

ADAPTIVE ROBUST INTEGRATION OF INTERNAL DATA WITH EXTERNAL SUMMARIES UNDER DISTRIBUTIONAL SHIFT

Anonymous authors

Paper under double-blind review

ABSTRACT

Integrating evidence from heterogeneous datasets is challenging when predictor spaces differ and data distributions shift. Large datasets such as biobanks offer substantial sample sizes but often lack in-depth information due to cost constraints. In contrast, internal datasets from smaller analytic studies provide richer, individual-level detail. We propose a general Distributionally Robust Optimization (DRO) framework for integrating internal individual-level data with external summary-level data under distributional shift. Our method minimizes Cressie-Read divergence between a full model (fit to internal data with many predictors) and a reduced model (estimated from external data with fewer predictors), using a specialized nested-iteration algorithm. While effective under moderate shift, standard DRO can degrade when the distributional shift is severe. To mitigate this, we introduce an Empirical Bayes DRO (EB-DRO), which stabilizes estimates by adaptively shrinking toward internal-only solutions. We further develop an ensemble EB-DRO method that aggregates across multiple divergence families to improve robustness without selecting a single best family. Our proposed methods preserve privacy by operating on external summary statistics, support robust integration under shift, and enable valid inference when no shift is present. Simulations show that DRO improves over internal-only estimates under light shifts, EB-DRO adds stability under greater shifts, and ensemble EB-DRO achieves the most consistent robustness overall.

1 INTRODUCTION

The rapid expansion of large-scale data resources—spanning biobank-scale cohorts (e.g., UK Biobank, All of Us), healthcare systems with electronic health records, and digital traces from wearables and online platforms, is transforming both machine learning and applied sciences (Lapatás et al., 2015; Ashley, 2016; Russ et al., 2019). Large *external* datasets provide unprecedented breadth (sample size) but often lack detailed measurements due to cost constraints. Smaller *internal* studies, in contrast, collect richer covariates on a limited scale. The challenge is to integrate breadth with depth in a principled way, enhancing efficiency, robustness and generalizability of inference.

However, direct sharing of individual-level data is often blocked by privacy, regulatory, and logistical barriers (Rothstein, 2010; Kisselburgh and Beaver, 2022), so large cohorts typically release only *summary statistics* (e.g., regression coefficients, odds ratios). These summaries, in practice, are often derived from *reduced covariate sets* and from populations that differ from the richly measured internal study, creating pervasive *distributional shift* between internal and external data sources. If ignored, such shifts can bias estimates and erode the benefits of integration (Zhai and Han, 2022; Taylor et al., 2023). To address this, we develop a *distributionally robust data integration* framework that combines individual-level data with external summary statistics, enabling efficient learning that is stable under diverse shift scenarios.

Related work. Existing methods typically mitigate distributional shift by targeting specific forms, such as covariate shift in which covariate distributions differ but the conditional model remain invariant, or conditional shift in which the conditional model itself differs across studies. For example, Han and Lawless (2019) employed a small supplementary sample to bridge marginal differences in covariates. When the supplementary sample is unavailable, Estes et al. (2018) proposed a matrix-weighted empirical Bayes estimator that guarantees prediction error no worse and often better than the internal-only estimate. Gu et al. (2023) extended this idea to settings with multiple external stud-

ies. In a similar spirit, Han et al. (2024) used a James-Stein-type shrinkage method to cast the shift as a linear constraint. Furthermore, Sheng et al. (2024) modeled the distributional shift via unknown parameters and developed shift-adjusted estimation procedures. Besides, selective-integration methods, such as Chen et al. (2021); Zhai and Han (2022), applied penalized shrinkage to many external sources, borrowing strength only when they align with the internal population and thus avoiding bias under shift. Sheng et al. (2022) proposed a two-step procedure: conduct a homogeneous population hypothesis testing to identify compatible studies and integrate information only from those retained. While effective in their settings, these approaches are not designed to adapt broadly across different shift mechanisms.

Our approach. We leverage distributionally robust optimization (DRO), which seeks models that perform well under worst-case perturbations of the distribution (Duchi and Namkoong, 2021). Although widely studied in statistics and machine learning, DRO has not, to our knowledge, been applied to data integration with external summaries under moment constraints. In related work, empirical likelihood (EL) has provided a popular framework for incorporating summary statistics by treating them as moment conditions (Qin, 2000; Chatterjee et al., 2016; Han and Lawless, 2019; Kundu et al., 2019). In particular, DRO’s dual reweighting form generalizes EL, which preserves the empirical-likelihood structure while explicitly accommodating distributional shifts between studies. This makes DRO a natural and flexible foundation for the summary-based data integration. Building on the DRO framework, we introduce an empirical Bayes stabilization step that adaptively balances efficiency and robustness, and further develop an ensemble across divergence families. Taken together, these innovations yield an integrative method that delivers robust performance without requiring assumptions on the type or extent of distributional shift.

In this paper, we address the challenge of integrating internal individual-level data with external summary statistics when both distributions and covariate sets differ. Our main contributions are:

- We propose a DRO formulation for integrating internal individual-level data with external summaries under moment constraints, based on the Cressie-Read divergence with a dual representation and a scalable nested optimization algorithm.
- We develop an empirical Bayes stabilized variant (EB-DRO) that adaptively shrinks toward the internal-only estimator to safeguard against severe shift, and further introduce an ensemble EB-DRO that aggregates across the Cressie-Read family of divergences.
- Theoretically, we establish asymptotic normality of the DRO estimator, showing it is at least as efficient as the internal-only estimate without shift and converges to a pseudo-true value with a leading bias term under shift, and prove that EB-DRO enjoys an oracle risk guarantee.
- We demonstrate through simulations that DRO excels under mild shifts, EB-DRO stabilizes under stronger shifts, and the ensemble achieves the most consistent robustness across regimes.

2 METHODS

Problem setup. We consider two heterogeneous data sources:

- **Internal (individual-level) data:** $\{(Y_i, \mathbf{X}_i, \mathbf{Z}_i)\}_{i=1}^{N_{\text{in}}}$, where $\mathbf{X}_i \in \mathbb{R}^q$ are covariates measured in both studies and $\mathbf{Z}_i \in \mathbb{R}^{d-q}$ are additional covariates observed only internally.
- **External (summary-level) data:** a large study with only \mathbf{X} observed, providing the maximum likelihood estimate (MLE) θ^* for the reduced model $f_{\theta}(Y | \mathbf{X})$. Due to the massive sample size $N_{\text{ext}} \gtrsim 10^5$ or 10^6 , we treat θ^* as the true reduced-model parameter θ .

We model Y via the following generalized linear models (GLMs):

$$Y | (\mathbf{X}, \mathbf{Z}) \sim \text{GLM}\{\eta = (\mathbf{X}^\top, \mathbf{Z}^\top)\beta\}, \quad Y | \mathbf{X} \sim \text{GLM}(\eta = \mathbf{X}^\top\theta),$$

where $\beta \in \mathbb{R}^d$ parameterizes the internal study and $\theta \in \mathbb{R}^q$ is the reduced-model parameter for the external study. Notably, θ is not simply the \mathbf{X} -subvector of β , since in GLMs the estimated effect of a feature depends on other features included in the model.

Distributional shift. The two studies may involve different populations or follow different protocols, leading to *distributional shift*. Such shifts may stem from changes in the marginal covariate

distribution, the conditional distribution, or both. Our DRO formulation is agnostic to the shift type and provides robustness in all cases, though the conditional shift poses greater methodological challenges. Let β^ω denote the external full-model parameter, which may differ from the internal β ; the reduced parameter θ^* is the projection of β^ω onto the shared covariates \mathbf{X} . Our goal is to robustly estimate β by integrating internal data with the external summary θ^* , leveraging the large external sample while accounting for possible shifts (see Fig. A1).

2.1 DISTRIBUTIONALLY ROBUST DATA INTEGRATION

2.1.1 FORMULATION AND DUAL REPRESENTATION

We estimate β via a distributionally robust optimization that augments the internal likelihood with a divergence penalty and a moment constraint. Let $p = (p_1, \dots, p_{N_{\text{in}}})$ be non-negative weights defining a reweighted empirical distribution over the N_{in} internal samples. Our objective is

$$\begin{aligned} \min_{\beta} & \left\{ -\frac{2}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \log f_{\beta}(Y_i | \mathbf{X}_i, \mathbf{Z}_i) + \min_{p_1, \dots, p_{N_{\text{in}}}} \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \phi_k \left(\frac{p_i}{1/N_{\text{in}}} \right) \right\}, \\ \text{s.t.} & \sum_{i=1}^{N_{\text{in}}} p_i = 1, \quad p_i \geq 0, \quad \sum_{i=1}^{N_{\text{in}}} p_i \mathbf{u}(\mathbf{X}_i, \mathbf{Z}_i; \beta, \theta^*) = \mathbf{0}. \end{aligned} \quad (1)$$

Here f_{β} is the GLM likelihood for observation i , and $\phi_k(\cdot)$ denotes a Cressie-Read ϕ -divergence, which regularizes reweighting by measuring the deviation of p from the uniform weights, with k tuning its strength. The vector-valued moment function \mathbf{u} encodes the external summary, defined as

$$\mathbf{u}(\mathbf{X}, \mathbf{Z}; \beta, \theta^*) = \mathbb{E}_{f_{\beta}(Y|\mathbf{X},\mathbf{Z})} \{ \mathbf{h}(Y, \mathbf{X}; \theta^*) | \mathbf{X}, \mathbf{Z} \} = \int \mathbf{h}(Y, \mathbf{X}; \theta^*) f_{\beta}(Y | \mathbf{X}, \mathbf{Z}) dY,$$

where $\mathbf{h}(Y, \mathbf{X}; \theta^*) \in \mathbb{R}^q$ is the reduced-model estimating function used to obtain summary θ^* from external data. Although individual-level external data are unavailable, evaluating $\mathbf{h}(Y, \mathbf{X}; \theta^*)$ based on internal Y and \mathbf{X} enables moment matching for cross-source compatibility.

Theorem 1 (Variational dual of DRO integration). *Let ϕ be a proper, closed, convex divergence with $\phi(1) = 0$ and convex conjugate ϕ^* . We write $\mathbf{u}_i(\beta) \equiv \mathbf{u}(\mathbf{X}_i, \mathbf{Z}_i; \beta, \theta^*)$, then the estimator in Eq. 1 admits the saddle formulation:*

$$\min_{\beta} \max_{\gamma \in \mathbb{R}, \lambda \in \mathbb{R}^q} \left[-\frac{2}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \log f_{\beta}(Y_i | \mathbf{X}_i, \mathbf{Z}_i) + \gamma - \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \phi^* \{ \gamma + \lambda^{\top} \mathbf{u}_i(\beta) \} \right].$$

Specializing to the Cressie-Read family, the inner maximization over γ in Theorem 1 can be solved in closed form, yielding a compact expression depending only on λ , stated below.

Corollary 1 (Cressie-Read dual form). *For the Cressie-Read divergence $\phi_k(x) = \frac{2}{k(k+1)}(x^{-k} - 1)$ with conjugate $\phi_k^*(y) = -\frac{2}{k} \left(-\frac{k+1}{2} y \right)^{k/(k+1)} + \frac{2}{k(k+1)}$, the dual DRO problem reduces to the following form (derivation in Appendix C):*

$$\min_{\beta} \max_{\lambda \in \mathbb{R}^q} \left\{ -\sum_{i=1}^{N_{\text{in}}} \log f_{\beta}(Y_i | \mathbf{X}_i, \mathbf{Z}_i) + H_k(\lambda, \beta) \right\}, \quad (2)$$

where

$$H_k(\lambda, \beta) = \frac{N_{\text{in}}}{k(k+1)} \left[\frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \{ 1 + \lambda^{\top} \mathbf{u}_i(\beta) \}^{\frac{k}{k+1}} \right]^{k+1}.$$

We define the optimum value of Eq. 2 as the DRO estimator $\hat{\beta}_{\text{DRO},k}$.

For a fixed β , maximizing H_k is equivalent to minimizing M_k , whose specific forms are reported in Table 1. Since the mapping $H_k \mapsto M_k$ is monotone, we have

$$\max_{\lambda \in \mathbb{R}^q} H_k(\lambda) \iff \min_{\lambda \in \mathbb{R}^q} M_k(\lambda).$$

Notably, the inner minimization in $M_k(\boldsymbol{\lambda})$ is convex for all k in the Cressie-Read family. For $k < 0$, convexity follows from the map $t \mapsto t^{k/(k+1)}$, while for $k > 0$, the sign adjustment in the last row of Table 1 ensures convexity. In addition, Table 1 indicates that the proposed framework naturally incorporates both the empirical likelihood (EL) and the exponential tilting (ET) methods, thereby unifying a broad class of classical estimators. These properties are formalized in Lemmas A1-A2 of Appendix D: Lemma A1 establishes the EL and ET limits, and Lemma A2 proves the convexity of the inner dual.

Table 1: Inner objectives across the Cressie-Read family.

Special Cases	k	$H_k(\boldsymbol{\lambda})$	$M_k(\boldsymbol{\lambda})$
Empirical likelihood	0	$\sum_i \log(1 + \boldsymbol{\lambda}^\top \mathbf{u}_i)$	$-\sum_i \log(1 + \boldsymbol{\lambda}^\top \mathbf{u}_i)$
Exponential tilting	-1	$-N_{\text{in}} \log\{\sum_i \exp(\boldsymbol{\lambda}^\top \mathbf{u}_i/2)/N_{\text{in}}\}$	$\sum_i \exp(\boldsymbol{\lambda}^\top \mathbf{u}_i/2)$
GMM	-2	$N_{\text{in}}/2\{\sum_i (1 + \boldsymbol{\lambda}^\top \mathbf{u}_i)^2/N_{\text{in}}\}^{-1}$	$\sum_i (1 + \boldsymbol{\lambda}^\top \mathbf{u}_i)^2$
Pearson/Neyman χ^2	1	$N_{\text{in}}/2\{\sum_i \sqrt{1 + \boldsymbol{\lambda}^\top \mathbf{u}_i}/N_{\text{in}}\}^2$	$-\sum_i \sqrt{1 + \boldsymbol{\lambda}^\top \mathbf{u}_i}$
Freeman-Tukey	-1/2	$-4N_{\text{in}}\{\sum_i (1 + \boldsymbol{\lambda}^\top \mathbf{u}_i)^{-1}/N_{\text{in}}\}^{1/2}$	$\sum_i (1 + \boldsymbol{\lambda}^\top \mathbf{u}_i)^{-1}$
General CR family	$k < 0$	$N_{\text{in}}/\{k(k+1)\}\{\sum_i (1 + \boldsymbol{\lambda}^\top \mathbf{u}_i)^{k/(k+1)}/N_{\text{in}}\}^{k+1}$	$\sum_i (1 + \boldsymbol{\lambda}^\top \mathbf{u}_i)^{k/(k+1)}$
	$k > 0$	(same as above)	$-\sum_i (1 + \boldsymbol{\lambda}^\top \mathbf{u}_i)^{k/(k+1)}$

Consequently, Eq. 1 provides a novel DRO framework for data integration with external summaries. It recovers classical estimators such as EL, ET, GMM, and χ^2 . More importantly, its scope is considerably broader: varying the Cressie-Read index k can generate a continuum of penalty forms, each shaping the bias-variance-robustness trade-off in distinct ways. This flexibility is central to our proposed framework: while classical methods arise as fixed choices of k , our approach enables systematic exploration across a broad continuum of Cressie-Read families beyond the traditional cases.

2.1.2 ROBUSTNESS PROPERTIES AND INSIGHTS

We now turn to the theoretical properties of the proposed DRO estimator. In particular, we study its asymptotic behavior under varying levels of distributional shift.

Theorem 2 (Asymptotic normality under no shift). *Assume the regularity conditions in Appendix E and the well-specified moments $\mathbb{E}\{\mathbf{u}_i(\boldsymbol{\beta}^*)\} = \mathbf{0}$. Let $(\hat{\boldsymbol{\beta}}_{\text{DRO},k}, \hat{\boldsymbol{\lambda}}_k)$ solve the first-order conditions of the DRO objective with Cressie-Read index k . We define $\ell_i(\boldsymbol{\beta}) := \log f_{\boldsymbol{\beta}}(Y_i | \mathbf{X}_i, \mathbf{Z}_i)$, then*

$$\sqrt{N_{\text{in}}} \begin{pmatrix} \hat{\boldsymbol{\beta}}_{\text{DRO},k} - \boldsymbol{\beta}^* \\ \hat{\boldsymbol{\lambda}}_k \end{pmatrix} \xrightarrow{d} \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{J} & -\mathbf{C}^\top \\ -\mathbf{C} & \boldsymbol{\Omega}_2 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{J} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Omega}_1 \end{bmatrix} \begin{bmatrix} \mathbf{J} & -\mathbf{C}^\top \\ -\mathbf{C} & \boldsymbol{\Omega}_2 \end{bmatrix}^{-1}\right),$$

where $\mathbf{J} := \mathbb{E}\{-\nabla_{\boldsymbol{\beta}\boldsymbol{\beta}}^2 \ell_i(\boldsymbol{\beta}^*)\}$, $\mathbf{C} := \mathbb{E}\{\nabla_{\boldsymbol{\lambda}\boldsymbol{\beta}}^2 H_k(\mathbf{0}, \boldsymbol{\beta}^*)\}$, $\boldsymbol{\Omega}_1 = \mathbb{E}\{\nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}^*) \nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}^*)^\top\}$ and $\boldsymbol{\Omega}_2 = -\mathbb{E}\{\nabla_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^2 H_k(\mathbf{0}, \boldsymbol{\beta}^*)\}$. Profiling out $\boldsymbol{\lambda}$ yields

$$\sqrt{N_{\text{in}}}(\hat{\boldsymbol{\beta}}_{\text{DRO},k} - \boldsymbol{\beta}^*) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1} \mathbf{V} \mathbf{A}^{-1}), \quad \mathbf{A} := \mathbf{J} + \mathbf{C}^\top \boldsymbol{\Omega}_2^{-1} \mathbf{C}, \quad \mathbf{V} := \mathbf{J} + \mathbf{C}^\top \boldsymbol{\Omega}_2^{-1} \boldsymbol{\Omega}_1 \boldsymbol{\Omega}_2^{-1} \mathbf{C}.$$

Denote $\mathbf{G} := \mathbb{E}\{\nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}^*)\}$ and $\mathbf{S} := \mathbb{E}\{\mathbf{u}_i(\boldsymbol{\beta}^*) \mathbf{u}_i(\boldsymbol{\beta}^*)^\top\}$, for the Cressie-Read family, Lemmas A3-A4 show

$$\mathbf{C} = \mathbf{G}/(k+1), \quad \boldsymbol{\Omega}_1 = \boldsymbol{\Omega}_2 = \mathbf{S}/(k+1)^2.$$

Clearly, all $(k+1)$ -factors cancel, so the first-order limit coincides with the efficient EL/ET/GMM variance under correct specification. Full Derivations are in Appendix E.

Remark 1. Note that the asymptotic variance $(\mathbf{J} + \mathbf{G}^\top \mathbf{S}^{-1} \mathbf{G})^{-1}$ is no larger than \mathbf{J}^{-1} . Thus, relative to the “naive” estimator with variance \mathbf{J}^{-1} that ignores the external summaries, the DRO estimator achieves strictly greater efficiency whenever $\mathbf{G} \neq \mathbf{0}$, under no distributional shift.

Challenges under distributional shift. Deriving the asymptotic distribution of the DRO estimator is considerably more involved under shift than in the no-shift case. The main difficulty is that the dual penalty $H_k(\boldsymbol{\lambda}, \boldsymbol{\beta})$ is not additively separable across samples: it is a nonlinear functional of the empirical mean of transformed scores. This prevents direct use of standard M-estimation theory. To proceed, Lemma A6 provides an influence-function expansion of H_k that linearizes the non-additive structure and represents first-order fluctuations as sums of i.i.d. per-sample terms.

Theorem 3 (Asymptotic normality under distributional shift). *Assume the regularity conditions in Appendix E. Let $(\hat{\beta}_{\text{DRO},k}, \hat{\lambda}_k)$ solve the DRO estimating equations with Cressie-Read index k , and let β^\dagger denote the population target. Then*

$$\sqrt{N_{\text{in}}}(\hat{\beta}_{\text{DRO},k} - \beta^\dagger) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{V}_\beta^\dagger),$$

where \mathbf{V}_β^\dagger is a profile sandwich variance of the form

$$\begin{aligned} \mathbf{V}_\beta^\dagger = & (\mathbf{A}_\beta^\dagger)^{-1} \{ \mathbf{K}^\dagger + (\mathbf{C}^\dagger)^\top (\boldsymbol{\Omega}_2^\dagger)^{-1} \boldsymbol{\Omega}_1^\dagger (\boldsymbol{\Omega}_2^\dagger)^{-1} \mathbf{C}^\dagger \\ & + (\mathbf{C}^\dagger)^\top (\boldsymbol{\Omega}_2^\dagger)^{-1} \mathbf{K}_{sg}^\dagger + (\mathbf{K}_{sg}^\dagger)^\top (\boldsymbol{\Omega}_2^\dagger)^{-1} \mathbf{C}^\dagger \} (\mathbf{A}_\beta^\dagger)^{-1}. \end{aligned}$$

Remark 2. The variance expression reflects both the usual score variability and the additional contribution from the dual multipliers $\hat{\lambda}_k$, which capture the effect of distributional shift. All detailed definitions of the block matrices $\mathbf{A}_\beta^\dagger, \mathbf{K}^\dagger, \mathbf{C}^\dagger, \boldsymbol{\Omega}_1^\dagger, \boldsymbol{\Omega}_2^\dagger, \mathbf{K}_{sg}^\dagger$, together with the linearization of H_k into summable influence functions, are provided in Lemma A6 and Theorem A6 of Appendix F.

The estimator $\hat{\beta}_{\text{DRO},k}$ converges to a pseudo-true value β^\dagger , which generally differs from β^* under distributional shift. Lemma A7 in Appendix F gives that a local expansion around β^* :

$$\hat{\beta}_{\text{DRO},k} - \beta^* = (\mathbf{J} - \mathbf{G}^\top \mathbf{S}^{-1} \mathbf{G})^{-1} \mathbf{G}^\top \mathbf{S}^{-1} \boldsymbol{\delta} + O_p(N_{\text{in}}^{-1/2}) + o(\|\boldsymbol{\delta}\|),$$

with shift term $\boldsymbol{\delta} = \mathbb{E}\{\mathbf{u}_i(\beta^*)\}$. Hence, under mild shifts ($\|\boldsymbol{\delta}\|$ small), the bias is negligible relative to sampling noise, while it dominates and persists asymptotically under large shifts. Importantly, the leading bias and variance are k -free, though k affects higher-order terms and finite-sample behavior. Thus, DRO estimators are robust under small misspecification but deteriorate as shifts grow, motivating the stabilized EB-DRO method in Section 2.2.

2.1.3 ALGORITHMIC DETAILS

The DRO problem is solved by a nested optimization (Algorithm 1), with an outer update for β and an inner update for λ . Both use Newton-Raphson steps scaled by the Hessian, with stability ensured via a backtracking line search: the update step is multiplied by a factor τ , halved until the descent condition is satisfied. In the outer loop, the objective is $\ell_H(\beta) = -\sum_i \log f_\beta(Y_i | \mathbf{X}_i, \mathbf{Z}_i) + H_k\{\hat{\lambda}(\beta), \beta\}$, *implicit differentiation* is applied to incorporate the dependence of the inner optimizer $\hat{\lambda}(\beta)$ on β . This approach extends the empirical likelihood algorithm of Han and Lawless (2019) to the full Cressie-Read family, providing a unified method for any divergence index k . Gradient/Hessian derivations and implicit differentiation details appear in Appendix H.

