsymbol{\theta}} D_G(g_n, f_{\boldsymbol{\theta}}) = 0$  and let  $\boldsymbol{\theta}_{m+1,n} \in M_n(\boldsymbol{\theta}_{m,n})$  with  $\boldsymbol{\theta}_{0,n} \in B_2(r; \boldsymbol{\theta}^*)$ . Assume (G1)–(G4) at  $\boldsymbol{\theta}^\dagger := \arg \min_{\boldsymbol{\theta}} D_G(g, f_{\boldsymbol{\theta}})$ , and the contraction in Theorem 1. If  $\sqrt{n} \kappa^{m_n} \rightarrow 0$ , then*

$$\sqrt{n}(\boldsymbol{\theta}_{m_n,n} - \boldsymbol{\theta}^\dagger) \Rightarrow \mathcal{N}(0, H^{-1} V H^{-1}),$$

*with  $H := \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}^\dagger; g)$  and  $V := \text{Var}_g[A'(\frac{g}{f_{\boldsymbol{\theta}^\dagger}} - 1) s_{\boldsymbol{\theta}^\dagger}(Y)]$ . In particular, under correct model specification and  $A'(0) = 1$ , the covariance reduces to  $I(\boldsymbol{\theta}^*)^{-1}$ .*

**Corollary 4** (Finite-step (Godambe-)Wilks). *Under the assumptions of Corollary 3 and  $\sqrt{n} \kappa^{m_n} \rightarrow 0$ ,*

$$2n \left\{ D_G(g_n, f_{\boldsymbol{\theta}^\dagger}) - D_G(g_n, f_{\boldsymbol{\theta}_{m_n,n}}) \right\} \Rightarrow \sum_{j=1}^{p(K_0)} \lambda_j \chi_{1,j}^2,$$

*where  $\{\lambda_j\}$  are eigenvalues of  $J := H^{-1/2} V H^{-1/2}$ . If  $g = f_{\boldsymbol{\theta}^*}$  and  $A'(0) = 1$  then  $H = V = I(\boldsymbol{\theta}^*)$  and the limit is  $\chi_{p(K_0)}^2$ .*

## 4.2 Robustness

For  $\epsilon \in [0, 1)$  and a contamination density  $\eta_n$ , define the  $\epsilon$ -contaminated model

$$f_{\epsilon,n}(y; \boldsymbol{\theta}) := (1 - \epsilon) f(y; \boldsymbol{\theta}) + \epsilon \eta_n(y), \quad \boldsymbol{\theta} \in \Theta.$$

Our first focus is on contraction and noisy contraction under small contamination. For  $\epsilon \in [0, \epsilon_0]$ , set  $g_\epsilon := (1 - \epsilon)g + \epsilon \eta$  and  $\boldsymbol{\theta}_\epsilon^\dagger := \arg \min_{\boldsymbol{\theta}} D_G(g_\epsilon, f_{\boldsymbol{\theta}})$ . Define

$$Q_\epsilon(\boldsymbol{\theta}' | \boldsymbol{\theta}) := Q_G(\boldsymbol{\theta}' | \boldsymbol{\theta}; g_\epsilon), \quad M_\epsilon(\boldsymbol{\theta}) := \arg \min_{\boldsymbol{\theta}'} Q_\epsilon(\boldsymbol{\theta}' | \boldsymbol{\theta}),$$

and analogously  $Q_{\epsilon,n}$  and  $M_{\epsilon,n}$  with  $g_{\epsilon,n}$ . Assume the fixed-order smoothness/curvature and FOS conditions of Theorem 1 hold uniformly for  $\epsilon \in [0, \epsilon_0]$  on a common ball  $B_2(r; \boldsymbol{\theta}^\star)$ .

**Corollary 5** (Population contraction under contamination). *There exists  $\bar{\kappa} \in (0, 1)$  and  $r > 0$  such that, for every  $\epsilon \in [0, \epsilon_0]$ ,*

$$\|M_\epsilon(\boldsymbol{\theta}) - \boldsymbol{\theta}_\epsilon^\dagger\|_2 \leq \bar{\kappa} \|\boldsymbol{\theta} - \boldsymbol{\theta}_\epsilon^\dagger\|_2, \quad \forall \boldsymbol{\theta} \in B_2(r; \boldsymbol{\theta}_\epsilon^\dagger).$$

The proof of the corollary is in Appendix L. We now turn to the sample based contraction whose proof is relegated to the Appendix L.

**Corollary 6** (Noisy contraction and opt-to-stat under contamination). *Let*

$$\mathcal{M}_{\text{unif}}(n, r; \epsilon) := \sup_{\boldsymbol{\theta} \in B_2(r; \boldsymbol{\theta}^\star)} \sup_{\boldsymbol{\xi} \in M_\epsilon(\boldsymbol{\theta})} \inf_{\boldsymbol{\eta} \in M_{\epsilon,n}(\boldsymbol{\theta})} \|\boldsymbol{\eta} - \boldsymbol{\xi}\|_2.$$

*Then any selection  $\boldsymbol{\theta}_{t+1,\epsilon,n} \in M_{\epsilon,n}(\boldsymbol{\theta}_{t,\epsilon,n})$  with  $\boldsymbol{\theta}_{t,\epsilon,n} \in B_2(r; \boldsymbol{\theta}_\epsilon^\dagger)$  obeys*

$$\|\boldsymbol{\theta}_{t+1,\epsilon,n} - \boldsymbol{\theta}_\epsilon^\dagger\|_2 \leq \bar{\kappa} \|\boldsymbol{\theta}_{t,\epsilon,n} - \boldsymbol{\theta}_\epsilon^\dagger\|_2 + \mathcal{M}_{\text{unif}}(n, r; \epsilon).$$

*If, in addition,  $\|g_{\epsilon,n} - g_\epsilon\|_{\mathcal{H}} = o_p(n^{-1/2})$  and  $A'(0) = 1$  with  $A'_{\max} := \sup_{\delta \geq -1} |A'(\delta)| < \infty$ , then  $\mathcal{M}_{\text{unif}}(n, r; \epsilon) = o_p(n^{-1/2})$  and any  $(\delta > 0)$*

$$m_n \geq \left\lceil \frac{(\frac{1}{2} + \delta) \log n + c_0}{|\log \bar{\kappa}|} \right\rceil$$

*yields the opt-to-stat bound  $\sqrt{n} \|\boldsymbol{\theta}_{\epsilon,n}^{(m_n)} - \hat{\boldsymbol{\theta}}_{\epsilon,n}\|_2 \xrightarrow{p} 0$*

**Remark (Robustness under contamination: bounded-RAF vs. KL).** The sample-level term in Corollary 6 is controlled by the uniform operator-deviation bound

$$\mathcal{M}_{\text{unif}}(n, r; \varepsilon) \lesssim \frac{C_{\text{fos}} A'_{\text{max}} \text{Env}(K)}{\lambda_\varepsilon} \|g_{\varepsilon, n} - g_\varepsilon\|_{\mathcal{H}},$$

on  $B_2(r; \boldsymbol{\theta}_\varepsilon^\dagger)$  (see Supplement, Corollary S.R.1). Hence **bounded-RAF** generators (NED, vNED), for which  $A'_{\text{max}} < \infty$ , yield  $\mathcal{M}_{\text{unif}}(n, r; \varepsilon) = o_p(n^{-1/2})$  under the usual plug-in rate, and any  $m_n = O(\log n)$  gives the opt-to-stat conclusion. For **KL**, the same conclusion requires a local density/score floor; without it one can have  $\mathcal{M}_{\text{unif}}(n, r; \varepsilon) \neq o_p(n^{-1/2})$  and the opt-to-stat step may fail (see Supplement, Proposition S.R.2).

We now turn to evaluate the influence function. Let  $T(\cdot)$  denote the DM population functional and write  $\boldsymbol{\theta}_\varepsilon^\star(n) := T(f_{\varepsilon, n}(\cdot; \boldsymbol{\theta}))$ , with  $\boldsymbol{\theta}^\star = \boldsymbol{\theta}_0^\star$ .

**Theorem 4.** *Let  $\{\boldsymbol{\theta}_{\varepsilon, m, n}\}$  be a sample-level DM optimal sequence at contamination level  $\varepsilon$ , and set  $\boldsymbol{\theta}_{\varepsilon, n}^\star := \lim_{m \rightarrow \infty} \boldsymbol{\theta}_{\varepsilon, m, n}$  whenever the limit exists. Assume (C2) and (O2) hold uniformly in  $n$  and  $\varepsilon \in [0, \varepsilon_0)$ .*

1. *For each fixed  $\varepsilon \in [0, \varepsilon_0)$ , suppose  $\boldsymbol{\theta}_{\varepsilon, n}^\star$  is unique for all  $n \geq 1$ . Then  $\{\boldsymbol{\theta}_{\varepsilon, n}^\star\}_{n \geq 1}$  is a bounded sequence and*

$$\lim_{n \rightarrow \infty} \boldsymbol{\theta}_{\varepsilon, n}^\star = \boldsymbol{\theta}_\varepsilon^\star.$$

2. *If (M1)–(M2) hold, then  $T$  is Gâteaux differentiable at  $f(\cdot; \boldsymbol{\theta}^\star)$  along the mixture direction  $\eta_n$ , and*

$$\lim_{\varepsilon \downarrow 0} \frac{\boldsymbol{\theta}_{\varepsilon, n}^\star - \boldsymbol{\theta}_\varepsilon^\dagger}{\varepsilon} = [I(\boldsymbol{\theta}^\star)]^{-1} \int_{\mathbb{R}} \eta_n(y) u(y; \boldsymbol{\theta}^\star) dy,$$

where  $u(y; \boldsymbol{\theta}) := \nabla_{\boldsymbol{\theta}} \log f(y; \boldsymbol{\theta})$  and  $I(\boldsymbol{\theta}^\star)$  is the Fisher information.

**Remark 3** (Influence function and gross-error sensitivity). *For point-mass contamination  $\eta_n = \delta_y$ , (ii) gives the influence function  $\text{IF}(y; T, f) = I(\boldsymbol{\theta}^\star)^{-1} u(y; \boldsymbol{\theta}^\star)$ , so the gross-error sensitivity is  $\sup_y \|I(\boldsymbol{\theta}^\star)^{-1} u(y; \boldsymbol{\theta}^\star)\|$ . (First-order efficiency also follows since  $I(\boldsymbol{\theta}^\star)$  matches the MLE information.)*

### 4.3 Breakdown Point Analysis

Roughly speaking, the breakdown point of a functional is the smallest proportion of contaminated value(s) that can produce arbitrary estimates. We consider the contamination path

$$g_{\epsilon,n}(y) = (1 - \epsilon)g(y) + \epsilon\eta_n(y), \quad \epsilon \in [0, 1].$$

Following [Simpson \(1987\)](#), the (asymptotic) *breakdown point* of a functional  $T$  is

$$\epsilon^* := \inf \left\{ \epsilon \in [0, 1] : \exists \{\eta_n\} \text{ with } \|T(g_{\epsilon,n}) - T(g)\| \rightarrow \infty \text{ as } n \rightarrow \infty \right\}.$$

**Definition 4** (Contaminated surrogate and update). *Define, for  $\theta', \theta \in \Theta$ ,*

$$Q_{\epsilon,n}(\theta' | \theta) = \mathbf{E}_Y \left[ \mathbf{E}_{Z|Y} \left\{ G \left( -1 + \frac{g_{\epsilon,n}(Y) w(Z | Y; \theta)}{f(Y; \theta') w(Z | Y; \theta')} \right) \right\} \right].$$

*A sample-level DM sequence at contamination level  $\epsilon$  satisfies  $\theta_{\epsilon,m+1,n} \in M_\epsilon(\theta_{\epsilon,m,n})$  with  $M_\epsilon(\theta) := \arg \min_{\theta'} Q_{\epsilon,n}(\theta' | \theta)$ . If  $\lim_{n \rightarrow \infty} \theta_{\epsilon,m,n} = \theta_{\epsilon,m}^*$  with  $\theta_{\epsilon,m}^* = T(g_{\epsilon,m})$ , we call  $\{\theta_{\epsilon,m}^*\}_{m \geq 0}$  a population DM optimal sequence at contamination level  $\epsilon$ .*

**Notation for contamination and finite steps.** For  $\epsilon \in [0, \epsilon_0]$ , let  $\theta_\epsilon^\dagger := \arg \min_{\theta} D_G(g_\epsilon, f_\theta)$ ,  $\hat{\theta}_{\epsilon,n} := \arg \min_{\theta} D_G(g_{\epsilon,n}, f_\theta)$ , and  $\theta_{\epsilon,n}^{(m)} \in M_{\epsilon,n}^{(m)}(\theta_{\epsilon,n}^{(0)})$  be the  $m$ -step DM iterate. We assume the contraction in [Theorem 1](#) holds uniformly on  $\mathbb{B}_2(r; \theta^*)$  with factor  $\kappa \in (0, 1)$ , and choose  $m_n$  so that  $\sqrt{n}\kappa^{m_n} \rightarrow 0$  (e.g.,  $m_n \geq \lceil (\frac{1}{2} \log n + c_0) / |\log \kappa| \rceil$ ).

**Theorem 5** (Finite-step Godambe CLT under contamination). *Assume (G1)–(G4) hold uniformly for  $\epsilon \in [0, \epsilon_0]$  at  $\theta_\epsilon^\dagger$  and the bandwidth/plug-in rate satisfies  $\|g_{\epsilon,n} - g_\epsilon\|_{\mathcal{H}} = o_p(n^{-1/2})$ . Then, with  $m_n$  as above,*

$$\sqrt{n} \left( \theta_{\epsilon,n}^{(m_n)} - \theta_\epsilon^\dagger \right) \Rightarrow \mathcal{N}(0, H_\epsilon^{-1} V_\epsilon H_\epsilon^{-1}),$$

where  $H_\epsilon := \nabla_{\theta} \Psi(\theta_\epsilon^\dagger; g_\epsilon)$  and  $V_\epsilon := \text{Var}_{g_\epsilon} \left[ A' \left( \frac{g_\epsilon(Y)}{f_{\theta_\epsilon^\dagger}(Y)} - 1 \right) s_{\theta_\epsilon^\dagger}(Y) \right]$ . In particular, at the model ( $\epsilon = 0$ ) and with  $A'(0) = 1$ , the covariance reduces to  $I(\theta^*)^{-1}$ .

**Theorem 6.** Assume the conditions of Theorem 5 and  $\sqrt{n} \kappa^{m_n} \rightarrow 0$ . Then the same limit holds with  $\boldsymbol{\theta}_{\epsilon,n}^{(m)}$  replaced by  $\boldsymbol{\theta}_{\epsilon,n}^{(m_n)}$ .

**Theorem 7** (Uniform breakdown lower bound for bounded RAF). Let  $g_\epsilon = (1 - \epsilon)g + \epsilon q$ , and let  $\hat{\boldsymbol{\theta}}_\epsilon \in \arg \min_{\boldsymbol{\theta}} D_G(g_\epsilon, f_{\boldsymbol{\theta}})$ . Assume: (i)  $D_G(g, \cdot)$  is  $\lambda$ -strongly convex on  $B(\boldsymbol{\theta}^*, r)$  and attains its minimum at  $\boldsymbol{\theta}^*$ ; (ii) for some  $S_K < \infty$  and all  $\boldsymbol{\theta} \in B(\boldsymbol{\theta}^*, r)$ ,  $\|\nabla_{\boldsymbol{\theta}} D_G(q, f_{\boldsymbol{\theta}})\| \leq S_K A_{\max}$  with  $A_{\max} := \sup_{\delta \geq -1} |A(\delta)| < \infty$  (e.g., NED/vNED); (iii)  $B(\boldsymbol{\theta}^*, r)$  contains no other local minimum. Then, for any  $\epsilon < \epsilon^\dagger := \lambda r / (S_K A_{\max})$ , we have  $\hat{\boldsymbol{\theta}}_\epsilon \in B(\boldsymbol{\theta}^*, r)$  and  $\|\hat{\boldsymbol{\theta}}_\epsilon - \boldsymbol{\theta}^*\| \leq (\epsilon S_K A_{\max}) / \lambda$ .

**Remark 4.** In regular mixtures one may take  $S_K \lesssim C_1 \sqrt{p(K)}$  or  $S_K \lesssim C_2 K / \pi_{\min}$ , yielding  $\epsilon^\dagger \gtrsim (\lambda r) / (A_{\max} C_1 \sqrt{p(K)})$  (resp.  $\lambda r \pi_{\min} / (A_{\max} C_2 K)$ ). For HD (unbounded RAF), the same bound holds under a bounded density-ratio condition on  $q$  in the neighborhood.

**Remark(Breakdown under contamination: bounded-RAF vs unbounded-RAF).**

The lower bound in Theorem 7 hinges on the uniform envelope  $\|\nabla_{\boldsymbol{\theta}} D_G(q, f_{\boldsymbol{\theta}})\| \leq S_K A_{\max}$  with  $A_{\max} := \sup_{\delta \geq -1} |A(\delta)| < \infty$  (bounded RAF). Thus generators with bounded RAF (e.g. NED, vNED) enjoy a *uniform* breakdown lower bound  $\epsilon^\dagger \asymp \lambda r / (S_K A_{\max})$  on  $B_2(r; \boldsymbol{\theta}^*)$ . For unbounded RAFs (e.g. Hellinger, KL),  $A_{\max} = +\infty$  and the uniform bound in Theorem 7 need not hold. If a local density-ratio/score floor is imposed on the contamination within the neighborhood (e.g.  $q(y) \leq (1 + \Gamma) f(y; \boldsymbol{\theta})$  for all  $y$  and  $\boldsymbol{\theta} \in B_2(r; \boldsymbol{\theta}^*)$ ), then the same proof yields a local lower bound with  $A_{\max}$  replaced by  $A_\Gamma := \sup_{\delta \in [-1, \Gamma]} |A(\delta)| < \infty$ ; see Supplement, Corollary in S.BD. 1 (for HD,  $A_\Gamma = 2(\sqrt{1 + \Gamma} - 1)$ ; for KL,  $A_\Gamma = \max\{1, \Gamma - 1\}$ ). Conversely, without any local floor one can construct contaminations for unbounded RAFs that violate any uniform lower bound; see Supplement, Proposition in S.BD.2.

## 5 Extension of Asymptotic Analysis when $K$ is unknown

We treat the case where the number of components  $K$  is unknown. Split the sample into two independent parts  $D_{1n}$  and  $D_{2n}$  of sizes  $n_1$  and  $n_2$  with  $n_1/n \rightarrow \tau \in (0, 1)$  and  $n_2/n \rightarrow 1 - \tau$ . We use  $D_{1n}$  to select  $\hat{K}_n$  (e.g., via a divergence-based information criterion), and  $D_{2n}$  to estimate the parameters conditional on  $\hat{K}_n$ . Let  $K_0$  denote the true number of components and let  $\boldsymbol{\theta}^*$  be the true parameter with  $K_0$  components.

**Dimension Matching:** Let  $\hat{K}_n$  be the estimator of the order  $K_0$  using the selection split  $D_{1n}$  of size  $n_1$ . On the estimation split  $D_{2n}$  of size  $n_2$ , let

$$\hat{\boldsymbol{\theta}}_{n_2} \in \arg \min_{\boldsymbol{\theta} \in \Theta_{\hat{K}_n}} D_{2n}(\boldsymbol{\theta}) \quad \text{with} \quad D_{2n}(\boldsymbol{\theta}) := D(g_{n_2}, f_{\boldsymbol{\theta}}).$$

Define the *dimension-matched estimator*  $\bar{\boldsymbol{\theta}}_{n_2}$  and the *dimension-matched truth*  $\boldsymbol{\theta}^*(\hat{K}_n)$  by

$$\bar{\boldsymbol{\theta}}_{n_2} = \begin{cases} (\hat{\boldsymbol{\theta}}_{n_2}, \mathbf{0}), & \hat{K}_n < K_0, \\ \hat{\boldsymbol{\theta}}_{n_2}, & \hat{K}_n \geq K_0, \end{cases} \quad \boldsymbol{\theta}^*(\hat{K}_n) = \begin{cases} \boldsymbol{\theta}^*, & \hat{K}_n \leq K_0, \\ (\boldsymbol{\theta}^*, \mathbf{0}), & \hat{K}_n > K_0, \end{cases}$$

so both live in  $\mathbb{R}^{p \max\{\hat{K}_n, K_0\}}$ . See [Khalili & Vidyashankar \(2018\)](#) for background on dimension-matching in mixtures.

**Model Selection.** To determine the number of unknown mixture components, we use the divergence-based mixture complexity estimator. Let  $\Delta_K := \{\pi \in [0, 1]^K : \sum_{k=1}^K \pi_k = 1\}$  and

$$\Theta_K := \left\{ (\pi, \phi_1, \dots, \phi_K) : \pi \in \Delta_K, \phi_k \in \Phi \ (k = 1, \dots, K) \right\}.$$

On the selection split  $D_{1n}$  of size  $n_1$ , define the generalized divergence information criterion

$$\text{GDIC}_{n_1}(K) := \widehat{\mathcal{R}}_{n_1}(K) + \frac{b_{n_1}}{n_1} p(K), \quad \widehat{\mathcal{R}}_{n_1}(K) := \inf_{\boldsymbol{\theta} \in \Theta_K} D_{1n}(\boldsymbol{\theta}),$$

where  $p(K)$  is the (identified) parameter dimension of  $\Theta_K$ , and  $b_{n_1}$  is a penalty weight with  $b_{n_1} \rightarrow \infty$  and  $\frac{b_{n_1}}{n_1} \rightarrow 0$ . We then choose  $\hat{K}_n \in \arg \min_{1 \leq K \leq K_{\max}} \text{GDIC}_{n_1}(K)$ .

*Default choice.* To match the usual BIC when  $D_{1n}$  is the average negative log-likelihood, take  $b_{n_1} = \frac{1}{2} \log n_1$ , so the penalty is  $(p(K) \log n_1)/(2n_1)$ .

**Theorem 8** (Consistency of the GDIC selector). *The following hold:*

1. **Uniform LLN:** For each fixed  $K \leq K_{\max}$ ,  $\sup_{\theta \in \Theta_K} |D_{1n}(\theta) - D(\theta)| \xrightarrow{p} 0$  as  $n_1 \rightarrow \infty$ .
2. **Identifiability gap:** Let  $K_0$  be the true component number and  $D_K^* := \inf_{\theta \in \Theta_K} D(\theta)$ . Then  $D_K^* > D_{K_0}^*$  for all  $K < K_0$ .
3. **Local regularity at  $K_0$ :**  $D$  is twice continuously differentiable at its (unique) minimizer  $\theta^* \in \Theta_{K_0}$  with positive definite Hessian  $H(\theta^*)$ , and the sample minimizer  $\hat{\theta}_{K_0,1n} \in \arg \min_{\theta \in \Theta_{K_0}} D_{1n}$  satisfies  $\hat{\theta}_{K_0,1n}$  converges in probability to  $\theta^*$  and  $D_{1n}(\hat{\theta}_{K_0,1n}) - D(\theta^*) = O_p(n_1^{-1})$ .
4. **Overfit control:** For each  $K > K_0$ , there exists a (possibly boundary) population minimizer  $\theta_K^\dagger \in \Theta_K$  with  $D(\theta_K^\dagger) = D(\theta^*)$  such that the sample minimum obeys  $D_{1n}(\hat{\theta}_{K,1n}) - D(\theta^*) = O_p(n_1^{-1})$ .
5. If  $b_{n_1} \rightarrow \infty$  and  $b_{n_1}/n_1 \rightarrow 0$ , then  $\mathbb{P}(\hat{K}_n = K_0) \rightarrow 1$ .

We now state the asymptotic theorem when  $K$  is unknown.

**Theorem 9** (Plug-in and post-selection CLTs on the estimation split). *Let the sample be split into independent parts  $\mathcal{D}_{1n}$  (size  $n_1$ ) and  $\mathcal{D}_{2n}$  (size  $n_2$ ), with  $n_1/n \rightarrow \tau \in (0, 1)$  and  $n_2/n \rightarrow 1 - \tau$ . On  $\mathcal{D}_{1n}$  select  $\hat{K}_n$  (e.g., by GDIC), and on  $\mathcal{D}_{2n}$  compute the minimum-divergence estimator*

$$\hat{\theta}_{n_2} \in \arg \min_{\theta \in \Theta_{\hat{K}_n}} D(g_{n_2}, f_\theta) \quad (\text{for the order } K \text{ indicated below}).$$

*Assume correct model specification at the true order  $K_0$  (i.e.,  $g = f_{\theta^*}$  for some  $\theta^* \in \Theta_{K_0}$ ),*

generator calibration  $A'(0) = 1$ , and the plug-in rate  $\|g_{n_2} - g\|_{\mathcal{H}} = o_p(n_2^{-1/2})$  (e.g., empirical pmf or discrete-kernel with  $h \rightarrow 0$  and  $n_2 h \rightarrow \infty$ ). Also suppose the fixed-order regularity conditions hold at  $K_0$  (e.g., (F1), (C1)–(C2), (K1)–(K2), (M1)–(M8)).

**(a) Fixed order.** For any fixed  $K$  (in particular  $K = K_0$ ),

$$\sqrt{n_2}(\hat{\boldsymbol{\theta}}_{n_2} - \boldsymbol{\theta}^*) \Rightarrow \mathcal{N}(0, I(\boldsymbol{\theta}^*)^{-1}).$$

**(b) Post-selection (unconditional).** If, in addition, the selector is consistent,  $\mathbb{P}(\widehat{K}_n = K_0) \rightarrow 1$  (e.g., under the conditions of the GDIC consistency theorem), define the dimension-matched estimator  $\bar{\boldsymbol{\theta}}_{n_2}$  and truth  $\boldsymbol{\theta}^*(\widehat{K}_n)$  by

$$\bar{\boldsymbol{\theta}}_{n_2} = \begin{cases} (\hat{\boldsymbol{\theta}}_{n_2}, \mathbf{0}), & \widehat{K}_n < K_0, \\ \hat{\boldsymbol{\theta}}_{n_2}, & \widehat{K}_n \geq K_0, \end{cases} \quad \boldsymbol{\theta}^*(\widehat{K}_n) = \begin{cases} \boldsymbol{\theta}^*, & \widehat{K}_n \leq K_0, \\ (\boldsymbol{\theta}^*, \mathbf{0}), & \widehat{K}_n > K_0, \end{cases}$$

so both live in  $\mathbb{R}^{p \max\{\widehat{K}_n, K_0\}}$ . Then

$$\sqrt{n_2} \{ \bar{\boldsymbol{\theta}}_{n_2} - \boldsymbol{\theta}^*(\widehat{K}_n) \} \Rightarrow \mathcal{N}(0, I(\boldsymbol{\theta}^*)^{-1}).$$

**Remark.** If the DM algorithm is truncated on  $\mathcal{D}_{2n}$ , it suffices to take  $m_{n_2} = O(\log n_2)$  iterations (Theorem 3(ii)) so that the optimization error is  $o_p(n_2^{-1/2})$ .

**Robust selection under bounded-RAF disparities.** For  $\epsilon \in [0, \epsilon_0)$  let  $g_\epsilon := (1-\epsilon)g + \epsilon\eta$  and write  $D_\epsilon(\boldsymbol{\theta}) := D_G(g_\epsilon, f_\boldsymbol{\theta})$  and  $D_{\epsilon, 1n}(\boldsymbol{\theta}) := D_G(g_{\epsilon, 1n}, f_\boldsymbol{\theta})$  on the selection split  $\mathcal{D}_{1n}$ . For  $K < K_0$  set the (contaminated) population gap

$$\Delta_\epsilon(K) := \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_K} D_\epsilon(\boldsymbol{\theta}) - \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_{K_0}} D_\epsilon(\boldsymbol{\theta}), \quad \Delta_0(K) > 0 \text{ by Theorem 8(2).}$$

We next turn to investigate the robustness features of the model selector described in Theorem 8.

**Theorem 10** (GDIC: finite-sample over/under-estimation bounds for bounded RAFs). *Assume (F1), (C0)–(C3), (K1)–(K2), (M1)–(M8) on a fixed neighborhood of  $\boldsymbol{\theta}^*$ , and*



that the RAF derivative is bounded and nonincreasing on  $[0, \infty)$ :  $A'(0) = 1$ ,  $A'_{\max} := \sup_{\delta \geq -1} |A'(\delta)| < \infty$ . Suppose the selection-split plug-in satisfies  $\|g_{\epsilon, 1n} - g_\epsilon\|_{\mathcal{H}} = o_p(n_1^{-1/2})$ . Let  $b_{n_1} \rightarrow \infty$  with  $b_{n_1}/n_1 \rightarrow 0$  (e.g.  $b_{n_1} = \frac{1}{2} \log n_1$ ).

**(Underestimation,  $K < K_0$ ).** There exist constants  $c_1, c_2 > 0$  (depending only on the local envelopes and curvature) such that for every fixed  $K < K_0$  and all sufficiently small  $\epsilon \leq \epsilon^b$ ,

$$\mathbb{P}(\widehat{K}_n = K) \leq \exp\left(-c_1 n_1 \Delta_\epsilon(K)^2\right) + c_2 \mathbb{P}\left(\|g_{\epsilon, 1n} - g_\epsilon\|_{\mathcal{H}} > \frac{1}{2} \Delta_\epsilon(K)\right),$$

where  $\Delta_\epsilon(K) \geq \Delta_0(K) - C\epsilon$  with  $C \lesssim A'_{\max} \text{Env}(K)$ . In particular,  $\mathbb{P}(\widehat{K}_n < K_0) \leq \sum_{K < K_0} \exp(-c_1 n_1 (\Delta_0(K) - C\epsilon)^2) + o(1)$ .

**(Overestimation,  $K > K_0$ ).** Let  $\nu_K := p(K) - p(K_0)$  and assume the (Godambe-)Wilks expansion for the selection-split DM contrast:

$$2n_1 \left\{ D_{1n}(\widehat{\boldsymbol{\theta}}_{K_0, 1n}) - D_{1n}(\widehat{\boldsymbol{\theta}}_{K, 1n}) \right\} \rightsquigarrow \chi_{\nu_K}^2 \quad (K > K_0).$$

Then, with  $b_{n_1} = \frac{1}{2} \log n_1$  (BIC-type GDIC),

$$\mathbb{P}(\widehat{K}_n \geq K_0 + 1) \leq \sum_{K=K_0+1}^{K_{\max}} \mathbb{P}\left(\chi_{\nu_K}^2 \geq \nu_K \log n_1 + o(1)\right) = O\left(\sum_{K=K_0+1}^{K_{\max}} n_1^{-\nu_K/2}\right).$$

Hence GDIC with BIC-penalty is (over-fit) consistent even under small contamination for the bounded-RAF class.

The proof of the Theorem is relegated to the appendix.

**Remark (Finite-sample robustness of GDIC vs BIC/AIC under contamination).**

On the selection split, each observation contributes  $O(n_1^{-1})$  to  $D_{1n}(\cdot)$ . For divergence generators with *bounded RAF*  $A$  (e.g. NED, vNED), the per-point decrement in  $D_{1n}$  achievable by inserting an additional component is uniformly bounded by  $A_{\max}/n_1$ . With a BIC-type penalty  $b_{n_1} = \frac{1}{2} \log n_1$ , a single “extraneous” point cannot overturn the penalty

once  $n_1$  is large, so GDIC does not overestimate  $K$  on the basis of a single contaminated datum. In contrast, for KL (likelihood disparity), the per-point contribution to  $D_{1n}$  is unbounded in the tail; absent a local density/score floor, one extreme contaminated point can reduce the KL-based contrast by more than the BIC penalty and spuriously trigger an extra component. The next two results make this precise. (Full proofs are in the Supplement.)

**Proposition 5** (GDIC overfit control with bounded RAFs). *Assume the setting of Theorem 8 on the selection split, and that the RAF is bounded:  $A_{\max} := \sup_{\delta \geq -1} |A(\delta)| < \infty$ . Fix  $K > K_0$  and let  $\Delta p := p(K) - p(K_0) \geq 1$ . Consider  $b_{n_1} \geq c \log n_1$  (e.g.  $b_{n_1} = \frac{1}{2} \log n_1$ ). Then for every data set of size  $n_1$ , set  $SS_K := \{\text{selection-split points that a } K\text{-only component captures}\}$*

$$\text{GDIC}_{n_1}(K) - \text{GDIC}_{n_1}(K_0) \geq \frac{\Delta p b_{n_1}}{n_1} - \frac{A_{\max}}{n_1} \#SS_K - R_{n_1},$$

with  $R_{n_1} = O_p(n_1^{-1/2})$  from sampling noise. In particular, for any threshold  $m_\star$

$$\#\{SS_K\} \leq m_\star \implies \text{GDIC}_{n_1}(K) - \text{GDIC}_{n_1}(K_0) \geq \frac{\Delta p b_{n_1} - A_{\max} m_\star}{n_1} - R_{n_1}.$$

Hence a single extraneous point ( $m_\star = 1$ ) cannot induce overestimation for all sufficiently large  $n_1$ , because  $\Delta p b_{n_1}/n_1 \gg A_{\max}/n_1$ . More generally, at least

$$m_{\min}(n_1, K) := \left\lceil (\Delta p) b_{n_1} / A_{\max} \right\rceil$$

contaminated points on the selection split are necessary (up to  $o_p(1)$ ) to force  $\widehat{K}_n \geq K$ .

The proof is in the Supplement, Proposition in S-GDIC, B1.

**Corollary 7** (Tail bound under  $\varepsilon$ -contamination; bounded RAFs). *Under the assumptions of Proposition 5, if the selection split has i.i.d.  $\varepsilon$ -contamination with rate  $\varepsilon \in (0, 1)$ , then for any fixed  $K > K_0$  and all large  $n_1$ ,*

$$\mathbb{P}(\widehat{K}_n \geq K) \leq \mathbb{P}(\text{Bin}(n_1, \varepsilon) \geq m_{\min}(n_1, K)) + o(1) \leq \exp\left(-n_1 \text{kl}\left(\frac{m_{\min}(n_1, K)}{n_1} \parallel \varepsilon\right)\right) + o(1),$$

where  $\text{kl}(x\|y)$  is the Bernoulli KL divergence. In particular, with  $b_{n_1} = \frac{1}{2} \log n_1$ ,  $m_{\min}(n_1, K) \asymp (\Delta p) \log n_1 / A_{\max}$  and  $\mathbb{P}(\widehat{K}_n \geq K) = O(n_1^{-c})$  for some  $c > 0$ .

