

# Distributionally Robust Optimization via Targeted Integral Probability Metrics for General Data Processes

Lanran Fang\*, Jianqiang Cheng†, Grani A. Hanasusanto\*, and Yijie Wang‡

Distributionally robust optimization (DRO) has been successful in addressing decision-making problems under uncertainty when the underlying distribution is unknown. Existing data-driven DRO frameworks, however, often impose restrictive assumptions on the data-generating process. We propose a new DRO framework based on targeted integral probability metrics. The ambiguity set is defined directly through the loss functions induced by feasible decisions, leading to an expected hinge-constrained formulation that is equivalent to an infinitely constrained ambiguity set. This targeted construction aligns the discrepancy measure with the downstream task and yields finite-sample guarantees that bypass the curse of dimensionality: whenever a scalar pointwise concentration inequality is available, the ambiguity radius can be calibrated at the canonical  $\tilde{O}(N^{-1/2})$  rate. As a result, the framework applies broadly to settings including heavier-tailed distributions, Markovian data, outlier-corrupted observations, incomplete data, and contextual optimization. We derive exact infinite-dimensional dual reformulations, establish out-of-sample and excess-risk guarantees, and develop a Monte Carlo approximation scheme that yields conservative sampled ambiguity sets together with convergence and suboptimality guarantees. For piecewise affine losses, the sampled problems admit tractable conic reformulations, and the Monte Carlo approximation converges at a provably fast rate. Numerical experiments in inventory management under heavy-tailed demand and regression with outlier corruption demonstrate the superiority of our framework over state-of-the-art approaches.

*Key words:* distributionally robust optimization, integral probability metrics, curse of dimensionality, heavy-tailed distributions, Markovian data, corrupted and missing data, contextual optimization

---

## 1. Introduction

Distributionally Robust Optimization (DRO) provides a principled mathematical framework for optimization under uncertainty. Rather than placing blind faith in a single nominal distribution (e.g., empirical distribution), DRO identifies a decision that minimizes the worst-case expected loss over an ambiguity set  $\mathcal{P}$ , which is defined as a family of probability measures that are statistically

\* Department of Industrial and Enterprise Systems Engineering, University of Illinois Urbana-Champaign, Urbana, IL 61801, USA. Email: lanranf2, gah@illinois.edu.

† Department of Systems and Industrial Engineering, University of Arizona, Tucson, AZ 85721. Email: jqcheng@arizona.edu.

‡ School of Economics and Management, Tongji University, Shanghai, China. Email: yijiewang@tongji.edu.cn.

plausible given the available data. This robust approach can be formulated as the following min-max problem:

$$\min_{\mathbf{x} \in \mathcal{X}} \sup_{\mathbb{Q} \in \mathcal{P}} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})], \quad (1)$$

where  $\ell(\mathbf{x}, \boldsymbol{\xi})$  represents the loss function parameterized by the decision vector  $\mathbf{x} \in \mathcal{X}$  and the uncertain parameter  $\boldsymbol{\xi} \in \Xi$ . By optimizing over the most adverse distribution  $\mathbb{Q}$  within  $\mathcal{P}$ , the DRO framework yields decisions that remain reliable even when deployed in unseen environments.

Early literature predominantly constructed these ambiguity sets using moment bounds or  $\phi$ -divergence, which, despite computational advantages, exhibited irreconcilable statistical and topological flaws. Moment ambiguity sets constrain the family of probability distributions via a finite number of generalized moment conditions. A fundamental drawback of moment ambiguity sets is the lack of asymptotic consistency: gathering infinite data will not cause the DRO solution to converge to the true stochastic programming optimal solution, as low-order moments cannot uniquely identify the true distribution. To overcome the failure of moment sets to converge to the true distribution, Ben-Tal et al. (2012) introduce discrepancy-based ambiguity sets. These sets center around the empirical distribution  $\hat{\mathbb{P}}_N$  and use  $\phi$ -divergences to define a neighborhood, capturing all distributions that are statistically close to the empirical one. While  $\phi$ -divergence sets exhibit good statistical consistency, they have severe topological limitations when applied to continuous support spaces. Specifically, if the adversarial distribution  $\mathbb{Q}$  assigns positive probability to any event that has zero probability under the reference distribution, the divergence instantly becomes infinite. This prohibits the  $\phi$ -divergence to provide robustness against unseen parameter realizations, severely crippling its out-of-sample generalization capability.

Recently, the Wasserstein distance has become the dominant metric for constructing ambiguity sets due to its nice geometric interpretation and asymptotic consistency. The  $p$ -Wasserstein distance between two probability distributions  $\mathbb{Q}$  and  $\mathbb{P}$  is defined as the  $p$ -th root of the minimum expected cost of 'transporting' one measure to match the other. Unlike  $\phi$ -divergences, even if two distributions have mutually exclusive support sets (e.g., one discrete, one continuous), the Wasserstein distance between them remains finite. Therefore, the Wasserstein ambiguity sets enable to smoothly transport probability mass from historical observations to adjacent unobserved regions in the support set  $\Xi$ , enabling genuine out-of-sample robustness.

Despite its conceptual elegance, Wasserstein DRO is hindered by the curse of dimensionality. To offer finite-sample guarantees, the radius of the Wasserstein ambiguity set must be calibrated so that the unknown true distribution  $\mathbb{P}^*$  is contained within the Wasserstein ball with high probability. Using the concentration of empirical Wasserstein distance theory derived by Fournier and Guillin (2015), Mohajerin Esfahani and Kuhn (2018) show that the empirical measure converges

to the true measure in the Wasserstein metric at a slow rate. For a random variable residing in a  $d$ -dimensional space, the concentration inequality yields a radius that scales as  $\mathcal{O}(N^{-1/\max\{2,d\}})$ . This indicates that the decay rate of the radius is fundamentally bottlenecked by the dimension of data. For high dimensional problems, achieving a tight confidence region requires an exponentially large and practically unattainable sample size.

To bypass the curse of dimensionality inherent in Wasserstein DRO, a stream of studies has proposed a variety of advanced methodologies. Blanchet et al. (2019), Blanchet and Kang (2021), Blanchet et al. (2022) propose a novel approach where the Wasserstein ambiguity set is not required to contain the true distribution  $\mathbb{P}^*$ , but rather at least one distribution that shares the true optimal solution. This targeted coverage enables the optimal Wasserstein radius to achieve an asymptotic decay rate of  $\mathcal{O}(1/\sqrt{N})$ . Shafieezadeh-Abadeh et al. (2019) show that 1-Wasserstein DRO applied to specific linear models (e.g., support vector machine and logistic regression) reduces to empirical risk minimization with an additive dual norm regularization penalty. Because the infinite-dimensional supremum over the Wasserstein ball analytically collapses into a simple parametric penalty, the convergence of the learning algorithm no longer depends on matching probability measures in high dimensions, allowing the radius to shrink at a non-asymptotic rate of  $\mathcal{O}(1/\sqrt{N})$ . Concurrently, Wu et al. (2025) extend this regularization equivalence for general risk measures and Wasserstein balls.

A crucial step toward non-asymptotic frameworks for Wasserstein DRO was achieved by Gao (2023), which offers the first finite-sample guarantees that avoid the curse of dimensionality for generic loss functions. Gao (2023) assumes that the underlying data-generating distribution satisfies a Talagrand transportation-information inequality (Villani et al. 2008, Gozlan and Léonard 2010). Through this lens, he establishes a novel variation-based concentration result, showing that the decay rate of the tail probability is controlled by the variation of the loss rather than the ambient dimension. By combining this with localized Rademacher complexity theory, the framework elegantly bridges Wasserstein DRO with variation regularization, yielding an optimal  $\mathcal{O}(1/\sqrt{N})$  generalization bound for Wasserstein robust learning. Gao et al. (2024) further develop a unified variation regularization theory that establishes non-asymptotic performance guarantees for non-convex and non-smooth losses. Additionally, Azizian et al. (2023) derive exact generalization guarantees for broad model classes in both standard and regularized Wasserstein DRO. Le and Mallick (2024) extend these exact guarantees to highly complex and non-smooth parametric families, covering arbitrary transport costs and modern deep learning objectives.

While state-of-the-art Wasserstein frameworks achieve optimal non-asymptotic rates, their reliance on complex analysis and certain assumptions bottlenecks extensions to a wide range of non-standard data regimes. To overcome these theoretical hurdles, this paper proposes a fundamentally different approach to constructing the ambiguity set. We situate our framework within the broader

class of Integral Probability Metrics (IPMs), following the taxonomy established by Kuhn et al. (2025). Formally, given a family of integrable functions  $f \in \mathcal{F}$ , an IPM measures the discrepancy between two distributions  $\mathbb{P}, \hat{\mathbb{P}} \in \mathcal{P}(\Xi)$  as

$$D_{\mathcal{F}}(\mathbb{P}, \hat{\mathbb{P}}) = \sup_{f \in \mathcal{F}} \left| \int_{\mathcal{Z}} f(z) d\mathbb{P}(z) - \int_{\mathcal{Z}} f(\hat{z}) d\hat{\mathbb{P}}(\hat{z}) \right|.$$

The IPM is a highly expressive topological framework. Depending on the specific choice of the function class  $\mathcal{F}$ , it can recover several foundational statistical distances. For instance, when  $\mathcal{F}$  consists of all indicator functions of Borel sets, then  $D_{\mathcal{F}}$  recovers the total variation distance. If  $\mathcal{F}$  represents the set of all 1-Lipschitz continuous functions, then IPM recovers the 1-Wasserstein distance via the Kantorovich-Rubinstein duality. Finally, if  $\mathcal{F}$  is defined as the unit ball in a Reproducing Kernel Hilbert Space (RKHS),  $D_{\mathcal{F}}$  yields the maximum mean discrepancy (MMD) distance, which is widely adopted in kernel methods.