Algorithm 1 Nested DRO solver (Cressie-Read index k) with Newton updates and backtracking

```

1: Input: initial  $\beta^{(0)}$ ; tolerances  $(\varepsilon_{\text{out}}, \varepsilon_{\text{in}})$ ; max iters  $(T_{\text{out}}, T_{\text{in}})$ ;
2: for  $t = 0, 1, \dots, T_{\text{out}}$  do # outer loop
3:   Inner solve (given  $\beta^{(t)}$ ): set  $\lambda^{(0)} = \mathbf{0}$ 
4:   for  $s = 0, 1, \dots, T_{\text{in}}$  do # inner loop
5:     if  $\|\nabla_\lambda M_k(\lambda^{(s)}, \beta^{(t)})\| \leq \varepsilon_{\text{in}}$  then break
6:     end if
7:     Update  $\lambda^{(s+1)} \leftarrow \lambda^{(s)} - \tau_{\text{in}} \{ \nabla_{\lambda\lambda}^2 M_k(\lambda^{(s)}, \beta^{(t)}) \}^{-1} \nabla_\lambda M_k(\lambda^{(s)}, \beta^{(t)})$  with backtracking ( $\tau_{\text{in}} \leftarrow \tau_{\text{in}}/2$ ) until descent holds
8:   end for
9:   Set  $\hat{\lambda}(\beta^{(t)}) \leftarrow \lambda^{(s)}$ 
10:  if  $\|\nabla_\beta \ell_H\{\beta^{(t)}, \hat{\lambda}(\beta^{(t)})\}\| \leq \varepsilon_{\text{out}}$  then break
11:  end if
12:  Update  $\beta^{(t+1)} \leftarrow \beta^{(t)} - \tau_{\text{out}} [ \nabla_{\beta\beta}^2 \ell_H\{\beta^{(t)}, \hat{\lambda}(\beta^{(t)})\} ]^{-1} \nabla_\beta \ell_H\{\beta^{(t)}, \hat{\lambda}(\beta^{(t)})\}$  with backtracking ( $\tau_{\text{out}} \leftarrow \tau_{\text{out}}/2$ ) until descent holds
13: end for
14: Return:  $\hat{\beta}_{\text{DRO},k} = \beta^{(t)}, \hat{\lambda}_k = \lambda^{(s)}$ 

```

2.2 EMPIRICAL BAYES-STABILIZED DRO (EB-DRO)

Although the DRO estimator $\widehat{\beta}_{\text{DRO},k}$ can substantially improve efficiency under mild distributional shift, its bias may grow quickly as the shift increases. Under severe shifts, $\widehat{\beta}_{\text{DRO},k}$ can even perform worse than the internal-only estimator $\widehat{\beta}_I$, taken as the maximum likelihood estimator (MLE) from the GLM fit using the internal sample. To guard against this deterioration, we introduce an Empirical Bayes (EB) stabilization step: an adaptive shrinkage procedure that combines the DRO and naive estimators. This construction preserves DRO’s efficiency gains under mild shifts while ensuring robustness under large shifts.

Let $\mathbf{Z} := \widehat{\beta}_{\text{DRO},k} - \widehat{\beta}_I$, $\widehat{\mathbf{A}} := \mathbf{Z}\mathbf{Z}^\top$ and $\widehat{\Sigma}_I := \text{Var}(\widehat{\beta}_I)$. The EB-DRO estimator is

$$\widehat{\beta}_{\text{EB},k} = (\mathbf{I} - \widehat{\mathbf{W}})\widehat{\beta}_{\text{DRO},k} + \widehat{\mathbf{W}}\widehat{\beta}_I, \quad \widehat{\mathbf{W}} := \widehat{\mathbf{A}}(\widehat{\mathbf{A}} + \widehat{\Sigma}_I)^{-1}.$$

This form arises naturally from a Gaussian hierarchical model in which $\widehat{\beta}_{\text{DRO},k}$ serves as the prior center (uncertainty $\widehat{\mathbf{A}}$) and $\widehat{\beta}_I$ is modeled as a Gaussian observation with covariance $\widehat{\Sigma}_I$, as detailed in Appendix I.1. Furthermore, the following proposition makes explicit the resulting shrinkage property, with proof given in Appendix I.2.

Proposition 1. *The EB-DRO estimator admits the equivalent representation*

$$\widehat{\beta}_{\text{EB},k} = \widehat{\beta}_{\text{DRO},k} + \alpha\mathbf{Z} = (1 - \alpha)\widehat{\beta}_{\text{DRO},k} + \alpha\widehat{\beta}_I, \quad \alpha = \mathbf{Z}^\top \widehat{\Sigma}_I^{-1} \mathbf{Z} / (1 + \mathbf{Z}^\top \widehat{\Sigma}_I^{-1} \mathbf{Z}) \in (0, 1). \quad (3)$$

This proposition shows that EB-DRO adaptively interpolates between the DRO and naive estimators depending on their Mahalanobis discrepancy. When $\widehat{\beta}_{\text{DRO},k}$ and $\widehat{\beta}_I$ are close, $\alpha \approx 0$ and EB-DRO coincides with DRO, preserving efficiency. When they diverge, $\alpha \rightarrow 1$ and EB-DRO reverts to the naive estimator, ensuring robustness.

Theorem 4 (Oracle risk guarantee and regret). *For $w \in [0, 1]$, define*

$$\widehat{\beta}(w) = (1 - w)\widehat{\beta}_{\text{DRO},k} + w\widehat{\beta}_I, \quad R(w) = \mathbb{E}\|\widehat{\beta}(w) - \beta^*\|_2^2.$$

Let $\Pi_{[0,1]}(t) = \min\{1, \max\{0, t\}\}$. *The oracle weight*

$$w^* = \Pi_{[0,1]}\{a - m/(a + c - 2m)\}, \quad a = \|\mathbf{b}_k\|_2^2 + \text{tr}(\mathbf{V}_{\text{DRO},k}), \quad c = \text{tr}(\widehat{\Sigma}_I), \quad m = \text{tr}(\mathbf{C}_k),$$

with $\mathbf{b}_k = \mathbb{E}(\widehat{\beta}_{\text{DRO},k}) - \beta^$, $\mathbf{V}_{\text{DRO},k} = \text{Var}(\widehat{\beta}_{\text{DRO},k})$ and $\mathbf{C}_k = \text{Cov}(\widehat{\beta}_{\text{DRO},k}, \widehat{\beta}_I)$, minimizes $R(w)$ and satisfies*

$$R(w^*) \leq \min\{R(0), R(1)\}.$$

Hence the oracle EB-DRO is never worse than DRO or naive. Moreover, for any $w \in [0, 1]$,

$$R(w) = R(w^*) + (a + c - 2m)(w - w^*)^2,$$

giving an exact quadratic regret identity.

Remark 3. The EB-DRO estimator uses the data-adaptive weight α from Eq. 3. Although $\alpha \neq w^*$ in general, Appendix I.3 shows that $\alpha \xrightarrow{P} \alpha_0$, where $\alpha_0 = \mathbf{b}_k^\top \widehat{\Sigma}_I^{-1} \mathbf{b}_k / (1 + \mathbf{b}_k^\top \widehat{\Sigma}_I^{-1} \mathbf{b}_k)$. In the extreme regimes of no shift or large shift, $\alpha_0 = w^*$, so EB-DRO inherits the oracle guarantee asymptotically. Moreover, Theorem 4 implies that any plug-in weight (including α) incurs at most a quadratic excess-risk term relative to w^* , which vanishes asymptotically. This motivates our ensemble EB-DRO strategy: by aggregating across divergence indices k , we further dampen plug-in variability and stabilize performance across a wider range of shifts.

2.3 ENSEMBLE EB-DRO

Different divergence indices k yield EB-DRO estimators $\{\widehat{\beta}_{\text{EB},k} : k \in \mathcal{K}\}$, each associated with a posterior covariance

$$\widehat{\Omega}_k = (\widehat{\Sigma}_I^{-1} + \widehat{\mathbf{A}}_k^+)^{-1},$$

where $\widehat{\mathbf{A}}_k^+$ denotes the Moore-Penrose pseudoinverse of $\widehat{\mathbf{A}}_k = (\widehat{\beta}_{\text{DRO},k} - \widehat{\beta}_I)(\widehat{\beta}_{\text{DRO},k} - \widehat{\beta}_I)^\top$.

To stabilize performance across various shifts, we aggregate the individual EB-DRO estimators into a precision-weighted ensemble:

$$\hat{\beta}_{\text{Ens}} = \left\{ \sum_{k \in \mathcal{K}} (\hat{\Omega}_k)^{-1} \right\}^{-1} \left\{ \sum_{k \in \mathcal{K}} (\hat{\Omega}_k)^{-1} \hat{\beta}_{\text{EB},k} \right\}, \quad (4)$$

This ensemble can be interpreted as the posterior mean of a higher-level Gaussian model in which $\{\hat{\beta}_{\text{EB},k}\}_{k \in \mathcal{K}}$ serve as complementary noisy signals about the same target.

Theorem 5 (Stability of ensemble EB-DRO). *Let $\Omega_{k\ell} = \text{Cov}(\hat{\beta}_{\text{EB},k}, \hat{\beta}_{\text{EB},\ell})$ and assume $\hat{\Omega}_k \xrightarrow{p} \Omega_{kk}$ for all k . Then $\mathbb{E}\|\hat{\beta}_{\text{Ens}} - \beta^*\|_2^2 \leq \min_{j \in \mathcal{K}} \mathbb{E}\|\hat{\beta}_{\text{EB},j} - \beta^*\|_2^2 + \Delta_{\text{cross}} + o(1)$, where Δ_{cross} depends only on the off-diagonal cross-covariances $\{\Omega_{k\ell}\}_{k \neq \ell}$ (explicit form in Appendix J).*

Theorem 5 shows that the precision-weighted ensemble achieves mean squared error no larger than the best EB-DRO up to a correlation penalty Δ_{cross} . Individual EB-DROs can be unstable: plug-in variability in α (Eq. 3) may underweight the naive estimator, leading to inflated risk and sometimes performing worse than naive. The ensemble mitigates this by averaging across k , dampening variability and pulling risk back toward a safe baseline, thereby ensuring stability across shift regimes.

Summary. We introduced three estimators: (i) *DRO*, which leverages external information to improve efficiency under mild shift but may suffer bias under severe shift; (ii) *EB-DRO*, which adaptively shrinks between DRO and the internal estimator to stabilize performance; and (iii) the *ensemble EB-DRO*, a precision-weighted aggregation across divergences k that dampens plug-in variability and ensures robustness across a wider range of shifts. We next investigate their empirical behavior through simulations.

3 SIMULATION EXPERIMENTS

3.1 SIMULATION SETUP

We simulate binary outcomes from a logistic GLM. Let $\mathbf{X} = (1, X_1, \dots, X_5)^\top$ include an intercept and five covariates. The true coefficient vector is $\beta^* = (1, -1, -1, 1, -1, -1)^\top = (\beta_0^*, \beta_1^*, \dots, \beta_5^*)^\top \in \mathbb{R}^6$, where β_0^* is the intercept. Internal outcomes ($N_{\text{in}} = 500$) follow

$$Y_{\text{in}} | \mathbf{X} \sim \text{Bernoulli}\{\text{expit}(\mathbf{X}^\top \beta^*)\}, \quad \text{expit}(t) = \frac{1}{1+e^{-t}}.$$

External outcomes ($N_{\text{ext}} = 100,000$) are generated under a *shifted* coefficient β_{ext} , defined as a perturbation of β^* . The external summaries θ are obtained by fitting a *reduced* GLM with covariates $(1, X_1, X_2, X_3)^\top \in \mathbb{R}^4$ to the external sample.

Shift scenarios. We consider three families of distributional shift:

- (A) **Intercept shift:** Replace β_0^* with $\beta_0^* + \Delta$, $\Delta \in \{0, 0.1, \dots, 1.0\}$ (11 settings), keeping slopes fixed. *Intuition:* Change baseline prevalence without altering covariate effects.
- (B) **Tail perturbation:** For a fraction $p \in \{0, 0.05, \dots, 1.0\}$ (21 settings) of external units in the upper tail of X_1 , reduce $(\beta_1^*, \beta_2^*, \beta_3^*)$ by 1.5; the remainder use β^* . *Intuition:* Induces a localized, subgroup-specific shift concentrated in high- X_1 regions.
- (C) **Angle (directional) shift:** Rotate true β^* toward a fixed orthogonal direction v via Gram-Schmidt by an angle $\varphi \in \{0, 0.0776, \dots, 0.6981\}$ (10 values, up to $\approx 40^\circ$), i.e., $\beta_{\text{ext}}(\varphi) = \cos(\varphi)\beta^* + \sin(\varphi)v$. *Intuition:* Induces a global misspecification that changes the direction of the signal while preserving its magnitude.

Estimators and evaluation. At each shift setting, we generate one external dataset ($N_{\text{ext}} = 100,000$) for $\hat{\theta}$ and repeatedly draw internal datasets ($N_{\text{in}} = 500$, $R = 1000$ replicates). On each replicate we compute: (i) the internal MLE $\hat{\beta}_I$ (naive); (ii) DRO estimators $\hat{\beta}_{\text{DRO},k}$ across 19 divergence indices k ; (iii) EB-DRO $\hat{\beta}_{\text{EB},k}$ for the same set of k values; and (iv) the ensemble $\hat{\beta}_{\text{Ens}}$. Performance is measured by mean ℓ_2 error $\mathbb{E}\|\hat{\beta} - \beta^*\|_2$, averaged over replicates.

3.2 SIMULATION RESULTS

We evaluate DRO and EB-DRO estimators across 19 divergence indices k (values shown in the legends of Figs. 1-2). Specially, Fig. 1 shows that robustness varies strongly with k : some families

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

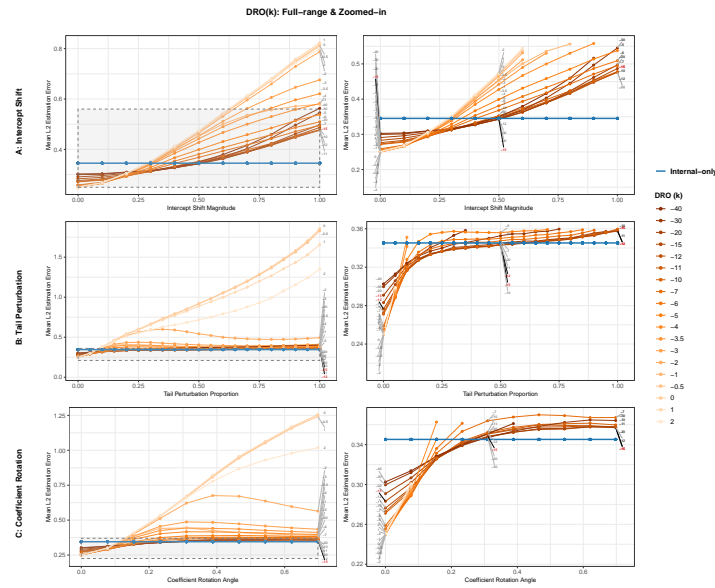


Figure 1: Performance of DRO across 19 divergence families. Left: full range of results. Right: zoomed-in view of the shaded region to highlight finer differences. Internal-only estimates are included for reference. Rows: (A) intercept shift, (B) tail perturbation, and (C) coefficient rotation.

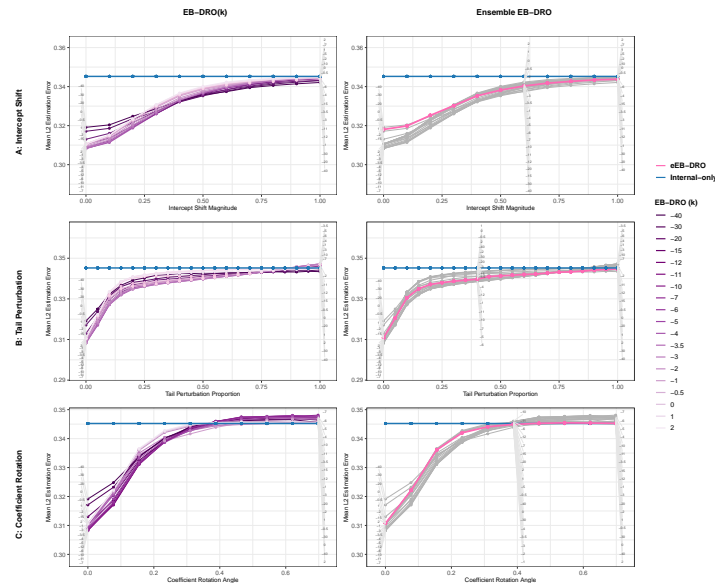


Figure 2: Performance of EB-DRO(k) and ensemble EB-DRO. Left: EB-DRO(k) for each family (colored). Right: per-family EB-DRO(k) (gray) with ensemble EB-DRO (pink) overlaid for comparison. Internal-only estimates are also shown.

(e.g., $k = -15$) remain competitive with or better than the naive estimator even under moderately large shifts, while others deteriorate quickly. This illustrates the k -dependent trade-off between bias, variance, and robustness: certain divergences yield stable performance across shifts, whereas others overshoot once bias accumulates. The pattern is consistent with Theorem 2 and Theorem 3, which shows DRO improves efficiency under mild shifts but becomes biased under severe ones. Exploring the full continuum of Cressie-Read divergences therefore extends beyond classical fixed cases (EL, ET, GMM, χ^2) to reveal a richer spectrum of robustness behaviors.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

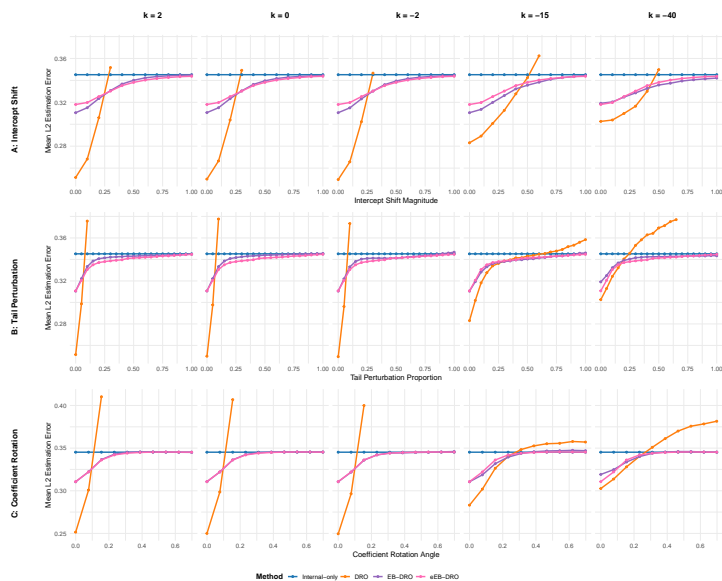


Figure 3: Performance of internal-only, DRO, EB-DRO(k), and ensemble EB-DRO for selected families ($k = 2, 0, -2, -15, -40$). Results for all 19 families are given in Figs A2-A4.

Fig. 2 displays that EB-DRO markedly improves over DRO, remaining close to or better than the naive estimator across most settings. In extreme cases (e.g., angle shifts > 0.4 or tail perturbations > 0.9), some EB-DRO(k) curves become slightly worse than naive due to noisy plug-in weights. Notably, the ensemble consistently pulls results back toward the naive baseline under severe shifts, highlighting its stability.

Fig. 3 directly compares internal-only, DRO, EB-DRO, and the ensemble. The pattern is clear: DRO can improve efficiency over naive under mild shifts but becomes unstable under severe ones; EB-DRO stabilizes performance by adaptively shrinking toward naive; and the ensemble consolidates these gains, avoiding DRO’s blow-ups while preserving EB-DRO’s efficiency under small shifts. These findings confirm the theoretical guarantees developed in Section 2.

4 DISCUSSION

We proposed a distributionally robust framework for integrating internal individual-level data with external summaries, introducing three estimators: (i) *DRO*, which improves efficiency under mild shifts but suffers bias under severe ones; (ii) *EB-DRO*, which stabilizes performance by adaptively shrinking between DRO and the internal estimator; and (iii) the *ensemble EB-DRO*, which aggregates across divergences via precision weighting to dampen plug-in variability.

This framework also unifies and extends classical estimators: the duality links DRO to the Cressie-Read family, recovering empirical likelihood, exponential tilting, GMM, and χ^2 as special cases. EB-DRO parallels empirical Bayes shrinkage, offering protection against noisy summaries, and the ensemble inherits the strengths of well-performing divergences while avoiding catastrophic failures.

Future directions. Several avenues remain open. First, extending the theory to high-dimensional covariates would broaden applicability. Second, exploring adaptive weighting schemes beyond precision weighting could further enhance stability. Third, combining our framework with conformal or post-selection inference may yield valid uncertainty quantification under severe distributional shifts. Finally, applying the framework to large-scale biomedical studies, such as those involving epigenetic profiling, where detailed measurements are rarely available in large cohorts, will be crucial for demonstrating practical impact.

REFERENCES

- 486
487
488 Ashley, E. A. (2016). Towards precision medicine. *Nature Reviews Genetics*, 17(9):507–522.
- 489
490 Chatterjee, N., Chen, Y.-H., Maas, P., and Carroll, R. J. (2016). Constrained maximum likelihood
491 estimation for model calibration using summary-level information from external big data sources.
492 *Journal of the American Statistical Association*, 111(513):107–117.
- 493
494 Chen, Z., Ning, J., Shen, Y., and Qin, J. (2021). Combining primary cohort data with external
495 aggregate information without assuming comparability. *Biometrics*, 77(3):1024–1036.
- 496
497 Duchi, J. C. and Namkoong, H. (2021). Learning models with uniform performance via distribu-
498 tionally robust optimization. *Annals of Statistics*, 49(3):1378–1406.
- 499
500 Estes, J. P., Mukherjee, B., and Taylor, J. M. (2018). Empirical bayes estimation and prediction
501 using summary-level information from external big data sources adjusting for violations of trans-
502 portability. *Statistics in Biosciences*, 10(3):568–586.
- 503
504 Gu, T., Taylor, J. M., and Mukherjee, B. (2023). A meta-inference framework to integrate multiple
505 external models into a current study. *Biostatistics*, 24(2):406–424.
- 506
507 Han, P. and Lawless, J. F. (2019). Empirical likelihood estimation using auxiliary summary infor-
508 mation with different covariate distributions. *Statistica Sinica*, 29(3):1321–1342.
- 509
510 Han, P., Li, H., Park, S. K., Mukherjee, B., and Taylor, J. M. (2024). Improving prediction of linear
511 regression models by integrating external information from heterogeneous populations: James-
512 stein estimators. *Biometrics*, 80(3):ujae072.
- 513
514 Kisselburgh, L. and Beever, J. (2022). The ethics of privacy in research and design: Principles,
515 practices, and potential. In *Modern Socio-Technical Perspectives on Privacy*, pages 395–426.
516 Springer International Publishing Cham.
- 517
518 Kundu, P., Tang, R., and Chatterjee, N. (2019). Generalized meta-analysis for multiple regression
519 models across studies with disparate covariate information. *Biometrika*, 106(3):567–585.
- 520
521 Lapatas, V., Stefanidakis, M., Jimenez, R. C., Via, A., and Schneider, M. V. (2015). Data integration
522 in biological research: An overview. *Journal of Biological Research-Thessaloniki*, 22(1):1–16.
- 523
524 Newey, W. K. and McFadden, D. (1994). Large sample estimation and hypothesis testing. In
525 *Handbook of Econometrics*, pages 2111–2245. Elsevier.
- 526
527 Qin, J. (2000). Miscellanea. combining parametric and empirical likelihoods. *Biometrika*,
528 87(2):484–490.
- 529
530 Rothstein, M. A. (2010). Is deidentification sufficient to protect health privacy in research? *The
531 American Journal of Bioethics*, 10(9):3–11.
- 532
533 Russ, T. C., Woelbert, E., Davis, K. A., Hafferty, J. D., Ibrahim, Z., Inkster, B., John, A., Lee,
534 W., Maxwell, M., McIntosh, A. M., et al. (2019). How data science can advance mental health
535 research. *Nature human behaviour*, 3(1):24–32.
- 536
537 Sheng, Y., Qin, J., and Huang, C.-Y. (2024). Sequential data integration under dataset shift. *Tech-
538 nometrics*, 66(4):662–670.
- 539
540 Sheng, Y., Sun, Y., Huang, C.-Y., and Kim, M.-O. (2022). Synthesizing external aggregated infor-
541 mation in the presence of population heterogeneity: A penalized empirical likelihood approach.
542 *Biometrics*, 78(2):679–690.
- 543
544 Taylor, J. M., Choi, K., and Han, P. (2023). Data integration: Exploiting ratios of parameter estimates
545 from a reduced external model. *Biometrika*, 110(1):119–134.
- 546
547 Van der Vaart, A. W. (1998). *Asymptotic Statistics*. Cambridge University Press.
- 548
549 Zhai, Y. and Han, P. (2022). Data integration with oracle use of external information from heteroge-
550 neous populations. *Journal of Computational and Graphical Statistics*, 31(4):1001–1012.