The Proof of the result standard.

**Theorem 11** (Sample-level breakdown lower bounds for GDIC with bounded RAF). *Work on the selection split  $\mathcal{D}_{1n}$  of size  $n_1$ . Assume (F1), (C0)–(C3), (K1)–(K2), (M1)–(M8) on a neighborhood of  $\boldsymbol{\theta}^*$ . Suppose the divergence generator has bounded RAF  $A$  (i.e.  $A_{\max} := \sup_{\delta \geq -1} |A(\delta)| < \infty$ ), and the plug-in satisfies  $\|g_{1n} - g\|_{\mathcal{H}} = o_p(n_1^{-1/2})$ . Let  $b_{n_1} \rightarrow \infty$  with  $b_{n_1}/n_1 \rightarrow 0$  and write  $\nu_K := p(K) - p(K_0)$ .*

1. **(Underestimation)** For each fixed  $K < K_0$  let  $\Delta_0(K) := \inf_{\boldsymbol{\theta} \in \Theta_K} D(\boldsymbol{\theta}) - D(\boldsymbol{\theta}^*) > 0$  (Theorem 8(2)). Then for any  $0 < \varepsilon < \varepsilon_{\text{under}}(K) := \Delta_0(K)/(4A_{\max})$ ,

$$\mathbb{P}(\widehat{K}_n = K) \leq \exp\left(-c_1 n_1 [\Delta_0(K) - 2A_{\max}\varepsilon]^2\right) + o(1),$$

for some  $c_1 > 0$  depending only on the local envelopes; hence  $\mathbb{P}(\widehat{K}_n < K_0) \rightarrow 0$  for any fixed  $\varepsilon < \min_{K < K_0} \varepsilon_{\text{under}}(K)$ .

2. **(Overestimation)** For each  $K > K_0$  define  $m_{\min}(K, n_1) := \lceil \nu_K b_{n_1} / A_{\max} \rceil$ . Let  $X_{n_1}$  be the number of contaminated points falling in  $\mathcal{D}_{1n}$  (e.g.  $X_{n_1} \sim \text{Bin}(n_1, \varepsilon)$  under i.i.d.  $\varepsilon$ -contamination). Then for every  $K > K_0$  and all large  $n_1$ ,

$$\mathbb{P}(\widehat{K}_n \geq K) \leq \mathbb{P}(X_{n_1} \geq m_{\min}(K, n_1)) + o(1).$$

In particular, with BIC-type penalty  $b_{n_1} = \frac{1}{2} \log n_1$  and any fixed  $\varepsilon < \varepsilon_{\text{over}}(K) := \frac{\nu_K}{2A_{\max}} \frac{\log n_1}{n_1}$ ,

$$\mathbb{P}(\widehat{K}_n \geq K) \leq \exp(-c_2 n_1) + o(1),$$

for some  $c_2 = c_2(\varepsilon, \nu_K, A_{\max}) > 0$ . In particular, a single extraneous point cannot force overestimation ( $\widehat{K}_n \geq K_0 + 1$ ) once  $n_1$  is large.

**Remark (Unbounded RAFs).** For generators with unbounded RAF (e.g. KL, HD),  $A_{\max} = +\infty$  and the uniform bounds in Theorem 11 do not apply. If a local density/score

floor holds on the selection split (e.g.  $q(y) \leq (1+\Gamma) f(y; \boldsymbol{\theta})$  and  $\inf_{\boldsymbol{\theta} \in B_2(r; \boldsymbol{\theta}^*)} f(y; \boldsymbol{\theta}) \geq c_* > 0$  on the relevant neighborhood), the same conclusions hold with  $A_{\max}$  replaced by  $A_\Gamma := \sup_{\delta \in [-1, \Gamma]} |A(\delta)| < \infty$ . Without such a floor, a single extreme point can overturn a BIC penalty for KL, and no sample-level breakdown lower bound holds (see Supplement, Prop. S.GDIC-KL).

**Remark (Positive breakdown and per-point influence for GDIC-DM).** On the selection split, each observation contributes  $O(n_1^{-1})$  to the GDIC contrast  $D_{1n}(\cdot)$ . For divergence generators with bounded RAF  $A$  (e.g. NED, vNED), the per-point decrement achievable by inserting an extra component is uniformly bounded by  $A_{\max}/n_1$ , whereas the BIC-type penalty contributes  $(\Delta p) b_{n_1}/n_1$  with  $\Delta p = p(K) - p(K_0) \geq 1$ . Thus, a single contaminated point cannot overturn the penalty once  $n_1$  is large; more generally, at least a nonvanishing fraction of contaminated points is required to force overestimation, as quantified in Theorem 11. In particular, GDIC-DM with bounded RAF enjoys a strictly positive sample-level breakdown bound, in sharp contrast to KL-based EM and BIC, whose unbounded RAF allows a single extreme point to destroy any uniform breakdown lower bound.

Our next result is the stable post-selection limit distribution of the minimum divergence estimator.

**Theorem 12** (Stable post-selection CLT). *Let  $\mathcal{D}_{1n}$  and  $\mathcal{D}_{2n}$  be independent splits with  $|\mathcal{D}_{2n}| = n_2$ . Let  $\widehat{K}_n = \widehat{K}_n(\mathcal{D}_{1n})$  be the order selected on  $\mathcal{D}_{1n}$  (e.g., GDIC), and let  $\bar{\boldsymbol{\theta}}_{n_2} = \bar{\boldsymbol{\theta}}_{n_2}(\mathcal{D}_{2n}; \widehat{K}_n)$  denote the dimension-matched MDE on  $\mathcal{D}_{2n}$  at order  $\widehat{K}_n$ . Assume that model is correctly specified at  $K_0$ ; that is,  $g = f_{\boldsymbol{\theta}^*}$  for some  $\boldsymbol{\theta}^* \in \boldsymbol{\Theta}_{K_0}$  and  $A'(0) = 1$ . Assume that  $\|g_{n_2} - g\|_{\mathcal{H}} = o_p(n_2^{-1/2})$  and that the conditions used for the fixed-order plug-in CLT hold at  $K_0$ ; that is, conditions (F1), (C1)–(C2), (K1)–(K2), (M1)–(M8) hold. Let  $\widehat{K}_n$*

be a consistent estimator of  $K_0$ . Then conditionally on  $\mathcal{D}_{1n}$ ,

$$\sqrt{n_2} \{ \bar{\boldsymbol{\theta}}_{n_2} - \boldsymbol{\theta}^*(K_0) \} \Rightarrow \mathcal{N}(0, I(\boldsymbol{\theta}^*)^{-1}) \quad \text{in probability,}$$

and hence the same convergence holds unconditionally.

**Remark 5** (Repeated splits / majority vote). *In practice, we may repeat the split  $C$  times and aggregate  $\widehat{K}_n$  by majority vote; if each split's selector is consistent, the aggregated selector is also consistent, and Theorem 8 continues to hold.*

## 6 Numerical Experiments

We evaluate DM instantiations with Hellinger (HD), and vNED divergences on synthetic finite mixtures (Poisson, Poisson–Gamma/negative binomial, Poisson–lognormal), focusing on contamination and model selection. We use 5000 Monte Carlo repetitions. Full experimental details, additional models (Poisson, Poisson-lognormal), and runtime comparisons are provided in the Supplementary materials 9.4. We first demonstrate robustness with known  $K$  and then evaluate the full split-select-estimate pipeline with unknown  $K$ . In the following tables, 'Ave' represents the average value of the estimates and 'StD' denotes the standard deviation." All methods are initialized using k-means clustering.

### 6.1 Robustness Simulation

#### PG Mixture

We simulate a two-component PG mixture ( $K = 2$ ) and true parameters  $(\pi_1, \alpha_1, \beta_1, \alpha_2, \beta_2) = (0.3, 10, 1, 1, 2)$  (details in supplementary materials 9.4). We inject point-mass contamination by replacing an  $\epsilon$ –fraction of values with the value 50. Figure 2 plots the average estimates versus  $\epsilon$ .

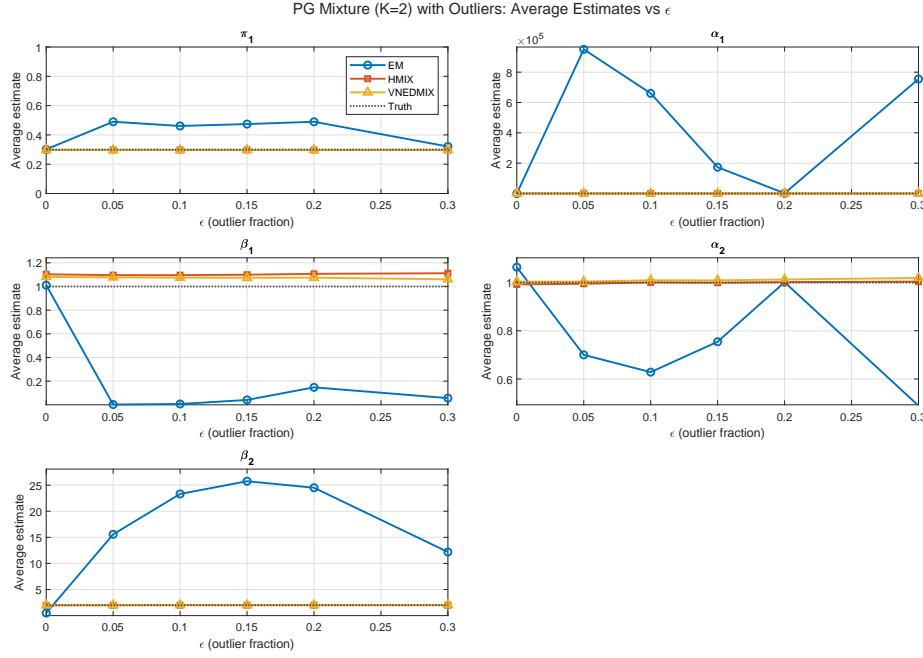


Figure 2: Average parameter estimates versus contamination level  $\epsilon$  in a two-component PG mixture with known  $K = 2$ .

As expected of likelihood-based methods, EM is sensitive to large injected counts and its averages drift with  $\epsilon$ . HMIX and vNEDMIX, which replace the likelihood by divergences against the empirical distribution, down-weight tail mismatches and remain closer to truth over a wider range of  $\epsilon$ . For readability,  $y$ -axes are trimmed using robust quantiles so occasional non-convergent runs at high contamination do not distort the display.

## 6.2 The number of mixtures $K$ is unknown

We now estimate both the number of components and the parameters via split-select-estimate with repeated splits. Let  $\mathcal{D}_{1n}$  and  $\mathcal{D}_{2n}$  be an even random split of the sample of sizes  $n_1$  and  $n_2$ . We adopt the following four-step process:

1. *Select  $K$  on  $\mathcal{D}_{1n}$* : compute  $\widehat{K}_n$  by GDIC with penalty  $b_{n_1} = \frac{1}{2} \log n_1$ .
2. *Estimate on  $\mathcal{D}_{2n}$* : fit the MDE at order  $\widehat{K}_n$  to obtain  $\widehat{\theta}$ .

3. *Repeat  $C$  times:* take a majority vote over  $\widehat{K}_n^{(1)}, \dots, \widehat{K}_n^{(C)}$ ; average parameters across runs that select the voted order.
4. *Repeat for  $R$  Monte Carlo runs* and report Ave/StD/MSE (we use  $R = 5,000$  in the tables).

All results in this subsection use the empirical kernel for  $g_n(\cdot)$ .

Table 1: Parameter Estimation with Unknown  $K$  with Data Splitting Method

PG model		$\hat{\pi}_1$	$\hat{\alpha}_1$	$\hat{\beta}_1$	$\hat{\alpha}_2$	$\hat{\beta}_2$
EM (100.0%)	Ave	0.297	10.62	1.054	0.937	1.826
	StD	0.004	0.675	0.063	0.078	0.192
HMIX (98.2%)	Ave	0.298	10.61	1.060	0.988	1.962
	StD	0.004	0.692	0.065	0.083	0.205
VNEDMIX (86.2%)	Ave	0.298	10.48	1.048	0.996	1.986
	StD	0.004	0.698	0.066	0.084	0.209

Table 1 summarizes a two-component PG mixture with truth  $(\pi_1, \alpha_1, \beta_1, \alpha_2, \beta_2) = (0.3, 10, 1, 1, 2)$  using  $C = 5$  repeats in the split-select-estimate pipeline. For each method (EM, HMIX, VNEDMIX) we first identify  $K$  and then report Ave/StD computed *only* on datasets where that method correctly identified  $K = 2$ . All three methods are essentially unbiased for component means. For example, the implied means  $\alpha_1/\beta_1$  are  $\approx 10.08$  (EM), 10.01 (HMIX), and 10.00 (VNEDMIX) versus the truth 10, and for component 2 they are approximately 0.51, 0.50, and 0.50 versus the truth 0.5. Dispersion parameters are close to the truth with small variability. When  $K$  is correctly identified, EM, HMIX, and VNEDMIX produce comparably accurate estimates in this PG setting.

## 7 Data Analysis: Image Segmentation

In this section, we apply our method to the image segmentation problem. Our methods will be tested on the Lena image (see original image [9.4](#)) which contains  $770 \times 776$  pixels. The gray-scale intensity values of the images range from 0 to 255 (all are integer-valued). We treat each pixel as one data point and fit a three-component Poisson mixture model for both images and apply EM algorithm, HMIX algorithm, and VNEDMIX algorithm. The empirical density estimate is used as the non-parametric density estimate. The gray-scale intensity values are labeled as 1, 100, and 200 to represent three classes.

### Robust Image Recovery

We study the robustness property of the proposed method. Specifically, we generate contaminated pixel data following the Poisson distribution with mean 250 with probability 0.3; if the generated value exceeds 255, then set it as 255 (since the maximum value for color is 255 under this format). We fit a three-component Poisson mixture model using EM algorithm, HMIX algorithm, and VNEDMIX algorithm. From Figure [3](#), we observe that the EM algorithm captures a lot of “noise” added to the image (the face has a lot of shadow area). In comparison, the HMIX algorithm and VNEDMIX algorithm are able to ignore the “outliers” well (the EM and VNEDMIX algorithm recovered figure are shown in supplementary material [9.4](#) and [9.4](#) respectively); in particular, the face contours are captured better. Furthermore, the point estimates based on the EM algorithm are highly affected while the other two algorithms give similar results to those from uncorrupted image.

One may argue that once we add 30% outliers, these outliers may be treated as another component. However, even if we fit a four-component Poisson mixture model and apply EM algorithm, it still produces as much shadow similar to using three-component Poisson mixture (image not displayed here). In contrast, HMIX and VNEDMIX will potentially



eliminate this additional component and produce more “reasonable” image. Additional image analysis including original image recovery, model selection tables can be found in supplementary analysis [9.4](#).



(a) Corrupted image for Lena



(b) Recovered image from HMIX algorithm

Figure 3: Lena image reconstruction after adding 30% outliers.

## 8 Concluding Remarks

We developed a divergence-minimization (DM) framework for models with latent structure that unifies classical algorithms (EM, HMIX, HELMIX) and yields robust, likelihood-free updates through a single operator. At the algorithmic level, we established a majorization-separation identity that links a complete-data surrogate to the observed-data divergence, and used it to prove monotone descent at the sample level together with global stationarity of limit points. At the population level, we provided a local contraction result under strong convexity and first-order stability, and derived a noisy-contraction bound for the sample operator.

For model selection, we proposed a divergence-based criterion (GDIC) with repeated sample-splitting and showed that it consistently recovers the mixture order under mild conditions;

the post-selection MDE enjoys the usual  $\sqrt{n}$ -asymptotic normality. On the empirical side, we studied discrete-kernel estimators of  $g_n$  in finite mixtures of counts, quantified bandwidth effects, and found that simple triangular kernels are competitive at small sample sizes. Finally, synthetic experiments and image-segmentation case studies demonstrate that HD/NED/vNED instantiations of DM are competitive with EM in well-specified regimes and deliver markedly greater stability under contamination and misspecification.

**Future work.** Natural extensions of our work include (i) scalable DM for high-dimensional latent structures (sparse or low-rank components), (ii) semiparametric DM with estimated nuisance  $g_n$  under weaker smoothness and discrete-support constraints, (iii) selection consistency beyond fixed  $K_{\max}$  (growing-model regimes), and (iv) tighter nonasymptotic guarantees for repeated sample-splitting and majority-vote selection.

## 9 Disclosure statement

The authors have no conflict of interest.

## SUPPLEMENTARY MATERIAL

### A: Background and Literature Review

Finite mixture models (FMM), a class of models for data with latent structure, have gained increasing attention for their flexibility in capturing multiple modes and unobserved heterogeneity. They are applied across diverse fields including astronomy, social sciences, biology, engineering, and medicine, and more recently have been adopted in neural network and deep learning research. Since the pioneering work of [Pearson \(1894\)](#) for mixtures of univariate normal distributions, FMM have been extended to a general class of distributions. Turning to fitting the models, starting with the work of [Radhakrishna Rao \(1948\)](#) who used Fisher’s method of scoring, the field has evolved into the routine use of EM algorithm and its variants as in the seminal paper of [Dempster et al. \(1977\)](#) (See also [Ganesalingam & McLachlan \(1978\)](#), [Ganesalingam & McLachlan \(1979a\)](#), [Ganesalingam & McLachlan \(1979b\)](#), [Ganesalingam & McLachlan \(1980\)](#), [O’Neill \(1978\)](#), and [Aitkin \(1980\)](#) for more details).

As is well known, the EM algorithm introduces a latent label for each observation, representing group membership, and treats it as missing data. This facilitates the formulation of “complete data” (missing and observed data) log-likelihood function, under the assumption of independence. The Expectation step (E-step) computes the conditional expectation of the complete-data log-likelihood given the current parameter estimates and observed data. The Maximization step (M-step) then updates the parameters by maximizing this expected log-likelihood. These steps are iterated until convergence. Regarding its convergence properties, [Dempster et al. \(1977\)](#) established that the log-likelihood function is non-decreasing after each iteration and [Wu \(1983\)](#) provided regularity conditions under which the limit

points of EM algorithm are stationary points of the likelihood function. Further convergence properties are classical; see [Wu \(1983\)](#) and [Vaida \(2005\)](#). Acceleration and variants are surveyed in [Louis \(1982\)](#), [McLachlan \(1995\)](#), [Tseng \(2004\)](#).

Despite the popularity of EM algorithm, it is well known that it has several limitations. First, the EM algorithm may converge slowly, and the situation is worse when the “incomplete information” dominates the likelihood. Various methods and modifications have been tried to improve the speed of convergence, such as Aitken’s acceleration method (see [McLachlan \(1995\)](#)), Louis’ method (see [Louis \(1982\)](#)), Conjugate Gradient method (see [Jamshidian & Jennrich \(1993\)](#)), and EM Gradient algorithm (see [Lange \(1995\)](#)). Second, it is not stable as it may often converge to the local optima. This problem is worse for contaminated data (see [Cutler & Cordero-Braña \(1996\)](#)). To address the issue, [Hu et al. \(2017\)](#) proposed a robust EM-type algorithm for log-concave mixtures regression models, where they use the trimmed least squares technique (see [Rousseeuw \(1985\)](#)) to achieve the robustness properties. However, the performance is not satisfactory when the percentage of outliers increase. Additionally, the use of trimmed least squares lead to loss of efficiency at the model. We address these challenges using the divergence method.

Divergence-based methods (see [Lindsay \(1994\)](#) and [Basu & Lindsay \(1994\)](#)) are appealing for parametric inference when the model is misspecified or when the data are contaminated. These methods have the property that they are first-order efficient when the model is correctly specified and are robust under model misspecification and the presence of outliers. Starting with the work of [Beran \(1977\)](#), who proposed minimum Hellinger distance (MHD) estimation for independent identically distributed (i.i.d.) data, extensions and variations of the theme have been studied, for example, in [Stather \(1981\)](#), [Donoho & Liu \(1988\)](#), [Eslinger & Woodward \(1991\)](#), [Basu & Harris \(1994\)](#), [Basu, Basu & Chaudhuri \(1997\)](#), [Cheng & Vidyashankar \(2006\)](#), [Simpson \(1987\)](#), [Simpson \(1989\)](#), and [Li et al. \(2019\)](#). Turning

to mixture models, [Woodward et al. \(1995\)](#) considered two- component normal mixture model and proposed to use MHD method to estimate the mixing proportion. [Cutler & Cordero-Braña \(1996\)](#) estimated all the unknown parameters using MHD method for normal mixture by the use of the so called HMIX algorithm. Note that the genesis of the HMIX algorithm and its connection to the EM algorithm are unknown. For discrete data, [Karlis & Xekalaki \(1998\)](#) considered the MHD estimation for Poisson mixtures and used the so called HELMIX algorithm, a variant of the HMIX algorithm. The HELMIX algorithm is specialized for Poisson models as it utilizes a recurrence relation for Poisson probabilities that do not generalize to other distributions. More recent references regarding the robust EM algorithm and divergences are referred [Nielsen & Sun \(2016\)](#), [Qin & Priebe \(2013\)](#), [Sammaknejad et al. \(2019\)](#), [Hu et al. \(2020\)](#), [Lücke & Forster \(2019\)](#), [Zhao et al. \(2020\)](#). Minimum-divergence (disparity) methods remain first-order efficient at the model yet damp the effect of large residuals via the RAF, which motivates our DM operator in Section 2 of the main text.

## B: Plot of Residual Adjustment Function (RAF)

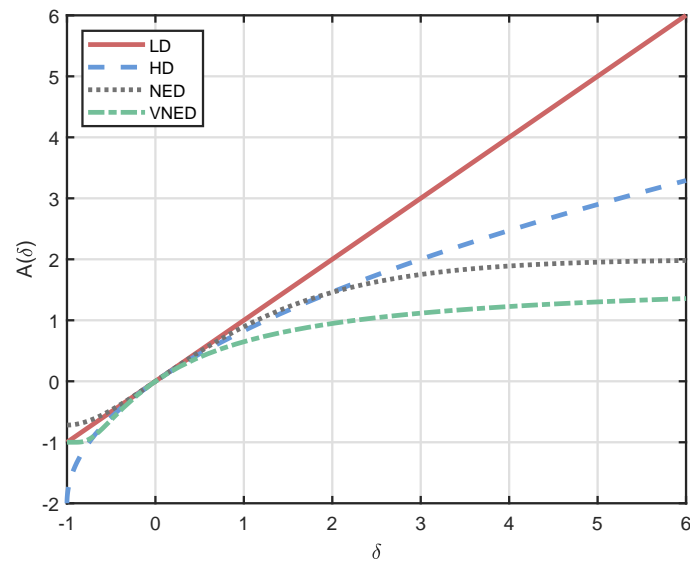


Figure 4: Plot of Residual Adjustment Function  $A(\delta)$  for LD, HD, NED, and vNED

## C: Derivation of updating mixing probability in (2.4)

*Proof of the update  $\pi$ .* Recall

$$Q(\boldsymbol{\theta}' \mid \boldsymbol{\theta}) = \sum_{k=1}^K \int_{\mathcal{Y}} \pi'_k h(y; \boldsymbol{\phi}'_k) G(\tau_k(y)) dy, \quad \tau_k(y) := \frac{g(y) w_k(y; \boldsymbol{\theta})}{\pi'_k h(y; \boldsymbol{\phi}'_k)},$$

where  $w_k(y; \boldsymbol{\theta}) = \frac{\pi_k h(y; \boldsymbol{\phi}_k)}{\sum_{\ell=1}^K \pi_\ell h(y; \boldsymbol{\phi}_\ell)}$  and the mixing weights satisfy  $\pi'_k > 0$  and  $\sum_{k=1}^K \pi'_k = 1$ .

For brevity write  $h'_k(y) := h(y; \boldsymbol{\phi}'_k)$  and  $\tau_k = \tau_k(y)$ .

We minimize  $Q(\boldsymbol{\theta}' \mid \boldsymbol{\theta})$  over  $\boldsymbol{\pi}'$  on the simplex via the Lagrangian

$$\mathcal{L}(\boldsymbol{\pi}', \lambda) = Q(\boldsymbol{\theta}' \mid \boldsymbol{\theta}) + \lambda \left( \sum_{k=1}^K \pi'_k - 1 \right).$$

**Step 1: Gradient with respect to  $\pi'_k$ .** For the integrand  $\pi'_k h'_k(y) G(\tau_k)$  with  $\tau_k = (g w_k)/(\pi'_k h'_k)$ , we have

$$\frac{\partial}{\partial \pi'_k} [\pi'_k h'_k G(\tau_k)] = h'_k G(\tau_k) + \pi'_k h'_k G'(\tau_k) \frac{\partial \tau_k}{\partial \pi'_k} = h'_k \left\{ G(\tau_k) - \tau_k G'(\tau_k) \right\},$$

because  $\frac{\partial \tau_k}{\partial \pi'_k} = -\frac{\tau_k}{\pi'_k}$ . Therefore

$$\frac{\partial \mathcal{L}}{\partial \pi'_k} = \int_{\mathcal{Y}} h'_k(y) \left\{ G(\tau_k(y)) - \tau_k(y) G'(\tau_k(y)) \right\} dy + \lambda.$$

Define

$$\Xi_k(\boldsymbol{\pi}', \boldsymbol{\phi}') := \int_{\mathcal{Y}} h'_k(y) \left\{ G(\tau_k(y)) - \tau_k(y) G'(\tau_k(y)) \right\} dy.$$

**Abbreviation.**  $B(u) := G(u) - u G'(u)$ . Then  $\frac{\partial Q^{(k)}(\boldsymbol{\theta}' \mid \boldsymbol{\theta})}{\partial \pi'_k} = \int h'_k(y) B(\tau_k(y)) dy$ .

Stationarity (KKT) gives, for all  $k$ ,

$$\Xi_k(\boldsymbol{\pi}', \boldsymbol{\phi}') = -\lambda. \tag{9.8}$$

**Step 2: Product identity and compact reparametrization.** Note that  $\pi'_k h'_k(y) \tau_k(y) \equiv g(y) w_k(y; \boldsymbol{\theta})$ . Multiplying (9.8) by  $\pi'_k$  and using this identity yields

$$\pi'_k \Xi_k(\boldsymbol{\pi}', \boldsymbol{\phi}') = \int_{\mathcal{Y}} \left\{ \pi'_k h'_k(y) G(\tau_k(y)) - g(y) w_k(y; \boldsymbol{\theta}) G'(\tau_k(y)) \right\} dy.$$

Introduce

$$\Phi_k(\boldsymbol{\pi}', \boldsymbol{\phi}') := \int_{\mathcal{Y}} \left\{ \pi'_k h'_k(y) G(\tau_k(y)) - g(y) w_k(y; \boldsymbol{\theta}) G'(\tau_k(y)) \right\} dy,$$

so that

$$\Phi_k(\boldsymbol{\pi}', \boldsymbol{\phi}') = \pi'_k \Xi_k(\boldsymbol{\pi}', \boldsymbol{\phi}') = -\lambda \pi'_k. \quad (9.9)$$

Moreover, since  $Q_k(\boldsymbol{\theta}' | \boldsymbol{\theta}) = \int_{\mathcal{Y}} \pi'_k h'_k(y) G(\tau_k(y)) dy$ , we have

$$\frac{\partial Q_k(\boldsymbol{\theta}' | \boldsymbol{\theta})}{\partial \pi'_k} = \Xi_k(\boldsymbol{\pi}', \boldsymbol{\phi}'), \quad \text{hence} \quad \Phi_k(\boldsymbol{\pi}', \boldsymbol{\phi}') = \pi'_k \frac{\partial Q_k(\boldsymbol{\theta}' | \boldsymbol{\theta})}{\partial \pi'_k}.$$

**Step 3: Normalization and the update.** Summing (9.9) over  $k$  and using  $\sum_k \pi'_k = 1$  gives

$$\sum_{k=1}^K \Phi_k(\boldsymbol{\pi}', \boldsymbol{\phi}') = -\lambda.$$

Therefore, for each  $k$ ,

$$\pi'_k = \frac{\Phi_k(\boldsymbol{\pi}', \boldsymbol{\phi}')}{\sum_{\ell=1}^K \Phi_{\ell}(\boldsymbol{\pi}', \boldsymbol{\phi}')} = \frac{\pi'_k \frac{\partial Q_k(\boldsymbol{\theta}' | \boldsymbol{\theta})}{\partial \pi'_k}}{\sum_{\ell=1}^K \pi'_{\ell} \frac{\partial Q_{\ell}(\boldsymbol{\theta}' | \boldsymbol{\theta})}{\partial \pi'_{\ell}}}.$$

By construction,  $\sum_k \pi'_k = 1$ . The integrability and finiteness of  $\Phi_k$  follow from the stated regularity conditions.  $\square$



# D: Relation of DM algorithm with other popular algorithms

## 1. Relation with the Proximal Point Algorithm

[Martinet \(1970\)](#), [Rockafellar \(1976b\)](#) and [Rockafellar \(1976a\)](#) proposed an iterative algorithm which is referred to as the proximal point algorithm and can be described as

$$\boldsymbol{\theta}_{m+1} = \operatorname{argmax}_{\boldsymbol{\theta}' \in \Theta} \left\{ \Psi(\boldsymbol{\theta}') - \frac{\beta_m}{2} \|\boldsymbol{\theta}' - \boldsymbol{\theta}_m\|_2^2 \right\},$$

where  $\Psi : \Theta \rightarrow \mathbb{R}$  and the quadratic penalty is relaxed using a positive sequence  $\{\beta_m\}$ .

[Tseng \(2004\)](#) introduced the entropy-like proximal point algorithm as

$$\boldsymbol{\theta}_{m+1} = \operatorname{argmax}_{\boldsymbol{\theta}' \in \Theta} \{ \Psi(\boldsymbol{\theta}') - H(\boldsymbol{\theta}' | \boldsymbol{\theta}_m) \},$$

where  $H : \Theta \times \Theta \rightarrow \mathbb{R}_+$  satisfies  $H(\boldsymbol{\theta}' | \boldsymbol{\theta}') = 0$  for all  $\boldsymbol{\theta}' \in \Theta$ . It is known that the EM algorithm can be considered as an example of proximal point algorithm (see [Chrétien & Hero \(2000\)](#)), and other extensions of proximal point algorithm are referred to [Tseng \(2004\)](#), [Chrétien & Hero \(2008\)](#), [Cunha et al. \(2010\)](#), and [Al Mohamad & Broniatowski \(2016\)](#).

The DM algorithm can also be treated as proximal point method by letting  $\Psi(\boldsymbol{\theta}') = -D(\boldsymbol{\theta}')$  and  $H(\boldsymbol{\theta}' | \boldsymbol{\theta}) = -D(\boldsymbol{\theta}') + Q(\boldsymbol{\theta}' | \boldsymbol{\theta})$ .

We now provide another relation via some examples between DM algorithm and the proximal point algorithm. In particular, we can view  $\Psi(\boldsymbol{\theta}')$  as the divergence information from the observed data and treat  $H(\boldsymbol{\theta}' | \boldsymbol{\theta}_m)$  as the divergence information for the latent variable. For example, for EM objective function, note that,

$$\begin{aligned} Q_{\text{KL}}(\boldsymbol{\theta}' | \boldsymbol{\theta}) &= - \int_{\mathcal{Y}} \int_{\mathcal{Z}} g(y) w(z | y; \boldsymbol{\theta}) \log \left( \frac{p(y, z; \boldsymbol{\theta}')}{g(y) w(z | y; \boldsymbol{\theta})} \right) dz dy \\ &= - \int_{\mathcal{Y}} \log \left( \frac{f(y; \boldsymbol{\theta}')}{g(y)} \right) g(y) dy - \int_{\mathcal{Y}} \left( \int_{\mathcal{Z}} w(z | y; \boldsymbol{\theta}) \log \left( \frac{w(Z | Y; \boldsymbol{\theta}')}{w(z | y; \boldsymbol{\theta})} \right) dz \right) g(y) dy. \end{aligned}$$

Let

$$\Psi_{\text{EM}}(\boldsymbol{\theta}') = - \int_{\mathcal{Y}} \log \left( \frac{f(y; \boldsymbol{\theta}')}{g(y)} \right) g(y) dy \quad \text{and} \quad (9.10)$$

$$H_{\text{EM}}(\boldsymbol{\theta}'|\boldsymbol{\theta}) = \int_{\mathcal{Y}} \left( \int_{\mathcal{Z}} w(z|y; \boldsymbol{\theta}) \log \left( \frac{w(Z|Y; \boldsymbol{\theta}')}{w(z|y; \boldsymbol{\theta})} \right) dz \right) g(y) dy, \quad (9.11)$$

then it becomes the proximal point algorithm. More importantly, the benefit for splitting  $Q_{\text{KL}}(\boldsymbol{\theta}'|\boldsymbol{\theta})$  into  $\Psi_{\text{EM}}(\boldsymbol{\theta}')$  and  $H_{\text{EM}}(\boldsymbol{\theta}'|\boldsymbol{\theta})$  is that  $Q_{\text{KL}}(\boldsymbol{\theta}'|\boldsymbol{\theta})$  is considered as the Kullback-Leibler divergence for complete data,  $\Psi_{\text{EM}}(\boldsymbol{\theta}')$  is exactly the Kullback-Leibler divergence for observed data, and  $H_{\text{EM}}(\boldsymbol{\theta}'|\boldsymbol{\theta})$  can be viewed as the Kullback-Leibler divergence for the latent variable given observed data (and parameters in previous step). As another example, taking  $\Psi_{\text{EM}}(\boldsymbol{\theta}')$  and  $H_{\text{EM}}(\boldsymbol{\theta}'|\boldsymbol{\theta})$  as (9.10) and (9.11) respectively, we have that

$$\Psi_{\text{EM}}(\boldsymbol{\theta}') \approx 2 \int_{\mathcal{Y}} \left( g^{\frac{1}{2}}(y) - f^{\frac{1}{2}}(y; \boldsymbol{\theta}') \right)^2 dy \equiv 2\Psi_{\text{HD}}(\boldsymbol{\theta}'), \quad \text{and} \quad (9.12)$$

$$\begin{aligned} H_{\text{EM}}(\boldsymbol{\theta}'|\boldsymbol{\theta}) &\approx -2 \int_{\mathcal{Y}} \left( \int_{\mathcal{Z}} \left( c^{\frac{1}{2}}(z|y; \boldsymbol{\theta}') - c^{\frac{1}{2}}(z|y; \boldsymbol{\theta}) \right)^2 dz \right) g(y) dy \\ &\equiv -2H_{\text{HD}}(\boldsymbol{\theta}'|\boldsymbol{\theta}). \end{aligned} \quad (9.13)$$

In addition, from (9.12) and (9.13) and using substitution principle, we have

$$Q_{\text{HD}}(\boldsymbol{\theta}'|\boldsymbol{\theta}) = \Psi_{\text{HD}}(\boldsymbol{\theta}') - H_{\text{HD}}(\boldsymbol{\theta}'|\boldsymbol{\theta}).$$

Similar as Kullback-Leibler divergence, the Hellinger distance divergence for complete data can also be divided into the “observed data divergence” part and the “latent data divergence” part, i.e.,  $Q_{\text{HD}}(\boldsymbol{\theta}'|\boldsymbol{\theta})$  is the Hellinger distance divergence for complete data,  $\Psi_{\text{HD}}(\boldsymbol{\theta}')$  is the Hellinger distance divergence for observed data, and  $H_{\text{HD}}(\boldsymbol{\theta}'|\boldsymbol{\theta})$  is the Hellinger distance divergence for latent variable given observed data (and current parameter). More generally, it is believed that for any divergence,  $Q(\boldsymbol{\theta}'|\boldsymbol{\theta})$  the divergence objective function for complete data,  $\Psi(\boldsymbol{\theta}')$  is the divergence objective function for observed data, and  $H(\boldsymbol{\theta}'|\boldsymbol{\theta})$  is the divergence objective function for latent variable given observed data (and current parameter). The detailed discussion is considered elsewhere.