Leveraging the IPM philosophy, we propose a generalized, hinge-constrained ambiguity set  $\mathcal{P}_{\epsilon}$  designed specifically to align with the optimization objective:

$$\mathcal{P}_{\epsilon} := \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : \begin{array}{l} \mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \right] \leq 0 \\ \mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega \end{array} \right\}. \quad (2)$$

Here,  $\epsilon$  serves as a tunable size parameter of the ambiguity set, and  $\mathbb{Z}$  is an auxiliary probability distribution with full support across the entire decision space  $\mathcal{X}$ . The reference distribution  $\hat{\mathbb{P}}_N$  is constructed from available historical data, and the second-moment bound ensures basic statistical regularity. The hinge constraint forces the inner term to be non-positive almost everywhere with respect to  $\mathbb{Z}$ . Since  $\mathbb{Z}$  has a strictly positive density on  $\mathcal{X}$ , this expected hinge constraint is equivalent to enforcing a semi-infinite constraint across the entirety of the decision space.

**PROPOSITION 1 (Informal).** *Under mild conditions, the ambiguity set (2) is equivalent to:*

$$\mathcal{P}_{\epsilon} := \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : \begin{array}{l} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq \epsilon \quad \forall \mathbf{z} \in \mathcal{X} \\ \mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega \end{array} \right\}.$$

This reformulation reveals the fundamental nature of our proposed set: It is a *targeted* Integral Probability Metric (TIPM) ambiguity set. Our TIPM departs from classical IPM constructions in two critical ways. First, by dropping the absolute value typically required to ensure a symmetric distance metric, our TIPM fundamentally shifts from a generic “statistical distance” to a “risk-aware discrepancy.” Because the inner supremum of the DRO problem seeks adversarial distributions that increase the expected loss, symmetric penalization is mathematically redundant. By employing a one-sided upper bound, we allow the ambiguity set to be more flexible: accommodating optimistic distributions while tightly constraining the exact pessimistic deviations that matter for

out-of-sample robustness. Second, unlike classical IPM ambiguity sets that define the test function class  $\mathcal{F}$  using arbitrary Lipschitz functions or Reproducing Kernel Hilbert Space (RKHS) balls, our TIPM strictly defines  $\mathcal{F}$  as the hypothesis class of the loss functions themselves, parameterized by the feasible decisions  $\mathbf{z} \in \mathcal{X}$ . By linking the discrepancy metric directly to the loss function of the specific optimization problem at hand, we create an ambiguity set that strictly monitors the discrepancies that actually matter to the decision-maker, ignoring irrelevant noise. As such, our proposed framework offers the following contributions:

- **A task-aware ambiguity set with exact characterization:** We introduce a new expected hinge-constrained ambiguity set based on targeted integral probability metrics (TIPMs), in which the discrepancy is defined directly through the loss functions induced by feasible decisions. We show that this ambiguity set is equivalent to an infinitely constrained loss-discrepancy ambiguity set, thereby aligning robustness with the downstream decision problem. The resulting ambiguity set is distributional in nature but task-aware in geometry. We further derive an exact infinite-dimensional dual reformulation of the resulting DRO problem.
- **Curse of dimensionality-free statistical guarantees under broad data regimes:** The TIPM construction yields finite-sample guarantees that bypass the curse of dimensionality. In particular, whenever a scalar pointwise concentration inequality is available, the ambiguity radius can be calibrated at the canonical  $\tilde{\mathcal{O}}(1/\sqrt{N})$  rate, without requiring stringent assumptions such as Talagrand-type transportation inequalities or problem-specific analyses. This allows the framework to accommodate a broad range of settings, including sub-Weibull losses, outlier-corrupted data, Markovian sequences, and incomplete data. Moreover, because the ambiguity set is defined directly in terms of the loss functions, it is agnostic to the representation of the underlying uncertainty and naturally handles settings with mixed discrete and continuous components.
- **A conservative Monte Carlo approximation with convergence and tractability guarantees:** To overcome the intractability caused by the continuous auxiliary distribution, we introduce a Monte Carlo approximation that replaces the ambiguity set with a conservative sampled counterpart. For piecewise affine losses, we establish convergence of optimal values and solutions together with a  $\tilde{\mathcal{O}}(1/\sqrt{M})$  sampling rate and suboptimality guarantees, and show that the resulting sampled problems admit tractable conic reformulations, including second-order cone approximations for two-stage linear losses.
- **Extension to unbounded support:** We extend the sampling and reformulation framework to unbounded support under a single-threshold tail bound implied by sub-Weibull tails. In this setting, we show that the same Monte Carlo principle remains valid up to explicit tail-dependent terms, while preserving tractability for piecewise affine losses.

### 1.1. Related Literature

**IPM ambiguity sets:** Among the various methods for constructing ambiguity sets, IPMs have gained increasing attention due to their rich topological properties and ability to unify classical statistical distances (Husain 2020, Birrell et al. 2022). By measuring the discrepancy using a specific class of test functions, the IPM framework flexibly recovers the Total Variation distance, the Wasserstein distance (Mohajerin Esfahani and Kuhn 2018, Blanchet and Murthy 2019, Gao and Kleywegt 2023), and the Maximum Mean Discrepancy distance (Zhu et al. 2021). Consequently, IPM-based and specifically MMD-based DRO formulations have been extensively studied for their computational tractability, generalization bounds, and deep theoretical connections to regularized empirical risk minimization (Shafieezadeh-Abadeh et al. 2019, Duchi and Namkoong 2019, Husain 2020, Zeng and Lam 2022, Iyengar et al. 2023). Building upon this framework, our work departs from generically defined test spaces by restricting the  $\mathcal{F}$  strictly to the hypothesis class of the loss functions induced by the feasible decisions  $\mathbf{z} \in \mathcal{X}$ . Benefiting from this targeted construction, the TIPM ambiguity sets achieve better structural compatibility for the downstream decision problem.

**Generalization bounds:** A central task in DRO is establishing generalization bounds that guarantee reliable out-of-sample performance. Early formulations primarily relied on moment-based ambiguity sets (Delage and Ye 2010, Goh and Sim 2010, Wiesemann et al. 2014) and  $\phi$ -divergence-based sets (Ben-Tal et al. 2013, Bertsimas et al. 2018). These approaches benefit from standard concentration inequalities, yielding favorable, dimension-independent  $\mathcal{O}(1/\sqrt{N})$  radius decay rate. Recently, Wasserstein DRO has become increasingly prominent due to its nice topological properties (Mohajerin Esfahani and Kuhn 2018, Blanchet and Murthy 2019, Gao and Kleywegt 2023). However, this topological richness comes at a statistical cost: standard concentration-of-measure results (Fournier and Guillin 2015) dictate that the convergence is at a slow rate  $\mathcal{O}(N^{-1/\max\{2,d\}})$ . To resolve this bottleneck, a line of research has developed novel frameworks to break the curse of dimensionality for Wasserstein DRO (Chen and Paschalidis 2018, Shafieezadeh-Abadeh et al. 2019, Blanchet et al. 2019, Blanchet and Kang 2021, Blanchet et al. 2022, Gao 2023, Gao et al. 2024, Wu et al. 2025), as detailed in our introduction. Besides, other studies have established generalization bounds for Wasserstein DRO using sample-size-independent fixed radii (Lee and Raginsky 2018, Sinha et al. 2017, Najafi et al. 2019), while parallel efforts have explored optimal radius calibration and bias-variance trade-offs within divergence-based frameworks (Duchi and Namkoong 2019, Lam 2019, Duchi et al. 2021). Furthermore, advanced mathematical tools, such as transportation-information inequalities, have also been leveraged to derive generalization bounds for general learning algorithms (Xu and Raginsky 2017, Lopez and Jog 2018, Russo and Zou 2019, Wang et al. 2019, 2021, Xie et al. 2022, He and Goldfeld 2025).

**Learning with non-structural data:** In practice, decision makers frequently operate in environments where the available data violate the standard independent and identically distributed (i.i.d.) assumption. For instance, time-series and sequential decision problems often involve Markovian data structures (Fan et al. 2021, Li et al. 2021, Nagaraj et al. 2020). Similarly, in selection problems such as hiring and college admissions, outcomes are only observable for accepted candidates, leading to incomplete data regimes (Najafi et al. 2019). Additionally, modern decision-making increasingly relies on auxiliary side information, necessitating optimization over conditional distributions in contextual data settings (Sadana et al. 2025). Furthermore, data collected in real-world environments are often contaminated by noise or adversarial perturbations, resulting in outlier-corrupted datasets (Nietert et al. 2023, 2024, Blanchet et al. 2024, Jiang and Xie 2024). While previous DRO models have attempted to tackle these complex data realities, existing solutions are predominantly problem-specific. They typically rely on specialized ambiguity set constructions or certain assumptions tailored to a single type of data anomaly. Our proposed framework bridges this gap. Benefited from our TIPM structure, the optimal radius convergence rate can be achieved immediately from any scalar, pointwise concentration inequality. This unified approach accommodates a remarkably broad spectrum of non-standard data regimes, offering a concise and unified DRO framework.

## 1.2. Paper Structure and Notation

**Paper structure.** The remainder of the paper is organized as follows. Section 2 formally introduces the expected hinge-constrained TIPM ambiguity set, establishes its equivalent semi-infinite loss-discrepancy representation, and derives coverage, out-of-sample, and excess-risk guarantees. It also shows how the same framework applies to several data-generating regimes through pointwise concentration inequalities, including sub-Weibull losses, Markovian samples, outlier-contaminated observations, incomplete data, and contextual information, and develops exact dual reformulations. Section 3 studies a Monte Carlo approximation of the ambiguity set, proves its conservative containment and convergence properties, and derives finite conic reformulations for piecewise affine losses. Section 4 extends these sampling and reformulation results to unbounded support under a single-threshold tail bound implied by sub-Weibull tails. Section 5 evaluates the proposed framework in multi-item newsvendor and outlier-corrupted regression experiments. The appendices collect all proofs, a two-stage linear decision-rule approximation, and detailed experimental specifications.