APPENDIX A SUPPLEMENTARY FIGURES

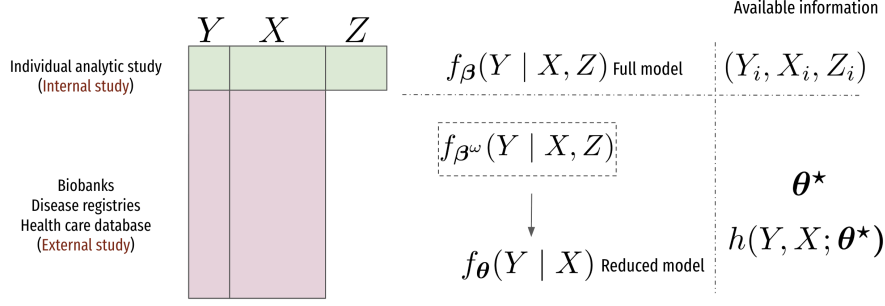


Figure A1: Problem Setup

APPENDIX B PROOF OF THEOREM 1 (VARIATIONAL DUAL FORMULATION)

Proof. For a given β , the inner problem in Eq. 1 over reweighting p is

$$\begin{aligned} \min_{p \geq 0 \in \mathbb{R}^N} \quad & \frac{1}{N} \sum_{i=1}^N \phi(Np_i) \\ \text{s.t.} \quad & \sum_{i=1}^N p_i = 1, \quad \sum_{i=1}^N p_i \mathbf{u}_i(\beta) = \mathbf{0}, \end{aligned}$$

where $N := N_{\text{in}}$. Introduce Lagrange multipliers $\gamma \in \mathbb{R}$ (normalization) and $\lambda \in \mathbb{R}^q$ (moments), and set $s_i := Np_i \geq 0$. The Lagrangian is

$$\mathcal{L}(s, \gamma, \lambda) = \frac{1}{N} \sum_{i=1}^N \phi(s_i) + \frac{1}{N} \sum_{i=1}^N s_i \{\gamma + \lambda^\top \mathbf{u}_i(\beta)\} - \gamma.$$

The dual function is the infimum over $s_i \geq 0$ (note separability in i):

$$g(\gamma, \lambda) = \sum_{i=1}^N \inf_{s_i \geq 0} \frac{1}{N} \{\phi(s_i) + s_i t_i\} - \gamma, \quad t_i := \gamma + \lambda^\top \mathbf{u}_i(\beta).$$

By the definition of the convex conjugate, $\inf_{s \geq 0} \{\phi(s) + st\} = -\phi^*(-t)$ (the domain restriction $s \geq 0$ is standard for ϕ -divergences and yields the same value for t in the effective domain). Therefore,

$$g(\gamma, \lambda) = -\frac{1}{N} \sum_{i=1}^N \phi^*(-t_i) - \gamma = -\frac{1}{N} \sum_{i=1}^N \phi^*\{-\gamma - \lambda^\top \mathbf{u}_i(\beta)\} - \gamma.$$

Maximizing g over (γ, λ) is equivalent by the sign change $(\gamma, \lambda) \mapsto (-\gamma, -\lambda)$ to

$$\max_{\gamma, \lambda} \gamma - \frac{1}{N} \sum_{i=1}^N \phi^*\{\gamma + \lambda^\top \mathbf{u}_i(\beta)\}.$$

By strong duality (Slater's condition holds whenever there exists a strictly feasible p with $p_i > 0$, $\sum p_i = 1$, and $\sum p_i \mathbf{u}_i(\beta) = \mathbf{0}$), the inner primal value equals this dual value. Since the negative log-likelihood term

$$-\frac{2}{N} \sum_{i=1}^N \log f_\beta(Y_i | \mathbf{X}_i, \mathbf{Z}_i)$$

does not depend on p , it passes unchanged through the dualization. Taking the outer minimization over β completes the proof. \square

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

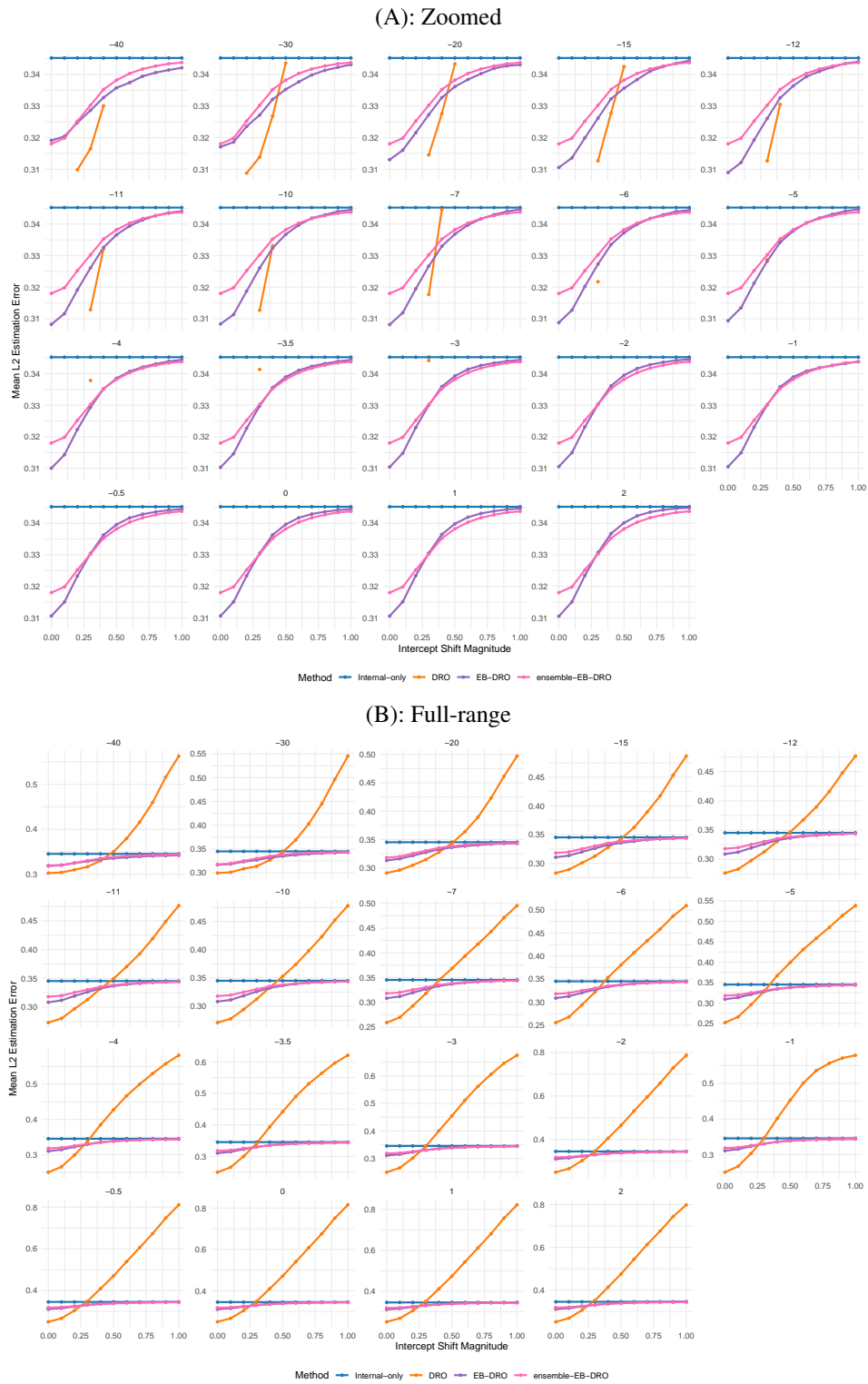


Figure A2: Theme A: Intercept shift magnitude

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

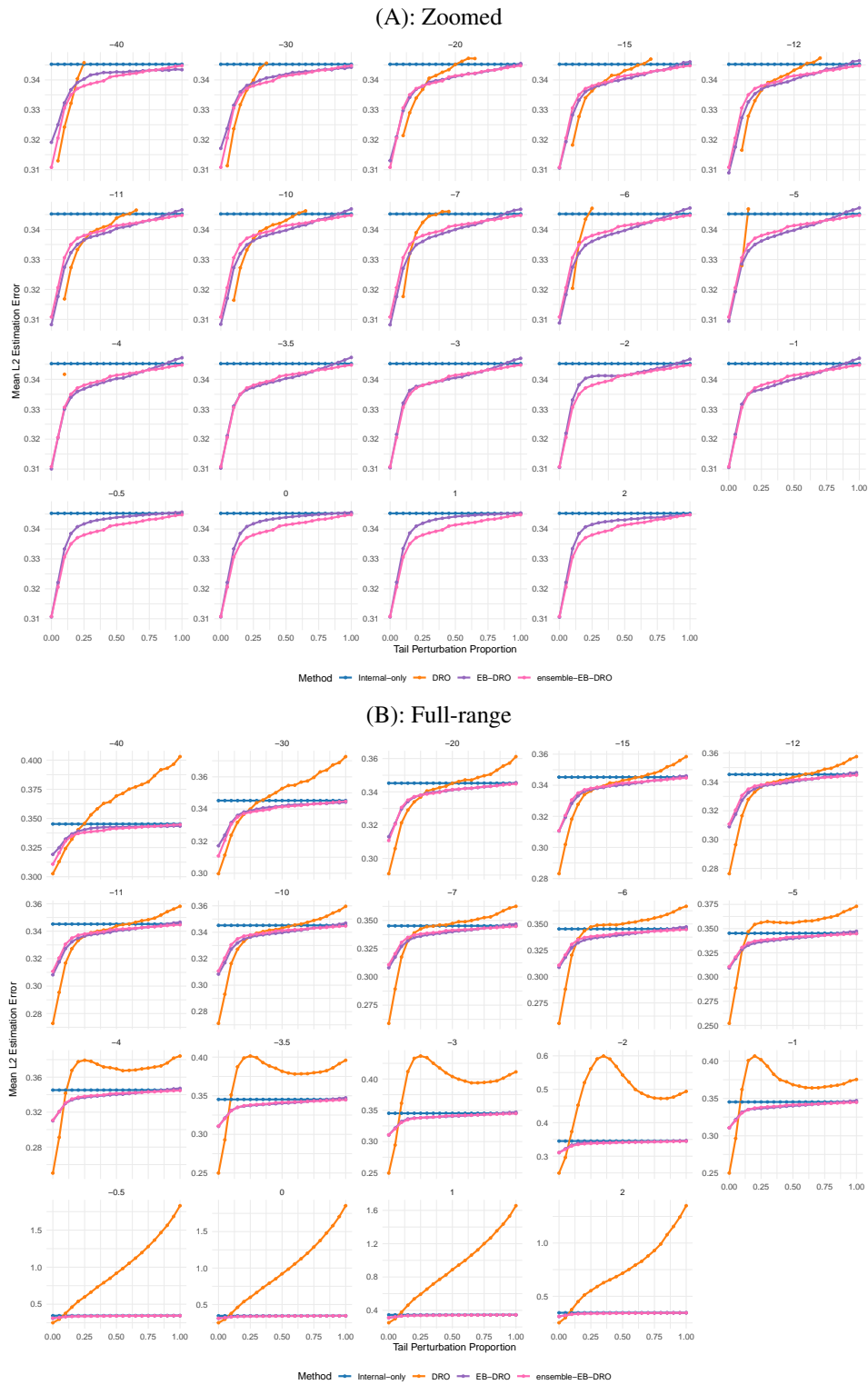


Figure A3: Theme B: Tail perturbation

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

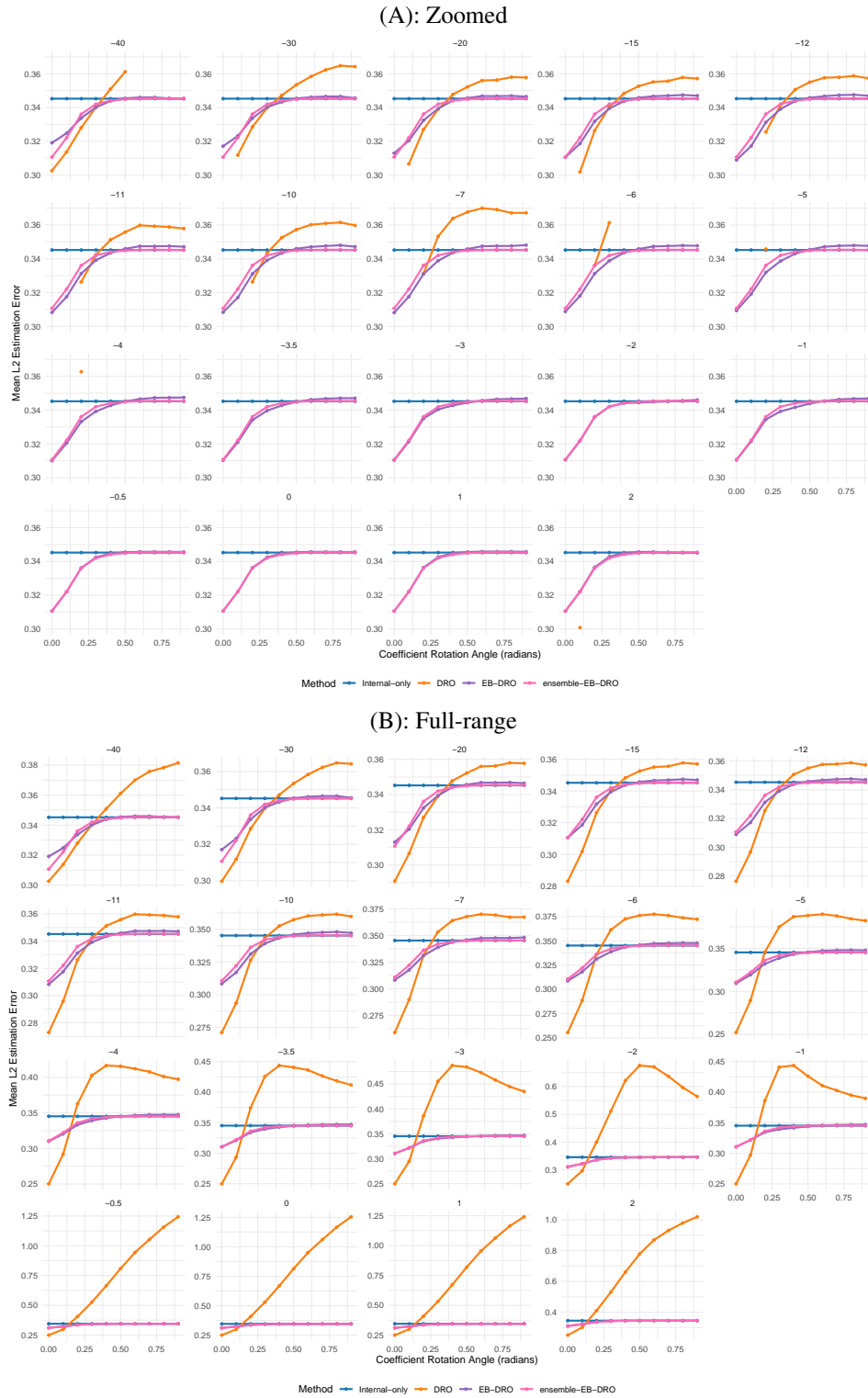


Figure A4: Theme C: Angle shift of the whole vector

APPENDIX C PROOF OF COROLLARY 1

Proof. By Theorem 1, for fixed $(\boldsymbol{\beta}, \boldsymbol{\lambda})$ the inner dual in $\gamma \in \mathbb{R}$ is

$$\max_{\gamma \in \mathbb{R}} \gamma - \frac{1}{N} \sum_{i=1}^N \phi_k^* \{\gamma + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta})\}, \quad N := N_{\text{in}}.$$

Set $t_i := 1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta})$ and assume $t_i > 0$ (the natural domain for CR). With the given ϕ_k^* , the objective equals (constants in ϕ_k^* cancel when summed and absorbed)

$$g(\gamma) = \gamma + \frac{2}{kN} \sum_{i=1}^N \left[-\frac{k+1}{2} \{\gamma + (t_i - 1)\} \right]^{\frac{k}{k+1}}.$$

Introduce $\alpha := -\frac{k+1}{2}\gamma$ ($\alpha \geq 0 \Leftrightarrow \gamma \leq 0$) so that $-\frac{k+1}{2}\{\gamma + (t_i - 1)\} = \alpha t_i$ and $\gamma = -\frac{2}{k+1}\alpha$. Then

$$g(\alpha) = -\frac{2}{k+1}\alpha + \frac{2}{kN} \sum_{i=1}^N (\alpha t_i)^{\frac{k}{k+1}} = -\frac{2}{k+1}\alpha + \frac{2}{k} \alpha^{\frac{k}{k+1}} \left(\frac{1}{N} \sum_{i=1}^N t_i^{\frac{k}{k+1}} \right).$$

Since $\frac{k}{k+1} \in (0, 1)$ for $k > -1$, the map $\alpha \mapsto g(\alpha)$ is strictly concave on $\alpha \geq 0$, so the maximizer is characterized by the first-order condition:

$$g'(\alpha) = -\frac{2}{k+1} + \frac{2}{k} \cdot \frac{k}{k+1} \alpha^{-\frac{1}{k+1}} \left(\frac{1}{N} \sum_{i=1}^N t_i^{\frac{k}{k+1}} \right) = 0.$$

Solving gives

$$\alpha^* = \left(\frac{1}{N} \sum_{i=1}^N t_i^{\frac{k}{k+1}} \right)^{k+1} \implies \gamma^* = -\frac{2}{k+1} \left(\frac{1}{N} \sum_{i=1}^N t_i^{\frac{k}{k+1}} \right)^{k+1}.$$

Evaluating $g(\alpha^*)$ yields

$$g(\alpha^*) = \frac{2}{k} \left(\frac{1}{N} \sum_{i=1}^N t_i^{\frac{k}{k+1}} \right)^{k+1} - \frac{2}{k+1} \left(\frac{1}{N} \sum_{i=1}^N t_i^{\frac{k}{k+1}} \right)^{k+1} = \frac{2}{k(k+1)} \left(\frac{1}{N} \sum_{i=1}^N t_i^{\frac{k}{k+1}} \right)^{k+1}.$$

Multiplying by N (recall the prefactor $1/N$ sits in front of the sum inside the dual) gives

$$H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{N}{k(k+1)} \left\{ \frac{1}{N} \sum_{i=1}^N (1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta}))^{\frac{k}{k+1}} \right\}^{k+1},$$

which depends only on $(\boldsymbol{\lambda}, \boldsymbol{\beta})$. Plugging back into the outer problem (with the likelihood term independent of γ) gives the stated min-max form. \square

APPENDIX D PROPERTIES OF THE INNER DRO PENALTY

D.1 SPECIAL CASES

Lemma A1 (Special cases of $k \rightarrow 0$ and $k \rightarrow -1$ limits). *Let*

$$H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{N_{\text{in}}}{k(k+1)} \left\{ \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} (1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta}))^{k/(k+1)} \right\}^{k+1},$$

with $t_i(\boldsymbol{\lambda}, \boldsymbol{\beta}) = 1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta}) > 0$. Then:

1. As $k \rightarrow 0$,

$$\lim_{k \rightarrow 0} \left(H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) - \frac{N_{\text{in}}}{k(k+1)} \right) = \sum_{i=1}^{N_{\text{in}}} \log(1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta})),$$

which coincides with the empirical likelihood (EL) dual objective.

810 2. As $k \rightarrow -1$,

$$811 \lim_{k \rightarrow -1} \left(H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) - \frac{N_{\text{in}}}{k(k+1)} \right) = -N_{\text{in}} \log \left(\frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \exp \left(\frac{1}{2} \boldsymbol{\theta}^\top \mathbf{u}_i(\boldsymbol{\beta}) \right) \right),$$

812 where $\boldsymbol{\theta}$ is a reparameterization of $\boldsymbol{\lambda}$. This corresponds to the exponential tilting (ET/KL)
813 dual form.

814 *Proof.* For the EL case $k \rightarrow 0$ (hence $\alpha = \frac{k}{k+1} \rightarrow 0$), expand

$$815 t_i^{k/(k+1)} = \exp \left(\frac{k}{k+1} \log t_i \right) = 1 + \frac{k}{k+1} \log t_i + O(k^2),$$

816 where $t_i = 1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta})$. Averaging gives

$$817 \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} t_i^{k/(k+1)} = 1 + \frac{k}{k+1} L(\boldsymbol{\lambda}, \boldsymbol{\beta}) + O(k^2), \quad L(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \log t_i.$$

818 Write $A_k := \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} t_i^{k/(k+1)}$, and $\delta_k := \frac{k}{k+1} L(\boldsymbol{\lambda}, \boldsymbol{\beta}) + O(k^2)$. Then

$$819 A_k^{k+1} = \exp \left((k+1) \log(1 + \delta_k) \right).$$

820 Since $\log(1 + \delta_k) = \delta_k + O(\delta_k^2)$ with $\delta_k = O(k)$,

$$821 (k+1) \log(1 + \delta_k) = kL(\boldsymbol{\lambda}, \boldsymbol{\beta}) + O(k^2).$$

822 Exponentiating gives

$$823 A_k^{k+1} = \exp \left(kL(\boldsymbol{\lambda}, \boldsymbol{\beta}) + O(k^2) \right) = 1 + kL(\boldsymbol{\lambda}, \boldsymbol{\beta}) + O(k^2).$$

824 Substituting into H_k gives

$$825 H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{N_{\text{in}}}{k(k+1)} \left(1 + kL(\boldsymbol{\lambda}, \boldsymbol{\beta}) + O(k^2) \right) = \frac{N_{\text{in}}}{k(k+1)} + \frac{N_{\text{in}}}{k+1} L(\boldsymbol{\lambda}, \boldsymbol{\beta}) + \frac{N_{\text{in}}}{k(k+1)} O(k^2)$$

826 where the last equality follows because $\frac{1}{k(k+1)} = O(1/k)$ as $k \rightarrow 0$ (indeed, $\frac{1}{k+1} = 1 + O(k)$),
827 hence $\frac{1}{k(k+1)} O(k^2) = O(k)$.

828 and thus

$$829 \lim_{k \rightarrow 0} \left(H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) - \frac{N_{\text{in}}}{k(k+1)} \right) = \sum_{i=1}^{N_{\text{in}}} \log \left(1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta}) \right).$$

830 For the ET case $k \rightarrow -1$, set $r = k+1 \rightarrow 0$ so that $\alpha = k/(k+1) = 1 - 1/r$. Reparameterize
831 $\boldsymbol{\lambda} = -\frac{r}{2} \boldsymbol{\vartheta}$ with $r = k+1 \rightarrow 0$. This scaling has two purposes: first, since $t_i(\boldsymbol{\lambda}, \boldsymbol{\beta}) = 1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta})$,
832 it guarantees $t_i(\boldsymbol{\lambda}, \boldsymbol{\beta}) = 1 + O(r) > 0$ for sufficiently small r , thus preserving the positivity
833 constraint; second, the factor $-\frac{1}{2}$ is chosen so that the expansion of $\alpha \log t_i$ with $\alpha = 1 - \frac{1}{r}$ yields
834 the exponential tilt form $\exp \left(\frac{1}{2} \boldsymbol{\vartheta}^\top \mathbf{u}_i(\boldsymbol{\beta}) \right)$ in the limit. A Taylor expansion shows

$$835 t_i^\alpha = \exp \left(\frac{1}{2} \boldsymbol{\vartheta}^\top \mathbf{u}_i(\boldsymbol{\beta}) \right) (1 + O(r)).$$

836 Hence

$$837 \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} t_i^\alpha = \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \exp \left(\frac{1}{2} \boldsymbol{\vartheta}^\top \mathbf{u}_i(\boldsymbol{\beta}) \right) + O(r).$$

838 Set

$$839 M_k := \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} t_i^\alpha = \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \exp \left(\frac{1}{2} \boldsymbol{\vartheta}^\top \mathbf{u}_i(\boldsymbol{\beta}) \right) + O(r) =: M_{-1} + O(r),$$

840 with $M_{-1} > 0$ bounded away from 0 in a neighborhood. Then

$$841 (M_k)^r = \exp \left(r \log M_k \right) = \exp \left(r \left[\log M_{-1} + O(r) \right] \right) = 1 + r \log M_{-1} + O(r^2),$$

864 i.e.

$$865 \left(\frac{1}{N_{\text{in}}} \sum t_i^\alpha \right)^r = 1 + r \log \left(\frac{1}{N_{\text{in}}} \sum \exp \left(\frac{1}{2} \boldsymbol{\vartheta}^\top \mathbf{u}_i(\boldsymbol{\beta}) \right) \right) + O(r^2).$$

866 Substituting into H_k gives

$$867 H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{N_{\text{in}}}{k(k+1)} \left[1 + r \log M_{-1} + O(r^2) \right] = \frac{N_{\text{in}}}{k(k+1)} + \frac{N_{\text{in}}}{k} \log M_{-1} + O(r),$$

868 where the last $O(r)$ term comes from $\frac{N_{\text{in}}}{k(k+1)} O(r^2) = \frac{N_{\text{in}}}{kr} O(r^2) = \frac{N_{\text{in}}}{k} O(r) = O(r)$ since $1/k =$
869 $O(1)$ as $k \rightarrow -1$. \square

870 D.2 CONVEXITY OF THE INNER OPTIMIZATION

871 **Lemma A2** (Inner dual equivalence and convexity). *For a fixed $\boldsymbol{\beta}$, we define*

$$872 H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{N_{\text{in}}}{k(k+1)} \left\{ \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} (1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta}))^{k/(k+1)} \right\}^{k+1},$$

873 with the conventions that $k = 0$ (empirical likelihood, EL) and $k = -1$ (exponential tilting, ET) are
874 taken as limits. Define the corresponding minimization form

$$875 M_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \begin{cases} -\sum_{i=1}^{N_{\text{in}}} \log(1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta})), & k = 0, \\ \sum_{i=1}^{N_{\text{in}}} \exp\left(\frac{1}{2} \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta})\right), & k = -1, \\ \sum_{i=1}^{N_{\text{in}}} (1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta}))^{k/(k+1)}, & k < 0, k \neq -1, \\ -\sum_{i=1}^{N_{\text{in}}} (1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta}))^{k/(k+1)}, & k > 0. \end{cases}$$

876 Then:

- 877 (i) **Monotone mapping and equivalence.** *There exists a strictly decreasing scalar map $h_k(\cdot)$*
878 *such that*

$$879 H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = h_k(M_k(\boldsymbol{\lambda}, \boldsymbol{\beta})),$$

880 where explicitly

$$881 h_k(m) = \begin{cases} -m, & k = 0, \\ -N_{\text{in}} \log(m/N_{\text{in}}) = -N_{\text{in}} \log m + \text{const}, & k = -1, \\ \frac{N_{\text{in}}^{-k}}{k(k+1)} m^{k+1}, & k < 0, k \neq -1, \\ \frac{N_{\text{in}}^{-k}}{k(k+1)} (-m)^{k+1}, & k > 0. \end{cases}$$

882 Hence,

$$883 \max_{\boldsymbol{\lambda}} H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) \iff \min_{\boldsymbol{\lambda}} M_k(\boldsymbol{\lambda}, \boldsymbol{\beta}).$$

- 884 (ii) **Convexity of the inner problem.** *For every $k \in \mathbb{R}$, the map $\boldsymbol{\lambda} \mapsto M_k(\boldsymbol{\lambda}, \boldsymbol{\beta})$ is convex on*
885 *its domain, and thus the inner program*

$$886 \min_{\boldsymbol{\lambda}} M_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) \quad (\text{equivalently, } \max_{\boldsymbol{\lambda}} H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}))$$

887 is a convex optimization problem in $\boldsymbol{\lambda}$.