## 2. Relation with MM algorithm

[Hunter & Lange \(2000b\)](#) proposed the general MM algorithm to construct optimization algorithms. Specifically, an MM algorithm creates a surrogate function that minorizes or majorizes the objective function and when the surrogate function is optimized, the objective function is forced to decrease or increase correspondingly. MM algorithms are widely used in a broad application areas, for instance, EM algorithm, robust regression (see [Huber \(1981\)](#)) quantile regression (see [Hunter & Lange \(2000a\)](#)), survival analysis (see [R Hunter & Lange \(2002\)](#)), paired and multiple comparisons (specifically on generalized Bradley-Terry models, see [Hunter \(2004\)](#)), etc. For a more detailed review of MM algorithm, see [Hunter & Lange \(2004\)](#).

The MM algorithm proceeds as follows. Let  $\boldsymbol{\theta}_m$  represent a fixed value of the parameter  $\boldsymbol{\theta}$ , and let  $\psi(\boldsymbol{\theta}|\boldsymbol{\theta}_m)$  denote a real-valued function of  $\boldsymbol{\theta}$  depending on  $\boldsymbol{\theta}_m$ . The function  $\psi(\boldsymbol{\theta}|\boldsymbol{\theta}_m)$  is said to majorize a real-valued function  $t(x)$  at the point  $\boldsymbol{\theta}_m$  provided

$$\psi(\boldsymbol{\theta}|\boldsymbol{\theta}_m) \geq t(\boldsymbol{\theta}) \quad \text{for all } \boldsymbol{\theta}, \text{ and } \psi(\boldsymbol{\theta}_m|\boldsymbol{\theta}_m) = t(\boldsymbol{\theta}_m). \quad (9.14)$$

Ordinarily,  $\boldsymbol{\theta}_m$  denotes the current iteration in a search surface of  $t(\boldsymbol{\theta})$ . In a majorize-minimize MM algorithm, we minimize the majorizing function  $\psi(\boldsymbol{\theta}|\boldsymbol{\theta}_m)$  rather than  $t(\boldsymbol{\theta})$ . If  $\boldsymbol{\theta}_{m+1}$  represents the minimizer of  $\psi(\boldsymbol{\theta}|\boldsymbol{\theta}_m)$ , then MM procedure forces  $t(\boldsymbol{\theta})$  downhill. To see this, note that

$$t(\boldsymbol{\theta}_{m+1}) = \psi(\boldsymbol{\theta}_{m+1}|\boldsymbol{\theta}_m) + t(\boldsymbol{\theta}_{m+1}) - \psi(\boldsymbol{\theta}_{m+1}|\boldsymbol{\theta}_m) \leq \psi(\boldsymbol{\theta}_m|\boldsymbol{\theta}_m) + t(\boldsymbol{\theta}_m) - \psi(\boldsymbol{\theta}_m|\boldsymbol{\theta}_m) = t(\boldsymbol{\theta}_m).$$

It turns out that DM algorithm also belongs to the class of MM algorithms by letting  $\psi(\boldsymbol{\theta}|\boldsymbol{\theta}_m) \equiv Q(\boldsymbol{\theta}|\boldsymbol{\theta}_m)$ ,  $t(\boldsymbol{\theta}) \equiv D(\boldsymbol{\theta})$ .

### Proximal Point Algorithm and MM algorithm

Note that the proximal point algorithm also belongs to the class of MM algorithms. Specifically, let  $\psi(\boldsymbol{\theta}|\boldsymbol{\theta}_m) = -(\Psi(\boldsymbol{\theta}) - H(\boldsymbol{\theta}|\boldsymbol{\theta}_m))$ ,  $t(\boldsymbol{\theta}) = -\Psi(\boldsymbol{\theta})$ . By definition of  $\Psi(\cdot)$

and  $H(\cdot|\cdot)$ , it follows that (9.14) is satisfied. Furthermore, if  $\boldsymbol{\theta}_{m+1}$  represents the minimizer of  $\psi(\boldsymbol{\theta}|\boldsymbol{\theta}_m)$ , then  $t(\boldsymbol{\theta}_{m+1}) \leq t(\boldsymbol{\theta}_m)$ .

**3. Relation of DM algorithm with coordinate descent.** Following Neal & Hinton (1998), EM can be viewed as a coordinate ascent algorithm. Similarly, DM can be viewed as a coordinate descent algorithm. Define

$$\mathcal{D}(\tilde{q}, \boldsymbol{\theta}') = \mathbf{E}_Y \left[ \mathbf{E}_{Z \sim \tilde{q}(\cdot|Y)} G \left( -1 + \frac{g(Y) \tilde{q}(Z|Y)}{f(Y; \boldsymbol{\theta}') w(Z|Y; \boldsymbol{\theta}')} \right) \right],$$

where  $\tilde{q}(\cdot|y)$  is any conditional density on  $Z$  given  $Y = y$ . By Lemma 5,  $\mathcal{D}(\tilde{q}, \boldsymbol{\theta}') \geq D(\boldsymbol{\theta}')$  with equality iff  $\tilde{q}(\cdot|y) = w(\cdot|y; \boldsymbol{\theta}')$ . Thus the DM algorithm alternates:

1. *D-step.* For fixed  $\boldsymbol{\theta}$ , minimize  $\mathcal{D}(\tilde{q}, \boldsymbol{\theta})$  over  $\tilde{q}$ , yielding  $\tilde{q}(z|y) = w(z|y; \boldsymbol{\theta})$ .
2. *M-step.* For this choice of  $\tilde{q}$ , minimize  $\mathcal{D}(\tilde{q}, \boldsymbol{\theta}')$  over  $\boldsymbol{\theta}'$ , which recovers the DM update.

Hence DM belongs to the class of two-block coordinate descent algorithms.

**Lemma 5.** *For any conditional density  $\tilde{q}(\cdot|y)$  and all  $\boldsymbol{\theta}' \in \boldsymbol{\Theta}$ ,*

$$\mathcal{D}(\tilde{q}, \boldsymbol{\theta}') \geq D(\boldsymbol{\theta}').$$

*Equality holds if and only if  $\tilde{q}(z|y) = w(z|y; \boldsymbol{\theta}')$  for almost every  $y \in \mathcal{Y}$ .*

**The generalized DM Algorithm:** It is possible that in the M-step, the minimizer  $\boldsymbol{\theta}_{m+1}$  is not unique. Let  $\boldsymbol{\Theta}_m$  denote the set of all minimizers at step  $m$ ; that is,  $\boldsymbol{\theta}_{m+1} \in \boldsymbol{\Theta}_{m+1}$ . Sometimes it may be difficult to perform M-step numerically; in this case, we can define a generalized DM algorithm (referred to as the G-DM algorithm) as follows: Let  $M : \boldsymbol{\Theta}_m \rightarrow \boldsymbol{\Theta}_{m+1}$  be a point to set map: then the G-DM algorithm is an iterative scheme such that

$$Q^{(G)}(\boldsymbol{\theta}'|\boldsymbol{\theta}_m) \leq Q^{(G)}(\boldsymbol{\theta}_m|\boldsymbol{\theta}_m) \quad \text{for all } \boldsymbol{\theta}' \in M(\boldsymbol{\theta}_m).$$

We notice here that for any G-DM sequence  $\{\boldsymbol{\theta}_m\}$ ,  $D^{(G)}(\boldsymbol{\theta}_{m+1}) \leq D^{(G)}(\boldsymbol{\theta}_m)$  and DM algorithm is a special case of G-DM algorithm. We emphasize here that by choosing

different divergences, we obtain many existing algorithms, including the EM, HMIX, and HELMIX algorithms.

## E: Special Cases of DM Algorithms for Various Choices of $G(\cdot)$

### Special case 1: EM Algorithm

Let  $G(\delta) = (\delta + 1) \log(\delta + 1)$ , which corresponds to the Kullback-Leibler divergence, we get the objective function obtained from the E-step in the EM algorithm. Specifically, since MLE can be obtained by minimizing the Kullback-Leibler divergence, it follows that

$$\begin{aligned}
 Q_{\text{KL}}(\boldsymbol{\theta}'|\boldsymbol{\theta}) &= \int_{\mathcal{Y}} \int_{\mathcal{Z}} \left( \frac{g(y)w(z|y; \boldsymbol{\theta})}{f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}')} \right) \log \left( \frac{g(y)w(z|y; \boldsymbol{\theta})}{f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}')} \right) f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}') dz dy \\
 &= \int_{\mathcal{Y}} \int_{\mathcal{Z}} g(y)w(z|y; \boldsymbol{\theta}) \log \left( \frac{g(y)w(z|y; \boldsymbol{\theta})}{f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}')} \right) dz dy \\
 &= \int_{\mathcal{Y}} \int_{\mathcal{Z}} g(y)w(z|y; \boldsymbol{\theta}) \log (g(y)w(z|y; \boldsymbol{\theta})) dz dy \\
 &\quad - \int_{\mathcal{Y}} \int_{\mathcal{Z}} g(y)w(z|y; \boldsymbol{\theta}) \log (f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}')) dz dy \\
 &\equiv \int_{\mathcal{Y}} \int_{\mathcal{Z}} g(y)w(z|y; \boldsymbol{\theta}) \log (g(y)w(z|y; \boldsymbol{\theta})) dz dy - \tilde{Q}(\boldsymbol{\theta}'|\boldsymbol{\theta}), \tag{9.15}
 \end{aligned}$$

where

$$\tilde{Q}(\boldsymbol{\theta}'|\boldsymbol{\theta}) = \int_{\mathcal{Y}} \left( \int_{\mathcal{Z}} w(z|y; \boldsymbol{\theta}) \log (f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}')) dz \right) g(y) dy.$$

The first term in (9.15) does not involve  $\boldsymbol{\theta}'$  and hence can be omitted in terms of optimization. The second term  $\tilde{Q}(\boldsymbol{\theta}'|\boldsymbol{\theta})$ , in the current literature, is referred to as the population level EM objective function, see [Balakrishnan et al. \(2017\)](#), [Dwivedi et al. \(2018\)](#). If we estimate  $g(\cdot)$  through empirical measure, then  $\tilde{Q}(\boldsymbol{\theta}'|\boldsymbol{\theta})$  will reduce to the sample level EM objective function, and is given by

$$\tilde{Q}_n(\boldsymbol{\theta}'|\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \int_{\mathcal{Z}} c(z|y_i; \boldsymbol{\theta}) \log (f(y_i; \boldsymbol{\theta}')c(z|y_i; \boldsymbol{\theta}')) dz.$$

### Special case 2: HMIX Algorithm

Next consider another example. By taking  $G(\delta) = 2[(\delta + 1)^{1/2} - 1]^2$ , which corresponds to the Hellinger distance divergence, we get the population level HMIX objective function

$$\begin{aligned} Q_{\text{HD}}(\boldsymbol{\theta}'|\boldsymbol{\theta}) &= 2 \int_{\mathcal{Y}} \int_{\mathcal{Z}} \left( \left( \frac{g(y)w(z|y; \boldsymbol{\theta})}{f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}')} \right)^{\frac{1}{2}} - 1 \right)^2 f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}') dz dy \\ &= 4 - 4 \int_{\mathcal{Y}} \int_{\mathcal{Z}} [g(y)w(z|y; \boldsymbol{\theta})p(y, z; \boldsymbol{\theta}')]^{\frac{1}{2}} dz dy, \\ &\equiv 4 - 4\mathcal{A}(\boldsymbol{\theta}'|\boldsymbol{\theta}). \end{aligned}$$

On the other hand, the sample level HMIX objective function is

$$\mathcal{A}_n(\boldsymbol{\theta}'|\boldsymbol{\theta}) = \int_{\mathcal{Y}} \int_{\mathcal{Z}} [g_n(y)w(z|y; \boldsymbol{\theta})f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}')]^{\frac{1}{2}} dz dy.$$

### Special case 3: Algorithm from the Negative Exponential Divergence

Next consider another example. By taking  $G(\delta) = [e^{-\delta} - 1 + \delta]$ , which corresponds to the negative exponential divergence, we get the population level objective function based on Negative exponential divergence. Specifically, it is given by

$$\begin{aligned} Q_{\text{NED}}(\boldsymbol{\theta}'|\boldsymbol{\theta}) &= \exp(1) \int_{\mathcal{Y}} \int_{\mathcal{Z}} \exp \left( -\frac{g(y)w(z|y; \boldsymbol{\theta})}{f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}')} \right) f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}') dz dy - 1 \\ &\equiv \exp(1)\text{NED}(\boldsymbol{\theta}'|\boldsymbol{\theta}) - 1. \end{aligned}$$

So  $Q_{\text{NED}}(\boldsymbol{\theta}'|\boldsymbol{\theta})$  is the population level objective function generated by Negative exponential divergence. The sample level  $Q_{\text{NED}_n}(\boldsymbol{\theta}'|\boldsymbol{\theta})$  is defined analogously, where  $g(\cdot)$  is replaced by  $g_n(\cdot)$ .

### Special case 4: Algorithm from the Variant Negative Exponential Divergence

We introduce a divergence similar to Negative exponential divergence called vNED. Specifically,  $G(\cdot)$  is given by

$$G_{\text{vNED}}(\delta) = \exp \left( -\frac{1}{1 + \delta} + 1 \right) (1 + \delta) - (2\delta + 1).$$

Then the corresponding population DM objective function  $Q(\cdot|\cdot)$  is as follows:

$$\begin{aligned} Q_{\text{vNED}}(\boldsymbol{\theta}'|\boldsymbol{\theta}) &= \exp(1) \int_{\mathcal{Y}} \int_{\mathcal{Z}} \exp\left(-\frac{f(y;\boldsymbol{\theta}')w(Z|Y;\boldsymbol{\theta}')}{g(y)w(z|y;\boldsymbol{\theta})}\right) g(y)w(z|y;\boldsymbol{\theta}) dz dy - 1 \\ &\equiv \exp(1)\text{vNED}(\boldsymbol{\theta}'|\boldsymbol{\theta}) - 1. \end{aligned}$$

### Special case 5: Algorithm from Blended Weighted Hellinger Distance Divergence

Basu & Lindsay (1994) proposed the blended weighted Hellinger distance (BWHD) divergence as

$$G_{\text{BWHD}}(\delta) = \frac{1}{2} \frac{\delta^2}{[\tau(\delta + 1)^{\frac{1}{2}} + 1 - \tau]^2}, \quad \text{where } \tau \in [0, 1].$$

Then the BWHD objective function from DM algorithm is given as

$$Q_{\text{BWHD}}(\boldsymbol{\theta}'|\boldsymbol{\theta}) = \int_{\mathcal{Y}} \int_{\mathcal{Z}} \left( \frac{g(y)w(z|y;\boldsymbol{\theta}) - f(y;\boldsymbol{\theta}')w(Z|Y;\boldsymbol{\theta}')}{\tau(g(y)w(z|y;\boldsymbol{\theta}))^{\frac{1}{2}} + (1 - \tau)(f(y;\boldsymbol{\theta}')w(Z|Y;\boldsymbol{\theta}'))^{\frac{1}{2}}} \right)^2 dz dy.$$

Apart from the examples described as above, one can get other DM algorithm objective functions following similar ideas (e.g., blended weighted Negative exponential divergence, blended weight chi-square divergence, etc), hence we omit here.

### Special case 6: Algorithm from the Cressie-Read Family

Next we consider an important subfamily of general divergence, the Cressie-Read (CR) family (see Cressie & Read (1984), Read & Cressie (1988)). For CR family,  $G(\cdot)$  is given by

$$G_{\text{CR}}(\delta) = \frac{(\delta + 1)^{\alpha+1} - 1}{\alpha(\alpha + 1)}, \quad \text{where } \alpha \in \mathbb{R}.$$

When  $\alpha = -1$ , it corresponds to the Kullback-Leibler divergence; when  $\alpha = -\frac{1}{2}$ , it corresponds to Hellinger distance divergence; when  $\alpha = 1$ , it corresponds to Pearson's  $\chi^2$  and when  $\alpha = -2$ , it corresponds to Neyman's  $\chi^2$ . Note that when  $\alpha = -1$ , the divergence is defined by continuity. The DM algorithm objective function for CR family is

$$Q_{\text{CR}}(\boldsymbol{\theta}'|\boldsymbol{\theta}) = \frac{1}{\alpha(1 + \alpha)} \int_{\mathcal{Y}} \int_{\mathcal{Z}} \left[ \left( \frac{g(y)w(z|y;\boldsymbol{\theta})}{f(y;\boldsymbol{\theta}')w(Z|Y;\boldsymbol{\theta}')} \right)^{\alpha+1} - 1 \right] f(y;\boldsymbol{\theta}')w(Z|Y;\boldsymbol{\theta}') dz dy.$$



## Special case 7: Algorithm from Power Divergence Family

Another important subfamily of general divergence is the power divergence (PD) family (see [Patra et al. \(2013\)](#)). For PD family,  $G(\cdot)$  is given by

$$G_{\text{PD}}(\delta) = \frac{1}{\alpha(1+\alpha)} \left[ (\delta+1)^{1+\alpha} - (\delta+1) \right] - \frac{\delta}{1+\alpha} \quad \text{where } \alpha \in \mathbb{R}.$$

The DM algorithm objective function for PD family is

$$Q_{\text{PD}}(\boldsymbol{\theta}'|\boldsymbol{\theta}) = \int_{\mathcal{Y}} \int_{\mathcal{Z}} \left\{ \frac{1}{\alpha(1+\alpha)} \left[ \left( \frac{g(y)w(z|y; \boldsymbol{\theta})}{f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}')} \right)^{1+\alpha} - \left( \frac{g(y)w(z|y; \boldsymbol{\theta})}{f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}')} \right) \right] \right. \\ \left. + \frac{1}{1+\alpha} \left[ 1 - \frac{g(y)w(z|y; \boldsymbol{\theta})}{f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}')} \right] \right\} f(y; \boldsymbol{\theta}')w(Z|Y; \boldsymbol{\theta}') dz dy.$$

Table 2: Representative special cases of DM algorithm for different  $G(\cdot)$ .

Algorithm	Generator $G(\delta)$	Population DM objective (up to constants)
KL (EM)	$(\delta+1) \log(\delta+1)$	$\iint g(y)w(z y; \theta) \log \frac{g(y)w(z y; \theta)}{f(y; \theta')w(z y; \theta')} dz dy$
Hellinger (HMIX)	$2((\delta+1)^{1/2} - 1)^2$	$4 - 4 \iint [g(y)w(z y; \theta) p(y, z; \theta')]^{1/2} dz dy$
NED	$e^{-\delta} - 1 + \delta$	$e \iint \exp\left(-\frac{g(y)w(z y; \theta)}{f(y; \theta')w(z y; \theta')}\right) f(y; \theta')w(z y; \theta') dz dy - 1$
vNED	$\exp\left(-\frac{1}{1+\delta} + 1\right) (1+\delta) - (2\delta+1)$	$e \iint \exp\left(-\frac{f(y; \theta')w(z y; \theta')}{g(y)w(z y; \theta)}\right) g(y)w(z y; \theta) dz dy - 1$
BWHD	$\frac{1}{2} \frac{\delta^2}{[\tau(\delta+1)^{1/2} + 1 - \tau]^2}$	$\iint \left( \frac{g(y)w(z y; \theta) - f(y; \theta')w(z y; \theta')}{\tau \sqrt{g(y)w(z y; \theta)} + (1-\tau) \sqrt{f(y; \theta')w(z y; \theta')}} \right)^2 dz dy$
CR family	$\frac{(\delta+1)^{\alpha+1} - 1}{\alpha(\alpha+1)}$	$\frac{1}{\alpha(1+\alpha)} \iint \left[ \left( \frac{g(y)w}{f(y; \theta')w'} \right)^{\alpha+1} - 1 \right] f(y; \theta')w' dz dy$
PD family	$\frac{1}{\alpha(1+\alpha)} \left[ (\delta+1)^{1+\alpha} - (\delta+1) \right] - \frac{\delta}{1+\alpha}$	$\iint \left\{ \frac{\left( \frac{g(y)w}{f(y; \theta')w'} \right)^{1+\alpha} - \frac{g(y)w}{f(y; \theta')w'}}{\alpha(1+\alpha)} + \frac{1}{1+\alpha} \left( 1 - \frac{g(y)w}{f(y; \theta')w'} \right) \right\} f(y; \theta')w' dz dy$

## F: DM Algorithm for Various Choices of $G(\cdot)$ with Application to FMM

In this subsection, we specialize the DM framework to finite mixture models (FMMs) described in Subsection 2.1. For each choice of divergence generator  $G(\cdot)$ , the corresponding population-level DM objective can be expressed in closed form. Recall that,

$$Q_G(\boldsymbol{\theta}' | \boldsymbol{\theta}) = \sum_{k=1}^K \int \pi'_k h(y; \boldsymbol{\phi}'_k) G(\tau_k(y)) dy \quad \text{and} \quad \tau_k(y) := \frac{g(y) w_k(y; \boldsymbol{\theta})}{\pi'_k h(y; \boldsymbol{\phi}'_k)}. \quad (9.16)$$

We list below the objectives for several important divergences; detailed update formulas for  $\boldsymbol{\theta}'$  will be given in Section 2.1.

For EM algorithm,  $\tilde{Q}(\cdot | \cdot)$  for finite mixture model is given by

$$Q_{\text{EM}}(\boldsymbol{\theta}' | \boldsymbol{\theta}) := \tilde{Q}(\boldsymbol{\theta}' | \boldsymbol{\theta}) = \sum_{k=1}^K \int_{\mathcal{Y}} g(y) w_k(y; \boldsymbol{\theta}) \log(\pi'_k h(y; \boldsymbol{\phi}'_k)) dy.$$

The HMX objective function  $\mathcal{A}(\cdot | \cdot)$  is given by

$$Q_{\text{HD}}(\boldsymbol{\theta}' | \boldsymbol{\theta}) := \mathcal{A}(\boldsymbol{\theta}' | \boldsymbol{\theta}) = \sum_{k=1}^K \int_{\mathcal{Y}} [g(y) w_k(y; \boldsymbol{\theta}) \pi'_k h(y; \boldsymbol{\phi}'_k)]^{\frac{1}{2}} dy.$$

The NED objective function for finite mixture model. It is given by

$$Q_{\text{NED}}(\boldsymbol{\theta}' | \boldsymbol{\theta}) = \sum_{k=1}^K \int_{\mathcal{Y}} \exp\left(-\frac{g(y) w_k(y; \boldsymbol{\theta})}{\pi'_k h(y; \boldsymbol{\phi}'_k)}\right) \pi'_k h(y; \boldsymbol{\phi}'_k) dy.$$

Next consider the vNED objective function for finite mixture model. It is given by

$$Q_{\text{vNED}}(\boldsymbol{\theta}' | \boldsymbol{\theta}) = \sum_{k=1}^K \int_{\mathcal{Y}} \exp\left(-\frac{\pi'_k h(y; \boldsymbol{\phi}'_k)}{g(y) w_k(y; \boldsymbol{\theta})}\right) g(y) w_k(y; \boldsymbol{\theta}) dy.$$

Next consider the Blended weight Hellinger distance (BWHD) objective function for finite mixture model, where for  $0 < \tau < 1$ ,  $Q_{\text{BWHD}}(\cdot | \cdot)$  is given by

$$Q_{\text{BWHD}}(\boldsymbol{\theta}' | \boldsymbol{\theta}) = \sum_{k=1}^K \int_{\mathcal{Y}} \left( \frac{g(y) w_k(y; \boldsymbol{\theta}) - \pi'_k h(y; \boldsymbol{\phi}'_k)}{\tau (g(y) w_k(y; \boldsymbol{\theta}))^{\frac{1}{2}} + (1 - \tau) (\pi'_k h(y; \boldsymbol{\phi}'_k))^{\frac{1}{2}}} \right)^2 dy.$$

Next consider the CR family, where  $\alpha \in \mathbb{R} - -1, 0$ .  $Q_{\text{CR}}(\cdot|\cdot)$  for finite mixture model is given by

$$Q_{\text{CR}}(\boldsymbol{\theta}'|\boldsymbol{\theta}) = \frac{1}{\alpha(1+\alpha)} \sum_{k=1}^K \int_{\mathcal{Y}} \left[ \left( \frac{g(y)w_k(y;\boldsymbol{\theta})}{\pi'_k h(y;\boldsymbol{\phi}'_k)} \right)^{\alpha+1} - 1 \right] \pi'_k h(y;\boldsymbol{\phi}'_k) dy.$$

Next consider the PD family (an equivalent alternative parametrization scaled to satisfy  $A'(0) = 1$ ).  $Q_{\text{PD}}(\cdot|\cdot)$  for finite mixture model is given by

$$\begin{aligned} Q_{\text{PD}}(\boldsymbol{\theta}'|\boldsymbol{\theta}) = \sum_{k=1}^K \int_{\mathcal{Y}} \left\{ \frac{1}{\alpha(1+\alpha)} \left[ \left( \frac{g(y)w_k(y;\boldsymbol{\theta})}{\pi'_k h(y;\boldsymbol{\phi}'_k)} \right)^{1+\alpha} - \left( \frac{g(y)w_k(y;\boldsymbol{\theta})}{\pi'_k h(y;\boldsymbol{\phi}'_k)} \right) \right] \right. \\ \left. + \frac{1}{1+\alpha} \left[ 1 - \frac{g(y)w_k(y;\boldsymbol{\theta})}{\pi'_k h(y;\boldsymbol{\phi}'_k)} \right] \right\} \pi'_k h(y;\boldsymbol{\phi}'_k) dy. \end{aligned}$$

We will provide specific algorithms for updating  $\boldsymbol{\theta}'$  in finite mixture models in Section [2.1](#).

# G: Discrete Kernels and K-means Clustering Algorithm

## Discrete Kernels

In general, let  $Y_1, \dots, Y_n$  be i.i.d. random variables with an unknown probability mass function (p.m.f.)  $f$  on  $\mathbb{Z}$ . A discrete kernel estimator of  $f$  can be defined as

$$g_n(y) = \frac{1}{n} \sum_{i=1}^n \mathcal{K}_{y,c}(Y_i), \quad y \in \mathbb{Z}, \quad (9.17)$$

where  $\mathcal{K}_{y,c}(\cdot)$  is a discrete kernel p.m.f. centered at  $y$ , and  $c_n$  is a sequence of smoothing bandwidths. More rigorously, following the notation from [Kokonendji & Kiesse \(2011\)](#), the discrete associated kernel is defined as follows.

**Definition 5.** Let  $\mathbb{T}$  be the discrete support of the p.m.f.  $f$  to be estimated,  $y$  a fixed target in  $\mathbb{T}$ , and  $c > 0$  a bandwidth. A p.m.f.  $\mathcal{K}_{y,c}(\cdot)$  on support  $\mathbb{S}_y$  (not depending on  $c$ ) is said to be an associated kernel if it satisfies:

$$y \in \mathbb{T}, \quad \lim_{c \rightarrow 0} \mathbf{E}[X_{y,c}] = y, \quad \text{and} \quad \lim_{c \rightarrow 0} \mathbf{Var}[X_{y,c}] = 0,$$

where  $X_{y,c}$  is a random variable with p.m.f.  $\mathcal{K}_{y,c}(\cdot)$ .

### 1. The empirical kernel:

$$\mathcal{K}_{y,c}(x) = I_{x=y}, \quad x \in \mathbb{T}, \quad c \geq 0, \quad (9.18)$$

where  $I_A$  is the indicator function.

### 2. Discrete triangular kernel: ([Kokonendji et al. \(2007\)](#))

For support  $\mathbb{T}$  (bounded or unbounded), bandwidth  $c > 0$ , and integer  $a > 0$ , define

$$\mathcal{K}_{y,c}(x) = \frac{(a+1)^c - |x-y|^c}{P(a,c)}, \quad x \in \{y-a, \dots, y+a\},$$

with normalizing constant  $P(a,c) = (2a+1)(a+1)^c - 2 \sum_{k=0}^a k^c$ .

3. **Poisson kernel:** For  $y \in \mathbb{N}, c > 0$ ,

$$\mathcal{K}_{y,c}(x) = \frac{(y+c)^x e^{-(y+c)}}{x!}, \quad x \in \mathbb{N}.$$

4. **Binomial kernel:** For  $y \in \mathbb{N}, c \in (0, 1]$ ,

$$\mathcal{K}_{y,c}(x) = \binom{y+1}{x} \left( \frac{y+c}{y+1} \right)^x \left( \frac{1-c}{y+1} \right)^{y+1-x}, \quad x \in \{0, 1, \dots, y+1\}.$$

5. **Negative binomial kernel:** For  $y \in \mathbb{N}, c > 0$ ,

$$\mathcal{K}_{y,c}(x) = \binom{y+x}{x} \left( \frac{y+c}{2y+1+c} \right)^x \left( \frac{y+1}{2y+1+c} \right)^{y+1}, \quad x \in \mathbb{N}.$$

In order to measure the performance of different discrete kernel estimators, we use the practical criterion integrated squared error (ISE) given by

$$\text{ISE} = \sum_{y \in \mathbb{T}} (g_n(y) - f(y; \boldsymbol{\theta}))^2,$$

## K-means clustering algorithm

Given samples  $y_1, \dots, y_n$  and fix the number of clusters  $K$ , place initial centroids  $c_1^{(0)}, \dots, c_k^{(0)}$  at random locations.

---

### Algorithm 2 The K-means Clustering Algorithm

---

1. Set  $m = 0$ ,

**repeat**

2. For each point  $y_i$ :

(i). Find nearest centroid, i.e.  $k^* = \underset{k}{\operatorname{argmin}} L(y_i, c_k^{(m)})$ , where  $L(\cdot)$  represents some distance metric.

(ii). Assign  $y_i$  to cluster  $k^*$ .

3. Find new centroids: new centroid  $c_k^{(m+1)} = \text{mean of all points } y_i \text{ assigned to cluster } k^*$ .

4. Set  $m = m + 1$ .

**until** None of the cluster assignments change.

---

**Remark 6.** In step 2 (i), the examples of distance metrics are Euclidean distance or  $L_1$  distance. Besides, the choice of distance metric depends on the model.

**Remark 7.** *This algorithm works well without any outliers or with few outliers. If there exists a considerable percentage of outliers in the sample, then the K-means clustering algorithm is more likely to treat the outlier group as a new cluster, and the initial points will be more likely to be away from the real values.*

## H: Additional Algorithms

If we choose the Hellinger distance divergence, then the algorithm (referred to as the DM-HMIX algorithm or simply HMIX algorithm) is provided in Algorithm 3.

Algorithm 4, we use vNED divergence. Specifically, the update for  $\pi'_k$  therefore can be written as

$$\pi'_k = \frac{\text{vNED}_k(\boldsymbol{\theta}'|\boldsymbol{\theta})}{\sum_{l=1}^K \text{vNED}_l(\boldsymbol{\theta}'|\boldsymbol{\theta})}, \quad \text{where} \quad \text{vNED}_k(\boldsymbol{\theta}'|\boldsymbol{\theta}) = \int_{\mathcal{Y}} \exp\left(-\frac{\pi_k h(y; \boldsymbol{\phi}'_k)}{g(y)w_k(y; \boldsymbol{\theta})}\right) \pi_k h(y; \boldsymbol{\phi}'_k) dy.$$

For Algorithm 5, by using recurrence relation in Poisson distribution, we can derive an algorithm similar to the HELMIX algorithm based on vNED. We refer this algorithm as the DM-NELMIX algorithm.

Let  $h(y; \lambda) = e^{-\lambda} \lambda^y / y!$ . For given initial values  $\boldsymbol{\theta}^{(0)} = (\pi_1^{(0)}, \dots, \pi_K^{(0)}, \lambda_1^{(0)}, \dots, \lambda_K^{(0)})$ .

---

**Algorithm 3** The DM-HMIX Algorithm

---

Set  $m = 0$ .