**Notation.** For a positive integer  $n$ , we write  $[n] := \{1, \dots, n\}$ . Random variables are designated by tilde signs (e.g.,  $\tilde{\xi}$ ), while their realizations are denoted by the same symbols without tildes (e.g.,  $\xi$ ). The true data-generating distribution is denoted by  $\mathbb{P}^*$ , the data-driven reference distribution by  $\hat{\mathbb{P}}_N$ , a generic distribution in the ambiguity set by  $\mathbb{Q}$ , and the sampling distribution over decisions by  $\mathbb{Z}$ . We use  $\mathcal{P}(\Xi)$  for the set of probability distributions on  $\Xi$ ,  $\mathbb{E}_{\mathbb{P}}[\cdot]$  and  $\text{Var}_{\mathbb{P}}[\cdot]$  for expectation

and variance under  $\mathbb{P}$ , and  $\mathbb{I}_A$  or  $\mathbb{I}[A]$  for the indicator of an event  $A$ . The feasible decision set is  $\mathcal{X} \subseteq \mathbb{R}^{D_x}$ , the uncertainty support is  $\Xi \subseteq \mathbb{R}^{D_\xi}$ , and  $R_x$  and  $R_\xi$  denote corresponding norm bounds when they are finite. The ambiguity radius is denoted by  $\epsilon$ , while  $\varepsilon(\cdot)$  denotes the nondecreasing function appearing in pointwise concentration bounds. The relaxation level is  $\eta$ , the second-moment budget is  $\Omega$ , the data and Monte Carlo sample sizes are  $N$  and  $M$  respectively, and  $\delta$  and  $\tau$  denote confidence parameters. We write  $\Delta^J$  for the probability simplex in  $\mathbb{R}_+^J$ ,  $\mathbf{e}$  for the all-ones vector, and  $\mathbf{e}_d$  for the  $d$ th unit vector. For a scalar  $u$ ,  $[u]_+ := \max\{u, 0\}$ .

The Euclidean norm is denoted by  $\|\cdot\|$  and the operator norm of a matrix by  $\|\cdot\|_{\text{op}}$ . For a real-valued random variable  $X$  and  $\vartheta > 0$ , we denote its Orlicz norm by

$$\|X\|_{\psi_\vartheta} := \inf \left\{ c > 0 : \mathbb{E} \left[ \exp \left( (|X|/c)^\vartheta \right) - 1 \right] \leq 1 \right\}.$$

For a set  $\Xi$ , its support function is  $\sigma_\Xi(\mathbf{z}) := \sup_{\xi \in \Xi} \mathbf{z}^\top \xi$ . We denote by  $\mathcal{C}(\mathcal{X})$  the space of continuous functions on  $\mathcal{X}$ , by  $\mathcal{C}_+(\mathcal{X}) := \{f \in \mathcal{C}(\mathcal{X}) : f(\mathbf{x}) \geq 0 \forall \mathbf{x} \in \mathcal{X}\}$  the cone of nonnegative continuous functions, and by  $\mathcal{M}_+(\mathcal{X})$  the cone of finite nonnegative Borel measures on  $\mathcal{X}$ . For a point  $\mathbf{x}$  and a set  $\mathcal{S}$ ,  $\text{dist}(\mathbf{x}, \mathcal{S})$  denotes the distance from  $\mathbf{x}$  to  $\mathcal{S}$ , and  $\xrightarrow{P}$  denotes convergence in probability. We use  $\mathcal{O}(\cdot)$  for upper bounds up to absolute constants, and  $\tilde{\mathcal{O}}(\cdot)$  when logarithmic factors are additionally suppressed. Model-specific notation, such as Markov spectral gaps, MoM block parameters, propensity scores, kernel bandwidths, and tail constants, is introduced locally where needed.

## 2. Distributionally Robust Optimization with Expected Hinge Constrained Ambiguity Sets

We consider the distributionally robust optimization problem

$$\min_{\mathbf{x} \in \mathcal{X}} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\xi})], \quad (3)$$

where the ambiguity set is given by

$$\mathcal{P}_\epsilon := \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : \begin{array}{l} \mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\tilde{\mathbf{z}}, \tilde{\xi})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\xi})] - \epsilon \right]_+ \right] \leq 0 \\ \mathbb{E}_{\mathbb{Q}}[\|\tilde{\xi}\|^2] \leq \Omega \end{array} \right\}, \quad (4)$$

and the feasible set  $\mathcal{X}$  is compact and convex. Here,  $\epsilon \geq 0$  is a size parameter and  $\mathbb{Z}$  is any probability distribution with full support on  $\mathcal{X}$ . In this data-driven ambiguity set, the reference distribution  $\hat{\mathbb{P}}_N$  is constructed from the available data, whose description depends on the specific applications. As we detail later, the hinge constraint ensures that the distribution  $\mathbb{Q}$  is close to  $\hat{\mathbb{P}}_N$  in terms of the expected loss. The second-moment constraint ensures regularity of the distributions in the ambiguity set, enabling rigorous theoretical development. This ambiguity set is designed to achieve both high-probability coverage of the underlying distribution  $\mathbb{P}^*$  and asymptotic convergence at a fast rate, thereby overcoming the curse of dimensionality of transportation-based ambiguity sets.

We consider a piecewise affine loss function of the form

$$\ell(\mathbf{x}, \boldsymbol{\xi}) := \max_{j \in [J]} (\mathbf{a}_j(\mathbf{x})^\top \boldsymbol{\xi} + b_j(\mathbf{x})), \quad (5)$$

where  $\mathbf{a}_j(\mathbf{x})$  and  $b_j(\mathbf{x})$ ,  $j \in [J]$ , are affine functions of  $\mathbf{x}$  given by

$$\mathbf{a}_j(\mathbf{x}) := \mathbf{A}_j \mathbf{x} + \bar{\mathbf{a}}_j, \quad b_j(\mathbf{x}) := \mathbf{b}_j^\top \mathbf{x} + \bar{b}_j.$$

This class is broad enough to capture many models arising in operations research and management science, and is standard in the distributionally robust optimization literature (Wiesemann et al. 2014, Mohajerin Esfahani and Kuhn 2018) because it often admits tractable conic reformulations.

It also includes two-stage linear loss functions as a special case:

$$\begin{aligned} \ell(\mathbf{x}, \boldsymbol{\xi}) &:= \mathbf{c}^\top \mathbf{x} + \inf \mathbf{q}^\top \mathbf{y} \\ \text{s.t. } &\mathbf{y} \in \mathbb{R}^{D_y} \\ &\mathbf{T}(\mathbf{x})\boldsymbol{\xi} + \mathbf{h}(\mathbf{x}) \leq \mathbf{W}\mathbf{y}, \end{aligned} \quad (6)$$

where  $\mathbf{T}(\mathbf{x})$  and  $\mathbf{h}(\mathbf{x})$  are a matrix and a vector that depend affinely on  $\mathbf{x}$ , and  $\mathbf{W}$  is a fixed recourse matrix. In this context,  $\mathbf{x}$  is a here-and-now decision made before the realization of uncertainty, while  $\mathbf{y}$  is a wait-and-see recourse decision that can be chosen after the uncertainty is realized.

The linear program in (6) can be dualized, yielding the representation

$$\begin{aligned} \ell(\mathbf{x}, \boldsymbol{\xi}) &= \mathbf{c}^\top \mathbf{x} + \sup (\mathbf{T}(\mathbf{x})\boldsymbol{\xi} + \mathbf{h}(\mathbf{x}))^\top \mathbf{p} \\ \text{s.t. } &\mathbf{p} \in \mathbb{R}_+^L \\ &\mathbf{q} = \mathbf{W}^\top \mathbf{p}. \end{aligned} \quad (7)$$

Strong linear programming duality holds under the standard *relatively complete recourse* assumption that the primal problem is feasible for any  $\mathbf{x} \in \mathcal{X}$  and  $\boldsymbol{\xi} \in \Xi$ . Since the feasible set of the dual problem is a polyhedron with finitely many extreme points  $\{\mathbf{p}_1, \dots, \mathbf{p}_Q\}$ , the loss function can be rewritten in the piecewise affine form (5), where  $J = Q$ ,  $\mathbf{a}_j(\mathbf{x}) = \mathbf{T}(\mathbf{x})^\top \mathbf{p}_j$ , and  $b_j(\mathbf{x}) = \mathbf{h}(\mathbf{x})^\top \mathbf{p}_j$ .

We make the following assumptions:

(A) *Linear growth in  $\boldsymbol{\xi}$* : There exist  $c_1, c_2 \in \mathbb{R}_+$  such that

$$|\ell(\mathbf{x}, \boldsymbol{\xi})| \leq c_1 + c_2 \|\boldsymbol{\xi}\| \quad \forall \mathbf{x} \in \mathcal{X} \forall \boldsymbol{\xi} \in \Xi.$$

(B) *Strict second-moment feasibility of  $\hat{\mathbb{P}}_N$* : We assume that  $\mathbb{E}_{\hat{\mathbb{P}}_N} [\|\tilde{\boldsymbol{\xi}}\|^2] < \Omega$ .

The first assumption is trivially satisfied under the piecewise linear loss function (5). The second assumption is needed to ensure strong duality. It can be relaxed to a high-confidence requirement, so that any probabilistic guarantees in the paper hold under the event that the strict feasibility requirement holds. For clarity of exposition, we do not pursue this extension.

We first derive a fundamental property of this ambiguity set that will be useful for developing the applications and the solution schemes.

PROPOSITION 2. *The ambiguity set (4) is equivalent to the infinitely-constrained ambiguity set:*

$$\mathcal{P}_\epsilon := \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : \begin{array}{l} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq \epsilon \quad \forall \mathbf{z} \in \mathcal{X} \\ \mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega \end{array} \right\}. \quad (8)$$

Proposition 2 clarifies the role of the expected hinge constraint in (4). The ambiguity set  $\mathcal{P}_\epsilon$  can be viewed as a feasibility region defined by an infinite family of loss-discrepancy inequalities together with a moment constraint. This semi-infinite perspective highlights a key advantage of our targeted formulation: rather than controlling a problem-agnostic statistical distance, it restricts discrepancies through the lens of the specific loss function. By enforcing  $\mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq \epsilon$  for every feasible decision  $\forall \mathbf{z} \in \mathcal{X}$ , the formulation rules out adversarial distributions that could induce unexpectedly large losses for any feasible decision. Thus, the model does not require the adversarial distribution to remain close to the empirical distribution in all statistical directions, but only in those directions that are operationally relevant.