888 **Proof.** (i) Substitute the definition of M_k into H_k . For $k \neq -1, 0$,

$$889 H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{N_{\text{in}}}{k(k+1)} \left(\frac{1}{N_{\text{in}}} (\pm M_k) \right)^{k+1} = \frac{N_{\text{in}}^{-k}}{k(k+1)} (\pm M_k)^{k+1},$$

890 where the \pm sign matches the definition of M_k above. For $k = 0$ and $k = -1$, the continuous limits
891 give

$$892 H_0(\boldsymbol{\lambda}, \boldsymbol{\beta}) = -M_0(\boldsymbol{\lambda}, \boldsymbol{\beta}), \quad H_{-1}(\boldsymbol{\lambda}, \boldsymbol{\beta}) = -N_{\text{in}} \log(M_{-1}(\boldsymbol{\lambda}, \boldsymbol{\beta})/N_{\text{in}}).$$

Thus in all cases $H_k = h_k(M_k)$, with derivatives

$$h'_k(m) = \begin{cases} -1, & k = 0, \\ -\frac{N_{\text{in}}}{m}, & k = -1, \\ \frac{N_{\text{in}}^{-k}}{k} m^k, & k < 0, k \neq -1, \\ \frac{N_{\text{in}}^{-k}}{k} (-m)^k, & k > 0, \end{cases}$$

which is strictly negative on the domain of m in each case ($M_k > 0$ if $k < 0$, $M_k < 0$ if $k > 0$, $M_0 \in \mathbb{R}$, $M_{-1} > 0$). Thus, each h_k is strictly decreasing, yielding the equivalence

$$\max_{\lambda} H_k(\lambda, \beta) \iff \min_{\lambda} M_k(\lambda, \beta).$$

(ii) For convexity, write $t_i(\lambda) = 1 + \lambda^\top \mathbf{u}_i(\beta)$, an affine function of λ . If $k < 0$, then $\alpha = \frac{k}{k+1} > 1$ or $\alpha < 0$, so $t \mapsto t^\alpha$ is convex on $t > 0$ and M_k is a sum of convex functions of affine maps. If $k > 0$, then $0 < \alpha < 1$, so t^α is concave, but the leading minus sign in M_k flips concavity to convexity. For $k = 0$, $-\log t$ is convex on $t > 0$; for $k = -1$, $\exp(t/2)$ is convex on \mathbb{R} . Summation preserves convexity, so M_k is convex in all cases. \square

APPENDIX E PROOF OF THEOREM 2

Assumption A1 (Regularity assumptions). Let $\{(Y_i, \mathbf{X}_i, \mathbf{Z}_i)\}_{i=1}^{N_{\text{in}}}$ be i.i.d. observations from the true parameter $\beta^* \in \mathbb{R}^d$. Define the log-likelihood score $\mathbf{s}_i(\beta) := \nabla_{\beta} \log f_{\beta}(Y_i | \mathbf{X}_i, \mathbf{Z}_i)$ and the auxiliary moment function $\mathbf{u}_i(\beta) \in \mathbb{R}^q$. Assume:

- (a) (Likelihood regularity) $\mathbb{E}\{\mathbf{s}_i(\beta^*)\} = \mathbf{0}$, the Fisher information $\mathbf{J} := \mathbb{E}\{-\nabla_{\beta\beta}^2 \log f_{\beta^*}(Y_i | \mathbf{X}_i, \mathbf{Z}_i)\}$ is finite and positive definite, and standard MLE regularity conditions (dominated differentiability, interchange of differentiation and expectation, finite variance) hold.
- (b) (Well-specified moments) $\mathbf{S} := \mathbb{E}\{\mathbf{u}_i(\beta^*)\mathbf{u}_i(\beta^*)^\top\}$ is finite and positive definite.
- (c) (Jacobian and rank) $\mathbf{u}_i(\beta)$ is continuously differentiable near β^* with $\mathbf{G} := \mathbb{E}\{\nabla_{\beta} \mathbf{u}_i(\beta^*)\} \in \mathbb{R}^{q \times d}$ full column rank.
- (d) (Positivity domain) There exists an open $\Lambda \subset \mathbb{R}^q$ containing $\mathbf{0}$ such that $t_i(\lambda, \beta) := 1 + \lambda^\top \mathbf{u}_i(\beta) > 0$ almost surely for all $(\lambda, \beta) \in \Lambda \times \mathcal{N}(\beta^*)$.
- (e) (Moments) Finite $(2+\delta)$ moments exist for some $\delta > 0$, so uniform LLNs and CLTs apply to the empirical means $\bar{\mathbf{s}}_N(\beta) := N_{\text{in}}^{-1} \sum_i \mathbf{s}_i(\beta)$, $\bar{\mathbf{u}}_N(\beta) := N_{\text{in}}^{-1} \sum_i \mathbf{u}_i(\beta)$, and their Jacobians.

Table A1 lists the notations for gradients, Hessians and cross Jacobian blocks that were used throughout the Appendix.

E.1 INNER MAXIMIZER

Write the population inner objective as a composition

$$H_k(\lambda, \beta) = h_k(M_k(\lambda, \beta)), \quad M_k(\lambda, \beta) = \mathbb{E}[m_k(1 + \lambda^\top \mathbf{u}(\beta))], \quad t = \lambda^\top \mathbf{u}(\beta)$$

where, for the Cressie-Read family (including the EL/ET limits),

$$m_k(t\lambda) = \begin{cases} -\log[1 + t_i(\lambda)], & k = 0 \\ \exp[t_i(\lambda)/2], & k = -1, \\ [1 + t_i(\lambda)]^\alpha, & k < 0, k \neq -1, \\ -[1 + t_i(\lambda)]^\alpha, & k > 0. \end{cases} \quad h_k(M) = \begin{cases} -M, & k = 0, \\ -\log M, & k = -1, \\ \frac{1}{k(k+1)} M^{k+1}, & k < 0, k \neq -1. \\ \frac{1}{k(k+1)} (-M)^{k+1}, & k > 0, \end{cases}$$

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Table A1: Notations for gradients, Hessians, and cross blocks. Aggregate versions are defined as sums over N_{in} observations. Per-observation versions are the per-sample contributions. Population limit denotes the expectation of the Per-observation.

Block	Aggregate	Per-observation	Population limit
Log-likelihood score	$\mathcal{S}_n := \sum_{i=1}^{N_{\text{in}}} \mathbf{s}_i(\boldsymbol{\beta}^*)$	$\mathbf{s}_i(\boldsymbol{\beta}^*)$	$\mathbb{E}\{\mathbf{s}_i(\boldsymbol{\beta}^*)\} = \mathbf{0}$
Log-likelihood Hessian	$\mathcal{J}_n := -\sum_{i=1}^{N_{\text{in}}} \nabla_{\boldsymbol{\beta}}^2 \ell_i(\boldsymbol{\beta}^*)$	$\mathbf{j}_i := -\nabla_{\boldsymbol{\beta}}^2 \ell_i(\boldsymbol{\beta}^*)$	$\mathbf{J} := \mathbb{E}\{\mathbf{j}_i\}$
Penalty gradient	$\mathcal{G}_n := \nabla_{\boldsymbol{\lambda}} H_k(\mathbf{0}, \boldsymbol{\beta}^*)$	$\mathbf{g}_i := \frac{1}{k+1} \mathbf{u}_i(\boldsymbol{\beta}^*)$	$\mathbb{E}[\mathbf{g}_i] = \mathbf{0}$
Penalty Hessian	$\mathcal{H}_n := \nabla_{\boldsymbol{\lambda}}^2 H_k(\mathbf{0}, \boldsymbol{\beta}^*)$	$\mathbf{h}_i := -\frac{1}{(k+1)^2} \mathbf{u}_i(\boldsymbol{\beta}^*) \mathbf{u}_i(\boldsymbol{\beta}^*)^\top$	$-\mathbb{E}[\mathbf{h}_i] = \frac{1}{(k+1)^2} \mathbf{S}$
Cross Jacobian	$\mathcal{C}_n := \nabla_{\boldsymbol{\lambda}}^2 H_k(\mathbf{0}, \boldsymbol{\beta}^*)$	$\mathbf{c}_i := \frac{1}{k+1} \nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}^*)$	$\mathbb{E}[\mathbf{c}_i] = \frac{1}{(k+1)} \mathbf{G}$
Penalty covariance	$\boldsymbol{\Omega}_1^{\text{agg}} := \mathbb{E}[\mathcal{G}_n \mathcal{G}_n^\top]$	$\boldsymbol{\Omega}_1 := \mathbb{E}[\mathbf{g}_i \mathbf{g}_i^\top]$	$\boldsymbol{\Omega}_1 = \frac{1}{(k+1)^2} \mathbf{S}$
Penalty curvature	$\boldsymbol{\Omega}_2^{\text{agg}} := -\mathbb{E}[\mathcal{H}_n]$	$\boldsymbol{\Omega}_2 := -\mathbb{E}[\mathbf{h}_i]$	$\boldsymbol{\Omega}_2 = \frac{1}{(k+1)^2} \mathbf{S}$

Note. Population limits are stated under correct model specification, i.e. $\mathbb{E}[\mathbf{u}_i(\boldsymbol{\beta}^*)] = \mathbf{0}$. Definitions. $\ell_i(\boldsymbol{\beta}) := \log f_{\boldsymbol{\beta}}(Y_i | \mathbf{X}_i, \mathbf{Z}_i)$, $\mathbf{S} := \mathbb{E}[\mathbf{u}_i(\boldsymbol{\beta}^*) \mathbf{u}_i(\boldsymbol{\beta}^*)^\top]$, $\mathbf{G} := \mathbb{E}[\nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}^*)]$. Here $H_k(\boldsymbol{\lambda}, \boldsymbol{\beta})$ denotes the inner objective for Cressie-Read index k , and ∇ denotes gradients/Hessians evaluated at $(\boldsymbol{\lambda}, \boldsymbol{\beta}) = (\mathbf{0}, \boldsymbol{\beta}^*)$.

Domain of the mapping h_k is : $M > 0$ for $k < 0$ and $M < 0$ for $k > 0$. By the chain rule,

$$\nabla_{\lambda} H_k(\lambda, \beta) = h'_k(M_k(\lambda, \beta)) \cdot \mathbb{E} \left[m'_k(1 + \lambda^{\top} \mathbf{u}(\beta)) \mathbf{u}(\beta) \right].$$

For every CR member (including EL/ET), $m'_k(t)$ and $h'_k(M)$ are finite. Hence, at $(\lambda, \beta) = (\mathbf{0}, \beta^*)$,

$$\nabla_{\lambda} H_k(\mathbf{0}, \beta^*) = h'_k(M_k(\mathbf{0}, \beta^*)) \mathbb{E} [m'_k(1) \mathbf{u}(\beta^*)] = h'_k(\cdot) m'_k(0) \mathbb{E} [\mathbf{u}(\beta^*)] = \mathbf{0},$$

by the well-specified moment assumption $\mathbb{E}[\mathbf{u}(\beta^*)] = \mathbf{0}$. Equivalently, the population M -form first-order condition is

$$\nabla_{\lambda} M_k(\lambda, \beta) = \mathbb{E} \left[m'_k(1 + \lambda^{\top} \mathbf{u}(\beta)) \mathbf{u}(\beta) \right] = \mathbf{0},$$

which is satisfied at $(\lambda, \beta) = (\mathbf{0}, \beta^*)$. Since $M_k(\cdot, \beta)$ is convex in λ (and strictly convex under standard conditions), $\lambda^*(\beta^*) = \mathbf{0}$ is the unique population minimizer of M_k , and thus the unique maximizer of H_k .

E.2 LEMMAS FOR THE GRADIENT, HESSIAN, JACOBIAN OF H UNDER THE POPULATION OPTIMUM

Lemma A3 (Gradient and Hessian of $H_k(\lambda, \beta)$ w.r.t. λ for $k \neq -1$). *Let $k \in \mathbb{R} \setminus \{-1\}$, and $\alpha = \frac{k}{k+1}$ and $t_i = 1 + \lambda^{\top} \mathbf{u}_i(\beta^*)$. At the population optimum $(\beta^*, \lambda^* = \mathbf{0})$, define the aggregate gradient and Hessian*

$$\mathcal{G}_n := \nabla_{\lambda} H_k(\mathbf{0}, \beta^*), \quad \mathcal{H}_n := \nabla_{\lambda \lambda}^2 H_k(\mathbf{0}, \beta^*).$$

Then

$$\mathcal{G}_n = \frac{N_{\text{in}}}{k+1} \bar{\mathbf{u}}, \quad \mathcal{H}_n = N_{\text{in}} \left[\alpha^2 \bar{\mathbf{u}} \bar{\mathbf{u}}^{\top} - \frac{1}{(k+1)^2} \bar{\mathbf{S}} \right],$$

where $\bar{\mathbf{u}} := \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \mathbf{u}_i(\beta^*)$ and $\bar{\mathbf{S}} := \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \mathbf{u}_i(\beta^*) \mathbf{u}_i(\beta^*)^{\top}$.

Per-sample scaling. Define the per-sample gradient and Hessian contributions

$$\mathbf{g}_i := \frac{1}{k+1} \mathbf{u}_i(\beta^*), \quad \mathbf{h}_i := -\frac{1}{(k+1)^2} \mathbf{u}_i(\beta^*) \mathbf{u}_i(\beta^*)^{\top},$$

so that $\mathcal{G}_n = \sum_{i=1}^{N_{\text{in}}} \mathbf{g}_i$, $\mathcal{H}_n = \sum_{i=1}^{N_{\text{in}}} \mathbf{h}_i$.

Then under correct specification, $\mathbb{E}[\mathbf{u}_i(\beta^*)] = \mathbf{0}$,

$$\frac{1}{N_{\text{in}}} \mathcal{G}_n \xrightarrow{p} \mathbf{0}, \quad -\frac{1}{N_{\text{in}}} \mathcal{H}_n \xrightarrow{p} \frac{1}{(k+1)^2} \mathbf{S}, \quad \mathbf{S} := \mathbb{E}[\mathbf{u}_i(\beta^*) \mathbf{u}_i(\beta^*)^{\top}].$$

Penalty covariance and curvature. At the per-sample level define

$$\Omega_1 := \mathbb{E}[\mathbf{g}_i \mathbf{g}_i^{\top}], \quad \Omega_2 := -\mathbb{E}[\mathbf{h}_i].$$

Then

$$\Omega_1 = \frac{1}{(k+1)^2} \mathbf{S}, \quad \Omega_2 = \frac{1}{(k+1)^2} \mathbf{S} + O\left(\frac{1}{N_{\text{in}}}\right).$$

Special case. In the limit $k \rightarrow 0$ (empirical likelihood), these reduce to

$$EL(k \rightarrow 0) : \quad \mathbf{g}_i = \mathbf{u}_i(\beta^*), \quad -\mathbf{h}_i = \mathbf{u}_i(\beta^*) \mathbf{u}_i(\beta^*)^{\top}.$$

Proof of Lemma A3. Let β^* be the true parameter and define, for $\alpha = \frac{k}{k+1}$,

$$t_i(\lambda) = 1 + \lambda^{\top} \mathbf{u}_i(\beta^*), \quad m_i(\lambda) = t_i(\lambda)^{\alpha}, \quad A(\lambda) = \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} m_i(\lambda),$$

$$H_k(\lambda) = \frac{N_{\text{in}}}{k(k+1)} (A(\lambda))^{k+1}.$$

All derivatives are w.r.t. $\boldsymbol{\lambda}$. Since t_i is affine,

$$\nabla_{\boldsymbol{\lambda}} m_i(\boldsymbol{\lambda}) = \alpha t_i(\boldsymbol{\lambda})^{\alpha-1} \mathbf{u}_i, \quad \nabla_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^2 m_i(\boldsymbol{\lambda}) = \alpha(\alpha-1) t_i(\boldsymbol{\lambda})^{\alpha-2} \mathbf{u}_i \mathbf{u}_i^\top,$$

hence

$$\nabla_{\boldsymbol{\lambda}} A(\boldsymbol{\lambda}) = \frac{1}{N_{\text{in}}} \sum_i \nabla_{\boldsymbol{\lambda}} m_i(\boldsymbol{\lambda}), \quad \nabla_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^2 A(\boldsymbol{\lambda}) = \frac{1}{N_{\text{in}}} \sum_i \nabla_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^2 m_i(\boldsymbol{\lambda}).$$

By the chain rule,

$$\begin{aligned} \nabla_{\boldsymbol{\lambda}} H_k(\boldsymbol{\lambda}) &= \frac{N_{\text{in}}}{k} A(\boldsymbol{\lambda})^k \nabla_{\boldsymbol{\lambda}} A(\boldsymbol{\lambda}), \\ \nabla_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^2 H_k(\boldsymbol{\lambda}) &= \frac{N_{\text{in}}}{k} [k A(\boldsymbol{\lambda})^{k-1} \nabla_{\boldsymbol{\lambda}} A(\boldsymbol{\lambda}) \nabla_{\boldsymbol{\lambda}} A(\boldsymbol{\lambda})^\top + A(\boldsymbol{\lambda})^k \nabla_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^2 A(\boldsymbol{\lambda})]. \end{aligned}$$

Evaluate at $\boldsymbol{\lambda} = \mathbf{0}$: since $t_i(\mathbf{0}) = 1$, we have $A(\mathbf{0}) = 1$ and

$$\nabla_{\boldsymbol{\lambda}} A(\mathbf{0}) = \alpha \bar{\mathbf{u}}, \quad \nabla_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^2 A(\mathbf{0}) = \alpha(\alpha-1) \bar{\mathbf{S}},$$

with $\bar{\mathbf{u}} := N_{\text{in}}^{-1} \sum_i \mathbf{u}_i(\boldsymbol{\beta}^*)$ and $\bar{\mathbf{S}} := N_{\text{in}}^{-1} \sum_i \mathbf{u}_i(\boldsymbol{\beta}^*) \mathbf{u}_i(\boldsymbol{\beta}^*)^\top$. Define the *aggregate* penalty gradient and Hessian

$$\mathcal{G}_n := \nabla_{\boldsymbol{\lambda}} H_k(\mathbf{0}, \boldsymbol{\beta}^*), \quad \mathcal{H}_n := \nabla_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^2 H_k(\mathbf{0}, \boldsymbol{\beta}^*).$$

Then

$$\begin{aligned} \mathcal{G}_n &= \frac{N_{\text{in}}}{k} \alpha \bar{\mathbf{u}} = \frac{N_{\text{in}}}{k+1} \bar{\mathbf{u}}, \\ \mathcal{H}_n &= \frac{N_{\text{in}}}{k} [k \alpha^2 \bar{\mathbf{u}} \bar{\mathbf{u}}^\top + \alpha(\alpha-1) \bar{\mathbf{S}}] = N_{\text{in}} \left[\alpha^2 \bar{\mathbf{u}} \bar{\mathbf{u}}^\top - \frac{1}{(k+1)^2} \bar{\mathbf{S}} \right]. \end{aligned}$$

Under $\mathbb{E}[\mathbf{u}_i(\boldsymbol{\beta}^*)] = \mathbf{0}$,

$$\mathbb{E}[\bar{\mathbf{u}} \bar{\mathbf{u}}^\top] = \frac{1}{N_{\text{in}}} \mathbf{S}, \quad \mathbb{E}[\bar{\mathbf{S}}] = \mathbf{S}, \quad \mathbf{S} := \mathbb{E}[\mathbf{u}_i(\boldsymbol{\beta}^*) \mathbf{u}_i(\boldsymbol{\beta}^*)^\top],$$

so

$$\mathbb{E}[\mathcal{H}_n] = \frac{N_{\text{in}}}{k} \left[k \alpha^2 \cdot \frac{1}{N_{\text{in}}} \mathbf{S} + \alpha(\alpha-1) \mathbf{S} \right] = \alpha^2 \mathbf{S} - \frac{N_{\text{in}}}{(k+1)^2} \mathbf{S} = -\frac{N_{\text{in}}}{(k+1)^2} \mathbf{S} + O(1).$$

Now introduce the average per-sample versions

$$\bar{\mathbf{g}}_n := \frac{1}{N_{\text{in}}} \mathcal{G}_n = \frac{1}{k+1} \bar{\mathbf{u}}, \quad \bar{\mathbf{h}}_n := \frac{1}{N_{\text{in}}} \mathcal{H}_n = \alpha^2 \bar{\mathbf{u}} \bar{\mathbf{u}}^\top - \frac{1}{(k+1)^2} \bar{\mathbf{S}}.$$

Then $\bar{\mathbf{g}}_n \xrightarrow{p} \mathbf{0}$ and $-\bar{\mathbf{h}}_n \xrightarrow{p} \frac{1}{(k+1)^2} \mathbf{S}$.

Finally, the per-observation penalty gradient and Hessian are

$$\mathbf{g}_i := \frac{1}{k+1} \mathbf{u}_i(\boldsymbol{\beta}^*), \quad \mathbf{h}_i := -\frac{1}{(k+1)^2} \mathbf{u}_i(\boldsymbol{\beta}^*) \mathbf{u}_i(\boldsymbol{\beta}^*)^\top.$$

Under correct specification, their population limits define the penalty covariance and curvature:

$$\boldsymbol{\Omega}_1 := \mathbb{E}[\mathbf{g}_i \mathbf{g}_i^\top] = \frac{1}{(k+1)^2} \mathbf{S}, \quad \boldsymbol{\Omega}_2 := -\mathbb{E}[\mathbf{h}_i] = \frac{1}{(k+1)^2} \mathbf{S}.$$

□

Lemma A4 (Gradient and Hessian of $H_k(\boldsymbol{\lambda}, \boldsymbol{\beta})$, w.r.t. $\boldsymbol{\beta}$). *Let $k \in \mathbb{R} \setminus \{-1\}$, set $\alpha = \frac{k}{k+1}$, and define*

$$\begin{aligned} t_i(\boldsymbol{\lambda}, \boldsymbol{\beta}) &= 1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta}), \quad A(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} t_i(\boldsymbol{\lambda}, \boldsymbol{\beta})^\alpha, \\ H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) &= \frac{N_{\text{in}}}{k(k+1)} (A(\boldsymbol{\lambda}, \boldsymbol{\beta}))^{k+1}. \end{aligned}$$

1134 Assume $\mathbf{u}_i(\boldsymbol{\beta}) \in \mathbb{R}^q$ is differentiable in $\boldsymbol{\beta} \in \mathbb{R}^d$. At the population optimum ($\boldsymbol{\beta}^*, \boldsymbol{\lambda}^* = \mathbf{0}$):
 1135
 1136 $\boldsymbol{\beta}$ -block vanishes at $\boldsymbol{\lambda} = \mathbf{0}$.

$$1137 \quad \nabla_{\boldsymbol{\beta}} H_k(\mathbf{0}, \boldsymbol{\beta}^*) = \mathbf{0}, \quad \nabla_{\boldsymbol{\beta}\boldsymbol{\beta}}^2 H_k(\mathbf{0}, \boldsymbol{\beta}^*) = \mathbf{0}.$$

1139 **Cross blocks.** Let $\bar{\mathbf{J}}_u := \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}^*) \in \mathbb{R}^{q \times d}$. Then

$$1141 \quad \nabla_{\boldsymbol{\lambda}\boldsymbol{\beta}}^2 H_k(\mathbf{0}, \boldsymbol{\beta}^*) = \frac{N_{\text{in}}}{k+1} \bar{\mathbf{J}}_u, \quad \nabla_{\boldsymbol{\beta}\boldsymbol{\lambda}}^2 H_k(\mathbf{0}, \boldsymbol{\beta}^*) = \left(\nabla_{\boldsymbol{\lambda}\boldsymbol{\beta}}^2 H_k(\mathbf{0}, \boldsymbol{\beta}^*) \right)^\top.$$

1144 **Limits of the per-observation cross jacobian** Define the aggregate cross block $\mathcal{C}_n :=$
 1145 $\nabla_{\boldsymbol{\lambda}\boldsymbol{\beta}}^2 H_k(\mathbf{0}, \boldsymbol{\beta}^*)$ and its per-sample version $\mathbf{c}_n := \frac{1}{N_{\text{in}}} \mathcal{C}_n$. Then

$$1147 \quad \mathbf{c}_n = \frac{1}{k+1} \bar{\mathbf{J}}_u \xrightarrow{p} \frac{1}{k+1} \mathbf{G}, \quad \mathbf{G} := \mathbb{E}[\nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}^*)] \in \mathbb{R}^{q \times d}.$$

1149 Equivalently, at the per-observation level,

$$1151 \quad \mathbf{c}_i = \frac{1}{k+1} \nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}^*), \quad \mathbb{E}[\mathbf{c}_i] = \frac{1}{k+1} \mathbf{G}.$$

1153 All block dimensions are consistent: $\nabla_{\boldsymbol{\lambda}\boldsymbol{\beta}}^2 H_k \in \mathbb{R}^{q \times d}$ and $\nabla_{\boldsymbol{\beta}\boldsymbol{\lambda}}^2 H_k \in \mathbb{R}^{d \times q}$.