**repeat**

1. Compute

$$HD_n(\boldsymbol{\theta}^{(m)}) = \sum_{k=1}^K \int_{\mathcal{Y}} \left( g_n(y) w_k(y; \boldsymbol{\theta}^{(m)}) \pi_k^{(m)} h(y; \boldsymbol{\phi}_k^{(m)}) \right)^{\frac{1}{2}} dy, \quad \text{where} \quad w_k(y; \boldsymbol{\theta}) = \frac{\pi_k h(y; \boldsymbol{\phi}_k)}{\sum_{l=1}^K \pi_l h(y; \boldsymbol{\phi}_l)}.$$

2. Update  $\boldsymbol{\phi}_k^{(m+1)}$ :

$$\boldsymbol{\phi}_k^{(m+1)} = \underset{\boldsymbol{\phi}'_k \in \boldsymbol{\Theta}}{\operatorname{argmax}} \int_{\mathcal{Y}} \left( g_n(y) w_k(y; \boldsymbol{\theta}^{(m)}) \pi_k^{(m)} h(y; \boldsymbol{\phi}'_k) \right)^{\frac{1}{2}} dy.$$

3. Update  $\pi_k^{(m+1)}$ :

$$\pi_k^{(m+1)} = \frac{HD_{n,k}^2(\boldsymbol{\phi}^{(m+1)} | \boldsymbol{\theta}^{(m)})}{\sum_{l=1}^K HD_{n,l}^2(\boldsymbol{\phi}^{(m+1)} | \boldsymbol{\theta}^{(m)})}, \quad (\text{HMIX type update})$$

or

$$\pi_k^{(m+1)} = \frac{\sqrt{\pi_k^{(m)}} HD_{n,k}(\boldsymbol{\phi}^{(m+1)} | \boldsymbol{\theta}^{(m)})}{\sum_{l=1}^K \sqrt{\pi_l^{(m)}} HD_{n,l}(\boldsymbol{\phi}^{(m+1)} | \boldsymbol{\theta}^{(m)})}, \quad (\text{DM-Mix type update (HDMIX)})$$

where

$$HD_{n,k}(\boldsymbol{\phi}^{(m+1)} | \boldsymbol{\theta}^{(m)}) = \int_{\mathcal{Y}} \left( g_n(y) w_k(y; \boldsymbol{\theta}^{(m)}) h(y; \boldsymbol{\phi}_k^{(m+1)}) \right)^{\frac{1}{2}} dy.$$

4. Update  $w_k(y; \boldsymbol{\theta}^{(m+1)})$ , compute  $HD(\boldsymbol{\theta}^{(m+1)})$  and the difference

$$\epsilon_{m+1} = |HD_n(\boldsymbol{\theta}^{(m+1)}) - HD_n(\boldsymbol{\theta}^{(m)})|.$$

5. Set  $m = m+1$ .

**until**  $\epsilon_m < \text{threshold}$ .

---



---

**Algorithm 4** The DM-vNEDMIX Algorithm

---

Set  $m = 0$ .

**repeat**

1. Compute

$$\text{vNED}_n(\boldsymbol{\theta}^{(m)}) = \sum_{k=1}^K \int_{\mathcal{Y}} \exp\left(-\frac{\pi_k^{(m)} h(y; \boldsymbol{\phi}_k^{(m)})}{g_n(y) w_k(y; \boldsymbol{\theta}^{(m)})}\right) g_n(y) w_k(y; \boldsymbol{\theta}^{(m)}) dy, \quad \text{where} \quad w_k(y; \boldsymbol{\theta}) = \frac{\pi_k h(y; \boldsymbol{\phi}_k)}{\sum_{l=1}^K \pi_l h(y; \boldsymbol{\phi}_l)}.$$

2. Update  $\boldsymbol{\phi}_k^{(m+1)}$ :

$$\boldsymbol{\phi}_k^{(m+1)} = \underset{\boldsymbol{\phi}'_k \in \Theta}{\operatorname{argmin}} \int_{\mathcal{Y}} \exp\left(-\frac{\pi_k^{(m)} h(y; \boldsymbol{\phi}'_k)}{g_n(y) w_k(y; \boldsymbol{\theta}^{(m)})}\right) g_n(y) w_k(y; \boldsymbol{\theta}^{(m)}) dy.$$

3. Update  $\pi_k^{(m+1)}$ :

$$\pi_k^{(m+1)} = \frac{\text{vNED}_{n,k}(\boldsymbol{\phi}_k^{(m+1)} | \boldsymbol{\theta}^{(m)})}{\sum_{l=1}^K \text{vNED}_{n,l}(\boldsymbol{\phi}_l^{(m+1)} | \boldsymbol{\theta}^{(m)})}, \quad \text{where}$$

$$\text{vNED}_{n,k}(\boldsymbol{\phi}_k^{(m+1)} | \boldsymbol{\theta}^{(m)}) = \int_{\mathcal{Y}} \exp\left(-\frac{\pi_k^{(m)} h(y; \boldsymbol{\phi}_k^{(m+1)})}{g_n(y) w_k(y; \boldsymbol{\theta}^{(m)})}\right) \pi_k^{(m)} h(y; \boldsymbol{\phi}_k^{(m+1)}) dy.$$

4. Update  $w_k(y; \boldsymbol{\theta}^{(m+1)})$ , compute  $\text{vNED}_n(\boldsymbol{\theta}^{(m+1)})$  and the difference

$$\epsilon_{m+1} = |\text{vNED}_n(\boldsymbol{\theta}^{(m+1)}) - \text{vNED}_n(\boldsymbol{\theta}^{(m)})|.$$

5. Set  $m = m+1$ .

**until**  $\epsilon_m < \text{threshold}$ .

---

---

**Algorithm 5** The DM-NELMIX Algorithm

---

Set  $m = 0$ .

**repeat**

1. Compute

$$\text{vNED}_n^*(\boldsymbol{\theta}^{(m)}) = \sum_{k=1}^K \sum_{y \in \mathcal{Y}} \exp \left( -\frac{\pi_k^{(m)} h(y; \boldsymbol{\phi}_k^{(m)})}{g_n(y) w_k(y; \boldsymbol{\theta}^{(m)})} \right) g_n(y) w_k(y; \boldsymbol{\theta}^{(m)}),$$

where

$$w_k(y; \boldsymbol{\theta}^{(m)}) = \frac{\pi_k^{(m)} h(y; \lambda_k^{(m)})}{\sum_{l=1}^K \pi_l^{(m)} h(y; \lambda_l^{(m)})}.$$

2. Find  $\lambda_k^{(m+1)}$ :

$$\lambda_k^{(m+1)} = \sum_{y \in \mathcal{Y}} y w_{yk}^{(m)} / \sum_{y \in \mathcal{Y}} w_{yk}^{(m)},$$

where

$$w_{yk}^{(m)} = \exp \left( -\frac{\pi_k^{(m)} h(y; \lambda_k^{(m)})}{g_n(y) w_k(y; \boldsymbol{\theta}^{(m)})} \right) \pi_k^{(m)} h(y; \lambda_k^{(m)}).$$

3. Update  $\pi_k^{(m+1)}$ :

$$\pi_k^{(m+1)} = \text{vNED}_{n,k} \left( \pi_k^{(m)}, \lambda_k^{(m+1)} \right) / \sum_{l=1}^K \text{vNED}_{n,l} \left( \pi_l^{(m)}, \lambda_l^{(m+1)} \right),$$

where

$$\text{vNED}_{n,k} \left( \pi_k^{(m)}, \lambda_k^{(m+1)} \right) = \sum_{y \in \mathcal{Y}} \exp \left( -\frac{\pi_k^{(m)} h(y; \lambda_k^{(m+1)})}{g_n(y) w_k(y; \boldsymbol{\theta}^{(m+1)})} \right) \pi_k^{(m)} h(y; \lambda_k^{(m+1)}).$$

4. Update  $w_k(y; \boldsymbol{\theta}^{(m+1)})$ , compute  $\text{vNED}_n^*(\boldsymbol{\theta}^{(m+1)})$  and the difference

$$\epsilon_{m+1} = |\text{vNED}_n^*(\boldsymbol{\theta}^{(m+1)}) - \text{vNED}_n^*(\boldsymbol{\theta}^{(m)})|.$$

5. Set  $m = m+1$ .

**until**  $\epsilon_m < \text{threshold}$ .

---

# I: Additional Simulation Results

## Discrete Models

In this section, we assume  $K$  is known. For discrete data we estimate the nonparametric distribution  $g_n(\cdot)$  using discrete kernels. We implemented both cross-validation and moment-based bandwidth selection; the results were similar, while the moment-based method was significantly faster, so we report only the moment-based results.

### Poisson Mixture Model

Suppose  $f(\cdot; \boldsymbol{\theta}) = 0.4 \text{Pois}(0.5) + 0.6 \text{Pois}(10)$ . Figure 5 shows the average estimates of  $\pi_1$  across kernels and sample sizes  $n \in \{20, 50, 100, 200\}$ ; the dashed line indicates the true value 0.4.

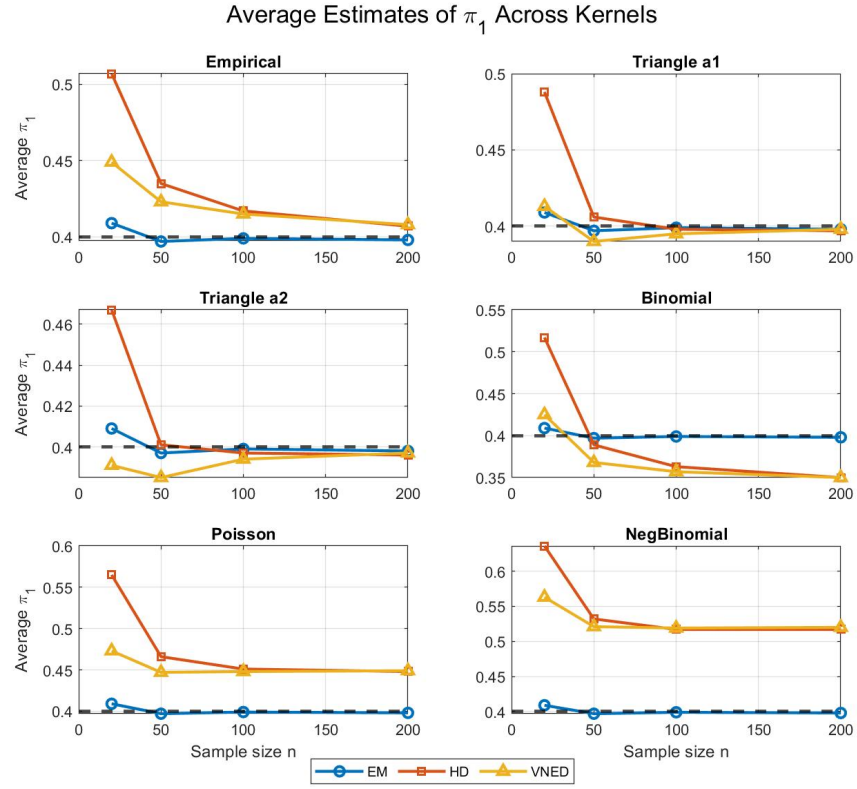
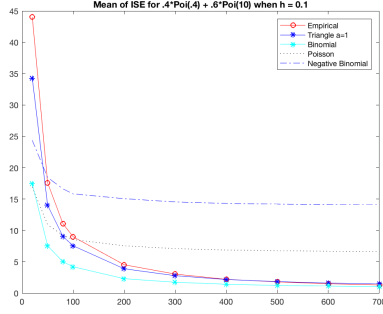
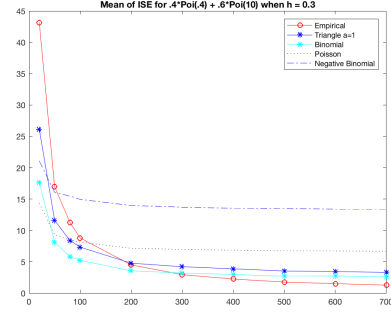


Figure 5: Average estimates of  $\pi_1$  across kernels for sample sizes  $n = 20, 50, 100, 200$ . Solid lines represent EM, HD, and vNED methods; dashed horizontal line indicates the true value 0.4.

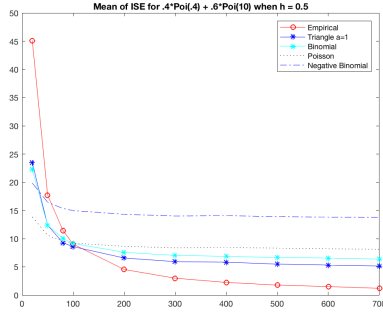
Figure 6 reports the mean ISE for the two-component Poisson mixture as  $n$  increases at several bandwidths  $c$ . The results indicate: (i) with small  $n$ , discrete kernels are beneficial; (ii) bandwidth choice is important; and (iii) the triangle kernel with  $a = 1$  performs well overall for this model.



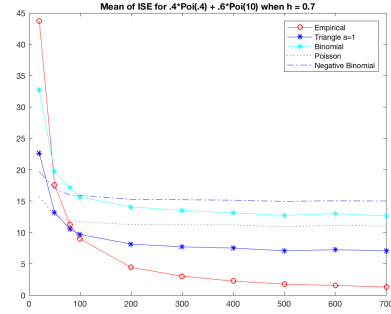
(a)  $c = 0.1$



(b)  $c = 0.3$



(c)  $c = 0.5$



(d)  $c = 0.7$

Figure 6: Mean of ISE for Model  $0.4Poi(0.5) + 0.6Poi(10)$  as  $n$  changes

In this section for discrete models, all simulation results use the empirical p.m.f.

$$g_n(y) = \frac{n_y}{n}, \quad n_y = \sum_{i=1}^n \mathbf{1}\{Y_i = y\}.$$

Table 3: Estimates Using Different Kernels for  $0.4Poi(0.5) + 0.6Poi(10)$

Kernel		Empirical			Triangle $a = 1$			Triangle $a = 2$		
$n$		$\hat{\pi}_1$	$\hat{\lambda}_1$	$\hat{\lambda}_2$	$\hat{\pi}_1$	$\hat{\lambda}_1$	$\hat{\lambda}_2$	$\hat{\pi}_1$	$\hat{\lambda}_1$	$\hat{\lambda}_2$
20	EM	0.409 (0.113)	0.544 (0.282)	9.904 (1.007)	0.409 (0.113)	0.544 (0.282)	9.904 (1.007)	0.409 (0.113)	0.544 (0.282)	9.904 (1.007)
	HD	0.507 (0.148)	0.410 (0.249)	9.605 (1.087)	0.488 (0.176)	0.487 (0.270)	9.604 (1.145)	0.467 (0.183)	0.498 (0.300)	9.578 (1.144)
	vNED	0.449 (0.135)	0.456 (0.268)	9.556 (1.120)	0.413 (0.155)	0.589 (0.292)	9.639 (1.146)	0.391 (0.157)	0.631 (0.349)	9.638 (1.135)
50	EM	0.397 (0.074)	0.490 (0.173)	10.01 (0.548)	0.397 (0.074)	0.490 (0.173)	10.01 (0.548)	0.397 (0.074)	0.490 (0.173)	10.01 (0.548)
	HD	0.435 (0.083)	0.442 (0.171)	9.788 (0.598)	0.406 (0.089)	0.570 (0.211)	9.772 (0.599)	0.401 (0.091)	0.572 (0.221)	9.768 (0.599)
	vNED	0.423 (0.079)	0.463 (0.175)	9.866 (0.598)	0.390 (0.083)	0.621 (0.216)	9.873 (0.582)	0.385 (0.085)	0.620 (0.227)	9.873 (0.583)
100	EM	0.399 (0.050)	0.502 (0.113)	10.02 (0.415)	0.399 (0.050)	0.502 (0.113)	10.02 (0.415)	0.399 (0.050)	0.502 (0.113)	10.02 (0.415)
	HD	0.417 (0.053)	0.478 (0.114)	9.849 (0.462)	0.398 (0.056)	0.584 (0.152)	9.841 (0.473)	0.397 (0.057)	0.584 (0.152)	9.835 (0.474)
	vNED	0.415 (0.052)	0.482 (0.110)	9.917 (0.463)	0.395 (0.054)	0.600 (0.147)	9.922 (0.464)	0.394 (0.055)	0.598 (0.145)	9.915 (0.464)
200	EM	0.398 (0.032)	0.502 (0.073)	10.07 (0.296)	0.398 (0.032)	0.502 (0.073)	10.07 (0.296)	0.398 (0.032)	0.502 (0.073)	10.07 (0.296)
	HD	0.407 (0.033)	0.489 (0.074)	9.879 (0.305)	0.397 (0.035)	0.551 (0.085)	9.875 (0.307)	0.396 (0.035)	0.553 (0.087)	9.869 (0.307)
	vNED	0.408 (0.033)	0.488 (0.070)	9.941 (0.299)	0.398 (0.035)	0.553 (0.082)	9.942 (0.300)	0.397 (0.035)	0.554 (0.084)	9.934 (0.300)
Kernel		Binomial			Poisson			Negative Binomial		
$n$		$\hat{\pi}_1$	$\hat{\lambda}_1$	$\hat{\lambda}_2$	$\hat{\pi}_1$	$\hat{\lambda}_1$	$\hat{\lambda}_2$	$\hat{\pi}_1$	$\hat{\lambda}_1$	$\hat{\lambda}_2$
20	EM	0.409 (0.113)	0.544 (0.282)	9.904 (1.007)	0.409 (0.113)	0.544 (0.282)	9.904 (1.007)	0.409 (0.113)	0.544 (0.282)	9.904 (1.007)
	HD	0.517 (0.176)	0.354 (0.225)	9.441 (1.210)	0.565 (0.183)	0.426 (0.252)	9.380 (1.221)	0.636 (0.171)	0.466 (0.268)	9.400 (1.270)
	vNED	0.425 (0.153)	0.408 (0.230)	9.453 (1.234)	0.473 (0.165)	0.519 (0.255)	9.396 (1.220)	0.563 (0.164)	0.565 (0.287)	9.366 (1.511)
50	EM	0.397 (0.074)	0.490 (0.173)	10.01 (0.548)	0.397 (0.074)	0.490 (0.173)	10.01 (0.548)	0.397 (0.074)	0.490 (0.173)	10.01 (0.548)
	HD	0.389 (0.090)	0.438 (0.153)	9.503 (0.635)	0.466 (0.095)	0.604 (0.213)	9.582 (0.727)	0.532 (0.093)	0.681 (0.255)	9.586 (0.805)
	vNED	0.368 (0.081)	0.463 (0.144)	9.554 (0.614)	0.447 (0.086)	0.651 (0.191)	9.649 (0.733)	0.521 (0.084)	0.732 (0.230)	9.618 (0.896)
100	EM	0.399 (0.050)	0.502 (0.113)	10.02 (0.415)	0.399 (0.050)	0.502 (0.113)	10.02 (0.415)	0.399 (0.050)	0.502 (0.113)	10.02 (0.415)
	HD	0.363 (0.053)	0.492 (0.111)	9.510 (0.522)	0.451 (0.058)	0.765 (0.196)	9.752 (0.576)	0.517 (0.060)	0.903 (0.260)	9.812 (0.630)
	vNED	0.357 (0.050)	0.502 (0.101)	9.562 (0.498)	0.448 (0.054)	0.777 (0.168)	9.818 (0.552)	0.519 (0.055)	0.912 (0.224)	9.832 (0.646)
200	EM	0.398 (0.032)	0.502 (0.073)	10.07 (0.296)	0.398 (0.032)	0.502 (0.073)	10.07 (0.296)	0.398 (0.032)	0.502 (0.073)	10.07 (0.296)
	HD	0.350 (0.031)	0.513 (0.072)	9.491 (0.322)	0.448 (0.035)	0.878 (0.125)	9.911 (0.374)	0.517 (0.035)	1.089 (0.177)	10.09 (0.445)
	vNED	0.350 (0.031)	0.515 (0.067)	9.547 (0.310)	0.449 (0.034)	0.861 (0.104)	9.942 (0.353)	0.520 (0.033)	1.055 (0.148)	10.05 (0.438)

Figure 7 and Figure 8 are figure illustrations of Table 3.

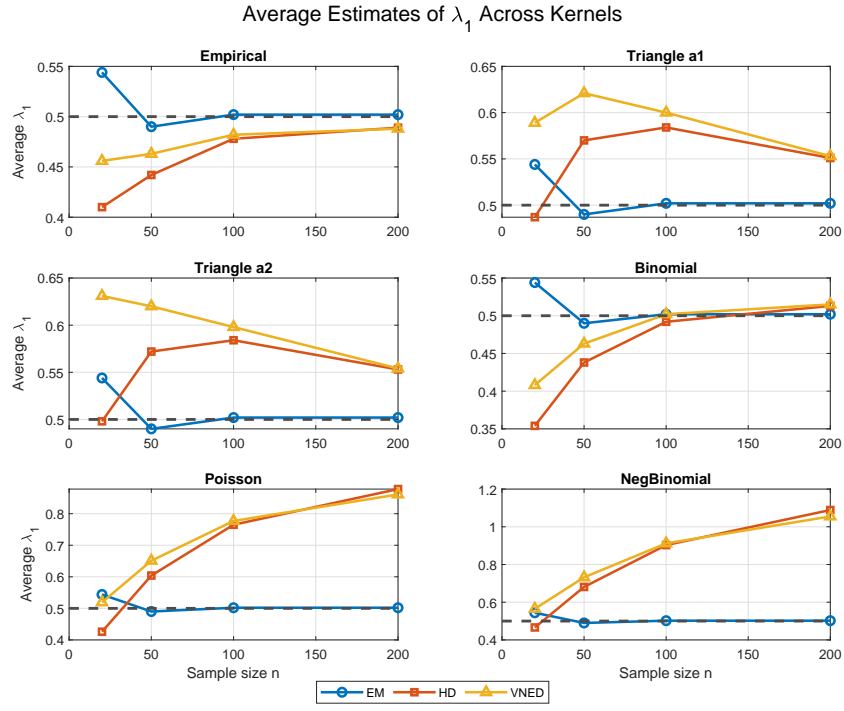


Figure 7: Average estimates of  $\lambda_1$  across kernels for sample sizes  $n = 20, 50, 100, 200$ . Solid lines represent EM, HD, and vNED methods; dashed horizontal line indicates the true value 0.4.

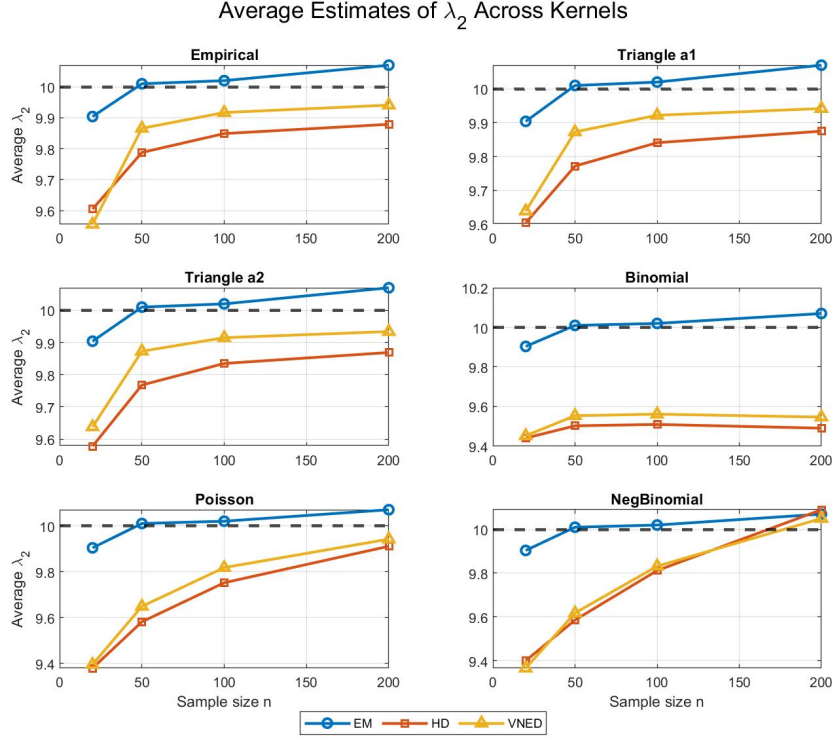


Figure 8: Average estimates of  $\lambda_2$  across kernels for sample sizes  $n = 20, 50, 100, 200$ . Solid lines represent EM, HD, and vNED methods; dashed horizontal line indicates the true value 0.4.

## PG Mixture Model with Discrete Kernel

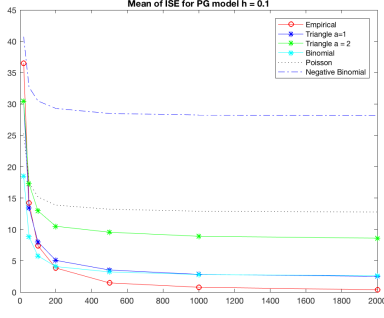
Now suppose the underlying model is two component PG mixture given by  $f(\cdot; \boldsymbol{\theta}) = 0.3PG(10, 1) + 0.7PG(1, 2)$ , where  $PG(\alpha, \beta)$  denotes a Poisson–Gamma mixture with  $\text{Gamma}(\alpha, \beta)$  mixing (shape  $\alpha$ , rate  $\beta$ ). For  $\epsilon \in \{0, 0.05, 0.10, 0.15, 0.20, 0.30\}$ , we inject outliers by replacing an  $\epsilon$  fraction of observations with the fixed value 50:  $Y_i^{(\epsilon)} = (1 - B_i)Y_i + 50 B_i$  with  $B_i \sim \text{Bernoulli}(\epsilon)$  independently. At each  $\epsilon$  we estimate by EM, HMIX, and vNEDMIX. For each parameter, the figure plots the Monte Carlo average (over 5,000 replications) versus  $\epsilon$ ; the dotted line marks the truth.



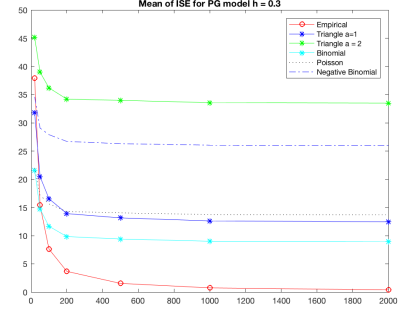
Table 4: Estimates Using Different Kernels for  $0.3PG(10, 1) + 0.7PG(1, 2)$

$n = 1000$		$\hat{\pi}_1$	$\hat{\alpha}_1$	$\hat{\beta}_1$	$\hat{\alpha}_2$	$\hat{\beta}_2$
EM		0.298 (0.018)	10.63 (2.827)	1.058 (0.269)	1.091 (0.378)	2.216 (0.938)
Empirical	HD	0.289 (0.018)	13.60 (4.089)	1.364 (0.391)	1.009 (0.315)	1.985 (0.760)
	vNED	0.289 (0.018)	12.82 (3.876)	1.287 (0.374)	1.043 (0.346)	2.081 (0.845)
Triangle $a = 1$	HD	0.292 (0.018)	13.54 (4.033)	1.358 (0.386)	1.079 (0.331)	2.072 (0.776)
	vNED	0.292 (0.018)	12.68 (3.766)	1.273 (0.363)	1.117 (0.363)	2.175 (0.862)
Triangle $a = 2$	HD	0.292 (0.018)	13.89 (4.270)	1.388 (0.407)	1.033 (0.281)	1.865 (0.624)
	vNED	0.293 (0.018)	12.93 (3.904)	1.293 (0.375)	1.064 (0.299)	1.945 (0.672)
Binomial	HD	0.272 (0.017)	14.21 (4.622)	1.420 (0.442)	1.226 (0.178)	1.826 (0.368)
	vNED	0.275 (0.017)	12.61 (3.787)	1.262 (0.363)	1.263 (0.182)	1.907 (0.376)
Poisson	HD	0.209 (0.016)	15.90 (4.911)	1.416 (0.423)	0.865 (0.044)	0.703 (0.069)
	vNED	0.216 (0.016)	12.76 (3.264)	1.146 (0.283)	0.890 (0.042)	0.745 (0.067)
Negative Binomial	HD	0.148 (0.014)	29.87 (16.50)	2.346 (1.290)	0.824 (0.029)	0.447 (0.036)
	vNED	0.152 (0.014)	23.22 (9.331)	1.845 (0.745)	0.834 (0.025)	0.459 (0.031)

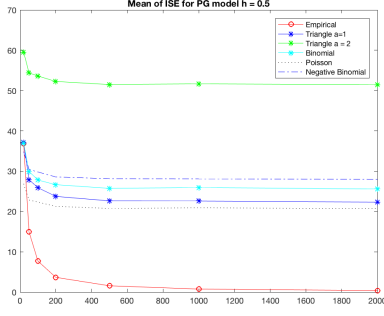
Figure 9 compare the mean of ISE for two component PG mixture model as  $n$  increases when choosing different bandwidth  $c$ . This indicates that (i) when the sample size is small, discrete kernels are beneficial, (ii) choosing proper discrete bandwidth is also crucial, and (iii) empirical kernel may be the best choice for this model.



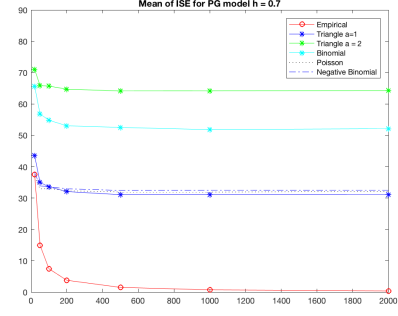
(a)  $c = 0.1$



(b)  $c = 0.3$



(c)  $c = 0.5$



(d)  $c = 0.7$

Figure 9: Mean of ISE for Model  $0.3PG(10, 1) + 0.7PG(1, 2)$  as  $n$  changes

## PL Mixture to Test Robustness

We simulate from a two-component PL mixture with  $\pi_1 = 0.3$ ,  $(\mu_1, \sigma_1) = (3, 0.5)$  and  $(\mu_2, \sigma_2) = (1, 0.5)$ . For each contamination level  $\epsilon \in \{0, 0.05, 0.10, 0.15, 0.20, 0.30\}$ , an  $\epsilon$ -fraction of observations is replaced by a large outlier value (50). Sample size  $n = 2000$ , replications  $R = 5000$ . Panels (a)–(e) plot the replication averages of the estimated parameters  $(\pi_1, \mu_1, \sigma_1, \mu_2, \sigma_2)$  versus  $\epsilon$  for three methods: EM (MLE), HMIX (Hellinger divergence), and vNEDMIX (vNED divergence). Horizontal dotted lines mark the true parameter values. For PL likelihood evaluations we use fixed Gauss–Hermite quadrature (common nodes/weights across methods). Warm starts are used across contamination levels

to stabilize fits, and axis limits are set robustly so occasional non-convergent runs do not distort the display.

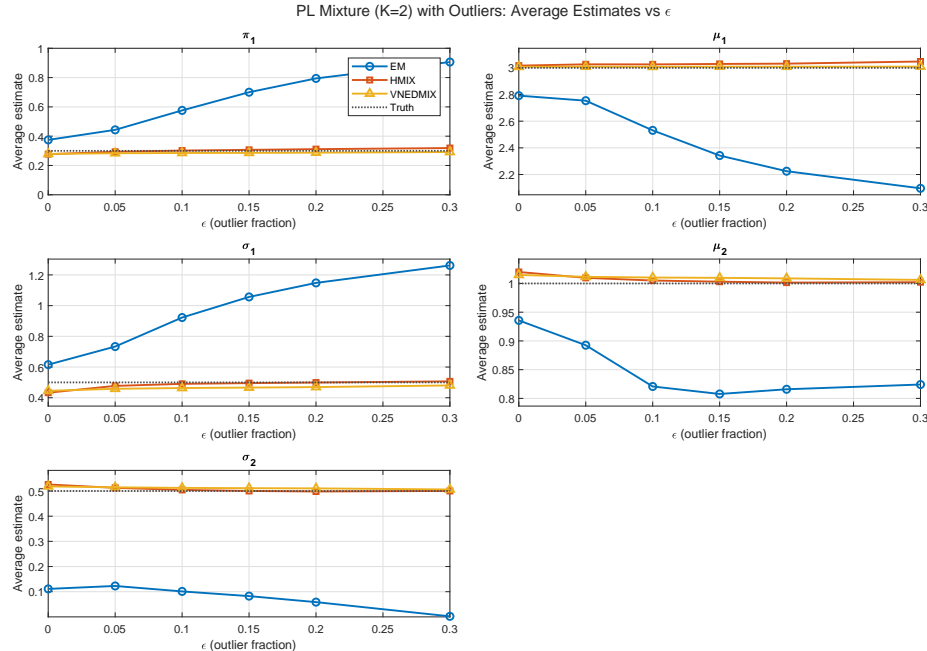


Figure 10: Outlier robustness for a Poisson-Lognormal (PL) mixture (known  $K = 2$ ).

This figure assesses robustness of three estimation strategies for the PL mixture under increasing contamination. All methods perform well at  $\epsilon = 0$ . As  $\epsilon$  grows, EM becomes progressively more sensitive to the injected large counts, typically inflating location/scale parameters for the high-mean component and, to a lesser extent, perturbing the mixture weight  $\pi_1$ . In contrast, HMIX and vNEDMIX downweight the influence of extreme observations through their discrepancy objectives, and therefore track the dotted truth lines more closely over a wider range of  $\epsilon$ . The divergence-based methods thus exhibit markedly improved stability for both location ( $\mu_k$ ) and dispersion ( $\sigma_k$ ) parameters under contamination, while maintaining competitive accuracy in the clean setting. The use of fixed Gauss-Hermite quadrature provides fast, consistent likelihood approximations for the PL components across all methods.

## Continuous Models

In this subsection, we consider the continuous model case. Specifically, we consider the two component mixture normal model, i.e.,

$$f(\cdot; \boldsymbol{\theta}) = \pi_1 N(\mu_1, \sigma_1^2) + (1 - \pi_1) N(\mu_2, \sigma_2^2).$$

In the following simulation, we let true parameters  $\pi_1 = 0.3, \mu_1 = 10, \sigma_1 = 1, \mu_2 = 0, \sigma_2 = 1$ , and the sample size is  $n = 200$ . For the nonparametric density estimate  $g_n(\cdot)$ , we use the Epanechnikov kernel and the optimal bandwidth is applied. Also, the number of components  $K$  is unknown and estimated based on the DIC criteria as described before.

Table 5: Two-Component Normal (Epanechnikov kernel; 100% chose  $K = 2$ )

		$\hat{\pi}_1$	$\hat{\mu}_1$	$\hat{\sigma}_1$	$\hat{\mu}_2$	$\hat{\sigma}_2$
EM	Ave	0.299	9.998	0.986	0.000	0.995
	StD	0.032	0.126	0.091	0.086	0.059
	MSE	0.001	0.016	0.009	0.007	0.004
HMIX	Ave	0.299	9.997	1.236	0.000	1.250
	StD	0.032	0.126	0.073	0.086	0.050
	MSE	0.001	0.016	0.061	0.007	0.065
vNEDMIX	Ave	0.299	9.997	1.243	0.000	1.253
	StD	0.032	0.126	0.073	0.086	0.050
	MSE	0.001	0.016	0.064	0.007	0.067

We observe that for continuous normal mixture models, the HMIX, vNEDMIX are competitive to the EM algorithm.

Table 6: Average Runtime (seconds) per simulation for mixture models with  $n = 20,000$  and true  $K = 2$ .