REMARK 1 (TASK-RELEVANT DISCREPANCY). Unlike classical ambiguity sets that rely on an a priori ambient-space ground metric, our targeted formulation captures task relevance natively. This structural advantage becomes particularly evident under heterogeneous uncertainty, such as when  $\tilde{\boldsymbol{\xi}}$  contains both continuous and discrete components, or when the loss is much more sensitive to some components than to others. A Wasserstein formulation can also reflect such structure, but only after specifying a suitable representation and ground cost. Our formulation instead measures the induced expected-loss deviations directly.

## 2.1. Performance Guarantees

In the next two subsections, we show that our DRO framework can be applied to a broad spectrum of problem classes, and establish that the radius  $\epsilon$  can be systematically adjusted with the number of samples  $N$  at the canonical  $\tilde{\mathcal{O}}(1/\sqrt{N})$  rate. The proofs of all subsequent coverage guarantees share the same two-step structure: (i) a *pointwise* concentration inequality for a fixed decision  $\mathbf{z} \in \mathcal{X}$ , and (ii) an extension to a *uniform* bound over all  $\mathbf{z} \in \mathcal{X}$  via a covering number argument. We capture this common structure in the following proposition, so that each specific application needs only establish the pointwise bound in step (i).

PROPOSITION 3. *Let  $\mathcal{X} \subseteq \mathbb{R}^{D_x}$  be compact with  $\sup_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\| \leq R_x$ . Assume that  $\mathbb{E}_{\mathbb{P}^*}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega$ . Suppose for every fixed  $\mathbf{z} \in \mathcal{X}$ , we have the pointwise concentration inequality*

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq \frac{1}{\sqrt{N}} \epsilon \left( \log \left( \frac{1}{\delta} \right) \right) \quad (9)$$

with probability at least  $1 - \delta$ , where  $\varepsilon : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  is nondecreasing and has at most polynomial growth. Let  $L := \max_{j \in [J]} \|\mathbf{A}_j\|_{\text{op}} \sqrt{\Omega} + \max_{j \in [J]} \|\mathbf{b}_j\|$  be the Lipschitz constant from Lemma 1. Then, by setting the radius to

$$\epsilon := \frac{1}{\sqrt{N}} \left( \varepsilon \left( D_{\mathbf{x}} \log \left( 1 + 2R_{\mathbf{x}} L \sqrt{N} \right) + \log \left( \frac{1}{\delta} \right) \right) + 2 \right), \quad (10)$$

one can ensure the distribution coverage

$$\mathbb{P}^* \in \mathcal{P}_\epsilon$$

with probability at least  $1 - \delta$ .

The proof combines the Lipschitz estimate in Lemma 1 with a covering argument; see Appendix A. The Lipschitz estimate controls nearby decisions, while monotonicity of  $\varepsilon$  lets the pointwise bound be lifted over the cover. Polynomial growth then keeps the extra covering penalty logarithmic in  $N$  and avoids a curse-of-dimensionality rate loss.

Proposition 3 provides the link between pointwise concentration and the applications below. Once (9) is established for any fixed decision  $\mathbf{z}$ , the covering argument in the proposition yields the uniform coverage condition  $\mathbb{P}^* \in \mathcal{P}_\epsilon$ . Thus, in the canonical-rate cases, each application only needs to specify the reference estimator  $\hat{\mathbb{P}}_N$  and the corresponding function  $\varepsilon(\cdot)$  in (9); the radius is then obtained from (10). The same pointwise-to-uniform argument applies to estimators with different rates, with the rate reflected in the radius calibration.

In the following, we derive out-of-sample performance and suboptimality guarantees for the solutions to the DRO problem (3). To this end, we denote  $\mathbf{x}^*$  and  $\hat{\mathbf{x}}$  as the optimal solutions of the true stochastic optimization problem and the DRO problem, respectively, i.e.,

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathcal{X}} \mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] \quad \text{and} \quad \hat{\mathbf{x}} \in \arg \min_{\mathbf{x} \in \mathcal{X}} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})]. \quad (11)$$

Note that the solution  $\mathbf{x}^*$  is attained since the mapping  $\mathbf{x} \rightarrow \mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})]$  is continuous, and the feasible set  $\mathcal{X}$  is compact. Furthermore, the objective function of the DRO problem constitutes a pointwise supremum of a family of lower semicontinuous (piecewise linear) functions. Therefore, it is also lower semicontinuous, and the set of optimal solutions is nonempty. That is,  $\hat{\mathbf{x}}$  is attained.

**COROLLARY 1.** *Let  $\hat{J}$  be the optimal value of the DRO problem (3). Then, setting the radius to (10), we can ensure the out-of-sample performance guarantee*

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})] \leq \hat{J}$$

with probability at least  $1 - \delta$ .

Remarkably, our DRO framework enables us to establish a guarantee on the excess risk of the distributionally robust solution, as stated in the following theorem.

THEOREM 1. *Suppose the pointwise concentration inequality in Proposition 3 holds in terms of absolute error, i.e.,*

$$\left| \mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \right| \leq \frac{1}{\sqrt{N}} \varepsilon \left( \log \left( \frac{1}{\delta} \right) \right) \quad (12)$$

with probability  $1 - \delta$  for any fixed  $\mathbf{z} \in \mathcal{X}$ . Then, setting the radius  $\epsilon$  to (10), we get

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})] \leq \mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{x}^*, \tilde{\boldsymbol{\xi}})] + \frac{1}{\sqrt{N}} \varepsilon \left( \log \left( \frac{1}{\delta} \right) \right) + \epsilon,$$

with probability  $1 - 2\delta$ .

This theorem shows that the excess risk is controlled by a single pointwise estimation term together with the ambiguity radius  $\epsilon$ . Therefore, whenever the radius calibration in Proposition 3 is of order  $\tilde{\mathcal{O}}(N^{-1/2})$ , the same is true for the excess risk. In this sense, the distributionally robust solution achieves guaranteed robustness while paying only a near-parametric excess-risk price, without the rate deteriorating in the ambient dimension  $D_{\boldsymbol{\xi}}$  of the uncertainty  $\tilde{\boldsymbol{\xi}}$ .

REMARK 2 (DECISION-DEPENDENT AMBIGUITY SETS). One could envisage an alternative decision-dependent ambiguity set

$$\mathcal{P}_\epsilon(\mathbf{x}) := \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : \begin{array}{l} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] \leq \epsilon \\ \mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega \end{array} \right\}.$$

Then, for any fixed  $\mathbf{x} \in \mathcal{X}$ , (9) implies the pointwise coverage guarantee

$$\text{Prob}(\mathbb{P}^* \in \mathcal{P}_\epsilon(\mathbf{x})) \geq 1 - \delta$$

However, this pointwise coverage does not generally imply the uniform coverage property

$$\text{Prob}(\mathbb{P}^* \in \mathcal{P}_\epsilon(\mathbf{x}) \forall \mathbf{x} \in \mathcal{X}) \geq 1 - \delta,$$

which is typically required to deduce an out-of-sample guarantee for the data-dependent optimizer  $\hat{\mathbf{x}}$ . Indeed,  $\hat{\mathbf{x}}$  depends on the same sample used to form  $\hat{\mathbb{P}}_N$ , so a guarantee that holds for each fixed  $\mathbf{x}$  need not hold at the random choice  $\mathbf{x} = \hat{\mathbf{x}}$  without a uniform bound.

## 2.2. Applications

We now present several pertinent applications of our DRO framework. Each result below verifies the pointwise concentration condition (9) under a different data regime or reference estimator. The corresponding function  $\varepsilon(\cdot)$  determines the radius in (10), which gives the coverage and out-of-sample guarantees through Proposition 3 and Corollary 1. When a two-sided version of the bound is available, Theorem 1 further yields the excess-risk guarantee.

### Sub-Weibull loss

In many applications such as financial risk, insurance claims, and inventory problems, the random loss  $\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})$  may exhibit substantial skewness and nontrivial tail behavior. At the same time, finite-sample DRO guarantees are often derived under light-tailed assumptions on the underlying uncertainty (Mohajerin Esfahani and Kuhn 2018, Gao 2023). It is therefore important to ask whether the same finite-sample guarantees can still be obtained under weaker tail assumptions.

Our framework extends naturally to losses with sub-Weibull tails, a broad family that unifies several tail regimes. It includes the sub-Gaussian and subexponential cases as special instances. The latter covers many familiar distributions, including the exponential, Gamma,  $\chi^2$ , Poisson, and geometric distributions, as well as the square of any sub-Gaussian random variable (Vershynin 2025, Example 2.8.7)). It also extends to the heavier-than-exponential regime when  $\vartheta \in (0, 1)$ . This allows us to treat light- and moderately heavy-tailed losses within a common framework.

Following Kuchibhotla and Chakraborty (2022), a centered random variable  $X$  is called sub-Weibull of order  $\vartheta \in (0, 2]$  if  $\|X\|_{\psi_\vartheta} < \infty$ . Under this parameterization,  $\vartheta = 2$  corresponds to the sub-Gaussian case,  $\vartheta = 1$  to the subexponential case, and  $0 < \vartheta < 1$  to heavier-than-exponential tails. For  $\vartheta < 1$ , the quantity  $\|\cdot\|_{\psi_\vartheta}$  is only a quasi-norm, which makes sharp Bernstein-type concentration substantially more delicate.

The following proposition provides a pointwise concentration bound for the empirical loss under a sub-Weibull assumption.

PROPOSITION 4. *Suppose that there exist  $\vartheta \in (0, 2]$  and  $K > 0$  such that*

$$\|\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) - \mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})]\|_{\psi_\vartheta} \leq K \quad \forall \mathbf{z} \in \mathcal{X}.$$

*Then there exist constants  $C_\vartheta, c_\vartheta > 0$ , depending only on  $\vartheta$ , such that the pointwise concentration (9) holds with  $\varepsilon(y) = C_\vartheta K \sqrt{1 + y}$ . More precisely,*

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq \frac{C_\vartheta K}{\sqrt{N}} \sqrt{1 + \log\left(\frac{1}{\delta}\right)} \quad (13)$$

*with probability at least  $1 - \delta$ , provided that*

$$\delta \in \begin{cases} [2 \exp(-c_\vartheta N^{\vartheta/(2-\vartheta)}), 1), & \vartheta \in (0, 1), \\ [2 \exp(-c_\vartheta N), 1), & \vartheta \in [1, 2), \\ (0, 1), & \vartheta = 2. \end{cases}$$

REMARK 3 (SUB-WEIBULL CONFIDENCE REGIME). Proposition 4 highlights several features of the sub-Weibull setting.