1155 *Proof of Lemma A4.* Write $A(\boldsymbol{\lambda}, \boldsymbol{\beta}) = N_{\text{in}}^{-1} \sum_i t_i(\boldsymbol{\lambda}, \boldsymbol{\beta})^\alpha$ with $t_i = 1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta})$ and $\alpha = \frac{k}{k+1}$,
 1156 where $\boldsymbol{\lambda} \in \mathbb{R}^q$, $\boldsymbol{\beta} \in \mathbb{R}^d$, and $\mathbf{u}_i(\boldsymbol{\beta}) \in \mathbb{R}^q$. Then

$$1158 \quad H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{N_{\text{in}}}{k(k+1)} A(\boldsymbol{\lambda}, \boldsymbol{\beta})^{k+1}.$$

1160 Differentiate w.r.t. $\boldsymbol{\beta}$:

$$1162 \quad \nabla_{\boldsymbol{\beta}} H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{N_{\text{in}}}{k} A(\boldsymbol{\lambda}, \boldsymbol{\beta})^k \nabla_{\boldsymbol{\beta}} A(\boldsymbol{\lambda}, \boldsymbol{\beta}).$$

1163 For each i , define the Jacobian and Hessian of \mathbf{u}_i with respect to $\boldsymbol{\beta}$:

$$1165 \quad \mathbf{J}_i(\boldsymbol{\beta}) := \nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}) \in \mathbb{R}^{q \times d}, \quad \mathbf{H}_{ij}(\boldsymbol{\beta}) := \nabla_{\boldsymbol{\beta}\boldsymbol{\beta}}^2 u_{ij}(\boldsymbol{\beta}) \in \mathbb{R}^{d \times d},$$

1166 for $j = 1, \dots, q$.

1168 Then

$$1169 \quad \nabla_{\boldsymbol{\beta}} t_i(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \mathbf{J}_i(\boldsymbol{\beta})^\top \boldsymbol{\lambda} \in \mathbb{R}^d, \quad \nabla_{\boldsymbol{\beta}\boldsymbol{\beta}}^2 t_i(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \sum_{j=1}^q \lambda_j \mathbf{H}_{ij}(\boldsymbol{\beta}) \in \mathbb{R}^{d \times d}.$$

1172 By the chain rule with $\phi(x) = x^\alpha$, $\phi'(x) = \alpha x^{\alpha-1}$, $\phi''(x) = \alpha(\alpha-1)x^{\alpha-2}$,

$$1174 \quad \nabla_{\boldsymbol{\beta}} A(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{\alpha}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} t_i(\boldsymbol{\lambda}, \boldsymbol{\beta})^{\alpha-1} \mathbf{J}_i(\boldsymbol{\beta})^\top \boldsymbol{\lambda} \in \mathbb{R}^d, \quad (5)$$

$$1177 \quad \nabla_{\boldsymbol{\beta}\boldsymbol{\beta}}^2 A(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \left\{ \phi''(t_i) \nabla_{\boldsymbol{\beta}} t_i \nabla_{\boldsymbol{\beta}} t_i^\top + \phi'(t_i) \nabla_{\boldsymbol{\beta}\boldsymbol{\beta}}^2 t_i \right\}$$

$$1181 \quad = \frac{\alpha}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \left\{ (\alpha-1) t_i(\boldsymbol{\lambda}, \boldsymbol{\beta})^{\alpha-2} (\mathbf{J}_i(\boldsymbol{\beta})^\top \boldsymbol{\lambda}) (\mathbf{J}_i(\boldsymbol{\beta})^\top \boldsymbol{\lambda})^\top + t_i(\boldsymbol{\lambda}, \boldsymbol{\beta})^{\alpha-1} \sum_{j=1}^q \lambda_j \mathbf{H}_{ij}(\boldsymbol{\beta}) \right\}.$$

1184 At $\boldsymbol{\lambda} = \mathbf{0}$, we have $t_i(\mathbf{0}, \boldsymbol{\beta}) \equiv 1$, so

$$1185 \quad \nabla_{\boldsymbol{\beta}} A(\mathbf{0}, \boldsymbol{\beta}) = \mathbf{0}, \quad \nabla_{\boldsymbol{\beta}\boldsymbol{\beta}}^2 A(\mathbf{0}, \boldsymbol{\beta}) = \mathbf{0}.$$

1187 Therefore,

$$\nabla_{\boldsymbol{\beta}} H_k(\mathbf{0}, \boldsymbol{\beta}) \equiv \mathbf{0}.$$

A second derivative gives

$$\nabla_{\beta\beta}^2 H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{N_{\text{in}}}{k} \left[k A^{k-1} (\nabla_{\beta} A) (\nabla_{\beta} A)^{\top} + A^k \nabla_{\beta\beta}^2 A \right],$$

so at $\boldsymbol{\lambda} = \mathbf{0}$ we also have $\nabla_{\beta\beta}^2 H_k(\mathbf{0}, \boldsymbol{\beta}) \equiv \mathbf{0}$.

Remark: The DRO objective $H_k(\boldsymbol{\lambda}, \boldsymbol{\beta})$ depends on $\boldsymbol{\beta}$ only through the moments $\mathbf{u}_i(\boldsymbol{\beta})$, which are multiplied by $\boldsymbol{\lambda}$. When $\boldsymbol{\lambda} = \mathbf{0}$, the moment conditions are effectively “switched off,” so that $A(\boldsymbol{\lambda}, \boldsymbol{\beta}) \equiv 1$ regardless of $\boldsymbol{\beta}$. This implies that $H_k(\mathbf{0}, \boldsymbol{\beta})$ is constant in $\boldsymbol{\beta}$, with no gradient or curvature in the $\boldsymbol{\beta}$ -block:

$$\nabla_{\beta} H_k(\mathbf{0}, \boldsymbol{\beta}) = \mathbf{0}, \quad \nabla_{\beta\beta}^2 H_k(\mathbf{0}, \boldsymbol{\beta}) = \mathbf{0}.$$

In other words, at the population optimum the $\boldsymbol{\beta}$ -block carries no direct information; all information about $\boldsymbol{\beta}$ is mediated through the cross block $\nabla_{\lambda\beta}^2 H_k$, which involves the score functions $\mathbf{u}_i(\boldsymbol{\beta})$.

Cross derivative $\nabla_{\lambda\beta}^2 A(\boldsymbol{\lambda}, \boldsymbol{\beta})$. Recall the per-observation notations used previously

$$\mathbf{J}_i(\boldsymbol{\beta}) := \nabla_{\beta} \mathbf{u}_i(\boldsymbol{\beta}) \in \mathbb{R}^{q \times d}, \quad \mathbf{H}_{ij}(\boldsymbol{\beta}) := \nabla_{\beta\beta}^2 u_{ij}(\boldsymbol{\beta}) \in \mathbb{R}^{d \times d}.$$

Then

$$\nabla_{\lambda} t_i(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \mathbf{u}_i(\boldsymbol{\beta}), \quad \nabla_{\beta} t_i(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \mathbf{J}_i(\boldsymbol{\beta})^{\top} \boldsymbol{\lambda}, \quad \nabla_{\lambda\beta}^2 t_i(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \mathbf{J}_i(\boldsymbol{\beta}).$$

With $\phi(x) = x^{\alpha}$ ($\phi'(x) = \alpha x^{\alpha-1}$, $\phi''(x) = \alpha(\alpha-1)x^{\alpha-2}$), and Eq. 5 the cross block of A is

$$\nabla_{\lambda\beta}^2 A(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \left\{ \phi''(t_i) \mathbf{u}_i(\boldsymbol{\beta}) (\mathbf{J}_i(\boldsymbol{\beta})^{\top} \boldsymbol{\lambda})^{\top} + \phi'(t_i) \mathbf{J}_i(\boldsymbol{\beta}) \right\}.$$

At $\boldsymbol{\lambda} = \mathbf{0}$, we have $t_i \equiv 1$ and hence

$$\nabla_{\lambda\beta}^2 A(\mathbf{0}, \boldsymbol{\beta}) = \frac{\alpha}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \mathbf{J}_i(\boldsymbol{\beta}).$$

Cross derivative $\nabla_{\lambda\beta}^2 H_k(\boldsymbol{\lambda}, \boldsymbol{\beta})$. Since

$$H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{N_{\text{in}}}{k(k+1)} (A(\boldsymbol{\lambda}, \boldsymbol{\beta}))^{k+1},$$

its gradient is

$$\nabla_{\lambda} H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{N_{\text{in}}}{k} A(\boldsymbol{\lambda}, \boldsymbol{\beta})^k \nabla_{\lambda} A(\boldsymbol{\lambda}, \boldsymbol{\beta}).$$

Its cross derivative (a $q \times d$ matrix) is

$$\nabla_{\lambda\beta}^2 H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{N_{\text{in}}}{k} \left[k A(\boldsymbol{\lambda}, \boldsymbol{\beta})^{k-1} (\nabla_{\lambda} A(\boldsymbol{\lambda}, \boldsymbol{\beta})) (\nabla_{\beta} A(\boldsymbol{\lambda}, \boldsymbol{\beta}))^{\top} + A(\boldsymbol{\lambda}, \boldsymbol{\beta})^k \nabla_{\lambda\beta}^2 A(\boldsymbol{\lambda}, \boldsymbol{\beta}) \right]. \quad (6)$$

At $\boldsymbol{\lambda} = \mathbf{0}$, for the Cressie-Read inner average

$$A(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} (1 + \boldsymbol{\lambda}^{\top} \mathbf{u}_i(\boldsymbol{\beta}))^{\alpha}, \quad \alpha = \frac{k}{k+1},$$

we have

$$A(\mathbf{0}, \boldsymbol{\beta}) = 1, \quad \nabla_{\beta} A(\mathbf{0}, \boldsymbol{\beta}) = \mathbf{0}, \quad \nabla_{\lambda\beta}^2 A(\mathbf{0}, \boldsymbol{\beta}) = \frac{\alpha}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \mathbf{J}_i(\boldsymbol{\beta}),$$

where $\mathbf{J}_i(\boldsymbol{\beta}) = \nabla_{\beta} \mathbf{u}_i(\boldsymbol{\beta}) \in \mathbb{R}^{q \times d}$.

Thus, the first term in Eq. 6 vanishes, and

$$\nabla_{\lambda\beta}^2 H_k(\mathbf{0}, \boldsymbol{\beta}) = \frac{N_{\text{in}}}{k} \nabla_{\lambda\beta}^2 A(\mathbf{0}, \boldsymbol{\beta}) = \frac{N_{\text{in}}}{k} \cdot \frac{\alpha}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \mathbf{J}_i(\boldsymbol{\beta}) = \frac{1}{k+1} \sum_{i=1}^{N_{\text{in}}} \mathbf{J}_i(\boldsymbol{\beta}) =: \mathbf{C}_n.$$

Per-sample scaling gives

$$\mathbf{c}_n := \frac{1}{N_{\text{in}}} \mathcal{C}_n = \frac{1}{k+1} \bar{\mathbf{J}}_u(\boldsymbol{\beta}), \quad \bar{\mathbf{J}}_u(\boldsymbol{\beta}) := \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \mathbf{J}_i(\boldsymbol{\beta}).$$

At $\boldsymbol{\beta} = \boldsymbol{\beta}^*$, a LLN gives

$$\mathbf{c}_n \xrightarrow{p} \frac{1}{k+1} \mathbf{G}, \quad \mathbf{G} := \mathbb{E}[\nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}^*)] \in \mathbb{R}^{q \times d}.$$

□

Lemma A5 (Exponential tilting (ET) blocks). *For the ET dual objective*

$$H_{-1}(\boldsymbol{\lambda}, \boldsymbol{\beta}) = -N_{\text{in}} \log \left\{ \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \exp \left(\frac{1}{2} \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta}) \right) \right\},$$

let $\mathbf{u}_i(\boldsymbol{\beta}) \in \mathbb{R}^q$ and $\mathbf{J}_i(\boldsymbol{\beta}) := \nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}) \in \mathbb{R}^{q \times d}$. Define

$$\bar{\mathbf{u}} := \frac{1}{N_{\text{in}}} \sum_i \mathbf{u}_i(\boldsymbol{\beta}), \quad \bar{\mathbf{S}} := \frac{1}{N_{\text{in}}} \sum_i \mathbf{u}_i(\boldsymbol{\beta}) \mathbf{u}_i(\boldsymbol{\beta})^\top, \quad \bar{\mathbf{J}}_u := \frac{1}{N_{\text{in}}} \sum_i \mathbf{J}_i(\boldsymbol{\beta}).$$

Then at $\boldsymbol{\lambda} = \mathbf{0}$,

$$\begin{aligned} \nabla_{\boldsymbol{\lambda}} H_{-1}(\mathbf{0}, \boldsymbol{\beta}) &= -\frac{N_{\text{in}}}{2} \bar{\mathbf{u}}, & \nabla_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^2 H_{-1}(\mathbf{0}, \boldsymbol{\beta}) &= -\frac{N_{\text{in}}}{4} (\bar{\mathbf{S}} - \bar{\mathbf{u}} \bar{\mathbf{u}}^\top), \\ \nabla_{\boldsymbol{\lambda}\boldsymbol{\beta}}^2 H_{-1}(\mathbf{0}, \boldsymbol{\beta}) &= -\frac{N_{\text{in}}}{2} \bar{\mathbf{J}}_u. \end{aligned}$$

Let

$$\mathbf{g}_n := \frac{1}{N_{\text{in}}} \nabla_{\boldsymbol{\lambda}} H_{-1}, \quad \mathbf{h}_n := \frac{1}{N_{\text{in}}} \nabla_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^2 H_{-1}, \quad \mathbf{c}_n := \frac{1}{N_{\text{in}}} \nabla_{\boldsymbol{\lambda}\boldsymbol{\beta}}^2 H_{-1}.$$

Then at $\boldsymbol{\lambda} = \mathbf{0}$,

$$\mathbf{g}_n = -\frac{1}{2} \bar{\mathbf{u}}, \quad \mathbf{h}_n = -\frac{1}{4} (\bar{\mathbf{S}} - \bar{\mathbf{u}} \bar{\mathbf{u}}^\top), \quad \mathbf{c}_n = -\frac{1}{2} \bar{\mathbf{J}}_u.$$

Under correct specification $\mathbb{E}[\mathbf{u}_i(\boldsymbol{\beta}^*)] = \mathbf{0}$,

$$\mathbf{g}_n \xrightarrow{p} \mathbf{0}, \quad -\mathbf{h}_n \xrightarrow{p} \frac{1}{4} \mathbf{S}, \quad \mathbf{c}_n \xrightarrow{p} -\frac{1}{2} \mathbf{G},$$

where $\mathbf{S} := \mathbb{E}[\mathbf{u}_i(\boldsymbol{\beta}^*) \mathbf{u}_i(\boldsymbol{\beta}^*)^\top]$ and $\mathbf{G} := \mathbb{E}[\nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}^*)]$.

Proof of Lemma A5. Let $\boldsymbol{\lambda} \in \mathbb{R}^q$ and $\boldsymbol{\beta} \in \mathbb{R}^d$. For each $i = 1, \dots, N_{\text{in}}$, let $\mathbf{u}_i(\boldsymbol{\beta}) \in \mathbb{R}^q$ and

$$w_i(\boldsymbol{\lambda}, \boldsymbol{\beta}) := \exp \left(\frac{1}{2} \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta}) \right).$$

Define the sample mean

$$m(\boldsymbol{\lambda}, \boldsymbol{\beta}) := \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} w_i(\boldsymbol{\lambda}, \boldsymbol{\beta}),$$

and the ET dual objective

$$H_{-1}(\boldsymbol{\lambda}, \boldsymbol{\beta}) := -N_{\text{in}} \log m(\boldsymbol{\lambda}, \boldsymbol{\beta}).$$

Let $D(\boldsymbol{\lambda}, \boldsymbol{\beta}) := \sum_{j=1}^{N_{\text{in}}} w_j(\boldsymbol{\lambda}, \boldsymbol{\beta}) = N_{\text{in}} m(\boldsymbol{\lambda}, \boldsymbol{\beta})$ and define softmax weights

$$p_i(\boldsymbol{\lambda}, \boldsymbol{\beta}) := \frac{w_i(\boldsymbol{\lambda}, \boldsymbol{\beta})}{D(\boldsymbol{\lambda}, \boldsymbol{\beta})}.$$

For any array a_i , write the p -weighted average $\mathbb{E}_p[a_i] := \sum_{i=1}^{N_{\text{in}}} p_i a_i$. Set

$$\boldsymbol{\mu}_u := \mathbb{E}[\mathbf{u}_i(\boldsymbol{\beta})] \in \mathbb{R}^q, \quad \boldsymbol{\Sigma}_u := \mathbb{E}[\mathbf{u}_i \mathbf{u}_i^\top] - \boldsymbol{\mu}_u \boldsymbol{\mu}_u^\top \in \mathbb{R}^{q \times q}, \quad \mathbf{J}_i(\boldsymbol{\beta}) := \nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}) \in \mathbb{R}^{q \times d}.$$

(i) **Gradient w.r.t. $\boldsymbol{\lambda}$.** By the chain rule,

$$\nabla_{\boldsymbol{\lambda}} H_{-1}(\boldsymbol{\lambda}, \boldsymbol{\beta}) = -N_{\text{in}} \frac{\nabla_{\boldsymbol{\lambda}} m}{m}.$$

1296 Because

$$1297 \nabla_{\lambda} m = \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \nabla_{\lambda} w_i = \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \left(\frac{1}{2} \mathbf{u}_i \right) w_i = \frac{D}{N_{\text{in}}} \cdot \frac{1}{2} \mathbb{E}_p[\mathbf{u}_i] = m \cdot \frac{1}{2} \boldsymbol{\mu}_u,$$

1300 we get

$$1301 \nabla_{\lambda} H_{-1}(\boldsymbol{\lambda}, \boldsymbol{\beta}) = -\frac{N_{\text{in}}}{2} \boldsymbol{\mu}_u.$$

1304 **(ii) Hessian w.r.t. $\boldsymbol{\lambda}$.** Differentiate the preceding display:

$$1305 \nabla_{\lambda\lambda}^2 H_{-1}(\boldsymbol{\lambda}, \boldsymbol{\beta}) = -\frac{N_{\text{in}}}{2} \nabla_{\lambda} \boldsymbol{\mu}_u.$$

1308 Since \mathbf{u}_i does not depend on $\boldsymbol{\lambda}$,

$$1309 \nabla_{\lambda} \boldsymbol{\mu}_u = \sum_i (\nabla_{\lambda} p_i) \mathbf{u}_i^{\top} = \sum_i p_i (\nabla_{\lambda} \log p_i) \mathbf{u}_i^{\top}.$$

1311 But $\log p_i = \frac{1}{2} \boldsymbol{\lambda}^{\top} \mathbf{u}_i - \log D$, hence

$$1312 \nabla_{\lambda} \log p_i = \frac{1}{2} \mathbf{u}_i - \frac{1}{D} \sum_j w_j \frac{1}{2} \mathbf{u}_j = \frac{1}{2} (\mathbf{u}_i - \boldsymbol{\mu}_u).$$

1316 Therefore

$$1317 \nabla_{\lambda} \boldsymbol{\mu}_u = \frac{1}{2} \left(\mathbb{E}_p[\mathbf{u}_i \mathbf{u}_i^{\top}] - \boldsymbol{\mu}_u \boldsymbol{\mu}_u^{\top} \right) = \frac{1}{2} \boldsymbol{\Sigma}_u,$$

1318 and thus

$$1319 \nabla_{\lambda\lambda}^2 H_{-1}(\boldsymbol{\lambda}, \boldsymbol{\beta}) = -\frac{N_{\text{in}}}{4} \boldsymbol{\Sigma}_u.$$

1322 **(iii) Cross block $\nabla_{\lambda\beta}^2 H_{-1}$.** From (i),

$$1323 \nabla_{\lambda\beta}^2 H_{-1}(\boldsymbol{\lambda}, \boldsymbol{\beta}) = -\frac{N_{\text{in}}}{2} \nabla_{\beta} \boldsymbol{\mu}_u = -\frac{N_{\text{in}}}{2} \left(\sum_i (\nabla_{\beta} p_i) \mathbf{u}_i^{\top} + \sum_i p_i \mathbf{J}_i \right).$$

1326 Differentiate $\log p_i$ w.r.t. $\boldsymbol{\beta}$:

$$1327 \nabla_{\beta} \log p_i = \frac{1}{2} \mathbf{J}_i^{\top} \boldsymbol{\lambda} - \frac{1}{D} \sum_j w_j \frac{1}{2} \mathbf{J}_j^{\top} \boldsymbol{\lambda} = \frac{1}{2} \left(\mathbf{J}_i^{\top} \boldsymbol{\lambda} - \mathbb{E}_p[\mathbf{J}^{\top} \boldsymbol{\lambda}] \right).$$

1330 Thus $\nabla_{\beta} p_i = p_i \nabla_{\beta} \log p_i$, and

$$1331 \sum_i (\nabla_{\beta} p_i) \mathbf{u}_i^{\top} = \frac{1}{2} \left(\mathbb{E}_p[\mathbf{u}_i (\mathbf{J}_i^{\top} \boldsymbol{\lambda})^{\top}]^{\top} - \boldsymbol{\mu}_u \mathbb{E}_p[\mathbf{J}^{\top} \boldsymbol{\lambda}]^{\top} \right).$$

1334 Let $\boldsymbol{\mu}_{J\lambda} := \mathbb{E}_p[\mathbf{J}_i^{\top} \boldsymbol{\lambda}] \in \mathbb{R}^d$. Then

$$1335 \nabla_{\lambda\beta}^2 H_{-1}(\boldsymbol{\lambda}, \boldsymbol{\beta}) = -\frac{N_{\text{in}}}{2} \mathbb{E}_p[\mathbf{J}_i] - \frac{N_{\text{in}}}{4} \left(\mathbb{E}_p[\mathbf{u}_i (\mathbf{J}_i^{\top} \boldsymbol{\lambda})^{\top}]^{\top} - \boldsymbol{\mu}_u \boldsymbol{\mu}_{J\lambda}^{\top} \right).$$

1338 The last bracket is $\text{Cov}_p(\mathbf{u}_i, \mathbf{J}_i^{\top} \boldsymbol{\lambda})$.

1339 **Specialization at $\boldsymbol{\lambda} = \mathbf{0}$.** When $\boldsymbol{\lambda} = \mathbf{0}$, $p_i = 1/N_{\text{in}}$, hence the weighted averages become simple averages:

$$1342 \bar{\mathbf{u}} := \frac{1}{N_{\text{in}}} \sum_i \mathbf{u}_i, \quad \bar{\mathbf{S}} := \frac{1}{N_{\text{in}}} \sum_i \mathbf{u}_i \mathbf{u}_i^{\top}, \quad \bar{\mathbf{J}}_u := \frac{1}{N_{\text{in}}} \sum_i \mathbf{J}_i, \quad \text{and} \quad \boldsymbol{\mu}_{J\lambda} = \mathbf{0}.$$

1344 Substituting into the general formulas yields the aggregate forms

$$1345 \nabla_{\lambda} H_{-1}(\mathbf{0}, \boldsymbol{\beta}) = -\frac{N_{\text{in}}}{2} \bar{\mathbf{u}}, \quad \nabla_{\lambda\lambda}^2 H_{-1}(\mathbf{0}, \boldsymbol{\beta}) = -\frac{N_{\text{in}}}{4} (\bar{\mathbf{S}} - \bar{\mathbf{u}} \bar{\mathbf{u}}^{\top}),$$

$$1348 \nabla_{\lambda\beta}^2 H_{-1}(\mathbf{0}, \boldsymbol{\beta}) = -\frac{N_{\text{in}}}{2} \bar{\mathbf{J}}_u.$$

1349 These agree with the per-sample versions defined in the lemma after dividing by N_{in} . \square

E.3 PROOF OF THEOREM 2

Proof of Theorem 2. Write the score $s_i(\boldsymbol{\beta}) = \nabla_{\boldsymbol{\beta}} \log f_{\boldsymbol{\beta}}(Y_i | \mathbf{X}_i, \mathbf{Z}_i)$, the moment vector $\mathbf{u}_i(\boldsymbol{\beta}) \in \mathbb{R}^q$, $\mathbf{J} := \mathbb{E}[-\nabla_{\boldsymbol{\beta}}^2 \log f_{\boldsymbol{\beta}^*}]$, $\mathbf{S} := \mathbb{E}[\mathbf{u}_i(\boldsymbol{\beta}^*)\mathbf{u}_i(\boldsymbol{\beta}^*)^\top]$, $\mathbf{G} := \mathbb{E}[\nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}^*)]$. Let the Cressie-Read inner dual be

$$H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \frac{N_{\text{in}}}{k(k+1)} \left\{ \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} (1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta}))^{k/(k+1)} \right\}^{k+1}.$$

Define the (outer) criterion

$$Q_{N_{\text{in}}}(\boldsymbol{\beta}) = - \sum_{i=1}^{N_{\text{in}}} \log f_{\boldsymbol{\beta}}(Y_i | \mathbf{X}_i, \mathbf{Z}_i) + \max_{\boldsymbol{\lambda} \in \Lambda} H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}), \quad \widehat{\boldsymbol{\beta}}_k \in \arg \min_{\boldsymbol{\beta}} Q_{N_{\text{in}}}(\boldsymbol{\beta}),$$

and let $\widehat{\boldsymbol{\lambda}}_k$ be a maximizer of $H_k(\cdot, \widehat{\boldsymbol{\beta}}_k)$. Assumptions (moved to Appendix C) ensure: (i) $t_i(\boldsymbol{\lambda}, \boldsymbol{\beta}) := 1 + \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta}) > 0$ on $\Lambda \times \mathcal{N}(\boldsymbol{\beta}^*)$; (ii) \mathbf{u}_i is C^1 in $\boldsymbol{\beta}$; (iii) moments of order $2 + \delta$; (iv) inner program is concave in $\boldsymbol{\lambda}$ (Table 1 / Lemma A2); (v) identification at $(\boldsymbol{\beta}^*, \boldsymbol{\lambda}^* = \mathbf{0})$ under correct specification $\mathbb{E}[\mathbf{u}_i(\boldsymbol{\beta}^*)] = \mathbf{0}$.