	$K$ known	$K$ unknown
<b>Pois</b>	0.2 s	16 s
<b>PG</b>	0.7 s	48 s
<b>PL</b>	46 s	66 s

## J: Additional Image Analysis

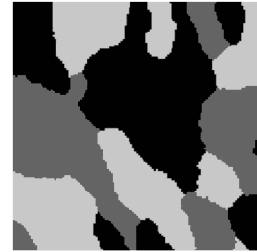
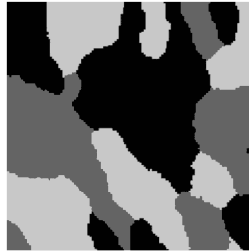
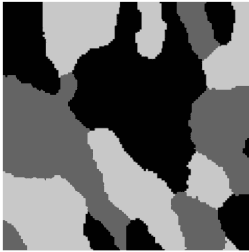
### Original Image Recovery

We observe from Table 7 that the BIC, HIC, and vNEDIC all prefer  $K = 5$ . However, when we recover the images the difference is negligible (figures not shown here). Thus, we choose  $K = 3$ . Figure 12 and Figure 13 show that the EM, HMIX, and vNED algorithms are comaparable in reconstructing the images.



(a) Original three-class phantom image

(b) Original lena image with 256 gray scale values



(a) EM algorithm

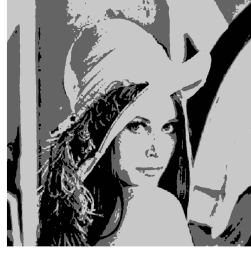
(b) HMIX algorithm

(c) vNEDMIX algorithm

Figure 12: Three-class phantom image segamentation  $K = 3$



(a) EM algorithm



(b) HMX algorithm



(c) vNEDMIX algorithm

Figure 13: Lena Image Segmentation  $K = 3$

The following tables are additional analysis for Phantom image and Lena image.

Table 7: Comparison of Different Model Selection Criteria (Phantom and Lena)

Image	Criterion	# of components = 1	2	3	4	5
Phantom	BIC	$2.327 \times 10^7$	$1.021 \times 10^7$	<b><math>5.584 \times 10^6</math></b>	$5.584 \times 10^6$	$5.584 \times 10^6$
	HIC	1.353	0.998	<b>0.804</b>	0.804	0.804
	vNEDIC	0.897	0.817	<b>0.768</b>	0.768	0.768
Lena	BIC	$1.599 \times 10^7$	$8.497 \times 10^6$	$6.965 \times 10^6$	$6.399 \times 10^6$	<b><math>6.290 \times 10^6</math></b>
	HIC	0.689	0.396	0.192	0.062	<b>0.022</b>
	vNEDIC	0.648	0.546	0.479	0.407	<b>0.382</b>

Table 8: Parameter Point Estimates for Phantom Image

		$\hat{\pi}_1$	$\hat{\lambda}_1$	$\hat{\pi}_2$	$\hat{\lambda}_2$	$\hat{\pi}_3$	$\hat{\lambda}_3$	$\hat{\pi}_4$	$\hat{\lambda}_4$	$\hat{\lambda}_5$
$K = 2$	EM	0.3665	28.99	0.6335	173.9					
	HMIX	0.5835	28.79	0.4165	218.1					
	vNEDMIX	0.6030	28.96	0.3970	218.7					
$K = 3$	EM	0.3636	28.65	0.3050	123.5	0.3314	219.3			
	HMIX	0.4097	28.79	0.2987	124.0	0.2916	218.2			
	vNEDMIX	0.3851	28.90	0.2857	124.2	0.3291	218.7			
$K = 4$	EM	0.1966	28.65	0.1670	28.65	0.3050	123.5	0.3314	219.3	
	HMIX	0.2215	28.79	0.1882	28.79	0.2987	124.0	0.3609	218.2	
	vNEDMIX	0.1991	28.90	0.1860	28.90	0.2857	124.2	0.3291	218.7	
$K = 5$	EM	0.1818	28.65	0.1818	28.65	0.3050	123.5	0.1657	219.3	219.3
	HMIX	0.2049	28.79	0.2049	28.79	0.2987	124.0	0.1458	218.2	218.2
	vNEDMIX	0.1926	28.90	0.1926	28.90	0.2857	124.2	0.1646	218.7	218.7



Table 9: Parameter Point Estimates for Lena Image

		$\hat{\pi}_1$	$\hat{\lambda}_1$	$\hat{\pi}_2$	$\hat{\lambda}_2$	$\hat{\pi}_3$	$\hat{\lambda}_3$	$\hat{\pi}_4$	$\hat{\lambda}_4$	$\hat{\lambda}_5$
$K = 2$	EM	0.3788	67.34	0.6212	151.1					
	HMIX	0.3747	91.62	0.6253	146.1					
	vNEDMIX	0.3119	48.28	0.6881	141.1					
$K = 3$	EM	0.2284	52.03	0.3731	108.8	0.3985	167.9			
	HMIX	0.2254	48.77	0.2977	96.55	0.4769	148.1			
	vNEDMIX	0.2470	48.01	0.5180	137.2	0.2350	189.9			
$K = 4$	EM	0.1907	48.16	0.2331	90.85	0.3935	138.0	0.1827	189.8	
	HMIX	0.1932	48.29	0.2370	92.36	0.3939	139.3	0.1760	191.0	
	vNEDMIX	0.1945	47.74	0.2401	94.13	0.3950	141.3	0.1705	192.8	
$K = 5$	EM	0.1794	47.08	0.1832	84.43	0.2162	118.7	0.2832	150.8	196.3
	HMIX	0.1792	47.48	0.1802	85.45	0.2187	119.3	0.2839	150.9	196.3
	vNEDMIX	0.1782	47.46	0.1807	87.02	0.2213	120.3	0.2823	151.2	196.5



(a) Recovered image from EM algorithm



(b) Recovered image from vNEDMIX algorithm

Figure 14: Lena image reconstruction after adding 30% outliers: EM vs. vNEDMIX.

## K: Assumptions

This section lays out all the necessary assumptions for the proofs of the main lemmas and theorems in Sections 3 and 4.

**Connection with Z-estimation.** The standing assumptions in this appendix are tailored so that the minimum-divergence estimator and its DM-based approximations can be analyzed along the lines of the classical Z-estimation scheme. In particular, for fixed  $K$  we can write the MDE as the solution of a system of estimating equations

$$\Psi_n(\boldsymbol{\theta}) := \nabla_{\boldsymbol{\theta}} D_G(g_n, f_{\boldsymbol{\theta}}) = 0, \quad \Psi(\boldsymbol{\theta}; g) := \nabla_{\boldsymbol{\theta}} D_G(g, f_{\boldsymbol{\theta}}),$$

with population target  $\boldsymbol{\theta}^* \in \arg \min_{\boldsymbol{\theta}} D_G(g, f_{\boldsymbol{\theta}})$  satisfying  $\Psi(\boldsymbol{\theta}^*; g) = 0$ . Assumptions (C0)–(C3), (F1), (K1)–(K2) and (M1)–(M8) guarantee: (i) identification and a local quadratic expansion of  $D_G$  around  $\boldsymbol{\theta}^*$ ; (ii) a uniform law of large numbers for  $\Psi_n(\boldsymbol{\theta})$  on a neighborhood of  $\boldsymbol{\theta}^*$ ; and (iii) Fréchet differentiability and nonsingularity of the Jacobian  $H(\boldsymbol{\theta}^*) := \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}^*; g)$ . These are exactly the ingredients that replace the abstract stochastic-equicontinuity conditions in the general Z-estimation theory, and they allow us to derive consistency and  $\sqrt{n}$ -normality for the solution of  $\Psi_n(\boldsymbol{\theta}) = 0$  in the usual way.

At the same time, the DM algorithm introduces two additional layers: the plug-in density  $g_n$  and the finite number of iterations  $m_n$ . The operator-level bounds in Section 3 (in particular the contraction and noisy-contraction results) are used to control the optimization error  $\|\boldsymbol{\theta}_{m_n, n} - \hat{\boldsymbol{\theta}}_n\|$ , where  $\hat{\boldsymbol{\theta}}_n$  solves  $\Psi_n(\boldsymbol{\theta}) = 0$ , while the assumptions in this appendix ensure that  $\hat{\boldsymbol{\theta}}_n$  itself behaves as a standard Z-estimator. Appendix N makes these connections explicit by working out the linearization of  $\Psi_n$  around  $\boldsymbol{\theta}^*$  and the associated Godambe information, and by showing how the DM-specific error terms can be absorbed into the usual Z-estimation expansions under our choice of  $m_n$ .

### Assumptions for Model Identifiability

(A0) Assume that  $G(0) = 0$ ,  $G'(0) = 0$  and  $G''(0) = 1$ . Note that the above conditions imply  $G'(\delta) < 0$  when  $\delta < 0$  and  $G'(\delta) > 0$  when  $\delta > 0$ .

We now state two identifiability conditions below which enable us to deduce results concerning the mixture model with random effects.

(A1)  $\lambda_1 \neq \lambda_2$  implies  $s(y|\lambda_1) \neq s(y|\lambda_2)$  for all  $y$ .

(A2)  $\phi_1 \neq \phi_2$  implies  $r(\lambda, \phi_1) \neq r(\lambda, \phi_2)$  on a set of positive Lebesgue measure.

Note that (D2) is satisfied if the condition (A1)-(A2) hold.

**Assumptions for  $Q_n(\cdot|\cdot)$ ,  $\Theta_{0,n}$ , and  $D_n(\cdot)$ .**

### Parameter Space, Stationary Points, and Divergence Function Assumptions

(B1)  $Q(\theta'|\theta)$  and  $Q_n(\theta'|\theta)$  are both continuous in  $\theta'$  and  $\theta$ . Also,  $Q(\theta'|\theta)$  is continuously differentiable w.r.t.  $\theta'$ .

(B2)  $\Theta_{\theta_0} = \{\theta' \in \Theta : D(\theta') \leq D(\theta_0)\}$  and  $\Theta_{\theta_{0,n}} = \{\theta'_n \in \Theta : D_n(\theta'_n) \leq D_n(\theta_{0,n})\}$  are compact.

(B3)  $D(\theta') < \infty$  and  $D_n(\theta') < \infty$  for  $\theta' \in \Theta$ .

(B4) All stationary points in  $S_D$  and  $S_{n,D}$  are isolated; that is, for any  $\theta \in S_D$  and  $\theta \in S_{n,D}$ , there exists a neighborhood of  $\theta$  which does not contain any other points of  $S_D$  and  $S_{n,D}$ .

### Assumptions for Consistency of the DM Algorithm Sequence

(C0) The density  $s(\cdot|\lambda)$  is bounded,  $r(\cdot; \phi_k)$  is continuous in  $\phi_k$  for each  $1 \leq k \leq K$ . Also, assume that there exists a function  $R(\cdot)$  such that  $r(\lambda; \phi_k) \leq R(\lambda)$  for all  $1 \leq k \leq K$  and  $\int_{\mathbb{R}} R(\lambda) d\lambda < \infty$ .

(C1) The parameter space  $\Theta$  is compact.

(C2)  $G(-1) + G'(\infty) < \infty$ , where  $G'(\infty) = \lim_{u \rightarrow \infty} G(u)/u$ .

(C3)  $g_n(\cdot)$  is strongly consistent to  $g(\cdot)$ .

(C4)  $c_n \rightarrow 0, nc_n \rightarrow \infty$  as  $n \rightarrow \infty$ .

### DM Algorithm Objective Function Assumptions

(D1) Closed graph / outer semicontinuity of the update map. The DM update correspondences  $M_n : \Theta \rightrightarrows \Theta$  and  $M : \Theta \rightrightarrows \Theta$  have closed graphs on  $\Theta$ ; that is, if  $\theta^{(j)} \rightarrow \bar{\theta}$  and  $\eta^{(j)} \in M_n(\theta^{(j)})$  with  $\eta^{(j)} \rightarrow \bar{\eta}$ , then  $\bar{\eta} \in M_n(\bar{\theta})$  (and analogously for  $M$ ). Equivalently,  $M_n$  and  $M$  are outer semicontinuous on  $\Theta$ . *Single-valued specialization.* When, for each  $\theta$  in a compact neighborhood, the minimizer of  $Q_n(\cdot | \theta)$  (resp.  $Q(\cdot | \theta)$ ) is unique,  $M_n$  (resp.  $M$ ) is single-valued and (D1) reduces to ordinary continuity of the function  $\theta \mapsto M_n(\theta)$  (resp.  $\theta \mapsto M(\theta)$ ) by standard argmin-continuity/Berge's maximum theorem.

(D2) For given  $\theta \in \Theta$ , both  $Q(\cdot | \theta)$  and  $Q_n(\cdot | \theta)$  have a unique global minimum.

Note that, under **B**, if (D2) holds, then (D1) is always satisfied. However, the converse may not hold. We can obtain the continuity of  $M(\cdot)$  by replacing conditions (B') in the proposition with conditions (B).

### Kernel Assumptions for Central Limit Theorem

(K1) The kernel  $\mathcal{K}(\cdot)$  is symmetric at 0, has second moment, and twice continuously differentiable on a compact support  $Supp(\mathcal{K})$ .

(K2) The bandwidth  $c_n$  satisfies  $c_n \rightarrow 0$ ,  $\sqrt{nc_n} \rightarrow \infty$ , and  $\sqrt{nc_n^4} \rightarrow 0$  as  $n \rightarrow \infty$ .

### Regularity Conditions for Central Limit Theorem

Below we state the main assumptions that are required in the proof of the CLT. These assumptions describe the smoothness properties of the postulated density.

(M1)  $A(\delta)$ ,  $A'(\delta)$ ,  $A'(\delta)(\delta + 1)$  and  $A''(\delta)(\delta + 1)$  are bounded on  $[-1, \infty)$ .

(M2) The following conditions ensure  $L^1$ -continuity of both the Hessian of the density and the score-quadratic component under parameter convergence. Specifically, as  $n \rightarrow \infty$ ,  $\phi_n \rightarrow \theta^*$ ,

$$\int_{\mathbb{R}} |\nabla^2 f(y; \phi_n) - \nabla^2 f(y; \theta^*)| dy = o_p(1) \quad \text{and} \\ \int_{\mathbb{R}} |u(y; \phi_n)u'(y; \phi_n)f(y; \phi_n) - u(y; \theta^*)u'(y; \theta^*)f(y; \theta^*)| dy = o_p(1).$$

(M3) The matrix  $I(\theta^*)$  given by

$$I(\theta^*) = \int_{\mathbb{R}} u(y; \theta^*)u'(y; \theta^*)f(y; \theta^*)dy$$

is finite.

(M4) An integrability constraint ensuring that the squared second derivative, scaled by the influence ratio, remains finite.

$$\int_{\mathbb{R}} \frac{|u(y; \theta^*)|}{f(y; \theta^*)} f''^2(y; \theta^*) dy < \infty.$$

In addition, let  $\{\alpha_n : n \geq 1\}$  be a sequence diverging to infinity, and

$$\int_{\mathbb{R}} \int_{|y| \leq \alpha_n} \frac{|u(y; \theta^*)|}{f(y; \theta^*)} \tilde{H}_n(t, y) dy dt < \infty,$$

where  $\tilde{H}_n(t, y) = \sup\{H_n(t, y), t \in \text{Supp}(\mathcal{K}), |y| \leq \alpha_n\}$ ,  $H_n(t, y) = \mathcal{K}(t)t^2 f''(\tilde{y}_n(t); \theta)$ , and  $\tilde{y}_n(t)$  is a point between  $(\min(y - c_n t, y), \max(y - c_n t, y))$ .

(M5) This condition controls the tail probability that the shifted observation  $Y_1 - c_n t$  falls outside the truncation region.

$$h_n \equiv n \sup_{t \in \text{Supp}(\mathcal{K})} \mathbb{P}(|Y_1 - c_n t| > \alpha_n) = O_p(1).$$

(M6) The following condition ensures that the truncated mass of the score component becomes negligible under the normalization  $1/(\sqrt{n} c_n)$ .

$$\lim_{n \rightarrow \infty} \frac{1}{\sqrt{n} c_n} \int_{|y| \leq \alpha_n} u(y; \theta^*) dy = 0.$$

(M7) This condition requires that the ratio of the shifted density to the original density remains uniformly bounded over the truncation region.

$$\sup_{|y| \leq \alpha_n} \sup_{t \in \text{Supp}(\mathcal{K})} \frac{f(y - c_n t; \boldsymbol{\theta}^*)}{f(y; \boldsymbol{\theta}^*)} = O_p(1).$$

(M8) The following condition ensures that the squared  $L^2$  distance between the shifted and unshifted score functions vanishes uniformly over  $t$ .

$$\lim_{n \rightarrow \infty} \sup_{t \in \text{Supp}(\mathcal{K})} \int_{\mathbb{R}} (u(y + c_n t; \boldsymbol{\theta}^*) - u(y; \boldsymbol{\theta}^*)) (u(y + c_n t; \boldsymbol{\theta}^*) - u(y; \boldsymbol{\theta}^*))' f(y; \boldsymbol{\theta}^*) dy = 0.$$

### Parametric Model Assumptions

(F1)  $f(\cdot; \boldsymbol{\theta})$  is twice continuously differentiable.

### Regularity Conditions for Breakdown Point Analysis

(O1)  $\int_{\mathbb{R}} |g(y) - \eta_n(y)| dy \rightarrow 2$  as  $n \rightarrow \infty$ .

(O2)  $\int_{\mathbb{R}} |f(y; \boldsymbol{\theta}) - \eta_n(y)| dy \rightarrow 2$  as  $n \rightarrow \infty$  uniformly for  $|\boldsymbol{\theta}| \leq c$  for every fixed  $c$ .

(O3)  $\int_{\mathbb{R}} |g(y) - f(y; \boldsymbol{\theta}_{\epsilon, n}^*)| dy \rightarrow 2$  as  $n \rightarrow \infty$  if  $|\boldsymbol{\theta}_{\epsilon, n}^*| \rightarrow \infty$  as  $n \rightarrow \infty$ , where  $\boldsymbol{\theta}_{\epsilon, n}^* = T(g_{\epsilon, n})$ .

# L: Proofs

## 9.1 Proofs of Results in Section 2 of the Main

*Detailed Proof of Lemma 5.* Let  $U(t) = tG(-1 + a/t)$ , where  $a \geq 0, t > 0$ . We first show that  $U(\cdot)$  is a convex function on  $(0, \infty)$ . Since  $G(\cdot)$  is thrice differentiable and convex, we take the second derivative with respect to  $t$ . Let  $\delta = (-1 + a/t)$ , so the second derivative of  $U(\cdot)$  w.r.t.  $t$  is

$$\begin{aligned} \frac{\partial^2 U}{\partial t^2} &= \frac{\partial G(\delta)}{\partial \delta} \left( -\frac{a}{t^2} \right) - \left( \frac{\partial^2 G(\delta)}{\partial \delta^2} \left( -\frac{a}{t^2} \right) \frac{a}{t} + \frac{\partial G(\delta)}{\partial \delta} \left( -\frac{a}{t^2} \right) \right) \\ &= \frac{\partial^2 G(\delta)}{\partial \delta^2} \frac{a^2}{t^3} \geq 0 \end{aligned}$$

as  $a \geq 0, t > 0$  and convexity of  $G(\cdot)$ . Hence  $U(\cdot)$  is convex. Next we will show that for almost all  $y, a \geq 0, \boldsymbol{\theta}' \in \boldsymbol{\Theta}$ , and every  $\tilde{q}(\cdot|y)$ ,

$$f(y; \boldsymbol{\theta}') G \left( -1 + \frac{a}{f(y; \boldsymbol{\theta}')} \right) \leq \int_{\mathcal{Z}} f(y; \boldsymbol{\theta}') w(Z|Y; \boldsymbol{\theta}') G \left( -1 + \frac{a \tilde{q}(z|y)}{f(y; \boldsymbol{\theta}') w(Z|Y; \boldsymbol{\theta}')} \right) dz. \quad (9.19)$$

First note that the left hand side (LHS) of (9.19) is  $U(f(y; \boldsymbol{\theta}'))$ ; for the right hand side (RHS) of (9.19), by Jensen's inequality, we have

$$\begin{aligned} U(f(y; \boldsymbol{\theta}')) &= U \left( \mathbf{E}_{\tilde{q}} \left[ \frac{f(y; \boldsymbol{\theta}') w(Z|y; \boldsymbol{\theta}')}{\tilde{q}(Z|y)} \right] \right) \\ &\leq \mathbf{E}_{\tilde{q}} \left[ U \left( \frac{f(y; \boldsymbol{\theta}') w(Z|y; \boldsymbol{\theta}')}{\tilde{q}(Z|y)} \right) \right] \\ &= \int_{\mathcal{Z}} f(y; \boldsymbol{\theta}') w(Z|y; \boldsymbol{\theta}') G \left( -1 + \frac{a \tilde{q}(z|y)}{f(y; \boldsymbol{\theta}') w(Z|y; \boldsymbol{\theta}')} \right) dz. \end{aligned}$$

Hence (9.19) holds. Now let  $a = g(y)$ , and taking another integral w.r.t.  $\mathcal{Y}$ , we have  $\mathcal{D}(q, \boldsymbol{\theta}') \geq D(\boldsymbol{\theta}')$ . In addition, the equality holds if and only if  $\tilde{q}(z|y) = w(z|y; \boldsymbol{\theta}')$  for almost every  $y$ . The proof is complete.  $\square$

*Proof of Lemma (finite stationary points at a fixed value).* Fix  $D_n^* \in \mathbb{R}$  and set

$$S := \{\boldsymbol{\theta} \in \boldsymbol{\Theta} : \nabla D_n(\boldsymbol{\theta}) = 0, D_n(\boldsymbol{\theta}) = D_n^*\}.$$

By (B1)  $D_n$  and  $\nabla D_n$  are continuous; hence  $S = \nabla D_n^{-1}(\{0\}) \cap D_n^{-1}(\{D_n^*\})$  is closed. Let  $c > D_n^*$ ; by (B2) the sublevel set  $\{\boldsymbol{\theta} : D_n(\boldsymbol{\theta}) \leq c\}$  is compact, so  $S$  is compact. By (B4) every stationary point is isolated: for each  $\boldsymbol{\theta} \in S$  there exists  $r_\theta > 0$  such that  $B(\boldsymbol{\theta}, r_\theta) \cap \{\nabla D_n = 0\} = \{\boldsymbol{\theta}\}$ . Then  $\{B(\boldsymbol{\theta}, r_\theta/2) : \boldsymbol{\theta} \in S\}$  is an open cover of the compact set  $S$ , so it admits a finite subcover, forcing  $S$  to be finite (each ball contains exactly one point of  $S$  by isolation).  $\square$

### Proof of Proposition 3

*Proof of Proposition 3 (limit-set structure: convergence or finite cycle).* Let  $\{\boldsymbol{\theta}_{m,n}\}_{m \geq 0}$  be a sample-level DM sequence with  $\boldsymbol{\theta}_{m+1,n} \in M_n(\boldsymbol{\theta}_{m,n})$ . By Proposition 1 (majorization descent),  $D_n(\boldsymbol{\theta}_{m+1,n}) \leq D_n(\boldsymbol{\theta}_{m,n})$  for all  $m$ , and  $D_n(\boldsymbol{\theta}_{m,n}) \downarrow D_n^*$  for some  $D_n^* \in \mathbb{R}$ . Denote the  $\omega$ -limit set  $\Omega := \{\bar{\boldsymbol{\theta}} : \exists m_j \uparrow \infty, \boldsymbol{\theta}_{m_j,n} \rightarrow \bar{\boldsymbol{\theta}}\}$ . By (B2), all iterates lie in a compact sublevel set, so  $\Omega \neq \emptyset$ . By continuity of  $D_n$  and the convergence of  $D_n(\boldsymbol{\theta}_{m,n})$ , every  $\bar{\boldsymbol{\theta}} \in \Omega$  satisfies  $D_n(\bar{\boldsymbol{\theta}}) = D_n^*$ . By Proposition 1, every  $\bar{\boldsymbol{\theta}} \in \Omega$  is stationary for  $D_n$ ; together with the previous lemma,  $\Omega$  is a *finite* set.

Next, we show  $\Omega$  is  $M_n$ -invariant and that  $M_n$  acts as a permutation on  $\Omega$ . Assume (D1) (continuity of  $M_n$ ). If  $\boldsymbol{\theta}_{m_j,n} \rightarrow \bar{\boldsymbol{\theta}}$ , then  $\boldsymbol{\theta}_{m_j+1,n} = M_n(\boldsymbol{\theta}_{m_j,n}) \rightarrow M_n(\bar{\boldsymbol{\theta}})$ ; hence  $M_n(\bar{\boldsymbol{\theta}}) \in \Omega$ , so  $M_n(\Omega) \subseteq \Omega$ . Conversely, given any  $z \in \Omega$ , pick  $\{m_j\}$  with  $\boldsymbol{\theta}_{m_j+1,n} \rightarrow z$ ; compactness yields a subsequence  $\boldsymbol{\theta}_{m_{j_k},n} \rightarrow y \in \Omega$  with  $M_n(y) = z$  by continuity, so  $M_n : \Omega \rightarrow \Omega$  is surjective, hence bijective on the finite set  $\Omega$ . Therefore  $M_n|_{\Omega}$  is a *permutation*, so  $\Omega$  is a finite union of cycles. Because  $\{\boldsymbol{\theta}_{m,n}\}$  has  $\Omega$  as its *entire* limit set, only one such cycle can occur: otherwise, invariance would force the orbit to eventually remain in a neighborhood of a single cycle and it could not accumulate on the others. Hence  $\Omega = \{\boldsymbol{\theta}_{1,n}^*, \dots, \boldsymbol{\theta}_{t,n}^*\}$  with

$$M_n(\boldsymbol{\theta}_{i,n}^*) = \boldsymbol{\theta}_{i+1,n}^* \quad (i = 1, \dots, t-1), \quad M_n(\boldsymbol{\theta}_{t,n}^*) = \boldsymbol{\theta}_{1,n}^*,$$



and each  $\boldsymbol{\theta}_{i,n}^*$  is stationary with  $D_n(\boldsymbol{\theta}_{i,n}^*) = D_n^*$ .

It remains to establish (iii): convergence of the parallel subsequences  $\{\boldsymbol{\theta}_{tm+i,n}\}_{m \geq 0}$  to  $\boldsymbol{\theta}_{i,n}^*$ . By continuity of  $M_n$  and isolation (B4), fix disjoint open balls  $U_i := B(\boldsymbol{\theta}_{i,n}^*, r)$  so small that  $M_n(\overline{U_i}) \subset U_{i+1}$  for all  $i$  (indices modulo  $t$ ); this is possible since  $M_n(\boldsymbol{\theta}_{i,n}^*) = \boldsymbol{\theta}_{i+1,n}^*$  and  $M_n$  is continuous. Because  $\text{dist}(\boldsymbol{\theta}_{m,n}, \Omega) \rightarrow 0$  and the  $U_i$  cover a neighborhood of  $\Omega$ , there exists  $m_0$  after which every iterate belongs to  $\bigcup_i U_i$  and, moreover, membership cycles deterministically:  $\boldsymbol{\theta}_{m,n} \in U_i \Rightarrow \boldsymbol{\theta}_{m+1,n} \in U_{i+1}$ . In particular, for each  $i$  there exists  $m_i \geq m_0$  with  $\boldsymbol{\theta}_{m_i,n} \in U_i$ , and then by  $t$ -step invariance  $M_n^t(\overline{U_i}) \subset U_i$  we have  $\boldsymbol{\theta}_{m_i+kt,n} \in U_i$  for all  $k \geq 0$ . Any limit point of  $\{\boldsymbol{\theta}_{m_i+kt,n}\}_{k \geq 0}$  lies in  $\Omega \cap U_i = \{\boldsymbol{\theta}_{i,n}^*\}$ , so the subsequence converges:  $\boldsymbol{\theta}_{tm+i,n} \rightarrow \boldsymbol{\theta}_{i,n}^*$  as claimed.  $\square$

## Proof of Proposition 4

*Proof of Proposition 4 (no cycles under uniqueness; strict descent).* Assume (B) and (D2). Uniqueness in (D2) implies that for each  $\boldsymbol{\theta}$ , the map  $\boldsymbol{\theta}' \mapsto Q_n(\boldsymbol{\theta}' \mid \boldsymbol{\theta})$  has a *unique* minimizer, so the update is single-valued and continuous (standard argmin continuity), i.e., (D1) holds.

*Strict descent away from stationarity.* If  $\boldsymbol{\theta}_{m+1,n} = M_n(\boldsymbol{\theta}_{m,n}) \neq \boldsymbol{\theta}_{m,n}$ , then by uniqueness

$$Q_n(\boldsymbol{\theta}_{m+1,n} \mid \boldsymbol{\theta}_{m,n}) < Q_n(\boldsymbol{\theta}_{m,n} \mid \boldsymbol{\theta}_{m,n}) = D_n(\boldsymbol{\theta}_{m,n}),$$

and by majorization  $D_n(\boldsymbol{\theta}_{m+1,n}) \leq Q_n(\boldsymbol{\theta}_{m+1,n} \mid \boldsymbol{\theta}_{m,n})$ , hence  $D_n(\boldsymbol{\theta}_{m+1,n}) < D_n(\boldsymbol{\theta}_{m,n})$ .

*Excluding cycles.* If  $\Omega$  contained a cycle of length  $t \geq 2$ , then along that cycle all points are stationary and have objective value  $D_n^*$ , so once the orbit enters a small invariant neighborhood of the cycle, the updates cannot be strictly decreasing (the values along one pass around the cycle remain equal), contradicting the strict descent above unless the update lands exactly at a fixed point. Therefore  $\Omega$  is a singleton  $\{\boldsymbol{\theta}_n^*\}$ , and the sequence converges:  $\boldsymbol{\theta}_{m,n} \rightarrow \boldsymbol{\theta}_n^*$ .