- (i) The restriction on  $\delta$  comes from expressing the sub-Weibull concentration result in a pure  $1/\sqrt{N}$  form. Over the confidence regime considered here, the generalized Bernstein inequality of Kuchibhotla and Chakraborty (2022, Theorem 3.1) reduces to the same pointwise concentration template used in Proposition 3. More generally, if the prescribed  $\delta \in (0, 1)$  lies outside this regime, an additional higher-order term appears, which decays at the faster  $\mathcal{O}(1/N)$  rate.
- (ii) A key advantage of Proposition 4 is that it remains applicable under substantially heavier-tailed uncertainty distributions than those covered by existing finite-sample DRO guarantees. In particular, Gao (2023) relies on a transportation-information condition on the data-generating law, which implies standard sub-Gaussian-type tail assumptions, while the classical Wasserstein guarantee of Mohajerin Esfahani and Kuhn (2018), via Fournier and Guillin (2015), requires the stronger stretched-exponential moment condition  $\mathbb{E}[\exp(\|\tilde{\boldsymbol{\xi}}\|^r)] < +\infty$  for some  $r > 1$  on the raw uncertainty. By contrast, Proposition 4 directly covers the broader sub-Weibull regime.

### Markovian data

Rather strikingly, our ambiguity set can easily accommodate circumstances when the samples are not i.i.d. Consider the setting where  $\{\hat{\boldsymbol{\xi}}_i\}_{i \in [N]}$  is a sample trajectory of a Markov chain with invariant distribution  $\mathbb{P}^*$ . In this case, our mean estimator  $\mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})]$  is no longer constructed using i.i.d. samples, which prohibit the use of classical concentration inequalities. Nevertheless, using Hoeffding’s inequality for general Markov chains (Fan et al. 2021), we can obtain the desired pointwise concentration inequality.

**PROPOSITION 5.** *Assume the loss function is bounded:  $\ell(\mathbf{z}, \boldsymbol{\xi}) \in [\underline{\ell}, \bar{\ell}]$  for all  $\mathbf{z} \in \mathcal{X}$  and  $\boldsymbol{\xi} \in \Xi$ . Suppose that the Markov chain is time-homogeneous with an absolute spectral gap of  $1 - \lambda > 0$ , where  $\lambda \in [0, 1)$  is the operator norm of the Markov transition kernel acting on the Hilbert space of square-integrable mean-zero functions under  $\mathbb{P}^*$ . Then, the pointwise concentration (9) holds with  $\varepsilon(y) = \frac{\bar{\ell} - \underline{\ell}}{2} \sqrt{\frac{2(1+\lambda)}{1-\lambda}} y$ , i.e.,*

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq \frac{1}{\sqrt{N}} \frac{\bar{\ell} - \underline{\ell}}{2} \sqrt{\frac{2(1+\lambda)}{1-\lambda}} \log\left(\frac{1}{\delta}\right) \quad (14)$$

*with probability at least  $1 - \delta$ . The same result holds for time-inhomogeneous chains with  $\lambda$  replaced by  $\max_{i \in [N]} \lambda_i$ , where  $1 - \lambda_i$  is the spectral gap of the Markov transition kernel for step  $i$ .*

**REMARK 4 (MARKOVIAN DATA WITH GENERAL STATE SPACES).** Our proposed DRO model applies to (time-inhomogeneous) Markovian data with general state spaces. We note that DRO with time-homogeneous Markovian data has been studied in (Li et al. 2021) in a discrete-state-space setting. Unlike ours, their statistical guarantee is asymptotic as it relies on large deviation principles; hence, it does not provide a prescription for adjusting the radius with respect to the sample size  $N$ .

## Outlier corrupted data

In many practical settings, the observed data may be contaminated by outliers, i.e., samples that deviate arbitrarily from the underlying distribution  $\mathbb{P}^*$ . Under such corruption, the standard empirical-mean reference that uses all data points is sensitive to these atypical observations and may no longer reliably approximate  $\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})]$ , potentially invalidating the coverage guarantee. To address this, we replace the empirical-mean reference with a *Median-of-Means* (MoM) reference functional (Laforgue et al. 2021), which retains a  $1/\sqrt{N}$  convergence rate while remaining robust to a constant fraction of outliers.

*MoM estimator.* Given  $N$  observations  $\hat{\boldsymbol{\xi}}_1, \dots, \hat{\boldsymbol{\xi}}_N$ , fix an integer  $K \geq 1$  and partition the index set  $[N]$  into  $K$  disjoint blocks  $\mathcal{B}_1, \dots, \mathcal{B}_K$  of equal size  $B = \lfloor N/K \rfloor$ , independently of the data. For each block  $k \in [K]$ , define the block mean

$$f_k(\mathbf{z}) := \frac{1}{B} \sum_{i \in \mathcal{B}_k} \ell(\mathbf{z}, \hat{\boldsymbol{\xi}}_i).$$

Then the MoM reference functional is defined as

$$\hat{\mu}_{\text{MoM}}(\mathbf{z}) := \text{median}(f_1(\mathbf{z}), \dots, f_K(\mathbf{z})).$$

In the outlier-corrupted setting, we use  $\hat{\mu}_{\text{MoM}}(\mathbf{z})$  in place of  $\mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})]$  as the reference functional in the loss constraints.

Suppose the sample points  $\{\hat{\boldsymbol{\xi}}_1, \dots, \hat{\boldsymbol{\xi}}_N\}$  contains  $N - n_o$  inliers drawn i.i.d. from  $\mathbb{P}^*$ , and  $n_o$  arbitrary outliers. We assume that the fraction of outliers is  $\omega := n_o/N \in [0, 1/2)$ . Following (Laforgue et al. 2021) with the arithmetic-mean calibration, we define the inflation constant

$$\Gamma(\omega) := \frac{\sqrt{2(1+2\omega)}}{(1-2\omega)^{3/2}},$$

which satisfies  $\Gamma(\omega) \rightarrow \infty$  as  $\omega \rightarrow 1/2$ .

**PROPOSITION 6 (Proposition 2 (eq. (3)) of Laforgue et al. (2021)).** *Suppose the inlier loss has bounded variance:  $\text{Var}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq \sigma^2$  for all  $\mathbf{z} \in \mathcal{X}$ . Set  $K = \left\lceil \frac{4(1+2\omega)}{(1-2\omega)^2} \log \frac{1}{\delta} \right\rceil$  and let  $\delta \in \left[ e^{-(1-2\omega)^2 N / (4(1+2\omega))}, e^{-(1-2\omega)^2 N / 8} \right]$ . Then, the MoM reference satisfies the pointwise concentration bound*

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \hat{\mu}_{\text{MoM}}(\mathbf{z}) \leq 4\sqrt{e} \sigma \Gamma(\omega) \sqrt{\frac{1 + \log(1/\delta)}{N}}.$$

with probability at least  $1 - \delta$ .

We note that the covering argument in Proposition 3 carries over under this replacement, because it only requires a pointwise concentration bound and a Lipschitz reference functional. The latter property holds for  $\hat{\mu}_{\text{MoM}}$  with the same constant  $L$ : the median satisfies  $|\text{median}(\mathbf{a}) - \text{median}(\mathbf{b})| \leq$

$\max_k |a_k - b_k|$  for any two vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^K$ , and each block mean satisfies  $|f_k(\mathbf{z}) - f_k(\mathbf{z}')| \leq L\|\mathbf{z} - \mathbf{z}'\|$  by the pathwise Lipschitz bound  $|\ell(\mathbf{z}, \boldsymbol{\xi}) - \ell(\mathbf{z}', \boldsymbol{\xi})| \leq L\|\mathbf{z} - \mathbf{z}'\|$  (see proof of Lemma 1).

This replacement does not require interpreting  $\hat{\mu}_{\text{MoM}}$  as an expectation under a new empirical distribution. When we later derive the exact dual reformulation, the same generalized moment problem argument applies to the MoM-based ambiguity set as long as this set satisfies the usual Slater condition for the chosen radius  $\epsilon$ . Under this condition, strong duality and dual attainment hold with  $\hat{\mu}_{\text{MoM}}(\mathbf{z})$  replacing  $\mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})]$  in the dual objective.

**REMARK 5 (OUTLIER ROBUSTNESS THROUGH ROBUST ESTIMATORS).** Distributionally robust optimization with outlier corrupted data has gained increasing attention. Nietert et al. (2023) and Nietert et al. (2024) employ mixed Wasserstein-Total Variation ambiguity sets to capture both local geometric and global adversarial corruptions. Meanwhile, Blanchet et al. (2024) utilize optimal transport with concave cost functions to automatically discount and rectify large-distance outliers. Additionally, Jiang and Xie (2024) propose a distributionally favorable optimization model that seeks the most favorable distribution to naturally ignore outlier-induced high losses. However, these approaches rely on problem-specific ambiguity sets or non-standard optimization paradigms tailored to specific outlier structures. In contrast, our Targeted IPM framework provides a unified framework. Since our fast convergence rate only depends on a scalar concentration inequality, we can seamlessly plug in any statistically robust estimator, such as MoM, to achieve outlier resistance while preserving tractability and bypassing the curse of dimensionality.

## Incomplete Data

Our framework is applicable to settings where a subset of the random parameters is not always observable. This is relevant to many selection problems, such as company hiring, college admissions, and loan approval, where outcomes (e.g., high-performing candidates, qualified students, or applicants who could repay the loan) are only observed for selected individuals. Formally, we define the random parameters as  $\tilde{\boldsymbol{\xi}} = (\tilde{\boldsymbol{\chi}}, \tilde{\boldsymbol{\omega}}) \in \mathcal{X} \times \Omega$ , where  $\tilde{\boldsymbol{\chi}}$  are entries that are always observable, while  $\tilde{\boldsymbol{\omega}}$  are those that will be observed only if the candidate is selected.