KKT system and identification. Define the stacked estimating map

$$\boldsymbol{\Psi}_{N_{\text{in}}}(\boldsymbol{\vartheta}) := \begin{pmatrix} \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} s_i(\boldsymbol{\beta}) - \frac{1}{N_{\text{in}}} \nabla_{\boldsymbol{\beta}} H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) \\ \frac{1}{N_{\text{in}}} \nabla_{\boldsymbol{\lambda}} H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) \end{pmatrix}, \quad \boldsymbol{\vartheta} := \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\lambda} \end{pmatrix}, \quad \widehat{\boldsymbol{\vartheta}}_k := \begin{pmatrix} \widehat{\boldsymbol{\beta}}_k \\ \widehat{\boldsymbol{\lambda}}_k \end{pmatrix}.$$

By construction, $\boldsymbol{\Psi}_{N_{\text{in}}}(\widehat{\boldsymbol{\vartheta}}_k) = \mathbf{0}$ (first line is the outer first-order condition (FOC); second line is the inner FOC). At the population optimum $\boldsymbol{\vartheta}^* := (\boldsymbol{\beta}^{*\top}, \mathbf{0}^\top)^\top$, correct specification implies

$$\mathbb{E} \left[\frac{1}{N_{\text{in}}} \sum_i s_i(\boldsymbol{\beta}^*) \right] = \mathbf{0}, \quad \mathbb{E} \left[\frac{1}{N_{\text{in}}} \nabla_{\boldsymbol{\lambda}} H_k(\mathbf{0}, \boldsymbol{\beta}^*) \right] = \mathbf{0},$$

hence $\mathbb{E}[\boldsymbol{\Psi}_{N_{\text{in}}}(\boldsymbol{\vartheta}^*)] = \mathbf{0}$. Inner concavity (Lemma A2) yields uniqueness/continuity of $\boldsymbol{\lambda}^*(\boldsymbol{\beta})$ in a neighborhood of $\boldsymbol{\beta}^*$, so identification holds.

Consistency. By the uniform law of large numbers and continuity, $Q_{N_{\text{in}}}(\boldsymbol{\beta}) \rightarrow_p Q(\boldsymbol{\beta})$ uniformly on compact sets, where

$$Q(\boldsymbol{\beta}) = -\mathbb{E}[\log f_{\boldsymbol{\beta}}(Y | \mathbf{X}, \mathbf{Z})] + \max_{\boldsymbol{\lambda} \in \Lambda} \mathbb{E}[H_k(\boldsymbol{\lambda}, \boldsymbol{\beta})].$$

Under $\mathbb{E}[\mathbf{u}_i(\boldsymbol{\beta}^*)] = \mathbf{0}$ and the inner FOC, $\boldsymbol{\lambda}^*(\boldsymbol{\beta}^*) = \mathbf{0}$ maximizes the population inner dual, and standard likelihood identification plus convexity of the inner problem ensure that $Q(\boldsymbol{\beta})$ is uniquely minimized at $\boldsymbol{\beta}^*$. The argmin theorem thus gives $\widehat{\boldsymbol{\beta}}_k \rightarrow_p \boldsymbol{\beta}^*$.

Joint asymptotic normality of $(\widehat{\boldsymbol{\beta}}_k, \widehat{\boldsymbol{\lambda}}_k)$. Apply a mean-value expansion of the stacked estimating equations at $\boldsymbol{\vartheta}^*$:

$$\mathbf{0} = \boldsymbol{\Psi}_{N_{\text{in}}}(\widehat{\boldsymbol{\vartheta}}_k) = \boldsymbol{\Psi}_{N_{\text{in}}}(\boldsymbol{\vartheta}^*) + [\nabla_{\boldsymbol{\vartheta}} \boldsymbol{\Psi}_{N_{\text{in}}}(\bar{\boldsymbol{\vartheta}})] (\widehat{\boldsymbol{\vartheta}}_k - \boldsymbol{\vartheta}^*),$$

for some $\bar{\boldsymbol{\vartheta}}$ on the line segment between $\widehat{\boldsymbol{\vartheta}}_k$ and $\boldsymbol{\vartheta}^*$. Rearrange and scale:

$$\sqrt{N_{\text{in}}}(\widehat{\boldsymbol{\vartheta}}_k - \boldsymbol{\vartheta}^*) = - [\nabla_{\boldsymbol{\vartheta}} \boldsymbol{\Psi}_{N_{\text{in}}}(\bar{\boldsymbol{\vartheta}})]^{-1} \sqrt{N_{\text{in}}} \boldsymbol{\Psi}_{N_{\text{in}}}(\boldsymbol{\vartheta}^*).$$

By the CLT and regularity,

$$\sqrt{N_{\text{in}}} \boldsymbol{\Psi}_{N_{\text{in}}}(\boldsymbol{\vartheta}^*) \rightarrow \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{J} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Omega}_1 \end{bmatrix} \right),$$

and by the LLN,

$$\nabla_{\boldsymbol{\vartheta}} \boldsymbol{\Psi}_{N_{\text{in}}}(\bar{\boldsymbol{\vartheta}}) \xrightarrow{p} - \begin{bmatrix} \mathbf{J} & -\mathbf{C}^\top \\ -\mathbf{C} & \boldsymbol{\Omega}_2 \end{bmatrix}.$$

Hence, by Slutsky,

$$\sqrt{N_{\text{in}}} \begin{pmatrix} \widehat{\beta}_k - \beta^* \\ \widehat{\lambda}_k \end{pmatrix} \rightarrow \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{J} & -\mathbf{C}^\top \\ -\mathbf{C} & \Omega_2 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{J} & \mathbf{0} \\ \mathbf{0} & \Omega_1 \end{bmatrix} \begin{bmatrix} \mathbf{J} & -\mathbf{C}^\top \\ -\mathbf{C} & \Omega_2 \end{bmatrix}^{-1} \right).$$

Profiled asymptotics for $\widehat{\beta}_k$. From the joint normal limit, profiling out λ (blockwise Schur complement) yields

$$\sqrt{N_{\text{in}}}(\widehat{\beta}_k - \beta^*) \rightarrow \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1}\mathbf{V}\mathbf{A}^{-1}),$$

$$\mathbf{A} := \mathbf{J} + \mathbf{C}^\top \Omega_2^{-1} \mathbf{C}, \quad \mathbf{V} := \mathbf{J} + \mathbf{C}^\top \Omega_2^{-1} \Omega_1 \Omega_2^{-1} \mathbf{C}.$$

CR inner dual: $\Omega_1 = \Omega_2$ and k -invariance. By Lemma A3 (gradient/Hessian of H_k at $(\beta^*, \lambda = \mathbf{0})$), under correct specification

$$\Omega_1 := \frac{1}{N_{\text{in}}} \mathbb{E} \left[\{ \nabla_{\lambda} H_k \} \{ \nabla_{\lambda} H_k \}^\top \right] = \frac{1}{(k+1)^2} \mathbf{S},$$

$$\Omega_2 := -\frac{1}{N_{\text{in}}} \mathbb{E} \left[\nabla_{\lambda\lambda}^2 H_k \right] = \frac{1}{(k+1)^2} \mathbf{S}.$$

$$\mathbf{C} := \frac{1}{N_{\text{in}}} \mathbb{E} \left[\nabla_{\lambda\beta}^2 H_k \right] = \frac{1}{(k+1)} \mathbf{G}.$$

Therefore $\Omega_1 = \Omega_2$, and the common factor $(k+1)^{-2}$ cancels in \mathbf{A} and \mathbf{V} :

$$\mathbf{A} = \mathbf{J} + \mathbf{G}^\top \mathbf{S}^{-1} \mathbf{G}, \quad \mathbf{V} = \mathbf{J} + \mathbf{G}^\top \mathbf{S}^{-1} \mathbf{G},$$

so

$$\sqrt{N_{\text{in}}}(\widehat{\beta}_k - \beta^*) \rightarrow \mathcal{N} \left(\mathbf{0}, (\mathbf{J} + \mathbf{G}^\top \mathbf{S}^{-1} \mathbf{G})^{-1} \right).$$

□

APPENDIX F PROOF OF THEOREM 3 UNDER DISTRIBUTIONAL SHIFT

Lemma A6 (Influence expansion of the CR dual and additivity of first-order terms). *Fix (λ, β) in the interior of the domain $\{(\lambda, \beta) : 1 + \lambda^\top \mathbf{u}_\beta(\mathbf{W}) > 0 \text{ a.s.}\}$. Let P denote the population law of $\mathbf{W} = (Y, \mathbf{X}, \mathbf{Z})$ and P_n the empirical measure. For the Cressie-Read dual,*

$$H_k(\lambda, \beta; P) := \frac{N_{\text{in}}}{k(k+1)} \left(A(P; \lambda, \beta) \right)^{k+1}, \quad A(P; \lambda, \beta) := \mathbb{E}_P \left[(1 + \lambda^\top \mathbf{u}_\beta(\mathbf{W}))^\alpha \right], \quad \alpha := \frac{k}{k+1},$$

assume \mathbf{u}_β is measurable, P -integrable to the needed orders, and that the maps $(\lambda, \beta) \mapsto \mathbf{u}_\beta(\mathbf{w})$ are smooth. Then:

- (i) **Gateaux derivative (influence function) of H_k .** Define $h_k(x) := \frac{N_{\text{in}}}{k(k+1)} x^{k+1}$ so that $H_k = h_k \circ A$. Then the influence function of H_k at P is

$$\text{IF}_H(\mathbf{W}; \lambda, \beta; P) = \dot{H}_{k,P}(\mathbf{W}) = h'_k(A(P)) \left(a(\mathbf{W}) - A(P) \right), \quad a(\mathbf{W}) := (1 + \lambda^\top \mathbf{u}_\beta(\mathbf{W}))^\alpha.$$

Hence the linearization in the empirical process direction is

$$H_k(\lambda, \beta; P_n) - H_k(\lambda, \beta; P) = \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \text{IF}_H(\mathbf{W}_i; \lambda, \beta; P) + o_p(N_{\text{in}}^{-1/2}).$$

- (ii) **Influence functions of the gradients.** Recall that subscripts denote partial derivatives with respect to components of $\vartheta = (\beta^\top, \lambda^\top)^\top$. Write

$$\mathbf{A}_\vartheta := \nabla_{\vartheta} A \in \mathbb{R}^{d+q}, \quad \mathbf{A}_{\vartheta\vartheta} := \nabla_{\vartheta}^2 A \in \mathbb{R}^{(d+q) \times (d+q)}.$$

By the chain rule,

$$\nabla_{\vartheta} H_k = h'_k(A) \mathbf{A}_\vartheta, \quad \nabla_{\vartheta}^2 H_k = h''_k(A) \mathbf{A}_\vartheta \mathbf{A}_\vartheta^\top + h'_k(A) \mathbf{A}_{\vartheta\vartheta}.$$

Moreover, the influence functions of the gradients are

$$\text{IF}_{\nabla_{\vartheta} H_k}(\mathbf{W}; \boldsymbol{\lambda}, \boldsymbol{\beta}; P) = h_k''(A)(a(\mathbf{W}) - A)\mathbf{A}_{\vartheta} + h_k'(A)(\mathbf{A}_{\vartheta}(\mathbf{W}) - \mathbf{A}_{\vartheta}),$$

where $\mathbf{A}_{\vartheta}(\mathbf{W})$ is obtained by differentiating $a(\mathbf{W})$ pointwise in ϑ , and $\mathbf{A}_{\vartheta} = \mathbb{E}_P[\mathbf{A}_{\vartheta}(\mathbf{W})]$. In particular, the centered linear terms $\text{IF}_{\nabla_{\vartheta} H_k}(\mathbf{W}) - \mathbb{E}_P[\text{IF}_{\nabla_{\vartheta} H_k}(\mathbf{W})]$ are additive in \mathbf{W} .

(iii) **Additive “meat” for the stacked score.** Consider the stacked KKT map

$$\Psi(P, \vartheta) := \begin{pmatrix} \mathbb{E}_P[\mathbf{s}(\boldsymbol{\beta})] - \nabla_{\boldsymbol{\beta}} H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}; P)/N_{\text{in}} \\ \nabla_{\boldsymbol{\lambda}} H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}; P)/N_{\text{in}} \end{pmatrix}, \quad \vartheta := (\boldsymbol{\beta}^{\top}, \boldsymbol{\lambda}^{\top})^{\top}.$$

Although H_k itself is not additively separable, its empirical fluctuations are linearized by the influence functions in (i)-(ii). Hence

$$\sqrt{N_{\text{in}}}(\Psi(P_n, \vartheta) - \Psi(P, \vartheta)) = \frac{1}{\sqrt{N_{\text{in}}}} \sum_{i=1}^{N_{\text{in}}} \underbrace{\left(\mathbf{s}_i(\boldsymbol{\beta}) - \mathbb{E}_P[\mathbf{s}(\boldsymbol{\beta})] - \frac{1}{N_{\text{in}}} \text{IF}_{\nabla_{\boldsymbol{\beta}} H_k}(\mathbf{W}_i) \right)}_{:= \varphi(\mathbf{W}_i; \vartheta, P)} + o_p(1).$$

Thus, the asymptotic “meat” of the sandwich covariance is

$$\mathcal{B}(\vartheta) := \text{Var}_P(\varphi(\mathbf{W}; \vartheta, P)),$$

which is a variance of per-sample contributions.

(iv) **CR specialization near $\boldsymbol{\lambda} = \mathbf{0}$.** At $\boldsymbol{\lambda} = \mathbf{0}$ we have $A(P) = 1$ and $h_k'(1) = N_{\text{in}}/k$. Writing $\mathbf{g}(\mathbf{W}) := \mathbf{u}_{\boldsymbol{\beta}}(\mathbf{W})$ and recalling $\alpha = \frac{k}{k+1}$,

$$\text{IF}_{\nabla_{\boldsymbol{\lambda}} H_k}(\mathbf{W}) = h_k'(1)(\mathbf{A}_{\boldsymbol{\lambda}}(\mathbf{W}) - \mathbf{A}_{\boldsymbol{\lambda}}) = \frac{N_{\text{in}}}{k} \alpha (\mathbf{g}(\mathbf{W}) - \mathbb{E}_P[\mathbf{g}(\mathbf{W})]) \propto \mathbf{g}(\mathbf{W}) - \mathbb{E}_P[\mathbf{g}(\mathbf{W})].$$

Consequently, the leading blocks of \mathcal{B}^{\dagger} reduce (up to known multiplicative constants that cancel in the small-shift CR expansions) to the familiar

$$\mathbf{K}^{\dagger} = \text{Var}(\mathbf{s}_i(\boldsymbol{\beta}^{\dagger})), \quad \boldsymbol{\Omega}_1^{\dagger} = \text{Var}(\mathbf{g}_i), \quad \mathbf{K}_{s_g}^{\dagger} = \text{Cov}(\mathbf{s}_i(\boldsymbol{\beta}^{\dagger}), \mathbf{g}_i),$$

with $\mathbf{g}_i := \mathbf{u}_i(\boldsymbol{\beta}^{\dagger})$.

Proof sketch. (i) View $H_k = h_k \circ a$ with $A(P) = \int a(\mathbf{W}) dP$. The Gateaux derivative in the direction $\delta_{\mathbf{w}} - P$ is $\dot{A}_P(\mathbf{W}) = a(\mathbf{W}) - A(P)$. Apply the chain rule to get $\dot{H}_{k,P}(\mathbf{W}) = h_k'(A)\dot{A}_P(\mathbf{W})$.

(ii) Differentiate H_k w.r.t. parameters via the chain rule; then take the Gateaux derivative in P again. The term $\mathbf{A}_{\vartheta}(\mathbf{W}) - \mathbf{A}_{\vartheta}$ comes from linearizing $\int \mathbf{A}_{\vartheta}(\mathbf{W}) dP$ and the factor $h_k''(A)(a(\mathbf{W}) - A)\mathbf{A}_{\vartheta}$ comes from differentiating $h_k'(A)$.

(iii) Stack the score for $\boldsymbol{\beta}$ with the dual score for $\boldsymbol{\lambda}$ and linearize $\Psi(P_n, \vartheta) - \Psi(P, \vartheta)$ using (i)-(ii). This yields a sum of i.i.d. mean-zero terms, which defines the additive influence vector $\varphi(\mathbf{W}_i; \vartheta, P)$.

(iv) Substitute $\boldsymbol{\lambda} = \mathbf{0}$ and simplify $a(\mathbf{W}) = (1 + \boldsymbol{\lambda}^{\top} \mathbf{g}(\mathbf{W}))^{\alpha}$ and its derivatives. The leading constants cancel in the CR small-shift profiles, giving the standard block forms used in the main theorem. \square

Theorem A6 (Asymptotic normality of DRO estimator under shift: full version). *Assume the regularity conditions in Appendix D. Let $(\hat{\boldsymbol{\beta}}_k, \hat{\boldsymbol{\lambda}}_k)$ solve the sample KKT system of the DRO objective with Cressie-Read index k . Define the stacked estimating map*

$$\Psi_{N_{\text{in}}}(\vartheta) := \begin{pmatrix} \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} \mathbf{s}_i(\boldsymbol{\beta}) - \frac{1}{N_{\text{in}}} \nabla_{\boldsymbol{\beta}} H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) \\ \frac{1}{N_{\text{in}}} \nabla_{\boldsymbol{\lambda}} H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}) \end{pmatrix}, \quad \vartheta := \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\lambda} \end{pmatrix},$$

where $\mathbf{s}_i(\boldsymbol{\beta}) = \nabla_{\boldsymbol{\beta}} \ell_i(\boldsymbol{\beta})$ and $\ell_i(\boldsymbol{\beta}) = \log f_{\boldsymbol{\beta}}(Y_i | \mathbf{X}_i, \mathbf{Z}_i)$. Let $\widehat{\boldsymbol{\vartheta}}_k$ be the solution to $\boldsymbol{\Psi}_{N_{\text{in}}}(\widehat{\boldsymbol{\vartheta}}_k) = \mathbf{0}$, and let $\boldsymbol{\vartheta}^\dagger = (\boldsymbol{\beta}^{\dagger\top}, \boldsymbol{\lambda}^{\dagger\top})^\top$ be the unique population solution satisfying $\mathbb{E}[\boldsymbol{\Psi}_{N_{\text{in}}}(\boldsymbol{\vartheta}^\dagger)] = \mathbf{0}$.

Bread and meat. Define the bread and meat as

$$\mathcal{A}^\dagger := \mathbb{E}[\nabla_{\boldsymbol{\vartheta}} \boldsymbol{\Psi}_{N_{\text{in}}}(\boldsymbol{\vartheta}^\dagger)], \quad \mathcal{B}^\dagger := \text{Var}(\boldsymbol{\Phi}_i(\boldsymbol{\vartheta}^\dagger)),$$

where the per-sample contribution $\boldsymbol{\Phi}_i$ is

$$\boldsymbol{\Phi}_i(\boldsymbol{\vartheta}^\dagger) := \begin{pmatrix} \mathbf{s}_i(\boldsymbol{\beta}^\dagger) + \mathbf{H}_{\boldsymbol{\beta}}^{(i)}(\boldsymbol{\lambda}^\dagger, \boldsymbol{\beta}^\dagger) \\ \mathbf{H}_{\boldsymbol{\lambda}}^{(i)}(\boldsymbol{\lambda}^\dagger, \boldsymbol{\beta}^\dagger) \end{pmatrix},$$

and (definition) the terms $\mathbf{H}_{\boldsymbol{\beta}}^{(i)}, \mathbf{H}_{\boldsymbol{\lambda}}^{(i)}$ are the scaled influence contributions of H_k at $\boldsymbol{\vartheta}^\dagger$, namely

$$\mathbf{H}_{\boldsymbol{\beta}}^{(i)} := -\frac{1}{N_{\text{in}}} \text{IF}_{\nabla_{\boldsymbol{\beta}} H_k}(\mathbf{W}_i; \boldsymbol{\vartheta}^\dagger; P), \quad \mathbf{H}_{\boldsymbol{\lambda}}^{(i)} := \frac{1}{N_{\text{in}}} \text{IF}_{\nabla_{\boldsymbol{\lambda}} H_k}(\mathbf{W}_i; \boldsymbol{\vartheta}^\dagger; P),$$

so that $\frac{1}{\sqrt{N_{\text{in}}}} \sum_i \boldsymbol{\Phi}_i$ is precisely the linearization of $\sqrt{N_{\text{in}}}(\boldsymbol{\Psi}(P_n, \boldsymbol{\vartheta}^\dagger) - \boldsymbol{\Psi}(P, \boldsymbol{\vartheta}^\dagger))$ (see Lemma A6).

Block form. Then \mathcal{A}^\dagger and \mathcal{B}^\dagger admit the block decompositions

$$\mathcal{A}^\dagger = \begin{bmatrix} \mathbf{J}^\dagger + \mathbf{C}_{\boldsymbol{\beta}\boldsymbol{\beta}}^\dagger & -(\mathbf{C}^\dagger)^\top \\ -\mathbf{C}^\dagger & \boldsymbol{\Omega}_2^\dagger \end{bmatrix}, \quad \mathcal{B}^\dagger = \begin{bmatrix} \mathbf{K}^\dagger & \mathbf{K}_{sg}^\dagger \\ (\mathbf{K}_{sg}^\dagger)^\top & \boldsymbol{\Omega}_1^\dagger \end{bmatrix},$$

with

$$\mathbf{J}^\dagger := \mathbb{E}[-\nabla_{\boldsymbol{\beta}\boldsymbol{\beta}}^2 \ell_i(\boldsymbol{\beta}^\dagger)],$$

$$\mathbf{C}_{\boldsymbol{\beta}\boldsymbol{\beta}}^\dagger := -\frac{1}{N_{\text{in}}} \mathbb{E}[\nabla_{\boldsymbol{\beta}\boldsymbol{\beta}}^2 H_k(\boldsymbol{\lambda}^\dagger, \boldsymbol{\beta}^\dagger)], \quad \mathbf{C}^\dagger := \frac{1}{N_{\text{in}}} \mathbb{E}[\nabla_{\boldsymbol{\lambda}\boldsymbol{\beta}}^2 H_k(\boldsymbol{\lambda}^\dagger, \boldsymbol{\beta}^\dagger)],$$

$$\boldsymbol{\Omega}_2^\dagger := -\frac{1}{N_{\text{in}}} \mathbb{E}[\nabla_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^2 H_k(\boldsymbol{\lambda}^\dagger, \boldsymbol{\beta}^\dagger)],$$

$$\mathbf{K}^\dagger := \text{Var}(\mathbf{s}_i(\boldsymbol{\beta}^\dagger) + \mathbf{H}_{\boldsymbol{\beta}}^{(i)}(\boldsymbol{\lambda}^\dagger, \boldsymbol{\beta}^\dagger)), \quad \mathbf{K}_{sg}^\dagger := \text{Cov}(\mathbf{s}_i(\boldsymbol{\beta}^\dagger) + \mathbf{H}_{\boldsymbol{\beta}}^{(i)}(\boldsymbol{\lambda}^\dagger, \boldsymbol{\beta}^\dagger), \mathbf{H}_{\boldsymbol{\lambda}}^{(i)}(\boldsymbol{\lambda}^\dagger, \boldsymbol{\beta}^\dagger)),$$

$$\boldsymbol{\Omega}_1^\dagger := \text{Var}(\mathbf{H}_{\boldsymbol{\lambda}}^{(i)}(\boldsymbol{\lambda}^\dagger, \boldsymbol{\beta}^\dagger)).$$

(Remarks.) (i) The signs and $1/N_{\text{in}}$ factors match $\boldsymbol{\Psi}_{N_{\text{in}}}$: the minus sign on $\nabla_{\boldsymbol{\beta}} H_k$ in $\boldsymbol{\Psi}_{N_{\text{in}}}$ is absorbed by the definition of $\mathbf{H}_{\boldsymbol{\beta}}^{(i)}$ above. (ii) H_k is the aggregate dual penalty from Eq. 2; its second derivatives scale like N_{in} , and the explicit $1/N_{\text{in}}$ factors above ensure each block of \mathcal{A}^\dagger is $O(1)$.