*Fixed point identity.* By continuity of  $M_n$  and the update relation  $\boldsymbol{\theta}_{m+1,n} = M_n(\boldsymbol{\theta}_{m,n})$ , passing to the limit yields  $\boldsymbol{\theta}_n^* = M_n(\boldsymbol{\theta}_n^*)$ . The final monotonicity statement follows directly from the strict descent shown above unless  $\boldsymbol{\theta}_{m,n}$  has already reached  $\boldsymbol{\theta}_n^*$ .  $\square$

## Proof of Theorem 1

*Proof of Theorem 1.* Since both  $M(\boldsymbol{\theta})$  and  $\boldsymbol{\theta}^*$  are in  $\Theta$ , we may apply condition

$$\langle \nabla Q(\boldsymbol{\theta}^* | \boldsymbol{\theta}^*), \boldsymbol{\theta}' - \boldsymbol{\theta}^* \rangle \geq 0 \quad \text{for all } \boldsymbol{\theta}' \in \Theta. \quad (\text{A.1})$$

with  $\boldsymbol{\theta}' = M(\boldsymbol{\theta})$ :

$$\langle \nabla Q(\boldsymbol{\theta}^* | \boldsymbol{\theta}^*), M(\boldsymbol{\theta}) - \boldsymbol{\theta}^* \rangle \geq 0 \quad \text{for all } \boldsymbol{\theta}' \in \Theta,$$

and apply condition

$$\langle \nabla Q(M(\boldsymbol{\theta}) | \boldsymbol{\theta}), \boldsymbol{\theta}' - M(\boldsymbol{\theta}) \rangle \geq 0 \quad \text{for all } \boldsymbol{\theta}' \in \Theta. \quad (\text{A.2})$$

with  $\boldsymbol{\theta}' = \boldsymbol{\theta}^*$ :

$$\langle \nabla Q(M(\boldsymbol{\theta}) | \boldsymbol{\theta}), \boldsymbol{\theta}^* - M(\boldsymbol{\theta}) \rangle \geq 0 \quad \text{for all } \boldsymbol{\theta}' \in \Theta.$$

Adding the above two inequalities and then perform some algebra yields the condition

$$\langle \nabla Q(M(\boldsymbol{\theta}) | \boldsymbol{\theta}^*) - \nabla Q(\boldsymbol{\theta}^* | \boldsymbol{\theta}^*), M(\boldsymbol{\theta}) - \boldsymbol{\theta}^* \rangle \leq \langle \nabla Q(M(\boldsymbol{\theta}) | \boldsymbol{\theta}^*) - \nabla Q(M(\boldsymbol{\theta}) | \boldsymbol{\theta}), M(\boldsymbol{\theta}) - \boldsymbol{\theta}^* \rangle. \quad (\text{A.3})$$

Now the  $\lambda$ -strong convexity condition implies that the left-hand side is lower bounded as (by letting  $\boldsymbol{\theta}_1 = M(\boldsymbol{\theta})$ ,  $\boldsymbol{\theta}_2 = \boldsymbol{\theta}^*$  and  $\boldsymbol{\theta}_1 = \boldsymbol{\theta}$ ,  $\boldsymbol{\theta}_2 = M(\boldsymbol{\theta})$  respectively)

$$\langle \nabla Q(M(\boldsymbol{\theta}) | \boldsymbol{\theta}^*) - \nabla Q(\boldsymbol{\theta}^* | \boldsymbol{\theta}^*), M(\boldsymbol{\theta}) - \boldsymbol{\theta}^* \rangle \geq \lambda \|\boldsymbol{\theta}^* - M(\boldsymbol{\theta})\|_2^2. \quad (\text{A.4})$$

On the other hand, the FOS( $\gamma$ ) condition together with the Cauchy-Schwarz inequality implies that the right-hand side upper bounded as

$$\langle \nabla Q(M(\boldsymbol{\theta}) | \boldsymbol{\theta}^*) - \nabla Q(M(\boldsymbol{\theta}) | \boldsymbol{\theta}), M(\boldsymbol{\theta}) - \boldsymbol{\theta}^* \rangle \leq \gamma \|\boldsymbol{\theta}^* - M(\boldsymbol{\theta})\|_2 \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2, \quad (\text{A.5})$$

Combining inequalities (A.4) and (A.5) with original bound (A.3) yields

$$\lambda \|\boldsymbol{\theta}^* - M(\boldsymbol{\theta})\|_2 \leq \gamma \|\boldsymbol{\theta}^* - M(\boldsymbol{\theta})\|_2 \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2, \quad (\text{A.6})$$

and canceling terms completes the proof.  $\square$

## Proof of Theorem 2

*Proof of Theorem 2.* From (3.6), for any fixed  $\boldsymbol{\theta} \in \mathbb{B}_2(r'; \boldsymbol{\theta}^*)$ ,

$$\|M_n(\boldsymbol{\theta}) - M(\boldsymbol{\theta})\|_2 \leq \epsilon_M(n, \rho)$$

with probability at least  $1 - \rho$ . It suffices to show that

$$\|\boldsymbol{\theta}_{m+1,n} - \boldsymbol{\theta}^*\|_2 \leq \kappa \|\boldsymbol{\theta}_{m,n} - \boldsymbol{\theta}^*\|_2 + \epsilon_M(n, \rho). \quad (\text{A.7})$$

Indeed, when this bound holds, we may iterate it to show that

$$\begin{aligned} \|\boldsymbol{\theta}_{m,n} - \boldsymbol{\theta}^*\|_2 &\leq \kappa \|\boldsymbol{\theta}_{m-1,n} - \boldsymbol{\theta}^*\|_2 + \epsilon_M(n, \rho) \\ &\leq \kappa \{ \kappa \|\boldsymbol{\theta}_{m-2,n} - \boldsymbol{\theta}^*\|_2 + \epsilon_M(n, \rho) \} + \epsilon_M(n, \rho) \\ &\leq \kappa^m \|\boldsymbol{\theta}_{0,n} - \boldsymbol{\theta}^*\|_2 + \left\{ \sum_{k=0}^{m-1} \kappa^k \right\} \epsilon_M(n, \rho) \\ &\leq \kappa^m \|\boldsymbol{\theta}_{0,n} - \boldsymbol{\theta}^*\|_2 + \frac{1}{1 - \kappa} \epsilon_M(n, \rho). \end{aligned}$$

It remains to prove (A.7), and we do so by induction. Beginning with  $m = 0$ , we have

$$\begin{aligned} \|\boldsymbol{\theta}_{1,n} - \boldsymbol{\theta}^*\|_2 &= \|M_n(\boldsymbol{\theta}_{0,n}) - \boldsymbol{\theta}^*\|_2 \leq \|M(\boldsymbol{\theta}_{0,n}) - \boldsymbol{\theta}^*\|_2 + \|M_n(\boldsymbol{\theta}_{0,n}) - M(\boldsymbol{\theta}_{0,n})\|_2 \\ &\leq \kappa \|\boldsymbol{\theta}_{0,n} - \boldsymbol{\theta}^*\|_2 + \epsilon_M(n, \rho) \\ &\leq \kappa r' + (1 - \kappa) r' = r'. \end{aligned}$$

In the induction from  $m \mapsto m + 1$ , suppose that  $\|\boldsymbol{\theta}_{m,n} - \boldsymbol{\theta}^*\|_2 \leq r$ , and the bound (A.7) at iteration  $m$ . Now at iteration  $m + 1$ ,

$$\begin{aligned} \|\boldsymbol{\theta}_{m+1,n} - \boldsymbol{\theta}^*\|_2 &= \|M_n(\boldsymbol{\theta}_{m,n}) - \boldsymbol{\theta}^*\|_2 \leq \|M(\boldsymbol{\theta}_{m,n}) - \boldsymbol{\theta}^*\|_2 + \|M_n(\boldsymbol{\theta}_{m,n}) - M(\boldsymbol{\theta}_{m,n})\|_2 \\ &\leq \kappa \|\boldsymbol{\theta}_{m,n} - \boldsymbol{\theta}^*\|_2 + \epsilon_M(n, \rho) \\ &\leq \kappa r' + (1 - \kappa)r' = r'. \end{aligned}$$

Thus  $\|\boldsymbol{\theta}_{m+1,n} - \boldsymbol{\theta}^*\|_2 \leq r'$ , and this completes the proof.  $\square$

### Proof of Theorem 3

(i) **Consistency.** Under (C0)–(C3) the DM surrogate decomposition  $Q_n(\boldsymbol{\theta}' \mid \boldsymbol{\theta}) = D_n(\boldsymbol{\theta}') + H_G(\boldsymbol{\theta}' \mid \boldsymbol{\theta})$  with  $H_G \geq 0$  yields  $D_n(\boldsymbol{\theta}_{m+1,n}) \leq D_n(\boldsymbol{\theta}_{m,n})$  (monotone descent). By sublevel compactness the sequence  $\{\boldsymbol{\theta}_{m,n}\}$  is tight for each  $n$ ; any limit point is stationary. Uniform convergence and identification (C0)–(C3) imply the stationary set concentrates at  $\boldsymbol{\theta}^*$  as  $n \rightarrow \infty$ ; hence  $\lim_n \lim_m \boldsymbol{\theta}_{m,n} = \boldsymbol{\theta}^*$  a.s. If  $m_n \rightarrow \infty$  then  $\boldsymbol{\theta}_{m_n,n} \xrightarrow{p} \boldsymbol{\theta}^*$ .

(ii)  **$\sqrt{n}$ –normality for truncated iterates.** Under Theorem 1 there exists  $\kappa \in (0, 1)$  such that  $M$  is  $\kappa$ –contractive on  $B_2(r; \boldsymbol{\theta}^*)$ . Let  $\hat{\boldsymbol{\theta}}_n := \arg \min_{\boldsymbol{\theta}} D_n(\boldsymbol{\theta})$  and note that  $\hat{\boldsymbol{\theta}}_n \in M_n(\hat{\boldsymbol{\theta}}_n)$  (unique by the fixed-order regularity). Under (F1), (C1)–(C2), (K1)–(K2), (M1)–(M8) the fixed-order MDE satisfies

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*) \Rightarrow \mathcal{N}(0, I(\boldsymbol{\theta}^*)^{-1}).$$

Moreover, the sample update map  $M_n$  inherits the contraction locally with factor  $\kappa_n \leq \kappa + o_p(1)$  (same FOS/strong-convexity argument with  $g$  replaced by  $g_n$  and uniform LLN). Hence for any initialization  $\boldsymbol{\theta}_{0,n} \in B_2(r; \boldsymbol{\theta}^*)$ ,

$$\|\boldsymbol{\theta}_{t,n} - \hat{\boldsymbol{\theta}}_n\|_2 \leq \kappa_n^t \|\boldsymbol{\theta}_{0,n} - \hat{\boldsymbol{\theta}}_n\|_2 \quad (\text{w.p.} \rightarrow 1).$$

Choose  $m_n \geq \lceil (\frac{1}{2} \log n + c_0) / |\log \kappa| \rceil$  so that  $\sqrt{n} \kappa^{m_n} \rightarrow 0$  and therefore  $\sqrt{n} \kappa_n^{m_n} \rightarrow 0$ .

Decompose

$$\sqrt{n}(\boldsymbol{\theta}_{m_n,n} - \boldsymbol{\theta}^*) = \underbrace{\sqrt{n}(\boldsymbol{\theta}_{m_n,n} - \hat{\boldsymbol{\theta}}_n)}_{\rightarrow 0 \text{ in prob.}} + \underbrace{\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*)}_{\Rightarrow \mathcal{N}(0, I(\boldsymbol{\theta}^*)^{-1})}.$$

The first term is  $o_p(1)$  by the geometric bound above, the second term is the fixed-order CLT. Slutsky's lemma gives  $\sqrt{n}(\boldsymbol{\theta}_{m_n,n} - \boldsymbol{\theta}^*) \Rightarrow \mathcal{N}(0, I(\boldsymbol{\theta}^*)^{-1})$ . Finally, the displayed choice implies  $m_n = \lceil (\frac{1}{2} \log n) / |\log \kappa| \rceil + O(1)$ .  $\square$

### Proof of Corollary 3 (Finite-step Godambe CLT)

Let  $\hat{\boldsymbol{\theta}}_n$  solve  $\Psi_n(\boldsymbol{\theta}) := \nabla_{\boldsymbol{\theta}} D_G(g_n, f_{\boldsymbol{\theta}}) = 0$  and let  $\boldsymbol{\theta}_{t+1,n} \in M_n(\boldsymbol{\theta}_{t,n})$  with  $\boldsymbol{\theta}_{0,n} \in B_2(r; \boldsymbol{\theta}^*)$ . Under (G1)–(G4) at  $\boldsymbol{\theta}^\dagger := \arg \min_{\boldsymbol{\theta}} D_G(g, f_{\boldsymbol{\theta}})$ , the Z-estimation CLT gives  $\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\dagger) \Rightarrow \mathcal{N}(0, H^{-1} V H^{-1})$  with  $H := \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}^\dagger; g)$ ,  $V := \text{Var}_g[A'(\frac{g}{f_{\boldsymbol{\theta}^\dagger}} - 1) s_{\boldsymbol{\theta}^\dagger}(Y)]$ . By Theorem 1 and the same uniformity argument as in Theorem 1 (ii), the sample update map contracts locally with factor  $\kappa_n \leq \kappa + o_p(1)$ , hence  $\|\boldsymbol{\theta}_{m_n,n} - \hat{\boldsymbol{\theta}}_n\|_2 \leq \kappa_n^{m_n} \|\boldsymbol{\theta}_{0,n} - \hat{\boldsymbol{\theta}}_n\|_2$ . With the threshold  $\sqrt{n} \kappa^{m_n} \rightarrow 0$ , we get  $\sqrt{n} \|\boldsymbol{\theta}_{m_n,n} - \hat{\boldsymbol{\theta}}_n\|_2 \rightarrow 0$  in probability. Decompose

$$\sqrt{n}(\boldsymbol{\theta}_{m_n,n} - \boldsymbol{\theta}^\dagger) = \sqrt{n}(\boldsymbol{\theta}_{m_n,n} - \hat{\boldsymbol{\theta}}_n) + \sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\dagger),$$

apply Slutsky's lemma, and conclude  $\sqrt{n}(\boldsymbol{\theta}_{m_n,n} - \boldsymbol{\theta}^\dagger) \Rightarrow \mathcal{N}(0, H^{-1} V H^{-1})$ . At the model ( $g = f_{\boldsymbol{\theta}^*}$ ) and with  $A'(0) = 1$ ,  $H = V = I(\boldsymbol{\theta}^*)$ , so the covariance is  $I(\boldsymbol{\theta}^*)^{-1}$ .  $\square$

### Proof of Corollary 4 (Finite-step (Godambe-Wilks))

Under the assumptions of Corollary 4 and  $\sqrt{n} \kappa^{m_n} \rightarrow 0$ , a quadratic expansion of  $D_G(g_n, f_{\boldsymbol{\theta}})$  at  $\hat{\boldsymbol{\theta}}_n$  yields

$$2n\{D_G(g_n, f_{\boldsymbol{\theta}^\dagger}) - D_G(g_n, f_{\boldsymbol{\theta}_{m_n,n}})\} = n(\boldsymbol{\theta}_{m_n,n} - \boldsymbol{\theta}^\dagger)^\top H (\boldsymbol{\theta}_{m_n,n} - \boldsymbol{\theta}^\dagger) + o_p(1),$$

since  $n\|\boldsymbol{\theta}_{m_n,n} - \hat{\boldsymbol{\theta}}_n\|_2^2 = (\sqrt{n} \kappa^{m_n})^2 = o_p(1)$ . Let  $J := H^{-1/2} V H^{-1/2}$  with eigenvalues  $\{\lambda_j\}_{j=1}^{p(K_0)}$ . By Corollary 3,  $n(\boldsymbol{\theta}_{m_n,n} - \boldsymbol{\theta}^\dagger)^\top H (\boldsymbol{\theta}_{m_n,n} - \boldsymbol{\theta}^\dagger) \Rightarrow \sum_{j=1}^{p(K_0)} \lambda_j \chi_{1,j}^2$ . If  $g = f_{\boldsymbol{\theta}^*}$  and

$A'(0) = 1$ , then  $H = V = I(\boldsymbol{\theta}^*)$  and the limit is  $\chi_{p(K_0)}^2$ .  $\square$

## Proof of Corollary 5 (Population contraction under contamination)

Define  $\Psi_\epsilon(\boldsymbol{\theta}; g) := \nabla_{\boldsymbol{\theta}} D_G(g, f_{\boldsymbol{\theta}}) = \int f_{\boldsymbol{\theta}}(y) s_{\boldsymbol{\theta}}(y) B(g(y)/f_{\boldsymbol{\theta}}(y)) dy$ ,  $s_{\boldsymbol{\theta}} = \nabla_{\boldsymbol{\theta}} \log f_{\boldsymbol{\theta}}$ ,  $B(u) = G(u) - uG'(u)$ . By (F1), (C1)–(C2), (K1)–(K2), (M1)–(M8), the maps  $\boldsymbol{\theta} \mapsto \nabla_{\boldsymbol{\theta}} \Psi_\epsilon(\boldsymbol{\theta}; g_\epsilon)$  and the Gâteaux derivative  $h \mapsto \partial_g \Psi_\epsilon(\boldsymbol{\theta}; g_\epsilon)[h]$  are continuous in  $(\boldsymbol{\theta}, g_\epsilon)$  on  $B_2(r; \boldsymbol{\theta}^*) \times \{g_\epsilon : \epsilon \in [0, \epsilon_0]\}$  (see Appendix N, Gâteaux lemma). At  $\epsilon = 0$ , Theorem 1 provides FOS( $\gamma$ ) and local curvature  $\lambda > 0$  on  $B_2(r_0; \boldsymbol{\theta}^*)$  with  $\kappa = \gamma/\lambda < 1$ . By continuity in  $\epsilon$ , there exist  $r \in (0, r_0]$  and  $\bar{\kappa} \in (0, 1)$  such that, for all  $\epsilon \in [0, \epsilon_0]$ , FOS( $\gamma_\epsilon$ ) and local  $\lambda_\epsilon$ -strong convexity of  $D_G(g_\epsilon, \cdot)$  hold on  $B_2(r; \boldsymbol{\theta}_\epsilon^\dagger)$  and  $\gamma_\epsilon/\lambda_\epsilon \leq \bar{\kappa} < 1$ . Let  $\boldsymbol{\theta}_+ \in M_\epsilon(\boldsymbol{\theta})$  be any minimizer of  $Q_\epsilon(\cdot | \boldsymbol{\theta})$ . By strong convexity of  $q_\epsilon(\cdot) := Q_\epsilon(\cdot | \boldsymbol{\theta}_\epsilon^\dagger)$  and its optimality at  $\boldsymbol{\theta}_\epsilon^\dagger$ ,

$$\lambda_\epsilon \|\boldsymbol{\theta}_+ - \boldsymbol{\theta}_\epsilon^\dagger\|_2^2 \leq \langle \nabla_{\boldsymbol{\theta}'} q_\epsilon(\boldsymbol{\theta}_+) - \nabla_{\boldsymbol{\theta}'} q_\epsilon(\boldsymbol{\theta}_\epsilon^\dagger), \boldsymbol{\theta}_+ - \boldsymbol{\theta}_\epsilon^\dagger \rangle.$$

Add and subtract  $\nabla_{\boldsymbol{\theta}'} Q_\epsilon(\boldsymbol{\theta}_+ | \boldsymbol{\theta})$  and use the optimality condition  $\langle \nabla_{\boldsymbol{\theta}'} Q_\epsilon(\boldsymbol{\theta}_+ | \boldsymbol{\theta}), \boldsymbol{\theta}_\epsilon^\dagger - \boldsymbol{\theta}_+ \rangle \geq 0$  to get

$$\lambda_\epsilon \|\boldsymbol{\theta}_+ - \boldsymbol{\theta}_\epsilon^\dagger\| \leq \left\| \nabla_{\boldsymbol{\theta}'} Q_\epsilon(\boldsymbol{\theta}_+ | \boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}'} Q_\epsilon(\boldsymbol{\theta}_+ | \boldsymbol{\theta}_\epsilon^\dagger) \right\| \leq \gamma_\epsilon \|\boldsymbol{\theta} - \boldsymbol{\theta}_\epsilon^\dagger\|,$$

where the last step is FOS( $\gamma_\epsilon$ ). Hence  $\|\boldsymbol{\theta}_+ - \boldsymbol{\theta}_\epsilon^\dagger\| \leq (\gamma_\epsilon/\lambda_\epsilon) \|\boldsymbol{\theta} - \boldsymbol{\theta}_\epsilon^\dagger\| \leq \bar{\kappa} \|\boldsymbol{\theta} - \boldsymbol{\theta}_\epsilon^\dagger\|$ , proving the contraction.  $\square$

## Proof of Corollary 6 (Noisy contraction and opt-to-stat under contamination)

Let  $\boldsymbol{\theta}_{t+1, \epsilon, n} \in M_{\epsilon, n}(\boldsymbol{\theta}_{t, \epsilon, n})$  with  $\boldsymbol{\theta}_{t, \epsilon, n} \in B_2(r; \boldsymbol{\theta}_\epsilon^\dagger)$ . Choose  $\boldsymbol{\eta}_t \in M_\epsilon(\boldsymbol{\theta}_{t, \epsilon, n})$  so that  $\|\boldsymbol{\theta}_{t+1, \epsilon, n} - \boldsymbol{\eta}_t\|_2 \leq \mathcal{M}_{\text{unif}}(n, r; \epsilon)$ . By Corollary 5 (population contraction under contamination),

$$\|\boldsymbol{\eta}_t - \boldsymbol{\theta}_\epsilon^\dagger\|_2 \leq \bar{\kappa} \|\boldsymbol{\theta}_{t, \epsilon, n} - \boldsymbol{\theta}_\epsilon^\dagger\|_2.$$

Therefore

$$\|\boldsymbol{\theta}_{t+1,\epsilon,n} - \boldsymbol{\theta}_\epsilon^\dagger\|_2 \leq \|\boldsymbol{\theta}_{t+1,\epsilon,n} - \boldsymbol{\eta}_t\|_2 + \|\boldsymbol{\eta}_t - \boldsymbol{\theta}_\epsilon^\dagger\|_2 \leq \bar{\kappa} \|\boldsymbol{\theta}_{t,\epsilon,n} - \boldsymbol{\theta}_\epsilon^\dagger\|_2 + \mathcal{M}_{\text{unif}}(n, r; \epsilon),$$

establishing the contaminated noisy recursion.

*Bounding  $\mathcal{M}_{\text{unif}}(n, r; \epsilon)$ .* Fix  $\boldsymbol{\theta} \in B_2(r; \boldsymbol{\theta}_\epsilon^\dagger)$ . By strong convexity of  $Q_\epsilon(\cdot \mid \boldsymbol{\theta})$  on the ball with modulus  $\lambda_\epsilon$  and the optimality conditions for  $M_\epsilon(\boldsymbol{\theta})$  and  $M_{\epsilon,n}(\boldsymbol{\theta})$ , the minimizer map is locally inverse-Lipschitz:

$$\text{dist}(M_{\epsilon,n}(\boldsymbol{\theta}), M_\epsilon(\boldsymbol{\theta})) \leq \frac{1}{\lambda_\epsilon} \sup_{\boldsymbol{\eta} \in B_2(r; \boldsymbol{\theta}_\epsilon^\dagger)} \|\nabla_{\boldsymbol{\theta}'} Q_{\epsilon,n}(\boldsymbol{\eta} \mid \boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}'} Q_\epsilon(\boldsymbol{\eta} \mid \boldsymbol{\theta})\|. \quad (\text{A.8})$$

Using the score representation and the mean-value inequality for the RAF derivative, for any  $\boldsymbol{\eta} \in B_2(r; \boldsymbol{\theta}_\epsilon^\dagger)$ ,

$$\begin{aligned} \|\nabla_{\boldsymbol{\theta}'} Q_{\epsilon,n}(\boldsymbol{\eta} \mid \boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}'} Q_\epsilon(\boldsymbol{\eta} \mid \boldsymbol{\theta})\| &\leq A'_{\max} \sup_{\boldsymbol{\zeta} \in B_2(r; \boldsymbol{\theta}_\epsilon^\dagger)} \|s_{\boldsymbol{\zeta}}\|_{\mathcal{H}^*} \|g_{\epsilon,n} - g_\epsilon\|_{\mathcal{H}} \\ &\leq C_{fos} \text{Env}(K) A'_{\max} \|g_{\epsilon,n} - g_\epsilon\|_{\mathcal{H}}, \end{aligned}$$

uniformly on the ball (the last inequality is the same envelope bound used in the main after Theorem 2, adapted to  $g_\epsilon$ ). Combining with (A.8) and taking the sup over  $\boldsymbol{\theta}$  gives the contaminated analogue of (3.7):

$$\mathcal{M}_{\text{unif}}(n, r; \epsilon) \lesssim \frac{C_{fos} A'_{\max} \text{Env}(K)}{\lambda_\epsilon} \|g_{\epsilon,n} - g_\epsilon\|_{\mathcal{H}}. \quad (\text{A.9})$$

Under the assumption  $\|g_{\epsilon,n} - g_\epsilon\|_{\mathcal{H}} = o_p(n^{-1/2})$  and the uniform boundedness of  $\lambda_\epsilon^{-1}$  for  $\epsilon \in [0, \epsilon_0]$ , we obtain  $\mathcal{M}_{\text{unif}}(n, r; \epsilon) = o_p(n^{-1/2})$ .

*Conclusion (opt-to-stat).* Iterating the noisy recursion and using Lemma 12 yields

$$\|\boldsymbol{\theta}_{\epsilon,n}^{(m_n)} - \hat{\boldsymbol{\theta}}_{\epsilon,n}\|_2 \leq \bar{\kappa}^{m_n} \|\boldsymbol{\theta}_{\epsilon,n}^{(0)} - \hat{\boldsymbol{\theta}}_{\epsilon,n}\|_2 + \frac{1 - \bar{\kappa}^{m_n}}{1 - \bar{\kappa}} \mathcal{M}_{\text{unif}}(n, r; \epsilon).$$

Multiplying by  $\sqrt{n}$  and using  $\sqrt{n} \bar{\kappa}^{m_n} \rightarrow 0$  together with  $\sqrt{n} \mathcal{M}_{\text{unif}}(n, r; \epsilon) = o_p(1)$  gives

$$\sqrt{n} \|\boldsymbol{\theta}_{\epsilon,n}^{(m_n)} - \hat{\boldsymbol{\theta}}_{\epsilon,n}\|_2 \rightarrow^p 0. \quad \square$$

## S.R.1 Robust contaminated noisy contraction for bounded-RAF class

**Corollary 8** (Robust class (NED/vNED)). *Assume  $A'_{\max} := \sup_{\delta \geq -1} |A'(\delta)| < \infty$  and  $A'$  is nonincreasing on  $[0, \infty)$  (e.g., NED, vNED). Then there exist  $r > 0$  and  $\bar{\kappa} \in (0, 1)$  such that, for every  $\varepsilon \in [0, \varepsilon_0]$  and any selection  $\boldsymbol{\theta}_{t+1, \varepsilon, n} \in M_{\varepsilon, n}(\boldsymbol{\theta}_{t, \varepsilon, n})$  with  $\boldsymbol{\theta}_{t, \varepsilon, n} \in B_2(r; \boldsymbol{\theta}_\varepsilon^\dagger)$ ,*

$$\|\boldsymbol{\theta}_{t+1, \varepsilon, n} - \boldsymbol{\theta}_\varepsilon^\dagger\|_2 \leq \bar{\kappa} \|\boldsymbol{\theta}_{t, \varepsilon, n} - \boldsymbol{\theta}_\varepsilon^\dagger\|_2 + \mathcal{M}_{\text{unif}}(n, r; \varepsilon),$$

and, if  $\|g_{\varepsilon, n} - g_\varepsilon\|_{\mathcal{H}} = o_p(n^{-1/2})$ , then  $\mathcal{M}_{\text{unif}}(n, r; \varepsilon) = o_p(n^{-1/2})$ ; consequently any  $m_n \geq \lceil (\frac{1}{2} \log n + c_0) / |\log \bar{\kappa}| \rceil$  satisfies  $\sqrt{n} \|\boldsymbol{\theta}_{\varepsilon, n}^{(m_n)} - \hat{\boldsymbol{\theta}}_{\varepsilon, n}\| \rightarrow^p 0$ .

*Proof.* (1) *Inverse-Lipschitz for the argmin map.* Strong convexity of  $Q_\varepsilon(\cdot | \boldsymbol{\theta})$  with modulus  $\lambda_\varepsilon$  on  $B_2(r; \boldsymbol{\theta}_\varepsilon^\dagger)$  yields, for every  $\boldsymbol{\theta}$  in the ball,

$$\text{dist}(M_{\varepsilon, n}(\boldsymbol{\theta}), M_\varepsilon(\boldsymbol{\theta})) \leq \frac{1}{\lambda_\varepsilon} \sup_{\boldsymbol{\eta} \in B_2(r; \boldsymbol{\theta}_\varepsilon^\dagger)} \|\nabla_{\boldsymbol{\theta}'} Q_{\varepsilon, n}(\boldsymbol{\eta} | \boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}'} Q_\varepsilon(\boldsymbol{\eta} | \boldsymbol{\theta})\|. \quad (\text{A.10})$$

(2) *Gradient perturbation with bounded RAF derivative.* Using the score representation and the mean-value form,

$$\sup_{\boldsymbol{\eta} \in B_2(r; \boldsymbol{\theta}_\varepsilon^\dagger)} \|\nabla_{\boldsymbol{\theta}'} Q_{\varepsilon, n}(\boldsymbol{\eta} | \boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}'} Q_\varepsilon(\boldsymbol{\eta} | \boldsymbol{\theta})\| \leq C_{\text{fos}} A'_{\max} \text{Env}(K) \|g_{\varepsilon, n} - g_\varepsilon\|_{\mathcal{H}},$$

uniformly on the ball (bounded, nonincreasing  $A'$  allows the envelope to be absorbed).

Combining with (A.10) and taking sup over  $\boldsymbol{\theta}$  yields

$$\mathcal{M}_{\text{unif}}(n, r; \varepsilon) \lesssim \frac{C_{\text{fos}} A'_{\max} \text{Env}(K)}{\lambda_\varepsilon} \|g_{\varepsilon, n} - g_\varepsilon\|_{\mathcal{H}}. \quad (\text{A.11})$$

Under the plug-in rate and a uniform curvature lower bound  $\inf_{\varepsilon \leq \varepsilon_0} \lambda_\varepsilon > 0$ , we get  $\mathcal{M}_{\text{unif}} = o_p(n^{-1/2})$ . Noisy recursion + Lemma (noisy linear recurrence) imply the opt-to-stat bound.

□



## S.R.2 KL requires a local floor: failure mode without it

**Proposition 6** (KL needs a local density/score floor). *Let  $G(u) = u \log u - u + 1$  (KL). Assume only the standing smoothness on  $B_2(r; \boldsymbol{\theta}_\varepsilon^\dagger)$  but no local lower density/score bound. Then there exist contamination sequences  $\{\eta_m\}$  with  $\|g_{\varepsilon,n} - g_\varepsilon\|_{\mathcal{H}} = O_p(n^{-1/2})$  for which  $\mathcal{M}_{\text{unif}}(n, r; \varepsilon) \neq o_p(n^{-1/2})$ ; hence the opt-to-stat step may fail. If, in addition, a local density/score floor holds (e.g.  $\inf_{\boldsymbol{\theta} \in B_2(r; \boldsymbol{\theta}_\varepsilon^\dagger)} \inf_{y \in \mathcal{Y}_r} f(y; \boldsymbol{\theta}) \geq c > 0$ ), then (A.11) applies with  $A'_{\max} = 1$  and the contaminated noisy contraction and opt-to-stat bounds hold for KL as well.*

*Proof.* For KL,  $A'(\delta) \equiv 1$ , so the bound (A.11) hinges on (i) a uniform curvature lower bound  $\lambda_\varepsilon > 0$  and (ii) a finite score envelope  $\text{Env}(K)$  on the ball. Without a local floor one can pick sets  $A_m$  with  $g_\varepsilon(A_m) > 0$  and  $\sup_{y \in A_m} f(y; \boldsymbol{\theta}) \rightarrow 0$  for all  $\boldsymbol{\theta}$  in the ball; letting  $\eta_m$  concentrate on  $A_m$  and taking  $g_{\varepsilon,n}^{(m)}$  with  $\|g_{\varepsilon,n}^{(m)} - g_\varepsilon^{(m)}\|_{\mathcal{H}} = O_p(n^{-1/2})$ , the score envelope diverges and the supremum in (A.10) cannot be  $O_p(n^{-1/2})$  uniformly in  $m$ . Equivalently  $\mathcal{M}_{\text{unif}}(n, r; \varepsilon) \neq o_p(n^{-1/2})$ , so the noisy-contraction/opt-to-stat step may fail. With a local floor, the same argument as in S.R.1 goes through with  $A'_{\max} = 1$ .  $\square$

## Proof of Theorem 4

*Proof of Theorem 4.* For clarity fix the target (contaminated) density at  $\boldsymbol{\theta}^*$ :

$$g_{\varepsilon,n}(y) := f_{\varepsilon,n}(y; \boldsymbol{\theta}^*) = (1 - \varepsilon) f(y; \boldsymbol{\theta}^*) + \varepsilon \eta_n(y),$$

and recall the DM population functional  $T(g) := \arg \min_{\boldsymbol{\theta} \in \Theta} D_G(g, f_{\boldsymbol{\theta}})$  (defined as a singleton under uniqueness). Thus  $\boldsymbol{\theta}_{\varepsilon,n}^* = T(g_{\varepsilon,n})$  and  $\boldsymbol{\theta}_\varepsilon^* = T(g_\varepsilon)$ , where  $g_\varepsilon := (1 - \varepsilon) f(\cdot; \boldsymbol{\theta}^*) + \varepsilon \eta$  is the  $n \rightarrow \infty$  limit when it exists.

**(1) Boundedness and convergence for fixed  $\varepsilon$ .** Assume (C2) and (O2) hold uniformly

in  $n$  and  $\epsilon \in [0, \epsilon_0)$ , and that for the fixed  $\epsilon$  the minimizer  $\boldsymbol{\theta}_{\epsilon,n}^*$  is unique for all  $n \geq 1$ . Set  $\phi_{\epsilon,n}(\boldsymbol{\theta}) := D_G(g_{\epsilon,n}, f_{\boldsymbol{\theta}})$  and  $\phi_{\epsilon}(\boldsymbol{\theta}) := D_G(g_{\epsilon}, f_{\boldsymbol{\theta}})$ . By linearity of  $D_G$  in its first argument and the regularity in **(C2)**–**(O2)**, we have

$$\sup_{\boldsymbol{\theta} \in \Theta} |\phi_{\epsilon,n}(\boldsymbol{\theta}) - \phi_{\epsilon}(\boldsymbol{\theta})| \longrightarrow 0 \quad (n \rightarrow \infty),$$

and the sublevel sets  $\{\boldsymbol{\theta} : \phi_{\epsilon,n}(\boldsymbol{\theta}) \leq c\}$  are compact uniformly in  $n$  (equicoercivity). Hence, by the argmin-continuity/Berge maximum theorem (or van der Vaart's Thm 5.7), the sequence of minimizers is bounded and

$$\boldsymbol{\theta}_{\epsilon,n}^* = \arg \min \phi_{\epsilon,n} \longrightarrow \arg \min \phi_{\epsilon} = \boldsymbol{\theta}_{\epsilon}^*.$$

**(2) Gâteaux derivative at the model along the mixture direction.** Let  $\Psi(\boldsymbol{\theta}; g) := \nabla_{\boldsymbol{\theta}} D_G(g, f_{\boldsymbol{\theta}})$  and note that  $T(g)$  is characterized by  $\Psi(T(g); g) = \mathbf{0}$ . Under **(M1)**–**(M2)** (differentiability in  $\boldsymbol{\theta}$  and Hadamard/Gâteaux differentiability in  $g$ ), the implicit function map  $g \mapsto T(g)$  is differentiable at  $(\boldsymbol{\theta}^*, g_0)$  with  $g_0 = f(\cdot; \boldsymbol{\theta}^*)$  and nonsingular  $H := \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}^*; g_0)$ :

$$DT_{g_0}[h] = -H^{-1} \partial_g \Psi(\boldsymbol{\theta}^*; g_0)[h].$$

Calibration  $A'(0) = 1$  and evaluation at the model give the well-known score form  $\partial_g \Psi(\boldsymbol{\theta}^*; g_0)[h] = -\int u(y; \boldsymbol{\theta}^*) h(y) dy$ , where  $u(y; \boldsymbol{\theta}) := \nabla_{\boldsymbol{\theta}} \log f(y; \boldsymbol{\theta})$ , and  $H = I(\boldsymbol{\theta}^*)$  (Fisher information). Consider the contamination path

$$g_{\epsilon,n} = (1 - \epsilon) g_0 + \epsilon \eta_n \quad \Rightarrow \quad \dot{g}_{\epsilon} \Big|_{\epsilon=0} = \eta_n - g_0.$$

Therefore

$$DT_{g_0}[\eta_n - g_0] = I(\boldsymbol{\theta}^*)^{-1} \int u(y; \boldsymbol{\theta}^*) (\eta_n(y) - g_0(y)) dy.$$

Since  $\int u(y; \boldsymbol{\theta}^*) g_0(y) dy = \mathbf{0}$ , we obtain

$$\lim_{\epsilon \downarrow 0} \frac{\boldsymbol{\theta}_{\epsilon,n}^* - \boldsymbol{\theta}^*}{\epsilon} = DT_{g_0}[\eta_n - g_0] = I(\boldsymbol{\theta}^*)^{-1} \int \eta_n(y) u(y; \boldsymbol{\theta}^*) dy.$$

This is the claimed influence-function expression. □

## Proof of Theorem 5

*Proof of Theorem 5.* Write  $\hat{\boldsymbol{\theta}}_{\epsilon,n} := \arg \min_{\boldsymbol{\theta}} D_G(g_{\epsilon,n}, f_{\boldsymbol{\theta}})$  and  $\boldsymbol{\theta}_{\epsilon}^{\dagger} := \arg \min_{\boldsymbol{\theta}} D_G(g_{\epsilon}, f_{\boldsymbol{\theta}})$ . Under (G1)–(G4) (uniformly in  $\epsilon \in [0, \epsilon_0]$ ), the Z-estimation CLT gives

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_{\epsilon,n} - \boldsymbol{\theta}_{\epsilon}^{\dagger}) \Rightarrow \mathcal{N}(0, H_{\epsilon}^{-1} V_{\epsilon} H_{\epsilon}^{-1}),$$

with  $H_{\epsilon} := \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}_{\epsilon}^{\dagger}; g_{\epsilon})$  and  $V_{\epsilon} := \text{Var}_{g_{\epsilon}}[A'(\frac{g_{\epsilon}}{f_{\boldsymbol{\theta}_{\epsilon}^{\dagger}}} - 1) s_{\boldsymbol{\theta}_{\epsilon}^{\dagger}}(Y)]$ . By the population contraction (Theorem 1) and the uniform plug-in rate  $\|g_{\epsilon,n} - g_{\epsilon}\|_{\mathcal{H}} = o_p(n^{-1/2})$ , the sample update map contracts locally with factor  $\kappa_n \leq \kappa + o_p(1)$  and the noisy deviation is  $o_p(n^{-1/2})$  (since  $\mathcal{M}_{\text{unif}}(n, r) \lesssim C_{fos} A'_{\max} \text{Env}(K) \|g_{\epsilon,n} - g_{\epsilon}\|_{\mathcal{H}}$ ). Hence, for any initialization in  $B_2(r; \boldsymbol{\theta}^{\star})$ ,

$$\|\boldsymbol{\theta}_{\epsilon,n}^{(m)} - \hat{\boldsymbol{\theta}}_{\epsilon,n}\|_2 \leq \kappa_n^m C_r \Rightarrow \sqrt{n} \|\boldsymbol{\theta}_{\epsilon,n}^{(m_n)} - \hat{\boldsymbol{\theta}}_{\epsilon,n}\|_2 \leq C_r \sqrt{n} \kappa_n^{m_n} \xrightarrow{p} 0,$$

because  $\sqrt{n} \kappa^{m_n} \rightarrow 0$  and  $\kappa_n \rightarrow \kappa$  in probability. Decompose

$$\sqrt{n}(\boldsymbol{\theta}_{\epsilon,n}^{(m_n)} - \boldsymbol{\theta}_{\epsilon}^{\dagger}) = \sqrt{n}(\boldsymbol{\theta}_{\epsilon,n}^{(m_n)} - \hat{\boldsymbol{\theta}}_{\epsilon,n}) + \sqrt{n}(\hat{\boldsymbol{\theta}}_{\epsilon,n} - \boldsymbol{\theta}_{\epsilon}^{\dagger}),$$

apply the previous display and the Z-estimation CLT, and conclude by Slutsky's lemma.