Let  $\tilde{S} \in \{0, 1\}$  denote the historical selection outcome with  $\tilde{S} = 1$  if the candidate is selected to advance and  $\tilde{S} = 0$  otherwise. Define  $\bar{\mathbb{P}}$  to be the joint distribution of  $(\tilde{S}, \tilde{\boldsymbol{\chi}}, \tilde{\boldsymbol{\omega}})$ . This binary random variable is governed by a logging policy

$$\pi(\boldsymbol{\chi}) := \bar{\mathbb{P}}(\tilde{S} = 1 | \tilde{\boldsymbol{\chi}} = \boldsymbol{\chi}),$$

which represents the probability that a candidate with covariate  $\boldsymbol{\chi}$  is selected. Based on the logging policy, the available dataset is given by

$$\{(\boldsymbol{\chi}_n, \boldsymbol{\omega}_n)\}_{n \in \mathcal{I}} \cup \{\boldsymbol{\chi}_n\}_{n \in [N] \setminus \mathcal{I}},$$

where  $\mathcal{I}$  is the set of candidates who were selected and have the outcome observed. Under the mild conditional exchangeability and positivity assumptions (Jia et al. 2024, Assumption 1 and 2), the inverse probability weighting estimator of the expected loss is given by

$$\mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] := \frac{1}{N} \sum_{n \in [N]} \frac{\mathbb{I}[S_n = 1]}{\pi(\boldsymbol{\chi}_n)} \ell(\mathbf{z}, \boldsymbol{\xi}_n), \quad (15)$$

where each full sample point is weighted by the reciprocal of its propensity score  $\pi(\boldsymbol{\chi}_n)$ . We use this IPW estimator as the reference functional in the loss constraints and obtain the following concentration bound.

**PROPOSITION 7.** *Suppose conditional exchangeability and positivity hold, with  $\pi(\boldsymbol{\chi}) \geq \underline{\pi} > 0$ , and suppose the losses are uniformly sub-Gaussian, i.e.,  $\|\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})\|_{\psi_2} \leq \sigma$  for all  $\mathbf{z} \in \mathcal{X}$ . The pointwise concentration (9) holds with  $\varepsilon(y) = C_{\text{IPW}} \frac{\sigma}{\underline{\pi}} \sqrt{y}$  for a universal constant  $C_{\text{IPW}} > 0$ , i.e., we have*

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq C_{\text{IPW}} \frac{\sigma}{\underline{\pi}} \sqrt{\frac{\log(1/\delta)}{N}}$$

with probability at least  $1 - \delta$ .

### Contextual optimization

Our framework can also be applied to stochastic optimization with side information. Suppose that a realization  $\mathbf{c} \in \mathbb{R}^{D_c}$  of some contextual covariates is observed, which changes the conditional distribution of  $\tilde{\boldsymbol{\xi}}$ . To exploit the side information, the decision-maker replaces the expectation in the objective with the conditional expectation

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) | \tilde{\mathbf{c}} = \mathbf{c}],$$

where  $\mathbb{P}^*$  denotes the joint distribution of  $(\tilde{\mathbf{c}}, \tilde{\boldsymbol{\xi}})$ . In a data driven setting, we observe historical data  $\{(\mathbf{c}_n, \boldsymbol{\xi}_n)\}_{n \in [N]}$  and approximate the conditional expectation using the weighted sum

$$\mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) | \tilde{\mathbf{c}} = \mathbf{c}] := \sum_{n \in [N]} w(\mathbf{c}, \mathbf{c}_n) \ell(\mathbf{z}, \boldsymbol{\xi}_n),$$

where the weights  $(w(\mathbf{c}, \mathbf{c}_n))_{n \in [N]}$  are designed such that they give higher values to data points that are close to the current realization  $\mathbf{c}$ . A popular scheme is the Nadaraya-Watson kernel regression, in which the weights are defined through a kernel function  $\mathcal{K}_h$  with bandwidth parameter  $h > 0$ :

$$w(\mathbf{c}, \mathbf{c}_n) := \frac{\mathcal{K}_h(\mathbf{c} - \mathbf{c}_n)}{\sum_{m \in [N]} \mathcal{K}_h(\mathbf{c} - \mathbf{c}_m)}$$

The resulting data-driven approximation is guaranteed to provide asymptotic consistency of the optimal solutions. The following proposition establishes a pointwise concentration for this data-driven conditional estimator.

PROPOSITION 8. Assume that the loss function takes values in the interval  $[0, 1]$ , the conditional expectation  $\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})|\tilde{\mathbf{c}} = \mathbf{c}]$  is  $L_c$ -Lipschitz continuous in  $\mathbf{c}$ , and the marginal density function  $f(\mathbf{c})$  satisfies  $f(\mathbf{c}) \geq \underline{f} > 0$  for all  $\mathbf{c}$  in its support. Then, the following pointwise concentration inequality holds:

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})|\tilde{\mathbf{c}} = \mathbf{c}] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})|\tilde{\mathbf{c}} = \mathbf{c}] \leq \left( \frac{2 \log(1/\delta)}{N c \underline{f}} \right)^{\frac{1}{D_c+2}} L_c^{\frac{D_c}{D_c+2}} \frac{D_c + 2}{2(D_c/2)^{\frac{D_c}{D_c+2}}}$$

with probability at least  $1 - \delta$ .

REMARK 6 (CONTEXTUAL RATES AND DIMENSIONALITY). Note that the bound suffers from the curse of dimensionality in  $D_c$ . This dependence in the dimension is unavoidable since learning the conditional expectation in  $D_c$ -dimensional nonparametric regression has the *minimax* squared-error rate of  $O(N^{-2/(D_c+2)})$  (Györfi et al. 2002), corresponding to the absolute-error rate  $O(N^{-1/(D_c+2)})$  in Proposition 8. Nevertheless, our DRO framework is applicable to provide a robustification with coverage guarantee  $\mathbb{P}^*(\cdot | \tilde{\mathbf{c}} = \mathbf{c}) \in \mathcal{P}_{\epsilon_N}$  where the radius can be set to

$$\epsilon_N := \frac{D_c + 2}{2(D_c/2)^{\frac{D_c}{D_c+2}}} L_c^{\frac{D_c}{D_c+2}} \left( \frac{2 \left( D_x \log(1 + 2R_x L \sqrt{N}) + \log(1/\delta) \right)}{N c \underline{f}} \right)^{\frac{1}{D_c+2}} + \frac{2}{\sqrt{N}}.$$

### 2.3. Dual Formulations

In this subsection, we derive an exact dual reformulation of the DRO problem (3) and establish convergence to empirical risk minimization as the ambiguity parameter  $\epsilon$  approaches 0. We start with the reformulation.

THEOREM 2. Fix  $\epsilon > 0$ . The DRO problem (3) is equivalent to the infinite linear program

$$\begin{aligned} \min \quad & \alpha + \beta \Omega + \int_{\mathcal{X}} \left( \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] + \epsilon \right) \nu(d\mathbf{z}) \\ \text{s.t.} \quad & \mathbf{x} \in \mathcal{X}, \alpha \in \mathbb{R}, \beta \in \mathbb{R}_+, \nu \in \mathcal{M}_+(\mathcal{X}) \\ & \ell(\mathbf{x}, \boldsymbol{\xi}) \leq \alpha + \beta \|\boldsymbol{\xi}\|^2 + \int_{\mathcal{X}} \ell(\mathbf{z}, \boldsymbol{\xi}) \nu(d\mathbf{z}) \quad \forall \boldsymbol{\xi} \in \Xi. \end{aligned} \tag{16}$$

The optimal solution  $(\mathbf{x}, \alpha, \beta, \nu) \in \mathcal{X} \times \mathbb{R} \times \mathbb{R}_+ \times \mathcal{M}_+(\mathcal{X})$  is attained.

The optimization problem remains intractable because it involves an infinite number of variables and constraints. In the following section, we will develop a tractable approximation scheme with guarantees.

We close the section by showing that the optimal value and solutions of the DRO problem converge to those of the (non-robust) empirical risk minimization problem as  $\epsilon$  approaches 0.

PROPOSITION 9. As  $\epsilon \rightarrow 0$ , the optimal value of (16) converges to that of the empirical risk minimization problem

$$\min_{\mathbf{x} \in \mathcal{X}} \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})]. \quad (17)$$

In addition, every cluster point  $\hat{\mathbf{x}}^*$  of a sequence  $\{\hat{\mathbf{x}}_\epsilon\}_{\epsilon \downarrow 0}$  of minimizers for (16) is a minimizer for (17).

### 3. A Monte Carlo Sampling Approach

Problem (16) remains difficult to solve as it involves a functional decision variable  $\nu$  and an expectation over a continuous distribution  $\mathbb{Z}$ . By Proposition 2, our DRO problem (3) constitutes an instance of DRO problems with an infinitely constrained ambiguity set, which have been studied by Chen et al. (2019) and enable a convergent cutting-plane algorithm for settings where the support  $\Xi$  is compact. In this paper, we propose a tractable approximation via Monte Carlo sampling.

Specifically, we draw  $M$  samples from  $\mathbb{Z}$  and obtain the approximate ambiguity set:

$$\mathcal{P}_\epsilon^M := \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : \begin{array}{l} \frac{1}{M} \sum_{m \in [M]} \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \leq 0 \\ \mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega \end{array} \right\}. \quad (18)$$

By construction, the constraint only enforces  $\mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] \leq \epsilon$  at finitely many sample points  $\{\mathbf{z}_m\}_{m=1}^M$ . Thus, we have a relaxed ambiguity set  $\mathcal{P}_\epsilon^M \supseteq \mathcal{P}_\epsilon$ , and the resulting distributionally robust optimization problem is a *conservative* approximation, and the out-of-sample performance guarantees provided in Corollary 1 remain valid. Unlike the classical sample-average approximation for stochastic programming that often introduces optimistic bias, our approximation always generates pessimistic bias and yields safe solutions. The conservativeness further endows the model with attractive theoretical properties, which we delineate in the following subsection.