Limit distribution. Then

$$\sqrt{N_{\text{in}}}(\widehat{\boldsymbol{\vartheta}}_k - \boldsymbol{\vartheta}^\dagger) \xrightarrow{d} \mathcal{N}(\mathbf{0}, (\mathcal{A}^\dagger)^{-1} \mathcal{B}^\dagger (\mathcal{A}^\dagger)^{-\top}),$$

where $(\mathcal{A}^\dagger)^{-\top} := ((\mathcal{A}^\dagger)^{-1})^\top$. In particular, the marginal limit for $\widehat{\boldsymbol{\beta}}_k$ is

$$\sqrt{N_{\text{in}}}(\widehat{\boldsymbol{\beta}}_k - \boldsymbol{\beta}^\dagger) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{V}_{\boldsymbol{\beta}}^\dagger),$$

with the profile sandwich variance

$$\mathbf{V}_{\boldsymbol{\beta}}^\dagger = (\mathbf{A}_{\boldsymbol{\beta}}^\dagger)^{-1} \left(\mathbf{K}^\dagger + (\mathbf{C}^\dagger)^\top (\boldsymbol{\Omega}_2^\dagger)^{-1} \boldsymbol{\Omega}_1^\dagger (\boldsymbol{\Omega}_2^\dagger)^{-1} \mathbf{C}^\dagger + (\mathbf{C}^\dagger)^\top (\boldsymbol{\Omega}_2^\dagger)^{-1} \mathbf{K}_{sg}^\dagger + (\mathbf{K}_{sg}^\dagger)^\top (\boldsymbol{\Omega}_2^\dagger)^{-1} \mathbf{C}^\dagger \right) (\mathbf{A}_{\boldsymbol{\beta}}^\dagger)^{-1},$$

where

$$\mathbf{A}_{\boldsymbol{\beta}}^\dagger := \mathbf{J}^\dagger + \mathbf{C}_{\boldsymbol{\beta}\boldsymbol{\beta}}^\dagger + (\mathbf{C}^\dagger)^\top (\boldsymbol{\Omega}_2^\dagger)^{-1} \mathbf{C}^\dagger.$$

Proof of Theorem A6. Let P denote the population law of $\mathbf{W} = (Y, \mathbf{X}, \mathbf{Z})$ and $P_n := N_{\text{in}}^{-1} \sum_{i=1}^{N_{\text{in}}} \delta_{\mathbf{W}_i}$ the empirical measure. Define the stacked KKT map (as a functional of $(P, \boldsymbol{\vartheta})$)

$$\boldsymbol{\Psi}(P, \boldsymbol{\vartheta}) := \begin{pmatrix} \mathbb{E}_P[\mathbf{s}_i(\boldsymbol{\beta})] - \frac{1}{N_{\text{in}}} \nabla_{\boldsymbol{\beta}} H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}; P) \\ \frac{1}{N_{\text{in}}} \nabla_{\boldsymbol{\lambda}} H_k(\boldsymbol{\lambda}, \boldsymbol{\beta}; P) \end{pmatrix}, \quad \boldsymbol{\Psi}_n(\boldsymbol{\vartheta}) := \boldsymbol{\Psi}(P_n, \boldsymbol{\vartheta}),$$

where $H_k(\lambda, \beta; P)$ is the CR dual objective with expectation under P replacing the empirical average. For the CR family, H_k is a smooth composition of linear functionals of P (e.g., $A(\lambda, \beta; P) = \int (1 + \lambda^\top \mathbf{u}_\beta)^\alpha dP$ with $\alpha = \frac{k}{k+1}$) and smooth maps (power, log), so $\Psi(\cdot, \cdot)$ is Hadamard differentiable jointly in (P, ϑ) in a neighborhood of (P, ϑ^\dagger) .

Consistency. By assumption, the population KKT system $\Psi(P, \vartheta) = \mathbf{0}$ has a unique solution $\vartheta^\dagger = (\beta^{\dagger\top}, \lambda^{\dagger\top})^\top$. A uniform law of large numbers for $\{\Psi(P, \vartheta) : \vartheta \in \Theta\}$ (together with identification; see Appendix D) implies $\widehat{\vartheta}_k \rightarrow_p \vartheta^\dagger$.

Joint linearization. A first-order Fréchet expansion of Ψ in both arguments around (P, ϑ^\dagger) gives

$$\mathbf{0} = \Psi(P_n, \widehat{\vartheta}_k) = \underbrace{\Psi(P, \vartheta^\dagger)}_{=\mathbf{0}} + \underbrace{\dot{\Psi}_\vartheta(P, \vartheta^\dagger)}_{=:\mathcal{A}^\dagger}(\widehat{\vartheta}_k - \vartheta^\dagger) + \underbrace{\dot{\Psi}_P(P, \vartheta^\dagger)[P_n - P]}_{=:\Delta_n} + r_n,$$

where $\dot{\Psi}_\vartheta$ is the Jacobian w.r.t. ϑ and $\dot{\Psi}_P(P, \vartheta^\dagger)[\cdot]$ is the Gâteaux derivative in the direction of a signed measure. Hadamard differentiability plus $\widehat{\vartheta}_k \rightarrow_p \vartheta^\dagger$ imply $r_n = o_p(N_{\text{in}}^{-1/2})$. Differentiating Ψ at ϑ^\dagger yields

$$\mathcal{A}^\dagger = \nabla_\vartheta \Psi(P, \vartheta)|_{\vartheta=\vartheta^\dagger} = \begin{bmatrix} \mathbf{J}^\dagger + \mathbf{C}_{\beta\beta}^\dagger & -(\mathbf{C}^\dagger)^\top \\ -\mathbf{C}^\dagger & \Omega_2^\dagger \end{bmatrix},$$

with the blocks defined exactly as in the theorem (note the explicit $1/N_{\text{in}}$ factors inherited from Ψ).

Influence representation. By the functional delta method (e.g., Van der Vaart (1998, Thm. 20.8)),

$$\sqrt{N_{\text{in}}}\Delta_n = \frac{1}{\sqrt{N_{\text{in}}}} \sum_{i=1}^{N_{\text{in}}} \varphi_i + o_p(1),$$

for a mean-zero influence vector φ_i obtained by applying Lemma A6 componentwise to Ψ . Writing the observation-level contributions to the H_k gradients as

$$\mathbf{H}_\beta^{(i)}(\lambda^\dagger, \beta^\dagger) := -\frac{1}{N_{\text{in}}} \text{IF}_{\nabla_\beta H_k}(\mathbf{W}_i; \vartheta^\dagger; P), \quad \mathbf{H}_\lambda^{(i)}(\lambda^\dagger, \beta^\dagger) := \frac{1}{N_{\text{in}}} \text{IF}_{\nabla_\lambda H_k}(\mathbf{W}_i; \vartheta^\dagger; P),$$

we can take

$$\varphi_i = \begin{pmatrix} \mathbf{s}_i(\beta^\dagger) + \mathbf{H}_\beta^{(i)}(\lambda^\dagger, \beta^\dagger) \\ \mathbf{H}_\lambda^{(i)}(\lambda^\dagger, \beta^\dagger) \end{pmatrix} - \mathbb{E} \left[\begin{pmatrix} \mathbf{s}_i(\beta^\dagger) + \mathbf{H}_\beta^{(i)}(\lambda^\dagger, \beta^\dagger) \\ \mathbf{H}_\lambda^{(i)}(\lambda^\dagger, \beta^\dagger) \end{pmatrix} \right].$$

Its covariance $\mathcal{B}^\dagger := \text{Var}(\varphi_i)$ has the block form stated in the theorem:

$$\mathcal{B}^\dagger = \begin{bmatrix} \mathbf{K}^\dagger & \mathbf{K}_{sg}^\dagger \\ (\mathbf{K}_{sg}^\dagger)^\top & \Omega_1^\dagger \end{bmatrix}.$$

Joint asymptotic normality. Rearranging the linearization,

$$\sqrt{N_{\text{in}}}(\widehat{\vartheta}_k - \vartheta^\dagger) = -(\mathcal{A}^\dagger)^{-1} \sqrt{N_{\text{in}}}\Delta_n + o_p(1) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathcal{A}^{\dagger-1} \mathcal{B}^\dagger \mathcal{A}^{\dagger-\top}).$$

Profiling to obtain \mathbf{V}_β^\dagger . Write $\mathcal{A}^\dagger = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{D} \end{bmatrix}$ with $\mathbf{A} := \mathbf{J}^\dagger + \mathbf{C}_{\beta\beta}^\dagger$, $\mathbf{B} := -\mathbf{C}^\dagger$, $\mathbf{D} := \Omega_2^\dagger$,

and $\mathcal{B}^\dagger = \begin{bmatrix} \mathbf{K}^\dagger & \mathbf{K}_{sg}^\dagger \\ (\mathbf{K}_{sg}^\dagger)^\top & \Omega_1^\dagger \end{bmatrix}$. The Schur complement of \mathbf{D} in \mathcal{A}^\dagger is

$$\mathbf{A}_\beta^\dagger := \mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{B}^\top = \mathbf{J}^\dagger + \mathbf{C}_{\beta\beta}^\dagger + (\mathbf{C}^\dagger)^\top (\Omega_2^\dagger)^{-1} \mathbf{C}^\dagger.$$

Standard block-matrix algebra gives the profiled covariance of the β -component:

$$\mathbf{V}_\beta^\dagger = (\mathbf{A}_\beta^\dagger)^{-1} \left(\mathbf{K}^\dagger + (\mathbf{C}^\dagger)^\top (\Omega_2^\dagger)^{-1} \Omega_1^\dagger (\Omega_2^\dagger)^{-1} \mathbf{C}^\dagger + (\mathbf{C}^\dagger)^\top (\Omega_2^\dagger)^{-1} \mathbf{K}_{sg}^\dagger + (\mathbf{K}_{sg}^\dagger)^\top (\Omega_2^\dagger)^{-1} \mathbf{C}^\dagger \right) (\mathbf{A}_\beta^\dagger)^{-1}.$$

Comment on non-additivity. Although $\nabla_{\beta}H_k$ and $\nabla_{\lambda}H_k$ are not simple sums over i , they are smooth functionals of empirical means (such as $A(P_n)$, $\log A(P_n)$, powers), hence admit linear influence representations in $(P_n - P)$. This is exactly what enters the preceding expansion via Lemma A6, so the standard Z-estimation machinery applies.

All steps are standard for Z-estimators with smooth (Hadamard differentiable) scores; see, e.g., Newey and McFadden (1994, Sec. 7) and Van der Vaart (1998, Ch. 5, 20). \square

APPENDIX G LOCAL EXPANSION UNDER DISTRIBUTIONAL SHIFT

Lemma A7 (First-order bias expansion under shift). *Let the assumptions of Theorem A6 hold and expand the KKT system locally around $\vartheta^* = (\beta^*, \mathbf{0})$. Define the shift $\delta := \mathbb{E}[\mathbf{u}_i(\beta^*)] \in \mathbb{R}^q$ and*

$$\mathbf{J} := \mathbb{E}[-\nabla_{\beta\beta}^2 \ell_i(\beta^*)], \quad \mathbf{G} := \mathbb{E}[\nabla_{\beta} \mathbf{u}_i(\beta^*)], \quad \mathbf{S} := \mathbb{E}[\mathbf{u}_i(\beta^*) \mathbf{u}_i(\beta^*)^{\top}].$$

Assume \mathbf{S} is positive definite and $\mathbf{A}_{\beta} := \mathbf{J} - \mathbf{G}^{\top} \mathbf{S}^{-1} \mathbf{G}$ is nonsingular. Then, for the Cressie-Read dual with $\|\widehat{\lambda}_k\| = o_p(1)$ (small shift),

$$\widehat{\beta}_k - \beta^* = (\mathbf{J} - \mathbf{G}^{\top} \mathbf{S}^{-1} \mathbf{G})^{-1} \mathbf{G}^{\top} \mathbf{S}^{-1} \delta + O_p(N_{\text{in}}^{-1/2}) + o(\|\delta\|).$$

In particular, the leading (deterministic) bias is k -free; dependence on k enters only through higher-order terms when the shift is not infinitesimal.

Proof of Lemma A7. Let P be the population law and P_n the empirical measure. Work with the scaled KKT map

$$\Psi_{N_{\text{in}}}(\vartheta) := \begin{pmatrix} \frac{1}{N_{\text{in}}} \sum_{i=1}^{N_{\text{in}}} s_i(\beta) - \frac{1}{N_{\text{in}}} \nabla_{\beta} H_k(\lambda, \beta; P_n) \\ \frac{1}{N_{\text{in}}} \nabla_{\lambda} H_k(\lambda, \beta; P_n) \end{pmatrix}, \quad \vartheta = \begin{pmatrix} \beta \\ \lambda \end{pmatrix}.$$

At the population level and at $\vartheta^* = (\beta^*, \mathbf{0})$, its population counterpart is

$$\Psi(P, \vartheta^*) := \begin{pmatrix} \mathbb{E}_P[s_i(\beta^*)] - \frac{1}{N_{\text{in}}} \nabla_{\beta} H_k(\mathbf{0}, \beta^*; P) \\ \frac{1}{N_{\text{in}}} \nabla_{\lambda} H_k(\mathbf{0}, \beta^*; P) \end{pmatrix}.$$

For the Cressie-Read (CR) family ($k \neq 0, -1$),

$$\frac{1}{N_{\text{in}}} \nabla_{\lambda} H_k(\mathbf{0}, \beta^*; P) = \frac{1}{k+1} \mathbb{E}_P[\mathbf{u}_i(\beta^*)] = \frac{1}{k+1} \delta, \quad \frac{1}{N_{\text{in}}} \nabla_{\beta} H_k(\mathbf{0}, \beta^*; P) = \mathbf{0},$$

so $\Psi(P, \vartheta^*) = (\mathbf{0}, c_k \delta)^{\top}$ with $c_k = \frac{1}{k+1}$.

Now take a first-order expansion of $\Psi(P, \cdot)$ around ϑ^* :

$$\mathbf{0} = \Psi(P, \widehat{\vartheta}_k) = \Psi(P, \vartheta^*) + \mathcal{A}^* \begin{pmatrix} \widehat{\beta}_k - \beta^* \\ \widehat{\lambda}_k - \mathbf{0} \end{pmatrix} + \mathbf{R}_1, \quad \mathbf{R}_1 = o(\|\widehat{\vartheta}_k - \vartheta^*\|),$$

where the Jacobian at $(\beta^*, \mathbf{0})$ is

$$\mathcal{A}^* := \mathbb{E}[\nabla_{\vartheta} \Psi_{N_{\text{in}}}(\vartheta^*)] = \begin{bmatrix} \mathbf{J} & -(\mathbf{C}^*)^{\top} \\ -\mathbf{C}^* & \Omega_2^* \end{bmatrix}, \quad \mathbf{J} := \mathbb{E}[-\nabla_{\beta\beta}^2 \ell_i(\beta^*)].$$

For CR at $\lambda = \mathbf{0}$ (under the same scaling by $1/N_{\text{in}}$ used in $\Psi_{N_{\text{in}}}$),

$$\mathbf{C}^* = \frac{1}{k+1} \mathbf{G}, \quad \Omega_2^* = \frac{1}{(k+1)^2} \mathbf{S},$$

$$\mathbf{G} := \mathbb{E}[\nabla_{\beta} \mathbf{u}_i(\beta^*)], \quad \mathbf{S} := \mathbb{E}[\mathbf{u}_i(\beta^*) \mathbf{u}_i(\beta^*)^{\top}].$$

1674 Ignoring sampling noise for the bias calculation and solving the linear system $\mathcal{A}^*(\hat{\beta}_k - \beta^*, \hat{\lambda}_k)^\top =$
 1675 $-(\mathbf{0}, c_k \delta)^\top$ by Schur complement yields

$$1676 \hat{\beta}_k - \beta^* = (\mathbf{J} - (\mathbf{C}^*)^\top (\boldsymbol{\Omega}_2^*)^{-1} \mathbf{C}^*)^{-1} (\mathbf{C}^*)^\top (\boldsymbol{\Omega}_2^*)^{-1} (c_k \delta) + (\text{higher order}).$$

1677 Now use

$$1678 (\mathbf{C}^*)^\top (\boldsymbol{\Omega}_2^*)^{-1} = \left(\frac{1}{k+1} \mathbf{G}^\top \right) \left((k+1)^2 \mathbf{S}^{-1} \right) = (k+1) \mathbf{G}^\top \mathbf{S}^{-1},$$

1679 so that the factor $c_k = \frac{1}{k+1}$ cancels $(k+1)$ exactly:

$$1680 (\mathbf{C}^*)^\top (\boldsymbol{\Omega}_2^*)^{-1} (c_k \delta) = \mathbf{G}^\top \mathbf{S}^{-1} \delta.$$

1681 Also,

$$1682 (\mathbf{C}^*)^\top (\boldsymbol{\Omega}_2^*)^{-1} \mathbf{C}^* = \mathbf{G}^\top \mathbf{S}^{-1} \mathbf{G}.$$

1683 Hence the deterministic first-order term simplifies to

$$1684 \hat{\beta}_k - \beta^* = (\mathbf{J} - \mathbf{G}^\top \mathbf{S}^{-1} \mathbf{G})^{-1} \mathbf{G}^\top \mathbf{S}^{-1} \delta + (\text{higher order}).$$

1685 Finally, (i) the empirical-process remainder contributes $O_p(N_{\text{in}}^{-1/2})$ to $\hat{\beta}_k - \beta^*$ by Theorem A6; (ii)
 1686 the Taylor remainder is $o(\|\delta\|)$ by smoothness (Hadamard differentiability in P and C^2 in (λ, β))
 1687 near $(\mathbf{0}, \beta^*)$. Collecting terms gives

$$1688 \hat{\beta}_k - \beta^* = (\mathbf{J} - \mathbf{G}^\top \mathbf{S}^{-1} \mathbf{G})^{-1} \mathbf{G}^\top \mathbf{S}^{-1} \delta + O_p(N_{\text{in}}^{-1/2}) + o(\|\delta\|),$$

1689 as claimed. \square

1690 **Interpretation.** The bias expansion (Lemma A7) shows that misspecification enters through the
 1691 shift term δ , with bias proportional to $(\mathbf{J} - \mathbf{G}^\top \mathbf{S}^{-1} \mathbf{G})^{-1} \mathbf{G}^\top \mathbf{S}^{-1} \delta$. Hence the estimator is stable
 1692 under mild shifts ($\|\delta\|$ small) but deteriorates when the shift is large.

1700 APPENDIX H INNER/OUTER DERIVATIVES AND IMPLICIT DIFFERENTIATION 1701 (GENERAL k)

1702 Let $N := N_{\text{in}}$. For a fixed β , define

$$1703 t_i(\lambda, \beta) := \lambda^\top \mathbf{u}_i(\beta), \quad m_i := \begin{cases} -\log(1 + t_i), & k = 0 \\ \exp\left(\frac{t_i}{2}\right), & k = -1 \\ (1 + t_i)^{\frac{k}{k+1}}, & k < 0, k \neq -1 \\ -(1 + t_i)^{\frac{k}{k+1}}, & k > 0, \end{cases} \quad M := \sum_{i=1}^N m_i.$$

1704 The dual objective can be written as $H_k(\lambda, \beta) = h_k(M)$ with

$$1705 h_k(m) = \begin{cases} -m, & k = 0 \\ -N \log(m/N), & k = -1 \\ \frac{N^{-k}}{k(k+1)} m^{k+1}, & k < 0, k \neq -1 \\ \frac{N^{-k}}{k(k+1)} (-m)^{k+1}, & k > 0. \end{cases}$$

1706 Hence

$$1707 h'_k(m) = \begin{cases} -1, & k = 0 \\ -\frac{N}{m}, & k = -1 \\ \frac{N^{-k}}{k} m^k, & k < 0, k \neq -1 \\ -\frac{N^{-k}}{k} (-m)^k, & k > 0 \end{cases} \quad h''_k(m) = \begin{cases} 0, & k = 0 \\ \frac{N}{m^2}, & k = -1 \\ N^{-k} m^{k-1}, & k < 0, k \neq -1 \\ N^{-k} (-m)^{k-1}, & k > 0. \end{cases}$$

Lemma A8 (Derivatives for the inner and outer problems). For each i , let $\mathbf{J}_i(\boldsymbol{\beta}) := \nabla_{\boldsymbol{\beta}} \mathbf{u}_i(\boldsymbol{\beta}) \in \mathbb{R}^{q \times d}$ and $\mathbf{H}_i(\boldsymbol{\beta}) := \nabla_{\boldsymbol{\beta}}^2 \mathbf{u}_i(\boldsymbol{\beta}) \in \mathbb{R}^{q \times d \times d}$.

(A) *Inner derivatives (w.r.t. λ).*

$$M_{\lambda} = \sum_{i=1}^N m'_i(t_i) \mathbf{u}_i, \quad M_{\lambda\lambda} = \sum_{i=1}^N m''_i(t_i) \mathbf{u}_i \mathbf{u}_i^{\top}.$$

Here

$$m'_i(t) = \begin{cases} -\frac{1}{1+t}, & k = 0 \\ \frac{1}{2} e^{t/2}, & k = -1 \\ \frac{k}{k+1} (1+t)^{-1/(k+1)}, & k < 0, k \neq -1 \\ -\frac{k}{k+1} (1+t)^{-1/(k+1)}, & k > 0 \end{cases}$$

$$m''_i(t) = \begin{cases} \frac{1}{(1+t)^2}, & k = 0 \\ \frac{1}{4} e^{t/2}, & k = -1 \\ -\frac{k}{(k+1)^2} (1+t)^{-(k+2)/(k+1)}, & k < 0, k \neq -1 \\ \frac{k}{(k+1)^2} (1+t)^{-(k+2)/(k+1)}, & k > 0. \end{cases}$$

(B) *Cross derivatives (w.r.t. (λ, β)).* For each coordinate $\beta^{[\ell]}$ (row index) and $\lambda^{[k]}$ (column index),

$$[M_{\lambda, \beta}]_{k\ell} = \sum_{i=1}^N \left(m'_i(t_i) [\mathbf{J}_i]_{k\ell} + m''_i(t_i) (\mathbf{J}_i^{\top} \boldsymbol{\lambda})_{\ell} u_i^{[k]} \right).$$

Equivalently, in matrix form

$$M_{\lambda, \beta} = \sum_{i=1}^N \left(m'_i(t_i) \mathbf{J}_i + m''_i(t_i) \mathbf{u}_i (\mathbf{J}_i^{\top} \boldsymbol{\lambda})^{\top} \right) \in \mathbb{R}^{q \times d}.$$

(C) *Implicit differentiation of the inner solution.* Let $\widehat{\boldsymbol{\lambda}}(\boldsymbol{\beta})$ solve the inner FOC $M_{\lambda}(\widehat{\boldsymbol{\lambda}}(\boldsymbol{\beta}), \boldsymbol{\beta}) = \mathbf{0}$. If $M_{\lambda\lambda}$ is nonsingular at $(\widehat{\boldsymbol{\lambda}}, \boldsymbol{\beta})$, then

$$\frac{d\widehat{\boldsymbol{\lambda}}}{d\boldsymbol{\beta}} = -\left(M_{\lambda\lambda}\right)^{-1} M_{\lambda, \beta}.$$

(D) *Total derivatives of $M(\widehat{\boldsymbol{\lambda}}(\boldsymbol{\beta}), \boldsymbol{\beta})$ w.r.t. $\boldsymbol{\beta}$.* Using $M_{\lambda} = 0$ at the inner optimum:

$$M_{\beta} = \frac{dM}{d\boldsymbol{\beta}} = \sum_{i=1}^N m'_i(t_i) \mathbf{J}_i^{\top} \boldsymbol{\lambda} \equiv \mathbf{A} \boldsymbol{\lambda}, \quad \mathbf{A} := \sum_{i=1}^N m'_i(t_i) \mathbf{J}_i^{\top} \in \mathbb{R}^{d \times q}.$$

The exact Hessian splits into a “direct” term and a Schur-complement term:

$$M_{\beta\beta} = \underbrace{\sum_{i=1}^N \left(m'_i(t_i) \mathcal{H}_i^{\top} \boldsymbol{\lambda} + m''_i(t_i) (\mathbf{J}_i^{\top} \boldsymbol{\lambda}) (\mathbf{J}_i^{\top} \boldsymbol{\lambda})^{\top} \right)}_{M_{\beta\beta}^{\text{dir}}} - \underbrace{M_{\lambda, \beta}^{\top} M_{\lambda\lambda}^{-1} M_{\lambda, \beta}}_{\text{Schur term}},$$

where \mathcal{H}_i stacks the slices of \mathbf{H}_i so that $\mathcal{H}_i^\top \boldsymbol{\lambda} \in \mathbb{R}^d$ has entries $\sum_k \lambda^{[k]} \partial^2 u_i^{[k]} / \partial \beta^{[l]} \partial \beta^{[a]}$.