At the model ( $\epsilon = 0$ ) with  $A'(0) = 1$ ,  $H_{\epsilon} = V_{\epsilon} = I(\boldsymbol{\theta}^{\star})$  and the covariance reduces to  $I(\boldsymbol{\theta}^{\star})^{-1}$ .  $\square$

## Proof of Theorem 6

*Proof of Theorem 6.* Under the assumptions of Theorem 5, the contraction bound yields

$$\|\boldsymbol{\theta}_{\epsilon,n}^{(m_n)} - \hat{\boldsymbol{\theta}}_{\epsilon,n}\|_2 \leq \kappa_n^{m_n} C_r \text{ with } \sqrt{n} \kappa^{m_n} \rightarrow 0, \text{ hence } \sqrt{n} \|\boldsymbol{\theta}_{\epsilon,n}^{(m_n)} - \hat{\boldsymbol{\theta}}_{\epsilon,n}\|_2 \rightarrow 0 \text{ in probability.}$$

Decompose

$$\sqrt{n}(\boldsymbol{\theta}_{\epsilon,n}^{(m_n)} - \boldsymbol{\theta}_{\epsilon}^{\dagger}) = \sqrt{n}(\boldsymbol{\theta}_{\epsilon,n}^{(m_n)} - \hat{\boldsymbol{\theta}}_{\epsilon,n}) + \sqrt{n}(\hat{\boldsymbol{\theta}}_{\epsilon,n} - \boldsymbol{\theta}_{\epsilon}^{\dagger}),$$

and invoke the Z-estimation CLT for the second term (Theorem 5 already assumes (G1)–(G4)). Slutsky's lemma yields the same Gaussian limit.  $\square$

## S.BD.1 Local breakdown lower bound for unbounded RAF under a density–ratio cap

**Corollary 9.** *Assume the setup of Theorem 7 except that (ii) is replaced by: (ii<sub>Γ</sub>) Local envelope under density–ratio cap. There exists  $\Gamma \in (0, \infty)$  such that, for all  $\boldsymbol{\theta} \in B_2(r; \boldsymbol{\theta}^*)$  and all contaminations  $q$  under consideration,*

$$0 \leq \frac{q(y)}{f(y; \boldsymbol{\theta})} - 1 \leq \Gamma \quad \text{for } f(\cdot; \boldsymbol{\theta})\text{-a.e. } y,$$

and  $\|\nabla_{\boldsymbol{\theta}} D_G(q, f_{\boldsymbol{\theta}})\| \leq S_K A_{\Gamma}$  with  $A_{\Gamma} := \sup_{\delta \in [-1, \Gamma]} |A(\delta)| < \infty$ .

Then the conclusion of Theorem 7 holds with  $\epsilon^{\dagger}$  replaced by

$$\epsilon^{\dagger}(\Gamma) := \frac{\lambda r}{S_K A_{\Gamma}}.$$

*Proof.* Repeat the proof of Theorem 7 replacing  $A_{\max}$  by  $A_{\Gamma}$  via (ii<sub>Γ</sub>). All constants are uniform on  $B_2(r; \boldsymbol{\theta}^*)$  by (i) and (iii), hence  $D_G((1 - \varepsilon)g + \varepsilon q, f_{\boldsymbol{\theta}}) - D_G(g, f_{\boldsymbol{\theta}}) \geq \varepsilon \|\nabla D_G(q, f_{\boldsymbol{\theta}})\| \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| - \frac{\lambda}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|^2$ , which yields  $\|\hat{\boldsymbol{\theta}}_{\varepsilon} - \boldsymbol{\theta}^*\| \leq \varepsilon S_K A_{\Gamma} / \lambda$  for any  $\varepsilon < \lambda r / (S_K A_{\Gamma})$  and ensures  $\hat{\boldsymbol{\theta}}_{\varepsilon} \in B_2(r; \boldsymbol{\theta}^*)$ .  $\square$

**Examples of  $A_{\Gamma}$ .** Hellinger:  $A(\delta) = 2(\sqrt{1 + \delta} - 1)$ , so  $A_{\Gamma} = 2(\sqrt{1 + \Gamma} - 1)$ . KL:  $A(\delta) = \delta$ , so  $A_{\Gamma} = \max\{1, \Gamma - 1\}$ .

## S.BD.2 Failure of any uniform breakdown lower bound for unbounded RAF without a local floor

**Proposition 7.** *Let  $G$  be an unbounded–RAF generator (e.g., KL with  $A(\delta) = \delta$  or HD with  $A(\delta) = 2(\sqrt{1 + \delta} - 1)$ ). Assume (i) and (iii) of Theorem 7 hold on  $B_2(r; \boldsymbol{\theta}^*)$ , but no local density/score floor is imposed; i.e., for every  $\Gamma < \infty$  there exists a probability  $q_{\Gamma}$  and  $\boldsymbol{\theta}_{\Gamma} \in B_2(r; \boldsymbol{\theta}^*)$  with  $\text{ess sup}_y \frac{q_{\Gamma}(y)}{f(y; \boldsymbol{\theta}_{\Gamma})} \geq \Gamma$ . Then for every  $\varepsilon_0 > 0$  and every  $C > 0$  there exists*

$\Gamma$  and a contamination  $g_\varepsilon^{(\Gamma)} = (1 - \varepsilon)g + \varepsilon q_\Gamma$  with some  $\varepsilon \in (0, \varepsilon_0]$  such that the  $\varepsilon$ -minimizer  $\hat{\boldsymbol{\theta}}_\varepsilon$  of  $D_G(g_\varepsilon^{(\Gamma)}, f_{\boldsymbol{\theta}})$  satisfies

$$\|\hat{\boldsymbol{\theta}}_\varepsilon - \boldsymbol{\theta}^*\| > C.$$

In particular, there is no uniform positive breakdown lower bound that depends only on  $(\lambda, S_K)$  when  $A_{\max} = +\infty$  and no density-ratio/score floor is assumed.

*Proof.* Fix  $\varepsilon_0, C > 0$ . By assumption, pick  $\Gamma$  so large that  $A_\Gamma > \lambda C / \varepsilon_0 S_K$  and select  $q_\Gamma$  with  $\text{ess sup}_y q_\Gamma(y) / f(y; \boldsymbol{\theta}) \geq \Gamma$  for all  $\boldsymbol{\theta} \in B_2(r; \boldsymbol{\theta}^*)$ . Let  $\varepsilon := \min\{\varepsilon_0, \lambda r / (2S_K A_\Gamma)\}$  and set  $g_\varepsilon^{(\Gamma)} := (1 - \varepsilon)g + \varepsilon q_\Gamma$ . Arguing as in the proof of Theorem 7 but with  $A_{\max}$  replaced by  $A_\Gamma$  (which is now arbitrarily large), the best bound one can obtain is  $\|\hat{\boldsymbol{\theta}}_\varepsilon - \boldsymbol{\theta}^*\| \leq \varepsilon S_K A_\Gamma / \lambda \geq C$ . Since  $\Gamma$  is arbitrary, no uniform lower bound  $\varepsilon^\dagger > 0$  (independent of  $\Gamma$ ) can be guaranteed.  $\square$

## Proof of Theorem 7

*Proof of Theorem 7.* Fix  $r > 0$  and consider any  $\boldsymbol{\theta}$  with  $\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2 = r$ . By (i) ( $\lambda$ -strong convexity of  $D_G(g, \cdot)$  on  $B(\boldsymbol{\theta}^*, r)$ ) and the optimality of  $\boldsymbol{\theta}^*$ ,

$$D_G(g, f_{\boldsymbol{\theta}}) - D_G(g, f_{\boldsymbol{\theta}^*}) \geq \lambda r^2.$$

By (ii), along the line segment between  $\boldsymbol{\theta}^*$  and  $\boldsymbol{\theta}$  (which lies in  $B(\boldsymbol{\theta}^*, r)$ ), the mean-value form of the fundamental theorem of calculus and the gradient envelope give

$$\left| D_G(q, f_{\boldsymbol{\theta}}) - D_G(q, f_{\boldsymbol{\theta}^*}) \right| \leq \sup_{\tilde{\boldsymbol{\theta}} \in B(\boldsymbol{\theta}^*, r)} \left\| \nabla_{\boldsymbol{\theta}} D_G(q, f_{\tilde{\boldsymbol{\theta}}}) \right\|_2 \cdot \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2 \leq S_K A_{\max} r.$$

Linearity of the disparity in its first argument yields, for the contaminated target  $g_\epsilon = (1 - \epsilon)g + \epsilon q$ ,

$$\Delta_\epsilon(\boldsymbol{\theta}) := D_G(g_\epsilon, f_{\boldsymbol{\theta}}) - D_G(g_\epsilon, f_{\boldsymbol{\theta}^*}) = (1 - \epsilon) \underbrace{\Delta_g(\boldsymbol{\theta})}_{\geq \lambda r^2} + \epsilon \Delta_q(\boldsymbol{\theta}),$$

with  $\Delta_q(\boldsymbol{\theta}) := D_G(q, f_{\boldsymbol{\theta}}) - D_G(q, f_{\boldsymbol{\theta}^*})$ . Using the worst-case (adversarial) sign for  $\Delta_q$ ,

$$\Delta_{\epsilon}(\boldsymbol{\theta}) \geq (1 - \epsilon) \lambda r^2 - \epsilon S_K A_{\max} r.$$

Consequently, for any  $\epsilon$  satisfying

$$(1 - \epsilon) \lambda r^2 - \epsilon S_K A_{\max} r > 0 \quad \Longleftrightarrow \quad \epsilon < \frac{\lambda r^2}{\lambda r^2 + S_K A_{\max} r},$$

we have  $\Delta_{\epsilon}(\boldsymbol{\theta}) > 0$  for all  $\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2 = r$ . By continuity,  $D_G(g_{\epsilon}, \cdot)$  cannot attain a minimum on or outside the sphere of radius  $r$ ; thus any minimizer  $\hat{\boldsymbol{\theta}}_{\epsilon}$  lies in  $B(\boldsymbol{\theta}^*, r)$ . The same display with  $\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2 \leq r$  yields the Lipschitz bound

$$D_G(g_{\epsilon}, f_{\boldsymbol{\theta}}) - D_G(g_{\epsilon}, f_{\boldsymbol{\theta}^*}) \geq \lambda \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2 - \epsilon S_K A_{\max} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2,$$

whose minimizer over the segment  $[0, r]$  is attained at  $\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2 \leq (\epsilon S_K A_{\max})/\lambda$ . Therefore, provided  $\epsilon < \epsilon^{\dagger} := \lambda r^2 / (\lambda r^2 + S_K A_{\max} r)$ , we have  $\hat{\boldsymbol{\theta}}_{\epsilon} \in B(\boldsymbol{\theta}^*, r)$  and  $\|\hat{\boldsymbol{\theta}}_{\epsilon} - \boldsymbol{\theta}^*\|_2 \leq (\epsilon S_K A_{\max})/\lambda$ .  $\square$

## Proof of Theorem 8

*Proof of Theorem 8.* We start with the proof of part (1). Fix  $K \leq K_{\max}$  and write, for  $\boldsymbol{\theta} \in \boldsymbol{\Theta}_K$ ,

$$D_{1n}(\boldsymbol{\theta}) = \int f_{\boldsymbol{\theta}}(y) G\left(\frac{g_{1n}(y)}{f_{\boldsymbol{\theta}}(y)}\right) dy, \quad D(\boldsymbol{\theta}) = \int f_{\boldsymbol{\theta}}(y) G\left(\frac{g(y)}{f_{\boldsymbol{\theta}}(y)}\right) dy,$$

where  $g_{1n}$  is the density/pmf estimator built on the selection split  $\mathcal{D}_{1n}$ . For each  $y$  and  $\boldsymbol{\theta}$ , set  $u_{1n}(y) := g_{1n}(y)/f_{\boldsymbol{\theta}}(y)$  and  $u(y) := g(y)/f_{\boldsymbol{\theta}}(y)$ . By the mean-value theorem,

$$\left| G(u_{1n}(y)) - G(u(y)) \right| \leq \left( \sup_{\delta \geq -1} |A'(\delta)| \right) \left| u_{1n}(y) - u(y) \right| = A'_{\max} \frac{|g_{1n}(y) - g(y)|}{f_{\boldsymbol{\theta}}(y)},$$

hence

$$\left| D_{1n}(\boldsymbol{\theta}) - D(\boldsymbol{\theta}) \right| \leq \int f_{\boldsymbol{\theta}}(y) A'_{\max} \frac{|g_{1n}(y) - g(y)|}{f_{\boldsymbol{\theta}}(y)} dy = A'_{\max} \|g_{1n} - g\|_{L^1}.$$

This bound is uniform in  $\boldsymbol{\theta} \in \boldsymbol{\Theta}_K$ , so

$$\sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_K} |D_{1n}(\boldsymbol{\theta}) - D(\boldsymbol{\theta})| \leq A'_{\max} \|g_{1n} - g\|_{L^1} \xrightarrow{p} 0,$$

provided  $\|g_{1n} - g\|_{L^1} \xrightarrow{p} 0$ . The latter holds for (i) the empirical pmf on a fixed finite alphabet (by the LLN/CLT in total variation), and (ii) kernel estimators on  $\mathbb{R}^d$  under standard bandwidth conditions ( $h \rightarrow 0$ ,  $n_1 h^d \rightarrow \infty$ ) ensuring  $L^1$ -consistency. This proves (1).

We next turn to the proof of (2), regarding the identifiability gap. Fix  $K < K_0$ . Under correct specification at  $K_0$  we have  $D_{K_0}^* = \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_{K_0}} D(\boldsymbol{\theta}) = D(\boldsymbol{\theta}^*) = 0$ . We will show  $D_K^* > 0$ . We do this in three steps.

*Step 1 (existence of a population minimizer on  $\boldsymbol{\Theta}_K$ ).* By the standing regularity (continuity of  $\boldsymbol{\theta} \mapsto f_{\boldsymbol{\theta}}$  and of  $D(\boldsymbol{\theta}) = \int f_{\boldsymbol{\theta}} G(g/f_{\boldsymbol{\theta}})$ ) and the compactness/coercivity used in §K for sublevel sets, the continuous map  $D : \boldsymbol{\Theta}_K \rightarrow \mathbb{R}_+$  attains its minimum on  $\boldsymbol{\Theta}_K$ : there exists  $\hat{\boldsymbol{\theta}}_K \in \boldsymbol{\Theta}_K$  with  $D(\hat{\boldsymbol{\theta}}_K) = \min_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_K} D(\boldsymbol{\theta}) =: D_K^*$ .

*Step 2 (strict propriety of the disparity).* Since  $G$  is convex with  $G(1) = 0$  and strictly convex at 1 (calibration  $A'(0) = 1$ ), the disparity is *strictly proper*:  $D(g, f_{\boldsymbol{\theta}}) \geq 0$  with equality iff  $f_{\boldsymbol{\theta}} = g$  a.e. (Equivalently,  $D(g, f_{\boldsymbol{\theta}}) = 0 \iff g/f_{\boldsymbol{\theta}} \equiv 1$  a.e.)

*Step 3 (minimality of  $K_0$  excludes  $g \in \boldsymbol{\Theta}_K$ ).* By the definition of  $K_0$  (true order is minimal) and identifiability of the mixture family,  $g = f(\cdot; \boldsymbol{\theta}^*) \notin \boldsymbol{\Theta}_K$  for all  $K < K_0$ ; i.e., there is no  $\boldsymbol{\theta} \in \boldsymbol{\Theta}_K$  with  $f_{\boldsymbol{\theta}} = g$ .

*Conclusion.* If, towards a contradiction,  $D_K^* = 0$ , then by Step 1 there exists  $\hat{\boldsymbol{\theta}}_K \in \boldsymbol{\Theta}_K$  with  $D(\hat{\boldsymbol{\theta}}_K) = 0$ . By Step 2 this forces  $f_{\hat{\boldsymbol{\theta}}_K} = g$  a.e., contradicting Step 3. Hence  $D_K^* > 0$  for every  $K < K_0$ , which proves the claimed identifiability gap.

We now turn to the proof of (3), namely the local regularity at  $K_0$ . To this end, work

under correct specification at order  $K_0$ :  $g = f_{\theta^*}$  with  $\theta^* \in \Theta_{K_0}$ . Let  $D(\theta) := D_G(g, f_\theta)$  and  $D_{1n}(\theta) := D_G(g_{1n}, f_\theta)$  on the selection split  $\mathcal{D}_{1n}$ . Write  $s_\theta(y) := \nabla_\theta \log f(y; \theta)$  and  $B(u) := G(u) - uG'(u)$ . Assume the fixed-order smoothness/identifiability you list (e.g., **(F1)**, **(K1)**–**(K2)**, **(M1)**–**(M8)**):  $f_\theta$  is  $C^2$  in a neighborhood  $\mathcal{N}(\theta^*)$  with integrable envelopes for  $s_\theta$  and  $\nabla s_\theta$ , and the Fisher information  $I(\theta^*) := \int f_{\theta^*} s_{\theta^*} s_{\theta^*}^\top dy$  is positive definite.

(a)  $C^2$  and curvature at  $\theta^*$ . By dominated differentiation,

$$\nabla_\theta D(\theta) = \int f_\theta(y) s_\theta(y) B\left(\frac{g(y)}{f_\theta(y)}\right) dy. \quad (*)$$

At  $\theta = \theta^*$  we have  $g/f_{\theta^*} \equiv 1$ , hence  $B(1) = G(1) - G'(1)$ , but  $\int f_{\theta^*} s_{\theta^*} dy = \mathbf{0}$ , so  $\nabla_\theta D(\theta^*) = \mathbf{0}$ . Differentiating once more and using  $B'(u) = G'(u) - (G'(u) + uG''(u)) = -uG''(u)$  gives

$$\nabla_\theta^2 D(\theta^*) = \int f_{\theta^*} s_{\theta^*} (-B'(1)) s_{\theta^*}^\top dy = G''(1) I(\theta^*) \succ 0,$$

so  $D$  is  $C^2$  and locally strongly convex at  $\theta^*$ . Therefore there exist  $r > 0$  and  $\lambda > 0$  such that

$$D(\theta) - D(\theta^*) \geq \frac{\lambda}{2} \|\theta - \theta^*\|_2^2 \quad \text{for all } \theta \in B_2(r; \theta^*). \quad (\text{A.12})$$

(b) *Consistency of the selection-split minimizer.* Let  $\hat{\theta}_{K_0, 1n} \in \arg \min_{\Theta_{K_0}} D_{1n}(\theta)$ . By part (1) of the Theorem (ULLN on  $\Theta_{K_0}$ ) and the uniqueness of the population minimizer at  $\theta^*$ , argmin continuity (Berge's maximum theorem / van der Vaart Thm 5.7) yields  $\hat{\theta}_{K_0, 1n} \xrightarrow{p} \theta^*$ .

(c) *Quadratic rate for the population risk.* A second-order Taylor expansion of  $D$  at  $\theta^*$  gives, for  $\hat{\theta}_{K_0, 1n}$  in  $B_2(r; \theta^*)$ ,

$$D(\hat{\theta}_{K_0, 1n}) - D(\theta^*) = \frac{1}{2} (\hat{\theta}_{K_0, 1n} - \theta^*)^\top \nabla_\theta^2 D(\tilde{\theta}_n) (\hat{\theta}_{K_0, 1n} - \theta^*),$$

for some  $\tilde{\theta}_n$  on the segment between  $\theta^*$  and  $\hat{\theta}_{K_0, 1n}$ . By continuity of the Hessian,  $\nabla_\theta^2 D(\tilde{\theta}_n) \rightarrow \nabla_\theta^2 D(\theta^*)$  in probability, hence the quadratic bound (A.12) implies

$$D(\hat{\theta}_{K_0, 1n}) - D(\theta^*) \asymp \|\hat{\theta}_{K_0, 1n} - \theta^*\|_2^2.$$



Thus it suffices to show  $\|\widehat{\boldsymbol{\theta}}_{K_0,1n} - \boldsymbol{\theta}^*\|_2 = O_p(n_1^{-1/2})$ , which we now verify under the fixed-order regularity.

(d) *Local expansion of the selection-split score and  $n_1^{-1/2}$  parameter rate.* The first-order condition  $\mathbf{0} = \nabla D_{1n}(\widehat{\boldsymbol{\theta}}_{K_0,1n})$  and a mean-value expansion around  $\boldsymbol{\theta}^*$  give

$$\mathbf{0} = \nabla D_{1n}(\boldsymbol{\theta}^*) + \left[ \nabla^2 D(\boldsymbol{\theta}^*) + o_p(1) \right] (\widehat{\boldsymbol{\theta}}_{K_0,1n} - \boldsymbol{\theta}^*),$$

where the  $o_p(1)$  term uses the uniform LLN of derivatives on a neighborhood (from your **(M)** and **(F)** conditions). By  $(*)$  with  $g$  replaced by  $g_{1n}$  and a first-order expansion of  $B$  at 1,

$$\nabla D_{1n}(\boldsymbol{\theta}^*) = \int f_{\boldsymbol{\theta}^*} s_{\boldsymbol{\theta}^*} B\left(\frac{g_{1n}}{f_{\boldsymbol{\theta}^*}}\right) dy = B'(1) \int s_{\boldsymbol{\theta}^*}(y) (g_{1n}(y) - g(y)) dy + R_{n,1},$$

where  $B'(1) = -G''(1)$  and the remainder  $R_{n,1} = O_p(\|g_{1n} - g\|_{L^2}^2) = O_p(n_1^{-1})$  for the discrete pmf case (finite support) and, more generally, under your bandwidth conditions. Hence  $\nabla D_{1n}(\boldsymbol{\theta}^*) = O_p(n_1^{-1/2})$ . Since  $\nabla^2 D(\boldsymbol{\theta}^*) = G''(1)I(\boldsymbol{\theta}^*)$  is nonsingular,

$$\widehat{\boldsymbol{\theta}}_{K_0,1n} - \boldsymbol{\theta}^* = -\left[ \nabla^2 D(\boldsymbol{\theta}^*) \right]^{-1} \nabla D_{1n}(\boldsymbol{\theta}^*) + o_p(n_1^{-1/2}) = O_p(n_1^{-1/2}).$$

Therefore  $D(\widehat{\boldsymbol{\theta}}_{K_0,1n}) - D(\boldsymbol{\theta}^*) = O_p(n_1^{-1})$  by the quadratic equivalence above.

Turning to (4), set  $\overline{\mathcal{R}}_{n_1}(K) := \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_K} D_{1n}(\boldsymbol{\theta})$  and  $\mathcal{R}(K) := \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_K} D(\boldsymbol{\theta})$ . By (1) (ULLN on each fixed  $\boldsymbol{\Theta}_K$ ),  $\overline{\mathcal{R}}_{n_1}(K) \xrightarrow{p} \mathcal{R}(K)$  for every  $K \leq K_{\max}$ . *Overfit* ( $K > K_0$ ): By (4), for each  $K > K_0$ ,

$$\overline{\mathcal{R}}_{n_1}(K) - \overline{\mathcal{R}}_{n_1}(K_0) = \{D_{1n}(\widehat{\boldsymbol{\theta}}_{K,1n}) - D(\boldsymbol{\theta}^*)\} - \{D_{1n}(\widehat{\boldsymbol{\theta}}_{K_0,1n}) - D(\boldsymbol{\theta}^*)\} = O_p(n_1^{-1}).$$

Finally, to prove (5), we consider the *Underfit* ( $K < K_0$ ). case. By part (2),  $\mathcal{R}(K) - \mathcal{R}(K_0) =: c_K > 0$ . Hence for any  $\delta \in (0, c_K)$ , with probability  $\rightarrow 1$ ,

$$\overline{\mathcal{R}}_{n_1}(K) - \overline{\mathcal{R}}_{n_1}(K_0) \geq c_K - \delta.$$

The penalty difference contributes  $(b_{n_1}/n_1)\{p(K) - p(K_0)\}$ , which is nonpositive and vanishes by  $b_{n_1}/n_1 \rightarrow 0$ . Therefore, eventually  $\text{GDIC}_{n_1}(K) > \text{GDIC}_{n_1}(K_0)$  for all  $K < K_0$ . The penalty difference is  $(b_{n_1}/n_1)\{p(K) - p(K_0)\}$ , which is strictly positive and dominates  $O_p(n_1^{-1})$  because  $b_{n_1} \rightarrow \infty$  while  $b_{n_1}/n_1 \rightarrow 0$  (e.g.,  $b_{n_1} = \frac{1}{2} \log n_1$ ). Hence  $\text{GDIC}_{n_1}(K) > \text{GDIC}_{n_1}(K_0)$  with probability  $\rightarrow 1$  for each  $K > K_0$ . Combining the two cases yields  $\mathbf{P}(\widehat{K}_n = K_0) \rightarrow 1$ .  $\square$

## Proof of Theorem 9

*Proof of Theorem 9. (a) Fixed order.* On the estimation split  $\mathcal{D}_{2n}$ , let  $\widehat{\boldsymbol{\theta}}_{n_2} \in \arg \min_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_K} D(g_{n_2}, f_{\boldsymbol{\theta}})$ . Under correct model specification at  $K_0$ , calibration  $A'(0) = 1$ , the plug-in rate  $\|g_{n_2} - g\|_{\mathcal{H}} = o_p(n_2^{-1/2})$ , and the fixed-order regularity **(F1)**, **(C1)**–**(C2)**, **(K1)**–**(K2)**, **(M1)**–**(M8)**, the standard Z-estimation CLT yields

$$\sqrt{n_2}(\widehat{\boldsymbol{\theta}}_{n_2} - \boldsymbol{\theta}^*) \Rightarrow \mathcal{N}(0, I(\boldsymbol{\theta}^*)^{-1}).$$

*(b) Post-selection (unconditional).* By Theorem 8,  $\mathbb{P}(\widehat{K}_n = K_0) \rightarrow 1$ . On the event  $\{\widehat{K}_n = K_0\}$ , the dimension-matched pair satisfies  $\bar{\boldsymbol{\theta}}_{n_2} = \widehat{\boldsymbol{\theta}}_{n_2}$  and  $\boldsymbol{\theta}^*(\widehat{K}_n) = \boldsymbol{\theta}^*$ . Hence

$$\sqrt{n_2}\{\bar{\boldsymbol{\theta}}_{n_2} - \boldsymbol{\theta}^*(\widehat{K}_n)\} = \sqrt{n_2}(\widehat{\boldsymbol{\theta}}_{n_2} - \boldsymbol{\theta}^*) \Rightarrow \mathcal{N}(0, I(\boldsymbol{\theta}^*)^{-1}),$$

and the same limit holds unconditionally by  $\mathbb{P}(\widehat{K}_n = K_0) \rightarrow 1$ .

*Truncated DM on  $\mathcal{D}_{2n}$ .* If  $\boldsymbol{\theta}_{n_2}^{(m_{n_2})}$  is the  $m_{n_2}$ -iterate with  $m_{n_2} = O(\log n_2)$  as in Theorem 3(ii), then  $\sqrt{n_2}\|\boldsymbol{\theta}_{n_2}^{(m_{n_2})} - \widehat{\boldsymbol{\theta}}_{n_2}\| = o_p(1)$  and the same limits follow by Slutsky.  $\square$

## Proof of Theorem 10 (GDIC over/under-estimation bounds; bounded-RAF)

*Proof. Underfit.* Write  $\Delta_\epsilon(K)$  as in the main text. Bound the empirical minimizer's value by quadratic expansion and operator-noise:  $\sup_{\theta \in \Theta_K} |D_{1n}(\theta) - D(\theta)| = O_p(\|g_{1n} - g\|_{\mathcal{H}}) = O_p(n_1^{-1/2})$  by bounded  $A'_{\max}$  and the score envelope. On the selection split, the value error at the (curved) minimizer is  $O_p(n_1^{-1})$ , so for  $t \in (0, \Delta_\epsilon(K))$ ,

$$\mathbb{P}\left(D_{1n}(\hat{\theta}_{K,1n}) - D_{1n}(\hat{\theta}_{K_0,1n}) \leq -t\right) \leq \exp(-c_1 n_1 t^2) + c_2 \mathbb{P}(\|g_{1n} - g\|_{\mathcal{H}} > t/2),$$

by a Bernstein (or Hoeffding) inequality for bounded-RAF contrasts. Plug  $t = \frac{1}{2}\Delta_\epsilon(K)$  and sum over  $K < K_0$ .

*Overfit.* Write  $\nu_K = p(K) - p(K_0)$ . By the (Godambe-)Wilks expansion on the selection split,  $2n_1\{D_{1n}(\hat{\theta}_{K,1n}) - D_{1n}(\hat{\theta}_{K_0,1n})\} \Rightarrow \chi_{\nu_K}^2$ . Thus, for BIC-type penalty ( $b_{n_1} = \frac{1}{2} \log n_1$ ),

$$\begin{aligned} \mathbb{P}\left(D_{1n}(\hat{\theta}_{K,1n}) + \frac{\log n_1}{2n_1} p(K) \leq D_{1n}(\hat{\theta}_{K_0,1n}) + \frac{\log n_1}{2n_1} p(K_0)\right) &\leq \mathbb{P}\left(\chi_{\nu_K}^2 \geq \nu_K \log n_1 + o(1)\right) \\ &= O(n_1^{-\nu_K/2}), \end{aligned}$$

using standard  $\chi^2$  tails. Union bound over  $K > K_0$  yields the claim.  $\square$

## Proof of Proposition 5 (BIC vs AIC; contamination and unbounded RAF)

*Proof.* (i)–(BIC) follow directly from Theorem 10. (ii) *AIC*: If  $b_{n_1} \equiv b \in (0, \infty)$ , the overfit test reduces to  $\mathbb{P}\{\chi_{\nu_K}^2 \geq \nu_K b + o(1)\}$ , which converges to a strictly positive limit  $c(\nu_K) \in (0, 1)$ . (iii) *KL*: With unbounded RAF  $A$ , the uniform operator-noise/curvature control used in Theorem S.GDIC.1 fails without a local density-ratio or score floor. Under such a floor, the same bounds hold (use  $A'_{\max} = 1$ ); without it, the failure mode in Prop. S.R.2 shows that neither underfit nor overfit probabilities need decay in general.  $\square$

The following proposition is a negative result describing how KL/BIC needs a local floor.

**Proposition 8** (KL/BIC needs a local floor; a one-point overfit failure mode). *Let  $G$  be the likelihood disparity (KL). Assume only the smoothness of Section 4 and no local density/score floor on  $B_2(r; \boldsymbol{\theta}^*)$ . Fix  $K_0$  and  $\Delta p \geq 1$ . Then for every  $n_1$  and every  $M > 0$  there exists a point  $y_M$  with  $\inf_{\boldsymbol{\theta} \in \Theta_{K_0}} \{-\log f(y_M; \boldsymbol{\theta})\} \geq M$ . For the  $\varepsilon$ -contaminated selection split (rate  $\varepsilon > 0$ ), choose  $M = (\Delta p)^{\frac{1}{2}} \log n_1 + 1$ . With probability at least  $1 - \exp(-c\varepsilon n_1)$  (some  $c > 0$ ), at least one observation falls in a neighbourhood of  $y_M$ ; refitting with  $K_0+1$  components that isolates this point reduces the average KL contrast by at least  $M/n_1$ , which exceeds the BIC penalty  $(\Delta p)^{\frac{\log n_1}{2n_1}}$ . Thus*

$$\liminf_{n_1 \rightarrow \infty} \mathbb{P}(\widehat{K}_n \geq K_0+1) \geq 1 - e^{-c\varepsilon} > 0.$$

*If a local density/score floor holds (e.g.,  $\inf_{\boldsymbol{\theta} \in B_2(r; \boldsymbol{\theta}^*)} \inf_{y \in \mathcal{Y}_r} f(y; \boldsymbol{\theta}) \geq c_* > 0$ ), then the overfit probability decays at least at the BIC rate (Wilks + penalty), and the GDIC–BIC conclusions of Theorem 8 and Theorem 10 carry over to KL.*

The proof is standard and is similar to the Proof of Theorem 10 and Proposition 5.

## Proof of Theorem 11: Overfit bound for bounded RAF

The proof of the Theorem relies on the proof of the following lemma.