Using the ambiguity set (18), we obtain the following dual reformulation:

$$\begin{aligned} \min \quad & \alpha + \beta\Omega + \frac{1}{M} \sum_{m \in [M]} \nu_m (\mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] + \epsilon) \\ \text{s.t.} \quad & \mathbf{x} \in \mathcal{X}, \alpha \in \mathbb{R}, \beta \in \mathbb{R}_+, \boldsymbol{\nu} \in \mathbb{R}_+^M \\ & \ell(\mathbf{x}, \boldsymbol{\xi}) \leq \alpha + \beta\|\boldsymbol{\xi}\|^2 + \frac{1}{M} \sum_{m \in [M]} \nu_m \ell(\mathbf{z}_m, \boldsymbol{\xi}) \quad \forall \boldsymbol{\xi} \in \Xi. \end{aligned} \quad (19)$$

This formulation constitutes the sample-average counterpart of Theorem 2. By replacing the expectation under  $\mathbb{Z}$  with an empirical average over  $\{\mathbf{z}_m\}_{m=1}^M$ , the functional multiplier  $\nu(\cdot)$  reduces to a finite-dimensional vector  $\boldsymbol{\nu} \in \mathbb{R}_+^M$ .

In the following, we develop the theoretical properties and tractable conic reformulations of this approximation scheme. We first consider the case where the support set  $\Xi$  is compact.

### 3.1. Theoretical Guarantees

Our first theoretical result establishes a high-confidence guarantee for the containment of the sample-average-based ambiguity set  $\mathcal{P}_\epsilon^M$  within a relaxed exact ambiguity set

$$\mathcal{P}_\epsilon(\eta) := \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : \begin{array}{l} \mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \right] \leq \eta \\ \mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega \end{array} \right\} \quad (20)$$

parameterized by  $\eta > 0$ .

**THEOREM 3.** *Assume  $\sup_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\| \leq R_x$  and  $\sup_{\boldsymbol{\xi} \in \Xi} \|\boldsymbol{\xi}\| \leq R_\xi$ , and let*

$$R := \max_{j \in [J]} \left\| \begin{bmatrix} \mathbf{A}_j^\top & \mathbf{b}_j \\ \mathbf{a}_j^\top & b_j \end{bmatrix} \right\|_{\text{op}} \sqrt{R_\xi^2 + 1} \sqrt{R_x^2 + 1}.$$

*Then, setting*

$$\eta := \mathcal{O} \left( R \sqrt{\frac{J \log J (\log M)^3}{M}} \right) + 2(c_1 + c_2 \sqrt{\Omega}) \sqrt{\frac{1}{2M} \log \left( \frac{1}{\tau} \right)}, \quad (21)$$

*we can ensure with probability at least  $1 - \tau$ :*

$$\mathcal{P}_\epsilon^M \subseteq \mathcal{P}_\epsilon(\eta).$$

This result shows that using the sample-average-based ambiguity set is akin to relaxing the expected hinge constraint by  $\eta := \tilde{\mathcal{O}}(1/\sqrt{M})$ . The result relies on a generalization bound for the sample-average approximation errors, stated as Proposition 11 in Appendix B. In the proof of that proposition, we bound the Rademacher complexity of the function class

$$\mathcal{H} := \left\{ \mathbf{z} \mapsto \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ : \mathbb{Q} \in \mathcal{P}(\Xi), \mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega \right\},$$

which contains functions parametrized by distributions  $\mathbb{Q} \in \mathcal{P}(\Xi)$  satisfying the second-moment constraint  $\mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega$ . Unlike in traditional settings, the parameter space is infinite-dimensional, rendering classical techniques inapplicable. We employ the contraction principle to reduce the analysis to the Rademacher complexity of the functions  $\mathbf{z} \mapsto \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})]$ . Using our loss function structure, we further reduce it to the complexity of  $J$ -fold maxima of hyperplanes (Attias and Kontorovich 2024, Corollary 5). This function-class argument is the key ingredient behind Theorem 3.

**REMARK 7 (ROLE OF FUNCTION-CLASS STRUCTURE).** The dimension-free  $\tilde{\mathcal{O}}(1/\sqrt{M})$  sampling rate in Proposition 11 is a key benefit of exploiting the piecewise-affine structure of the loss class, rather than a generic consequence of Lipschitz continuity. Indeed, for a generic class of bounded  $L$ -Lipschitz functions on a metric space with dimension  $D_x$  of the decision  $\mathbf{x}$ , the empirical Rademacher complexity can scale as  $\mathcal{O}(L/M^{1/(D_x+1)})$ ; see (Gottlieb et al. 2016, Theorem 4.3). In contrast, the piecewise-affine loss class considered here reduces to the class of  $J$ -fold maxima of hyperplanes, for which the sharper Rademacher bound in Proposition 11 is available.

We next establish the convergence of the optimal value and optimal solution of the approximate DRO problem. Let  $\hat{v}$  and  $\hat{v}_M$  be the optimal values of the exact problem and the approximate problem, respectively, i.e.,

$$\hat{v} := \min_{\mathbf{x} \in \mathcal{X}} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] \quad \text{and} \quad \hat{v}_M := \min_{\mathbf{x} \in \mathcal{X}} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon^M} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})].$$

We define  $\mathcal{S}$  as the set of minimizers of the exact problem. The minimizers  $\hat{\mathbf{x}} \in \mathcal{S}$  and  $\hat{\mathbf{x}}_M$  of the sampled problem are attained because the corresponding objectives are lower semicontinuous on the compact set  $\mathcal{X}$ .

**THEOREM 4.** *As  $M \rightarrow \infty$ , the following convergences in probability hold:*

$$\hat{v}_M \xrightarrow{P} \hat{v} \quad \text{and} \quad \text{dist}(\hat{\mathbf{x}}_M, \mathcal{S}) \xrightarrow{P} 0.$$

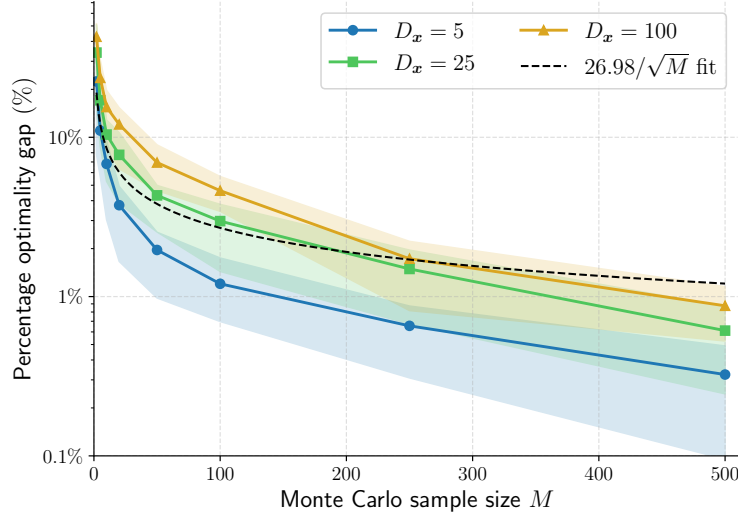
Furthermore, assume there exists an optimal dual measure  $\nu$  in Theorem 2 such that  $\gamma\mathbb{Z} - \nu \in \mathcal{M}_+(\mathcal{X})$  for some  $\gamma \in \mathbb{R}_+$ . Let  $\gamma^* := \inf\{\gamma \in \mathbb{R} : \gamma\mathbb{Z} - \nu \in \mathcal{M}_+(\mathcal{X})\}$ . Then, the following suboptimality bound holds with probability at least  $1 - \tau$ :

$$\begin{aligned} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}_M, \tilde{\boldsymbol{\xi}})] &\leq \min_{\mathbf{x} \in \mathcal{X}} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] \\ &+ \gamma^* \left( \mathcal{O} \left( R \sqrt{\frac{J \log J (\log M)^3}{M}} \right) + 2(c_1 + c_2 \sqrt{\Omega}) \sqrt{\frac{1}{2M} \log \left( \frac{1}{\tau} \right)} \right). \end{aligned}$$

**REMARK 8 (EXISTENCE OF THE DOMINANCE CONSTANT).** The domination condition used in the suboptimality bound of Theorem 4 is a condition on the auxiliary sampling distribution, not on the underlying data-generating model. Since  $\mathbb{Z}$  is chosen by the modeler, it can be enriched, if needed, to dominate the relevant optimal dual measure, as we now show. Specifically, if  $\nu$  is any optimal dual measure in Theorem 2, then for the case  $\nu(\mathcal{X}) = 0$ , any  $\gamma \geq 0$  works; otherwise, for any  $\lambda \in (0, 1)$ , the full-support mixture  $\mathbb{Z}_\lambda := (1 - \lambda)\mathbb{Z} + \lambda\nu/\nu(\mathcal{X})$  satisfies  $(\nu(\mathcal{X})/\lambda)\mathbb{Z}_\lambda - \nu \in \mathcal{M}_+(\mathcal{X})$ , and hence  $\gamma_\lambda^* \leq \nu(\mathcal{X})/\lambda < \infty$ . Thus the required dominance constant can always be made finite through the modeling choice of  $\mathbb{Z}$ . By Proposition 2, all full-support choices of the auxiliary distribution induce the same semi-infinite ambiguity set, so replacing  $\mathbb{Z}$  with  $\mathbb{Z}_\lambda$  does not alter the DRO problem itself.

**REMARK 9 (NUMERICAL ILLUSTRATION OF THE MONTE CARLO APPROXIMATION).** We provide a simple experiment to illustrate the sampling behavior in Theorem 4. The detailed experimental design is reported in Appendix E.1. To match the bounded-support assumptions, both the decision set and the uncertainty support are Euclidean unit balls. We use a shifted piecewise-affine loss and a symmetric empirical distribution so that the exact optimizer  $x_*$  is known and nonzero. This allows us to evaluate the sampled-dual solution  $\hat{\mathbf{x}}_M$  directly, without computing a large-reference Monte Carlo solution. We report the relative optimality gap

$$\frac{\mathbb{E}_{\hat{\mathcal{P}}_N}[\ell(\hat{\mathbf{x}}_M, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathcal{P}}_N}[\ell(x_*, \tilde{\boldsymbol{\xi}})]}{\mathbb{E}_{\hat{\mathcal{P}}_N}[\ell(x_*, \tilde{\boldsymbol{\xi}})]} \times 100\%$$



**Figure 1** Percentage optimality gap of the sampled-dual solution as a function of the Monte Carlo sample size  $M$ . The dashed curve is a fitted  $C/\sqrt{M}$  reference line.

across different Monte Carlo sample sizes  $M$ . Figure 1 shows that the sampled-dual solution improves rapidly as  $M$  increases. The fitted  $C/\sqrt{M}$  curve provides a visual benchmark for the sampling rate in Theorem 4. All curves exhibit a similar decay pattern. At the largest tested Monte Carlo sizes, the gaps are reduced to roughly the one-percent level across all dimensions. This provides numerical evidence that the sampled-dual approximation improves consistently with  $M$  and does not suffer from the curse of dimensionality.