(E) **Total derivatives of the composed dual** $H_k(\widehat{\boldsymbol{\lambda}}(\boldsymbol{\beta}), \boldsymbol{\beta}) = h_k(M)$. By the chain rule,

$$H_{\boldsymbol{\beta}} = h'_k(M)M_{\boldsymbol{\beta}}, \quad H_{\boldsymbol{\beta}\boldsymbol{\beta}} = h''_k(M)M_{\boldsymbol{\beta}}M_{\boldsymbol{\beta}}^\top + h'_k(M)M_{\boldsymbol{\beta}\boldsymbol{\beta}}.$$

(F) **Near- $\boldsymbol{\lambda} = \mathbf{0}$ simplification.** If $\|\widehat{\boldsymbol{\lambda}}\|$ is small (e.g., under mild shift), then $m'_i(t_i) = m'_i(0) + o(1)$ and the quadratic/second-derivative terms in $M_{\boldsymbol{\beta}\boldsymbol{\beta}}^{\text{dir}}$ are $o(1)$, giving the useful approximation

$$M_{\boldsymbol{\beta}\boldsymbol{\beta}} \approx -M_{\boldsymbol{\lambda},\boldsymbol{\beta}}^\top M_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^{-1} M_{\boldsymbol{\lambda},\boldsymbol{\beta}}, \quad H_{\boldsymbol{\beta}\boldsymbol{\beta}} \approx h''_k(M)M_{\boldsymbol{\beta}}M_{\boldsymbol{\beta}}^\top - h'_k(M)M_{\boldsymbol{\lambda},\boldsymbol{\beta}}^\top M_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^{-1} M_{\boldsymbol{\lambda},\boldsymbol{\beta}}.$$

Proof (sketch). (A)-(B) follow by the chain rule with $t_i = \boldsymbol{\lambda}^\top \mathbf{u}_i(\boldsymbol{\beta})$. (C) applies the implicit function theorem to the inner FOC $M_{\boldsymbol{\lambda}}(\widehat{\boldsymbol{\lambda}}(\boldsymbol{\beta}), \boldsymbol{\beta}) = \mathbf{0}$. (D) uses the total derivative

$$\frac{dM}{d\boldsymbol{\beta}} = M_{\boldsymbol{\beta}} + M_{\boldsymbol{\lambda}} \frac{d\widehat{\boldsymbol{\lambda}}}{d\boldsymbol{\beta}},$$

and the FOC $M_{\boldsymbol{\lambda}} = 0$ at the inner optimum; differentiating once more and substituting the implicit-diff formula in (C) yields the decomposition into the ‘‘direct’’ term and the Schur complement term for $M_{\boldsymbol{\beta}\boldsymbol{\beta}}$. (E) is a direct application of the chain rule to $H_k = h_k \circ M$. (F) drops terms that are $O(\|\widehat{\boldsymbol{\lambda}}\|)$ (or quadratic in $\widehat{\boldsymbol{\lambda}}$), which is justified under mild shift where the dual multipliers concentrate near $\mathbf{0}$.

□

APPENDIX I EMPIRICAL BAYES DERIVATION OF EB-DRO

I.1 EMPIRICAL BAYES DERIVATION OF EB-DRO

Hierarchical model. We model $\boldsymbol{\beta}$ hierarchically as

$$\boldsymbol{\beta} \mid \widehat{\boldsymbol{\beta}}_{\text{DRO},k} \sim \mathcal{N}(\widehat{\boldsymbol{\beta}}_{\text{DRO},k}, \mathbf{A}), \quad \widehat{\boldsymbol{\beta}}_I \mid \boldsymbol{\beta} \sim \mathcal{N}(\boldsymbol{\beta}, \boldsymbol{\Sigma}),$$

where $\widehat{\boldsymbol{\beta}}_{\text{DRO},k}$ is the DRO estimator, $\widehat{\boldsymbol{\beta}}_I$ is the internal-only GLM MLE, \mathbf{A} quantifies prior uncertainty around the DRO solution, and $\boldsymbol{\Sigma}$ is the covariance of the internal estimator.

Posterior mean. By Gaussian conjugacy,

$$\mathbb{E} \left[\boldsymbol{\beta} \mid \widehat{\boldsymbol{\beta}}_{\text{DRO},k}, \widehat{\boldsymbol{\beta}}_I \right] = \left(\mathbf{A}^{-1} + \boldsymbol{\Sigma}^{-1} \right)^{-1} \left(\mathbf{A}^{-1} \widehat{\boldsymbol{\beta}}_{\text{DRO},k} + \boldsymbol{\Sigma}^{-1} \widehat{\boldsymbol{\beta}}_I \right),$$

yielding the EB estimator in Section 2.2. Equivalently, using matrix identities,

$$\widehat{\boldsymbol{\beta}}_{\text{EB},k} = (\mathbf{I} - \mathbf{W})\widehat{\boldsymbol{\beta}}_{\text{DRO},k} + \mathbf{W}\widehat{\boldsymbol{\beta}}_I, \quad \mathbf{W} = \mathbf{A}(\mathbf{A} + \boldsymbol{\Sigma})^{-1},$$

which shows EB-DRO as a weighted average of the DRO and naive estimators.

Empirical Bayes estimates. In practice, we estimate the unknown quantities as follows:

- $\boldsymbol{\Sigma}$ is replaced by $\widehat{\boldsymbol{\Sigma}}_I$, the covariance estimator of $\widehat{\boldsymbol{\beta}}_I$ from the internal GLM.
- \mathbf{A} is taken as

$$\widehat{\mathbf{A}} = (\widehat{\boldsymbol{\beta}}_I - \widehat{\boldsymbol{\beta}}_{\text{DRO},k})(\widehat{\boldsymbol{\beta}}_I - \widehat{\boldsymbol{\beta}}_{\text{DRO},k})^\top,$$

a rank-one matrix encoding the discrepancy between the two estimators.

Substituting $(\widehat{\mathbf{A}}, \widehat{\boldsymbol{\Sigma}}_I)$ into the posterior mean yields the implementable EB estimator $\widehat{\boldsymbol{\beta}}_{\text{EB},k}$. Since $\widehat{\mathbf{A}}$ is rank-one, we interpret $\widehat{\mathbf{A}}^{-1}$ in the pseudoinverse sense (or equivalently as the $\varepsilon \downarrow 0$ limit of $\widehat{\mathbf{A}} + \varepsilon \mathbf{I}$), while $\widehat{\mathbf{A}} + \widehat{\boldsymbol{\Sigma}}_I$ is positive definite and invertible by construction.

I.2 PROOF OF PROPOSITION 1

In this section, we show that the matrix-valued weight $\widehat{\mathbf{W}}$ from the definition of the EB-DRO estimator can be simplified to a scalar multiplication by α .

Let $\mathbf{Z} := \widehat{\beta}_I - \widehat{\beta}_{\text{DRO},k}$, $\widehat{\mathbf{A}} := \mathbf{Z}\mathbf{Z}^\top$, and $\widehat{\Sigma}_I \succ 0$ denote the covariance estimator of $\widehat{\beta}_I$. Recall the EB-DRO estimator in matrix-average form

$$\widehat{\beta}_{\text{EB},k} = (\mathbf{I} - \widehat{\mathbf{W}})\widehat{\beta}_{\text{DRO},k} + \widehat{\mathbf{W}}\widehat{\beta}_I, \quad \widehat{\mathbf{W}} := \widehat{\mathbf{A}}(\widehat{\mathbf{A}} + \widehat{\Sigma}_I)^{-1}.$$

(Equivalently, this follows from Gaussian conjugacy together with $(\mathbf{A}^{-1} + \Sigma^{-1})^{-1}\mathbf{A}^{-1} = \Sigma(\mathbf{A} + \Sigma)^{-1}$ and $(\mathbf{A}^{-1} + \Sigma^{-1})^{-1}\Sigma^{-1} = \mathbf{A}(\mathbf{A} + \Sigma)^{-1}$.)

Set $s := \mathbf{Z}^\top \widehat{\Sigma}_I^{-1} \mathbf{Z} \geq 0$. By the Sherman-Morrison-Woodbury identity for a rank-one update,

$$(\widehat{\Sigma}_I + \mathbf{Z}\mathbf{Z}^\top)^{-1} = \widehat{\Sigma}_I^{-1} - \widehat{\Sigma}_I^{-1} \mathbf{Z} (1 + s)^{-1} \mathbf{Z}^\top \widehat{\Sigma}_I^{-1}.$$

Multiplying both sides by \mathbf{Z} gives

$$(\widehat{\Sigma}_I + \mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z} = \widehat{\Sigma}_I^{-1} \mathbf{Z} \left(1 - \frac{s}{1+s}\right) = \frac{1}{1+s} \widehat{\Sigma}_I^{-1} \mathbf{Z}. \quad (7)$$

Using $\widehat{\mathbf{W}} = \mathbf{Z}\mathbf{Z}^\top (\widehat{\Sigma}_I + \mathbf{Z}\mathbf{Z}^\top)^{-1}$ and Eq. 7,

$$\widehat{\mathbf{W}}\mathbf{Z} = \mathbf{Z}\mathbf{Z}^\top \frac{1}{1+s} \widehat{\Sigma}_I^{-1} \mathbf{Z} = \frac{s}{1+s} \mathbf{Z}.$$

Define

$$\alpha := \frac{s}{1+s} = \frac{\mathbf{Z}^\top \widehat{\Sigma}_I^{-1} \mathbf{Z}}{1 + \mathbf{Z}^\top \widehat{\Sigma}_I^{-1} \mathbf{Z}} \in [0, 1).$$

Since $\widehat{\beta}_{\text{EB},k} = \widehat{\beta}_{\text{DRO},k} + \widehat{\mathbf{W}}\mathbf{Z}$, we obtain

$$\widehat{\beta}_{\text{EB},k} = \widehat{\beta}_{\text{DRO},k} + \alpha \mathbf{Z} = (1 - \alpha)\widehat{\beta}_{\text{DRO},k} + \alpha \widehat{\beta}_I,$$

which is Eq. 3. Moreover, $\alpha = 0$ iff $\mathbf{Z} = \mathbf{0}$ (the two estimators coincide), and $\alpha \uparrow 1$ as the Mahalanobis discrepancy s grows, proving the stated range $\alpha \in (0, 1)$ for $\mathbf{Z} \neq \mathbf{0}$. \square

I.3 PROOF OF THEOREM 4

Proof of Theorem 4. Write

$$\widehat{\beta}(w) = (1 - w)\widehat{\beta}_{\text{DRO},k} + w\widehat{\beta}_I, \quad R(w) = \mathbb{E} \|\widehat{\beta}(w) - \beta^*\|_2^2.$$

Let $\mathbf{b}_k = \mathbb{E}[\widehat{\beta}_{\text{DRO},k}] - \beta^*$, $\mathbf{V}_{\text{DRO},k} = \text{Var}(\widehat{\beta}_{\text{DRO},k})$, $\Sigma_I = \text{Var}(\widehat{\beta}_I)$, and $\mathbf{C}_k = \text{Cov}(\widehat{\beta}_{\text{DRO},k}, \widehat{\beta}_I)$. Then

$$\mathbb{E}[\widehat{\beta}(w) - \beta^*] = (1 - w)\mathbf{b}_k, \quad \text{Var}(\widehat{\beta}(w)) = (1 - w)^2 \mathbf{V}_{\text{DRO},k} + w^2 \Sigma_I + 2w(1 - w)\mathbf{C}_k.$$

Hence

$$\begin{aligned} R(w) &= \|(1 - w)\mathbf{b}_k\|_2^2 + \text{tr} \left[(1 - w)^2 \mathbf{V}_{\text{DRO},k} + w^2 \Sigma_I + 2w(1 - w)\mathbf{C}_k \right] \\ &= (1 - w)^2 a + w^2 c + 2w(1 - w)m, \end{aligned}$$

with $a = \|\mathbf{b}_k\|_2^2 + \text{tr}(\mathbf{V}_{\text{DRO},k})$, $c = \text{tr}(\Sigma_I)$, $m = \text{tr}(\mathbf{C}_k)$. Expanding gives the quadratic

$$R(w) = a + (-2a + 2m)w + (a + c - 2m)w^2.$$

Let $\kappa := a + c - 2m = \|\mathbf{b}_k\|_2^2 + \text{tr} \text{Var}(\widehat{\beta}_{\text{DRO},k} - \widehat{\beta}_I) \geq 0$. The unconstrained minimizer is $w_{\mathbb{R}}^* = (a - m)/\kappa$ (if $\kappa = 0$, $R(w)$ is constant and any w is minimizer). Completing the square yields the **exact identity**

$$R(w) = R(w_{\mathbb{R}}^*) + \kappa(w - w_{\mathbb{R}}^*)^2, \quad (\kappa \geq 0). \quad (*)$$

Projecting onto $[0, 1]$ gives the oracle weight $w^* = \Pi_{[0,1]}(w_{\mathbb{R}}^*)$, which minimizes $R(w)$ over $[0, 1]$ by convexity. Therefore

$$R(w^*) \leq \min\{R(0), R(1)\},$$

establishing the dominance claim. Finally, using $R(w_{\mathbb{R}}^*) \leq R(w^*) \leq \min\{R(0), R(1)\}$ in (*) gives the **regret bound**

$$R(w) = R(w_{\mathbb{R}}^*) + \kappa(w - w_{\mathbb{R}}^*)^2 \leq \min\{R(0), R(1)\} + \kappa(w - w_{\mathbb{R}}^*)^2.$$

If $w^* \in (0, 1)$ (no projection), then $w^* = w_{\mathbb{R}}^*$ and the regret identity simplifies to $R(w) = R(w^*) + \kappa(w - w^*)^2$. \square

APPENDIX J PROOF OF THEOREM 5

Setup and notation. Stack the K EB-DRO estimators as

$$\widehat{\mathbf{B}} = \begin{bmatrix} \widehat{\beta}_{\text{EB},1} \\ \vdots \\ \widehat{\beta}_{\text{EB},K} \end{bmatrix} \in \mathbb{R}^{Kd}, \quad \mathbb{E}[\widehat{\mathbf{B}}] = (\mathbf{1}_K \otimes \mathbf{I}_d)(\beta^* + \mathbf{b}).$$

Let $\mathbf{V} = \text{Cov}(\widehat{\mathbf{B}})$ with $d \times d$ blocks $\Omega_{k\ell} = \text{Cov}(\widehat{\beta}_{\text{EB},k}, \widehat{\beta}_{\text{EB},\ell})$, and write the block-diagonal of marginal covariances as

$$\mathbf{D} = \text{diag}(\Omega_{11}, \dots, \Omega_{KK}).$$

Let $\mathbf{H} = \mathbf{1}_K \otimes \mathbf{I}_p \in \mathbb{R}^{Kd \times d}$, so that $\mathbb{E}[\widehat{\mathbf{B}}] = \mathbf{H}(\beta^* + \mathbf{b})$. The *oracle* GLS estimator of $\beta^* + \mathbf{b}$ from $\widehat{\mathbf{B}}$ is

$$\widetilde{\beta}_{\text{Ens}} = (\mathbf{H}^\top \mathbf{V}^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{V}^{-1} \widehat{\mathbf{B}}, \quad \text{Cov}(\widetilde{\beta}_{\text{Ens}}) = \mathbf{G}(\mathbf{V}) := (\mathbf{H}^\top \mathbf{V}^{-1} \mathbf{H})^{-1}.$$

The *precision-weighted* (PW) ensemble replaces \mathbf{V} by the block-diagonal precision:

$$\widehat{\beta}_{\text{PW}} = (\mathbf{H}^\top \widehat{\mathbf{D}}^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \widehat{\mathbf{D}}^{-1} \widehat{\mathbf{B}}, \quad \widehat{\mathbf{D}} = \text{diag}(\widehat{\Omega}_1, \dots, \widehat{\Omega}_K).$$

We assume $\widehat{\Omega}_k \rightarrow_p \Omega_{kk}$ for all k .

Throughout, mean-squared error (MSE) under squared loss decomposes as

$$\mathbb{E}\|\widehat{\theta} - \beta^*\|_2^2 = \underbrace{\|\mathbb{E}[\widehat{\theta}] - \beta^*\|_2^2}_{\text{(squared bias)}} + \underbrace{\text{tr}\{\text{Cov}(\widehat{\theta})\}}_{\text{(variance)}}.$$

Because all signals share the same mean $\beta^* + \mathbf{b}$, any linear unbiased combination has squared bias $\|\mathbf{b}\|_2^2$. Hence MSE comparisons reduce to comparing the *trace* of their covariance matrices.

Step 1: GLS dominates any component (Loewner order). Consider the class of linear unbiased estimators $\widehat{\beta} = \mathbf{A}^\top \widehat{\mathbf{B}}$ with $\mathbf{A} \in \mathbb{R}^{Kd \times d}$ satisfying the unbiasedness constraint $\mathbf{A}^\top \mathbf{H} = \mathbf{I}_p$. The GLS choice $\mathbf{A}_* = \mathbf{V}^{-1} \mathbf{H} (\mathbf{H}^\top \mathbf{V}^{-1} \mathbf{H})^{-1}$ minimizes $\text{Cov}(\mathbf{A}^\top \widehat{\mathbf{B}}) = \mathbf{A}^\top \mathbf{V} \mathbf{A}$ in the Loewner order; hence for any feasible \mathbf{A} ,

$$\mathbf{G}(\mathbf{V}) = \text{Cov}(\widetilde{\beta}_{\text{Ens}}) \preceq \mathbf{A}^\top \mathbf{V} \mathbf{A}.$$

Take the special feasible choice that picks a single component: for fixed j , set $\mathbf{A}_j = \mathbf{e}_j \otimes \mathbf{I}_p$ (so $\mathbf{A}_j^\top \mathbf{H} = \mathbf{I}_p$). Then $\mathbf{A}_j^\top \mathbf{V} \mathbf{A}_j = \Omega_{jj}$ and thus

$$\mathbf{G}(\mathbf{V}) \preceq \Omega_{jj} \quad \forall j \in \{1, \dots, K\}. \quad (8)$$

Taking traces yields the componentwise variance dominance:

$$\text{tr}\{\mathbf{G}(\mathbf{V})\} \leq \min_j \text{tr}(\Omega_{jj}).$$

Step 2: Comparing PW to GLS via a block-diagonal perturbation. Write $\mathbf{V} = \mathbf{D}^{1/2}(\mathbf{I} + \mathbf{E})\mathbf{D}^{1/2}$ with

$$\mathbf{E} := \mathbf{D}^{-1/2}(\mathbf{V} - \mathbf{D})\mathbf{D}^{-1/2},$$

so \mathbf{E} contains only the (scaled) off-diagonal blocks. Assume $\|\mathbf{E}\|_{\text{op}} < 1$, so $(\mathbf{I} + \mathbf{E})^{-1} = \sum_{r \geq 0} (-\mathbf{E})^r$. Then

$$\mathbf{V}^{-1} = \mathbf{D}^{-1/2}(\mathbf{I} + \mathbf{E})^{-1}\mathbf{D}^{-1/2} = \mathbf{D}^{-1} - \mathbf{D}^{-1/2}\mathbf{E}\mathbf{D}^{-1/2} + \mathbf{R},$$

with remainder $\mathbf{R} = \mathbf{D}^{-1/2} \sum_{r \geq 2} (-\mathbf{E})^r \mathbf{D}^{-1/2}$. Define

$$\mathbf{G}(\mathbf{M}) := (\mathbf{H}^\top \mathbf{M}^{-1} \mathbf{H})^{-1} \quad (\mathbf{M} \succ 0).$$

Then

$$\mathbf{G}(\mathbf{V})^{-1} - \mathbf{G}(\mathbf{D})^{-1} = \mathbf{H}^\top (\mathbf{V}^{-1} - \mathbf{D}^{-1}) \mathbf{H} = \mathbf{H}^\top \mathbf{D}^{-1/2} \left(-\mathbf{E} + \sum_{r \geq 2} (-\mathbf{E})^r \right) \mathbf{D}^{-1/2} \mathbf{H}.$$

Using the identity $\mathbf{A}^{-1} - \mathbf{B}^{-1} = \mathbf{A}^{-1}(\mathbf{B} - \mathbf{A})\mathbf{B}^{-1}$ and the positive definiteness of $\mathbf{G}(\cdot)$, one obtains the Loewner bound

$$\mathbf{G}(\mathbf{D}) - \mathbf{G}(\mathbf{V}) \succeq \mathbf{G}(\mathbf{D}) \mathbf{H}^\top \mathbf{D}^{-1/2} \left(\mathbf{E} - \sum_{r \geq 2} (-\mathbf{E})^r \right) \mathbf{D}^{-1/2} \mathbf{H} \mathbf{G}(\mathbf{D}) \succeq \mathbf{0}. \quad (9)$$

Taking traces and the norm bound $\sum_{r \geq 1} \|\mathbf{E}\|_{\text{op}}^r = \|\mathbf{E}\|_{\text{op}} / (1 - \|\mathbf{E}\|_{\text{op}})$,

$$\begin{aligned} \text{tr}\{\mathbf{G}(\mathbf{D})\} &\leq \text{tr}\{\mathbf{G}(\mathbf{V})\} + \frac{1}{1 - \|\mathbf{E}\|_{\text{op}}} \text{tr} \left[\mathbf{G}(\mathbf{D}) \mathbf{H}^\top \mathbf{D}^{-1} (\mathbf{V} - \mathbf{D}) \mathbf{D}^{-1} \mathbf{H} \mathbf{G}(\mathbf{D}) \right] \\ &=: \text{tr}\{\mathbf{G}(\mathbf{V})\} + \Delta_{\text{cross}}. \end{aligned} \quad (10)$$

Thus the variance of the block-diagonal precision estimator (with true \mathbf{D}) is within Δ_{cross} of GLS. By Slutsky and continuity of the map $(\cdot) \mapsto (\mathbf{H}^\top (\cdot)^{-1} \mathbf{H})^{-1}$ on \mathbb{S}_{++} , replacing \mathbf{D} by $\hat{\mathbf{D}} \rightarrow_p \mathbf{D}$ yields the same relation with an $o(1)$ term.

Step 3: Assemble the bound. By Step 1 and Eq. 10,

$$\text{tr}\{\text{Cov}(\hat{\beta}_{\text{PW}})\} \leq \text{tr}\{\mathbf{G}(\mathbf{D})\} + o(1) \leq \text{tr}\{\mathbf{G}(\mathbf{V})\} + \Delta_{\text{cross}} + o(1) \leq \min_j \text{tr}(\boldsymbol{\Omega}_{jj}) + \Delta_{\text{cross}} + o(1).$$

Adding the common squared bias $\|\mathbf{b}\|_2^2$ to both sides yields

$$\mathbb{E}\|\hat{\beta}_{\text{PW}} - \beta^*\|_2^2 \leq \min_{j \in \mathcal{K}} \mathbb{E}\|\hat{\beta}_{\text{EB},j} - \beta^*\|_2^2 + \Delta_{\text{cross}} + o(1),$$

which is Theorem 5.

Explicit form of the penalty. From Eq. 10, one convenient expression is

$$\Delta_{\text{cross}} \leq \frac{1}{1 - \|\mathbf{E}\|_{\text{op}}} \text{tr} \left[(\mathbf{H}^\top \mathbf{D}^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{D}^{-1} (\mathbf{V} - \mathbf{D}) \mathbf{D}^{-1} \mathbf{H} (\mathbf{H}^\top \mathbf{D}^{-1} \mathbf{H})^{-1} \right],$$

with $\mathbf{E} = \mathbf{D}^{-1/2}(\mathbf{V} - \mathbf{D})\mathbf{D}^{-1/2}$. This depends only on the off-diagonal cross-covariances $\{\boldsymbol{\Omega}_{k\ell}\}_{k \neq \ell}$. \square