**Lemma 6** (Per-point leverage, bounded RAF). *Fix  $K_0$  and  $K > K_0$ . On the selection split of size  $n_1$ , for any subset  $S \subset \{1, \dots, n_1\}$  with  $|S| = m$ , the maximal possible decrease of the empirical contrast  $D_{1n}$  achievable by refitting with  $K$  components that dedicate  $(K - K_0)$  new components to absorb exactly the  $m$  points in  $S$  is at most  $m A_{\max}/n_1$ .*

*Proof.* Write  $D_{1n}(\boldsymbol{\theta}) = \frac{1}{n_1} \sum_{i=1}^{n_1} \phi_{\boldsymbol{\theta}}(Y_i)$  with  $\phi_{\boldsymbol{\theta}}(y) := f_{\boldsymbol{\theta}}(y) G\left(\frac{1}{n_1 f_{\boldsymbol{\theta}}(y)}\right)$  when  $g_{1n}$  is the empirical pmf. For bounded RAF  $A$ , the pointwise decrement in  $\phi_{\boldsymbol{\theta}}(y)$  from changing  $f_{\boldsymbol{\theta}}$  arbitrarily at a single  $y$  is bounded by  $A_{\max}$ , so the total decrease contributed by the  $m$  indices in

$S$  is at most  $mA_{\max}$ ; dividing by  $n_1$  gives the claim. (A formal proof can be obtained by interpolating between the two fitted models and using the mean-value formula with  $A$  as the directional derivative of  $u \mapsto uG(1/u)$ .)  $\square$

**Proof of the Theorem:** Let  $\hat{\theta}_{K,1n}$  and  $\hat{\theta}_{K_0,1n}$  be selection-split minimizers. Then

$$\text{GDIC}_{n_1}(K) - \text{GDIC}_{n_1}(K_0) = \left\{ D_{1n}(\hat{\theta}_{K,1n}) - D_{1n}(\hat{\theta}_{K_0,1n}) \right\} + \frac{\nu_K b_{n_1}}{n_1}.$$

On the event  $\{X_{n_1} \leq m\}$ , Lemma 6 yields  $D_{1n}(\hat{\theta}_{K,1n}) - D_{1n}(\hat{\theta}_{K_0,1n}) \geq -mA_{\max}/n_1 - R_{n_1}$  with  $R_{n_1} = O_p(n_1^{-1/2})$  from sampling error (bounded RAF + ULLN). Thus for  $m < m_{\min}(K, n_1) = \lceil \nu_K b_{n_1}/A_{\max} \rceil$  and large  $n_1$ ,  $\text{GDIC}_{n_1}(K) - \text{GDIC}_{n_1}(K_0) > 0$  on  $\{X_{n_1} \leq m\}$ , which implies

$$\mathbb{P}(\widehat{K}_n \geq K) \leq \mathbb{P}(X_{n_1} \geq m_{\min}(K, n_1)) + o(1).$$

For  $X_{n_1} \sim \text{Bin}(n_1, \varepsilon)$ , Chernoff's bound gives  $\mathbb{P}(X_{n_1} \geq m_{\min}) \leq \exp\{-n_1 \text{kl}(m_{\min}/n_1 \parallel \varepsilon)\}$ , yielding the stated decay for BIC.  $\square$

## Proof of Theorem 11: Underfit bound for bounded RAF

For each fixed  $K < K_0$ , let  $\Delta_0(K) = \inf_{\theta \in \Theta_K} D(\theta) - D(\theta^*) > 0$ . Bound the contaminated population gap by linear response with bounded RAF:  $\Delta_\varepsilon(K) \geq \Delta_0(K) - 2A_{\max}\varepsilon$  (two contrasts shift by at most  $A_{\max}\varepsilon$  each). Selection-split ULLN gives  $\sup_{\theta \in \Theta_K} |D_{1n}(\theta) - D(\theta)| = O_p(n_1^{-1/2})$  uniformly for bounded-RAF class. Hence, for any  $t \in (0, \Delta_\varepsilon(K))$ ,

$$\mathbb{P}\left(D_{1n}(\hat{\theta}_{K,1n}) - D_{1n}(\hat{\theta}_{K_0,1n}) \leq -t\right) \leq \exp(-c_1 n_1 t^2) + o(1),$$

by a Bernstein/Hoeffding bound for bounded contrasts. Choosing  $t = \frac{1}{2}\Delta_\varepsilon(K)$  yields

$$\mathbb{P}(\widehat{K}_n = K) \leq \exp\left(-c_1 n_1 [\Delta_0(K) - 2A_{\max}\varepsilon]^2\right) + o(1),$$

and summing over  $K < K_0$  gives the claim.  $\square$

## One-point overfit failure mode for KL

Construct  $y_M$  with  $\inf_{\theta \in \Theta_{K_0}} [-\log f(y_M; \theta)] \geq M$  and let  $\varepsilon > 0$  be fixed. Under  $\varepsilon$ -contamination on  $\mathcal{D}_{1n}$ , with probability  $\rightarrow 1 - e^{-c\varepsilon}$  some  $Y_i$  lies in a small neighborhood of  $y_M$ . Refitting with  $K_0+1$  by dedicating one component to that  $Y_i$  reduces the average KL contrast at least by  $M/n_1$ , while the BIC penalty is  $(\nu_{K_0+1} \log n_1)/(2n_1)$ . Choosing  $M = (\nu_{K_0+1} \log n_1) + 1$  forces  $\text{GDIC}_{n_1}(K_0+1) < \text{GDIC}_{n_1}(K_0)$  on that event, so  $\liminf_{n_1 \rightarrow \infty} \mathbb{P}(\widehat{K}_n \geq K_0+1) \geq 1 - e^{-c\varepsilon} > 0$ . If a local density/score floor holds, the per-point leverage is bounded and S.GDIC-B1 applies with  $A_{\max}$  replaced by  $A_\Gamma$ .  $\square$

## Proof of Theorem 12

*Proof of Theorem 12.* Condition on  $\mathcal{D}_{1n}$ ; then  $\widehat{K}_n = \widehat{K}_n(\mathcal{D}_{1n})$  is fixed and independent of  $\mathcal{D}_{2n}$ . On  $\mathcal{D}_{2n}$ , by part (a) of Theorem 9 at  $K = \widehat{K}_n$  (and the fixed-order regularity at  $K_0$ ),

$$\sqrt{n_2} \{\bar{\boldsymbol{\theta}}_{n_2} - \boldsymbol{\theta}^*(\widehat{K}_n)\} \Rightarrow \mathcal{N}(0, I(\boldsymbol{\theta}^*)^{-1}) \quad \text{in probability.}$$

Since  $\mathbb{P}(\widehat{K}_n = K_0) \rightarrow 1$ , we have  $\sqrt{n_2} \{\bar{\boldsymbol{\theta}}_{n_2} - \boldsymbol{\theta}^*(K_0)\} \Rightarrow \mathcal{N}(0, I(\boldsymbol{\theta}^*)^{-1})$  in probability (stable convergence), and hence unconditionally as well. If a truncated iterate  $\boldsymbol{\theta}_{n_2}^{(m_{n_2})}$  is used, the additional  $o_p(1)$  term  $\sqrt{n_2} \|\boldsymbol{\theta}_{n_2}^{(m_{n_2})} - \bar{\boldsymbol{\theta}}_{n_2}\|$  vanishes by Theorem 4(ii), and the same limit holds by Slutsky.  $\square$

# M: Reference Results Used in the Proofs of the Main Theorems

We collect standard facts used throughout. Proofs are omitted; see the cited references.

**Lemma 7** (Uniform LLN on fixed-order classes). *For each fixed  $K \leq K_{\max}$ ,*

$$\sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_K} |D_{1n}(\boldsymbol{\theta}) - D(\boldsymbol{\theta})| \xrightarrow{p} 0 \quad (n_1 \rightarrow \infty),$$

*whenever  $G$  is convex with  $A'_{\max} < \infty$  and  $g_{1n}$  is  $L^1$ -consistent (empirical pmf on a finite alphabet or KDE with  $h \rightarrow 0$ ,  $n_1 h^d \rightarrow \infty$ ). Used in: Theorem 8 proof of (1) and (3).*

**Lemma 8** (Argmin consistency / continuity). *Let  $M_n(\boldsymbol{\theta}) := \arg \min_{\vartheta} D_{1n,K}(\vartheta)$  and  $M(\boldsymbol{\theta}) := \arg \min_{\vartheta} D_K(\vartheta)$  on a compact set. If  $\sup_{\vartheta} |D_{1n,K}(\vartheta) - D_K(\vartheta)| \rightarrow 0$  and  $\vartheta_K^*$  is the unique minimizer of  $D_K$ , then every measurable selection  $\hat{\vartheta}_{K,1n} \in M_n(\cdot)$  satisfies  $\hat{\vartheta}_{K,1n} \xrightarrow{p} \vartheta_K^*$ . Used in: Theorem 8, Part (3).*

**Lemma 9** (Berge maximum theorem (single-valued specialization)). *If  $Q_n(\cdot \mid \boldsymbol{\theta})$  has a unique minimizer for each  $\boldsymbol{\theta}$  and sublevel sets are compact, then the argmin map  $\boldsymbol{\theta} \mapsto M_n(\boldsymbol{\theta})$  is continuous in a neighborhood of any fixed point. Used in: Propositions 3 and 4 for local continuity, and cycle exclusion under uniqueness.*

**Lemma 10** (Closed graph / outer semicontinuity). *If  $Q_n(\cdot \mid \boldsymbol{\theta})$  is continuous and sublevel sets are compact, then the update correspondence  $M_n : \Theta \rightrightarrows \Theta$  has a closed graph (outer semicontinuous):  $\theta^{(j)} \rightarrow \bar{\theta}$ ,  $\eta^{(j)} \in M_n(\theta^{(j)})$ ,  $\eta^{(j)} \rightarrow \bar{\eta} \Rightarrow \bar{\eta} \in M_n(\bar{\theta})$ . Used in: Proposition 3 for verifying the invariance of the limit set.*

**Lemma 11** (Local quadratic expansion and curvature). *Under the fixed-order smoothness (F1), (K1)–(K2), (M1)–(M8),*

$$\nabla_{\boldsymbol{\theta}} D(\boldsymbol{\theta}^*) = \mathbf{0}, \quad \nabla_{\boldsymbol{\theta}}^2 D(\boldsymbol{\theta}^*) = G''(1) I(\boldsymbol{\theta}^*) \succ 0,$$

and for some  $r, \lambda > 0$ ,  $D(\boldsymbol{\theta}) - D(\boldsymbol{\theta}^*) \geq \frac{\lambda}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2$  on  $B_2(r; \boldsymbol{\theta}^*)$ . Used in: Theorem 8, part (3) and contraction radius in Theorem 3.

**Lemma 12** (Noisy linear recursion). *If  $x_{t+1} \leq \kappa x_t + b$  with  $\kappa \in (0, 1)$  then  $x_t \leq \kappa^t x_0 + \frac{1-\kappa^t}{1-\kappa} b \leq \kappa^t x_0 + \frac{b}{1-\kappa}$ .* Used in: Theorem 1, Theorem 3, finite-step bounds in Theorem 1.

**Lemma 13** (Stable (conditional) convergence under sample splitting). *If  $\mathcal{D}_{1n} \perp\!\!\!\perp \mathcal{D}_{2n}$  and, conditionally on  $\mathcal{D}_{1n}$ ,  $Z_{n_2}(\mathcal{D}_{2n}) \Rightarrow \mathcal{N}(0, \Sigma)$  in probability, then the same holds unconditionally; if  $\mathbb{P}(\widehat{K}_n = K_0) \rightarrow 1$ , post-selection versions follow by Slutsky.* Used in: Theorem 10.

**Lemma 14** (Weighted  $\chi^2$  limit for quadratic forms). *If  $\sqrt{n}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\dagger) \Rightarrow \mathcal{N}(0, H^{-1}VH^{-1})$  with  $H \succ 0$  and  $V \succeq 0$ , then*

$$n(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\dagger)^\top H (\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\dagger) \Rightarrow \sum_{j=1}^{p(K_0)} \lambda_j \chi_{1,j}^2,$$

where  $\{\lambda_j\}$  are the eigenvalues of  $J := H^{-1/2}VH^{-1/2}$ ; if  $H = V$ , the limit is  $\chi_{p(K_0)}^2$ . Used in: Theorem finite-step Theorem 6.



# N: Reference asymptotics: Z-estimation CLTs and Godambe–Wilks Theorems

## 9.2 Notation and standing assumptions

Let  $\Theta \subset \mathbb{R}^{p(K)}$  be the parameter space,  $f(\cdot; \theta)$  the observed-data model, and  $s_\theta(y) = \nabla_\theta \log f(y; \theta)$  the score. For a convex generator  $G : [-1, \infty) \rightarrow \mathbb{R}$  with  $G(0) = G'(0) = 0$  and  $G''(0) = 1$  (calibrated), define the residual-adjustment function

$$A(\delta) = (1 + \delta)G'(\delta) - G(\delta), \quad A'(0) = G''(0) = 1.$$

For any density/pmf  $g$ , define the divergence and DM score map

$$D_G(g, f_\theta) = \int G\left(\frac{g}{f_\theta} - 1\right) f_\theta dy, \quad \Psi(\theta; g) := \nabla_\theta D_G(g, f_\theta) = - \int A\left(\frac{g}{f_\theta} - 1\right) s_\theta f_\theta dy.$$

Let  $g_n$  be a plug-in estimate of  $g$  (empirical or a discrete-kernel estimator). We use the following assumptions; we quote them where needed.

**(A1) Local well-posedness.** There is a neighborhood  $\mathcal{N}$  of the target  $\theta^\dagger$  such that  $\Psi(\cdot; g)$  is continuously Fréchet-differentiable on  $\mathcal{N}$ , and the Jacobian

$$H := \nabla_\theta \Psi(\theta^\dagger; g) = \nabla_\theta^2 D_G(g, f_\theta) \Big|_{\theta=\theta^\dagger}$$

is nonsingular; further,  $D_G(g, \cdot)$  is  $\lambda$ -strongly convex on  $\mathcal{N}$  (eigenvalues of  $H$  bounded below by  $\lambda > 0$ ).

**(A2) Nuisance differentiability.** For any signed  $h$  with  $\int h = 0$ , the pathwise (Gâteaux) derivative  $\partial_g \Psi(\theta^\dagger; g)[h]$  exists and is continuous in  $h$ .

**(A3) Second moments and entropy control.** There is an envelope  $F \in L^2(g)$  such that  $\sup_{\theta \in \mathcal{N}} \|A'(\frac{g}{f_\theta} - 1)s_\theta\| \leq F$ , and the class  $\{A'(\frac{g}{f_\theta} - 1)s_\theta : \theta \in \mathcal{N}\}$  admits a Donsker/bracketing bound ensuring a uniform  $O_p(n^{-1/2})$  empirical process.

**(A4) Plug-in rate.**  $\|g_n - g\|_{\mathcal{H}} = o_p(n^{-1/2})$  in a norm that implies  $\sup_{\boldsymbol{\theta} \in \mathcal{N}} \|\Psi_n(\boldsymbol{\theta}) - \Psi(\boldsymbol{\theta}; g)\| = o_p(n^{-1/2})$ , where  $\Psi_n(\boldsymbol{\theta}) := \nabla_{\boldsymbol{\theta}} D_G(g_n, f_{\boldsymbol{\theta}})$ .

**(A5) FOS and operator noise (for contraction).** On  $B_2(r'; \boldsymbol{\theta}^*)$  the population DM map  $M(\boldsymbol{\theta}) := \arg \min_{\boldsymbol{\theta}'} \mathcal{Q}_G(\boldsymbol{\theta}' | \boldsymbol{\theta})$  satisfies the first-order stability (FOS)

$$\|M(\boldsymbol{\theta}) - M(\boldsymbol{\theta}^*)\| \leq \frac{\gamma_K}{\lambda} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|, \quad \gamma_K \leq \frac{C_{\text{fos}}}{\pi_{\min}} \cdot K \cdot L_{\text{comp}},$$

and the sample map  $M_n$  obeys

$$\sup_{\boldsymbol{\theta} \in B_2(r'; \boldsymbol{\theta}^*)} \|M_n(\boldsymbol{\theta}) - M(\boldsymbol{\theta})\| \leq C_{\text{op}} A'_{\max} \sqrt{\frac{p(K) + \log(1/\rho)}{n}}, \quad A'_{\max} := \sup_{\boldsymbol{\theta} \geq -1} |A(\boldsymbol{\theta})|.$$

**(A6) Robustness envelope.** For any density  $q$ ,  $\sup_{\boldsymbol{\theta} \in \mathcal{N}} \|\nabla_{\boldsymbol{\theta}} D_G(q, f_{\boldsymbol{\theta}})\| \leq S_K A_{\max}$ , with  $S_K \lesssim C_1 \sqrt{p(K)}$  or  $C_2 K / \pi_{\min}$  in regular mixtures.

At the correctly specified model  $g = f_{\boldsymbol{\theta}^*}$ , calibration gives  $H = I(\boldsymbol{\theta}^*)$  (Fisher information).

For Hellinger (unbounded RAF), (A6) holds if  $q/f_{\boldsymbol{\theta}}$  is uniformly bounded on  $\mathcal{N}$ .

### 9.3 Gâteaux derivative in the nuisance (orthogonality)

**Lemma 15** (Gâteaux derivative). *For any signed perturbation  $h$  with  $\int h = 0$ ,*

$$\partial_g \Psi(\boldsymbol{\theta}; g)[h] = - \int A' \left( \frac{g}{f_{\boldsymbol{\theta}}} - 1 \right) s_{\boldsymbol{\theta}} h \, dy.$$

*In particular, at  $(\boldsymbol{\theta}, g) = (\boldsymbol{\theta}^*, f_{\boldsymbol{\theta}^*})$ ,  $\partial_g \Psi(\boldsymbol{\theta}^*; g)[h] = - \int s_{\boldsymbol{\theta}^*} h \, dy$  (calibrated orthogonality).*

*Proof.* Write  $\delta(y; \boldsymbol{\theta}, g) = g/f_{\boldsymbol{\theta}} - 1$ . For  $g_t = g + th$ ,  $t \in \mathbb{R}$ ,

$$\frac{\Psi(\boldsymbol{\theta}; g_t) - \Psi(\boldsymbol{\theta}; g)}{t} = - \int \frac{A(\delta_t) - A(\delta)}{t} s_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}} \, dy \rightarrow - \int A'(\delta) \frac{h}{f_{\boldsymbol{\theta}}} s_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}} \, dy$$

by dominated convergence, using (A3). Evaluating at  $\delta \equiv 0$  gives  $A'(0) = 1$  and the claim at the model.  $\square$

## 9.4 DM–Wilks and Godambe–Wilks (with robust pivot)

**Theorem 13** (DM–Wilks and Godambe–Wilks). *Let  $\hat{\boldsymbol{\theta}}_n = \arg \min_{\boldsymbol{\theta}} D_G(g_n, f_{\boldsymbol{\theta}})$  and  $\boldsymbol{\theta}^\dagger = \arg \min_{\boldsymbol{\theta}} D_G(g, f_{\boldsymbol{\theta}})$ . Under (A1)–(A4),*

$$2n\{D_G(g_n, f_{\boldsymbol{\theta}^\dagger}) - D_G(g_n, f_{\hat{\boldsymbol{\theta}}_n})\} = n(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\dagger)^\top H(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\dagger) + o_p(1),$$

hence

$$2n\{D_G(g_n, f_{\boldsymbol{\theta}^\dagger}) - D_G(g_n, f_{\hat{\boldsymbol{\theta}}_n})\} \Rightarrow \sum_{j=1}^{p(K_0)} \lambda_j \chi_{1,j}^2,$$

where  $\{\lambda_j\}$  are the eigenvalues of  $J := H^{-1/2} V H^{-1/2}$ . Under correct specification and calibration,  $H = V = I(\boldsymbol{\theta}^*)$  and the limit is  $\chi_{p(K_0)}^2$ .

*Proof.* Let  $\tilde{\boldsymbol{\theta}}_n$  lie on the segment between  $\hat{\boldsymbol{\theta}}_n$  and  $\boldsymbol{\theta}^\dagger$ . A second–order expansion gives

$$D_G(g_n, f_{\boldsymbol{\theta}^\dagger}) - D_G(g_n, f_{\hat{\boldsymbol{\theta}}_n}) = \frac{1}{2}(\boldsymbol{\theta}^\dagger - \hat{\boldsymbol{\theta}}_n)^\top \nabla_{\boldsymbol{\theta}}^2 D_G(g_n, f_{\tilde{\boldsymbol{\theta}}_n})(\boldsymbol{\theta}^\dagger - \hat{\boldsymbol{\theta}}_n).$$

By (A1)–(A4),  $\sup_{\boldsymbol{\theta} \in \mathcal{N}} \|\nabla_{\boldsymbol{\theta}}^2 D_G(g_n, f_{\boldsymbol{\theta}}) - H\| = o_p(1)$ , so the right–hand side equals  $\frac{1}{2}(\boldsymbol{\theta}^\dagger - \hat{\boldsymbol{\theta}}_n)^\top H(\boldsymbol{\theta}^\dagger - \hat{\boldsymbol{\theta}}_n) + o_p(\|\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\dagger\|^2)$ . Multiplying by  $2n$  and using  $\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\dagger) \Rightarrow \mathcal{N}(0, H^{-1} V H^{-1})$  (Corollary 3 to Theorem 3) yields the weighted  $\chi^2$  limit.  $\square$

**Corollary 10** (Godambe–calibrated deviance). *Let  $\widehat{H} = \nabla_{\boldsymbol{\theta}}^2 D_G(g_n, f_{\boldsymbol{\theta}})|_{\hat{\boldsymbol{\theta}}_n}$  and*

$$\widehat{V} = \frac{1}{n} \sum_{i=1}^n \left[ A' \left( \frac{g_n(Y_i)}{f_{\hat{\boldsymbol{\theta}}_n}(Y_i)} - 1 \right) s_{\hat{\boldsymbol{\theta}}_n}(Y_i) \right] \left[ A' \left( \frac{g_n(Y_i)}{f_{\hat{\boldsymbol{\theta}}_n}(Y_i)} - 1 \right) s_{\hat{\boldsymbol{\theta}}_n}(Y_i) \right]^\top.$$

*Then the Wald–type pivot*

$$\Lambda_n^{\text{GW}} := n(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\dagger)^\top \widehat{V}^{-1}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\dagger) \Rightarrow \chi_{p(K_0)}^2,$$

*even under misspecification. At the model (calibrated), the raw DM deviance has a  $\chi_{p(K_0)}^2$  limit (DM–Wilks).*

*Proof.* By Corollary 3 to Theorem 3 and consistency of  $\widehat{V}$ ,  $\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\dagger) \Rightarrow \mathcal{N}(0, H^{-1} V H^{-1})$  and  $n(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\dagger)^\top \widehat{V}^{-1}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\dagger) \Rightarrow \chi_{p(K_0)}^2$ .  $\square$

# References

- Aitkin, M. (1980), ‘Mixture applications of the em algorithm in glim’, *COMPSTAT 1980* pp. 537–541.
- Al Mohamad, D. & Broniatowski, M. (2016), ‘A proximal point algorithm for minimum divergence estimators with application to mixture models’, *Entropy* **18**(8), 277.
- Balakrishnan, S., Wainwright, M. J., Yu, B. et al. (2017), ‘Statistical guarantees for the em algorithm: From population to sample-based analysis’, *The Annals of Statistics* **45**(1), 77–120.
- Basu, A., Basu, S. & Chaudhuri, G. (1997), ‘Robust minimum divergence procedures for count data models’, *Sankhyā Ser. B* **59**(1), 11–27.
- Basu, A. & Harris, I. R. (1994), ‘Robust predictive distributions for exponential families’, *Biometrika* **81**(4), 790–794.
- Basu, A. & Lindsay, B. G. (1994), ‘Minimum disparity estimation for continuous models: efficiency, distributions and robustness’, *Ann. Inst. Statist. Math.* **46**(4), 683–705.
- Basu, A., Sarkar, S. & Vidyashankar, A. N. (1997), ‘Minimum negative exponential disparity estimation in parametric models’, *J. Statist. Plann. Inference* **58**(2), 349–370.
- Beran, R. (1977), ‘Minimum Hellinger distance estimates for parametric models’, *The Annals of Statistics* **5**(3), 445–463.
- Cheng, A. & Vidyashankar, A. N. (2006), ‘Minimum hellinger distance estimation for randomized play the winner design’, *Journal of Statistical Planning and Inference* **136**(6), 1875–1910.

- Chrétien, S. & Hero, A. O. (2008), ‘On EM algorithms and their proximal generalizations’, *ESAIM Probab. Stat.* **12**, 308–326.
- Chrétien, S. & Hero, III, A. O. (2000), ‘Kullback proximal algorithms for maximum-likelihood estimation’, *IEEE Trans. Inform. Theory* **46**(5), 1800–1810. Information-theoretic imaging.
- Cressie, N. & Read, T. R. C. (1984), ‘Multinomial goodness-of-fit tests’, *J. Roy. Statist. Soc. Ser. B* **46**(3), 440–464.
- Cunha, F., da Cruz Neto, J. & Oliveira, P. (2010), ‘A proximal point algorithm with a  $\phi$ -divergence for quasiconvex programming’, *Optimization* **59**(5), 777–792.
- Cutler, A. & Cordero-Braña, O. I. (1996), ‘Minimum Hellinger distance estimation for finite mixture models’, *J. Amer. Statist. Assoc.* **91**(436), 1716–1723.
- Dempster, A. P., Laird, N. M. & Rubin, D. B. (1977), ‘Maximum likelihood from incomplete data via the EM algorithm’, *J. Roy. Statist. Soc. Ser. B* **39**(1), 1–38. With discussion.
- Donoho, D. L. & Liu, R. C. (1988), ‘The “automatic” robustness of minimum distance functionals’, *Ann. Statist.* **16**(2), 552–586.
- Doukhan, P. & León, J. R. (1990), ‘Déviation quadratique d’estimateurs de densité par projections orthogonales’, *C. R. Acad. Sci. Paris Sér. I Math.* **310**(6), 425–430.
- Dwivedi, R., Ho, N., Khamaru, K., Jordan, M. I., Wainwright, M. J. & Yu, B. (2018), ‘Singularity, misspecification, and the convergence rate of em’, *arXiv preprint arXiv:1810.00828*.
- Eslinger, P. W. & Woodward, W. A. (1991), ‘Minimum hellinger distance estimation for normal models’, *Journal of Statistical Computation and Simulation* **39**(1-2), 95–114.

- Fan, J. & Gijbels, I. (1996), *Local polynomial modelling and its applications*, Vol. 66 of *Monographs on Statistics and Applied Probability*, Chapman & Hall, London.
- Ganesalingam, S. & McLachlan, G. (1979a), ‘Small sample results for a linear discriminant function estimated from a mixture of normal populations’, *Journal of Statistical Computation and Simulation* **9**(2), 151–158.
- Ganesalingam, S. & McLachlan, G. (1980), ‘A comparison of the mixture and classification approaches to cluster analysis’, *Communications in Statistics-Theory and Methods* **9**(9), 923–933.
- Ganesalingam, S. & McLachlan, G. J. (1978), ‘The efficiency of a linear discriminant function based on unclassified initial samples’, *Biometrika* **65**(3), 658–662.
- Ganesalingam, S. & McLachlan, G. J. (1979b), ‘A case study of two clustering methods based on maximum likelihood’, *Statist. Neerlandica* **33**(2), 81–90.
- Hu, H., Yao, W. & Wu, Y. (2017), ‘The robust EM-type algorithms for log-concave mixtures of regression models’, *Comput. Statist. Data Anal.* **111**, 14–26.
- Hu, S., Pei, Y., Liang, P. P. & Liang, Y.-C. (2020), ‘Deep neural network for robust modulation classification under uncertain noise conditions’, *IEEE Transactions on Vehicular Technology* **69**(1), 564–577.
- Huber, P. J. (1981), *Robust statistics*, John Wiley & Sons, Inc., New York. Wiley Series in Probability and Mathematical Statistics.
- Hunter, D. R. (2004), ‘MM algorithms for generalized Bradley-Terry models’, *Ann. Statist.* **32**(1), 384–406.
- Hunter, D. R. & Lange, K. (2000a), ‘Quantile regression via an MM algorithm’, *J. Comput. Graph. Statist.* **9**(1), 60–77.

- Hunter, D. R. & Lange, K. (2000b), ‘Rejoinder’, *Journal of Computational and Graphical Statistics* **9**(1), 52–59.
- Hunter, D. R. & Lange, K. (2004), ‘A tutorial on MM algorithms’, *Amer. Statist.* **58**(1), 30–37.
- Jamshidian, M. & Jennrich, R. I. (1993), ‘Conjugate gradient acceleration of the EM algorithm’, *J. Amer. Statist. Assoc.* **88**(421), 221–228.
- Karlis, D. & Xekalaki, E. (1998), ‘Minimum hellinger distance estimation for poisson mixtures’, *Computational Statistics & Data Analysis* **29**(1), 81–103.
- Khalili, A. & Vidyashankar, A. N. (2018), ‘Hypothesis testing in finite mixture of regressions: Sparsity and model selection uncertainty’, *Canadian Journal of Statistics* **46**(3), 429–457.
- Kokonendji, C. C. & Kiese, T. S. (2011), ‘Discrete associated kernels method and extensions’, *Statistical Methodology* **8**(6), 497–516.
- Kokonendji, C., Senga Kiessé, T. & Zocchi, S. S. (2007), ‘Discrete triangular distributions and non-parametric estimation for probability mass function’, *Journal of Nonparametric Statistics* **19**(6-8), 241–254.
- Lange, K. (1995), ‘A gradient algorithm locally equivalent to the EM algorithm’, *J. Roy. Statist. Soc. Ser. B* **57**(2), 425–437.
- Li, L., Vidyashankar, A. N., Diao, G. & Ahmed, E. (2019), ‘Robust inference after random projections via hellinger distance for location-scale family’, *Entropy* **21**(4), 348.
- Lindsay, B. G. (1994), ‘Efficiency versus robustness: the case for minimum Hellinger distance and related methods’, *The Annals of Statistics* **22**(2), 1081–1114.

- Louis, T. A. (1982), ‘Finding the observed information matrix when using the EM algorithm’, *J. Roy. Statist. Soc. Ser. B* **44**(2), 226–233.
- Lücke, J. & Forster, D. (2019), ‘k-means as a variational em approximation of gaussian mixture models’, *Pattern Recognition Letters* **125**, 349–356.
- Martinet, B. (1970), ‘Régularisation d’inéquations variationnelles par approximations successives’, *Rev. Française Informat. Recherche Opérationnelle* **4**(Sér. R-3), 154–158.
- McLachlan, G. (1995), On aitken’s method and other approaches for accelerating convergence of the em algorithm, *in* ‘Proceedings of the AC Aitken centenary conference’, pp. 201–209.
- Neal, R. M. & Hinton, G. E. (1998), A view of the em algorithm that justifies incremental, sparse, and other variants, *in* ‘Learning in graphical models’, Springer, pp. 355–368.
- Nielsen, F. & Sun, K. (2016), ‘Guaranteed bounds on the kullback-leibler divergence of univariate mixtures using piecewise log-sum-exp inequalities’, *arXiv preprint arXiv:1606.05850*.
- O’Neill, T. J. (1978), ‘Normal discrimination with unclassified observations’, *J. Amer. Statist. Assoc.* **73**(364), 821–826.
- Patra, S., Maji, A., Basu, A. & Pardo, L. (2013), ‘The power divergence and the density power divergence families: the mathematical connection’, *Sankhya B* **75**(1), 16–28.
- Pearson, K. (1894), ‘Contributions to the mathematical theory of evolution’, *Philosophical Transactions of the Royal Society of London. A* **185**, 71–110.
- Qin, Y. & Priebe, C. E. (2013), ‘Maximum l q-likelihood estimation via the expectation-maximization algorithm: a robust estimation of mixture models’, *Journal of the American Statistical Association* **108**(503), 914–928.



- R Hunter, D. & Lange, K. (2002), ‘Computing estimates in the proportional odds model’, *Ann. Inst. Statist. Math.* **54**(1), 155–168.
- Radhakrishna Rao, C. (1948), ‘The utilization of multiple measurements in problems of biological classification’, *J. Roy. Statist. Soc. Ser. B.* **10**, 159–193; discussion, 194–203.
- Read, T. R. C. & Cressie, N. A. C. (1988), *Goodness-of-fit statistics for discrete multivariate data*, Springer Series in Statistics, Springer-Verlag, New York.
- Rockafellar, R. (1976a), ‘Monotone operators and the proximal point algorithm’, *SIAM J. Control Optimization* **14**(5), 877–898.
- Rockafellar, R. T. (1976b), ‘Augmented Lagrangians and applications of the proximal point algorithm in convex programming’, *Math. Oper. Res.* **1**(2), 97–116.
- Rousseeuw, P. (1985), Multivariate estimation with high breakdown point, in ‘Mathematical statistics and applications, Vol. B (Bad Tatzmannsdorf, 1983)’, Reidel, Dordrecht, pp. 283–297.
- Sammaknejad, N., Zhao, Y. & Huang, B. (2019), ‘A review of the expectation maximization algorithm in data-driven process identification’, *Journal of process control* **73**, 123–136.
- Simpson, D. G. (1987), ‘Minimum Hellinger distance estimation for the analysis of count data’, *Journal of the American Statistical Association* **82**(399), 802–807.
- Simpson, D. G. (1989), ‘Hellinger deviance tests: efficiency, breakdown points, and examples’, *Journal of the American Statistical Association* **84**(405), 107–113.
- Stather, C. R. (1981), Robust statistical inference using Hellinger distance methods, PhD thesis, La Trobe University.

- Tseng, P. (2004), ‘An analysis of the EM algorithm and entropy-like proximal point methods’, *Math. Oper. Res.* **29**(1), 27–44.
- Vaida, F. (2005), ‘Parameter convergence for EM and MM algorithms’, *Statist. Sinica* **15**(3), 831–840.
- Woodward, W. A., Whitney, P. & Eslinger, P. W. (1995), ‘Minimum Hellinger distance estimation of mixture proportions’, *J. Statist. Plann. Inference* **48**(3), 303–319.
- Wu, C. J. (1983), ‘On the convergence properties of the em algorithm’, *The Annals of statistics* pp. 95–103.
- Zhao, R., Li, Y. & Sun, Y. (2020), ‘Statistical convergence of the em algorithm on gaussian mixture models’.