### 3.2. Conic Programming Reformulations

In this subsection, we show that the approximate problem admits a tractable conic programming reformulation. Tractable approximations for two-stage problems are provided in Appendix D.

**THEOREM 5.** *Suppose that  $\mathcal{X}$  and  $\Xi$  are convex. Then problem (19) is equivalent to the following finite convex reformulation involving the support function  $\sigma_{\Xi}$ :*

$$\begin{aligned}
& \min \alpha + \beta\Omega + \frac{1}{M} \sum_{m \in [M]} \nu_m (\mathbb{E}_{\hat{\mathbb{P}}_N} [\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] + \epsilon) \\
& \text{s.t. } \mathbf{x} \in \mathcal{X}, \alpha \in \mathbb{R}, \beta \in \mathbb{R}_+, \boldsymbol{\nu} \in \mathbb{R}_+^M \\
& \quad \boldsymbol{\lambda}_{jm} \in \mathbb{R}_+^J \quad \forall m \in [M], \quad \boldsymbol{\theta}_j \in \mathbb{R}^{D_{\boldsymbol{\xi}}}, \zeta_j \in \mathbb{R}_+ \quad \forall j \in [J] \\
& \quad \mathbf{e}^\top \boldsymbol{\lambda}_{jm} = \nu_m \quad \forall m \in [M] \quad \forall j \in [J] \\
& \quad \zeta_j + \sigma_{\Xi}(\boldsymbol{\theta}_j) + b_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k b_k(\mathbf{z}_m) \leq \alpha \quad \forall j \in [J] \\
& \quad \left\| \left[ \mathbf{a}_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k \mathbf{a}_k(\mathbf{z}_m) - \boldsymbol{\theta}_j \right] \right\| \leq \zeta_j + \beta \quad \forall j \in [J].
\end{aligned} \tag{22}$$

In particular, if  $\mathcal{X}$  and  $\Xi$  are second-order cone representable, then the reformulation above is a second-order cone program.

#### 4. Extensions to Unbounded Support

We now extend the previous results to the setting where the distributional support is  $\Xi := \mathbb{R}^D \xi$ . We assume that there exist constants  $K_{\text{sw}}, \vartheta > 0$  such that

$$\mathbb{P}^*(\|\tilde{\xi}\| \leq t) \geq 1 - 2 \exp\left(- (t/K_{\text{sw}})^\vartheta\right) \quad \forall t > 0.$$

This is a standard sub-Weibull-type tail bound (Kuchibhotla and Chakraborty 2022, Definition 2.2). For a given threshold  $t > 0$ , we incorporate this tail information into the ambiguity set:

$$\mathcal{P}'_\epsilon := \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : \begin{array}{l} \mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\tilde{\mathbf{z}}, \tilde{\xi})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\xi})] - \epsilon \right]_+ \right] \leq 0 \\ \mathbb{E}_{\mathbb{Q}}[\|\tilde{\xi}\|^2] \leq \Omega, \quad \mathbb{Q}(\|\tilde{\xi}\| \leq t) \geq 1 - 2 \exp\left(- (t/K_{\text{sw}})^\vartheta\right) \end{array} \right\}. \quad (23)$$

In the sequel, we will apply this bound at thresholds  $t = t_M$  that vary with the Monte Carlo sample size  $M$ . Under the standing second-moment and tail assumptions on  $\mathbb{P}^*$ , Proposition 3 implies that  $\mathbb{P}^* \in \mathcal{P}'_\epsilon$  with probability at least  $1 - \delta$ .

We next show that the proposed framework extends beyond compact support. Under the single-threshold tail bound implied by sub-Weibull tails above, we derive an exact dual reformulation, establish that the Monte Carlo approximation remains conservative, and show that the same  $\tilde{\mathcal{O}}(1/\sqrt{M})$  sampling principle continues to hold up to explicit tail-dependent terms. We also obtain a tractable conic reformulation.

The following proposition provides the dual formulation of the DRO problem under this ambiguity set.

**PROPOSITION 10.** *The DRO problem (3) with the ambiguity set (23) is equivalent to the infinite linear program*

$$\begin{aligned} \min \quad & \alpha + \beta\Omega - \kappa \left( 1 - 2 \exp\left(- (t/K_{\text{sw}})^\vartheta\right) \right) + \int_{\mathcal{X}} \left( \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\xi})] + \epsilon \right) \nu(d\mathbf{z}) \\ \text{s.t.} \quad & \mathbf{x} \in \mathcal{X}, \alpha \in \mathbb{R}, \beta, \kappa \in \mathbb{R}_+, \nu \in \mathcal{M}_+(\mathcal{X}) \\ & \ell(\mathbf{x}, \xi) + \kappa \mathbb{I}_{\{\|\xi\| \leq t\}} \leq \alpha + \beta \|\xi\|^2 + \int_{\mathcal{X}} \ell(\mathbf{z}, \xi) \nu(d\mathbf{z}) \quad \forall \xi \in \Xi \end{aligned} \quad (24)$$

An optimal solution  $(\mathbf{x}, \alpha, \beta, \kappa, \nu) \in \mathcal{X} \times \mathbb{R} \times \mathbb{R}_+^2 \times \mathcal{M}_+(\mathcal{X})$  is attained.

Drawing  $M$  samples from  $\mathbb{Z}$ , we obtain an outer approximation of the ambiguity set:

$$\mathcal{P}'_\epsilon{}^M := \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : \begin{array}{l} \frac{1}{M} \sum_{m \in [M]} \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}_m, \tilde{\xi})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}_m, \tilde{\xi})] - \epsilon \right]_+ \leq 0 \\ \mathbb{E}_{\mathbb{Q}}[\|\tilde{\xi}\|^2] \leq \Omega, \quad \mathbb{Q}(\|\tilde{\xi}\| \leq t) \geq 1 - 2 \exp\left(- (t/K_{\text{sw}})^\vartheta\right) \end{array} \right\}. \quad (25)$$

The corresponding dual formulation is given by

$$\begin{aligned}
& \min \alpha + \beta\Omega - \kappa \left( 1 - 2 \exp \left( - (t/K_{\text{sw}})^\vartheta \right) \right) + \frac{1}{M} \sum_{m \in [M]} \nu_m (\mathbb{E}_{\hat{\mathbb{P}}_N} [\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] + \epsilon) \\
& \text{s.t. } \mathbf{x} \in \mathcal{X}, \alpha \in \mathbb{R}, \beta, \kappa \in \mathbb{R}_+, \boldsymbol{\nu} \in \mathbb{R}_+^M \\
& \ell(\mathbf{x}, \boldsymbol{\xi}) + \kappa \mathbb{I}_{\{\|\boldsymbol{\xi}\| \leq t\}} \leq \alpha + \beta \|\boldsymbol{\xi}\|^2 + \frac{1}{M} \sum_{m \in [M]} \nu_m \ell(\mathbf{z}_m, \boldsymbol{\xi}) \quad \forall \boldsymbol{\xi} \in \Xi.
\end{aligned} \tag{26}$$

REMARK 10 (CALIBRATION OF TAIL PARAMETERS). The results in this section provide a theoretical ambiguity-set construction with sampling error of order  $\tilde{\mathcal{O}}(1/\sqrt{M})$ . In practice, however, we use the ambiguity set (18) from Section 3 regardless of whether the support is bounded or unbounded. If one nevertheless wishes to work with the ambiguity set in this section, then, as with the constant  $\Omega$ , the tail parameters  $(K_{\text{sw}}, \vartheta)$  must be chosen or calibrated in advance. The threshold  $t$  governs the tradeoff in the bound: increasing  $t$  enlarges the localized radius, while shrinking the tail correction term  $\exp(- (t/K_{\text{sw}})^\vartheta)$ .

#### 4.1. Theoretical Guarantees

Throughout this subsection, we set  $t = t_M := K_{\text{sw}}(\log M)^{1/\vartheta}$ , so that  $2 \exp(- (t_M/K_{\text{sw}})^\vartheta) = 2/M$ . For notational clarity, we denote by  $\mathcal{P}'_{\epsilon, M}$  the exact ambiguity set (23) with  $t = t_M$ . The corresponding relaxed moving-threshold set is

$$\mathcal{P}'_{\epsilon, M}(\eta) := \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : \begin{array}{l} \mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \right] \leq \eta \\ \mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega, \mathbb{Q}(\|\tilde{\boldsymbol{\xi}}\| \leq t_M) \geq 1 - 2 \exp(- (t_M/K_{\text{sw}})^\vartheta) \end{array} \right\}, \tag{27}$$

parameterized by  $\eta > 0$ . The first theoretical result establishes a high-confidence guarantee for the containment of the sample-average-based ambiguity set  $\mathcal{P}'_{\epsilon}{}^M$  within this relaxed exact ambiguity set.

THEOREM 6. Assume  $\sup_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\| \leq R_{\mathbf{x}}$ , and let

$$R(t) := \left( \max_{j \in [J]} \left\| \begin{bmatrix} \mathbf{A}_j^\top & \mathbf{b}_j \\ \bar{\mathbf{a}}_j^\top & \bar{b}_j \end{bmatrix} \right\|_{\text{op}} \sqrt{t^2 + 1} \right) \sqrt{R_{\mathbf{x}}^2 + 1}, \quad t > 0.$$

Then, setting

$$\eta_M := \mathcal{O} \left( R(t_M) \sqrt{\frac{J \log J (\log M)^3}{M}} \right) + 2(c_1 + c_2 \sqrt{\Omega}) \sqrt{\frac{1}{2M} \log \left( \frac{1}{\tau} \right)} + \frac{4c_1}{M} + 2c_2 \sqrt{\frac{2\Omega}{M}},$$

we can ensure with probability at least  $1 - \tau$ :

$$\mathcal{P}'_{\epsilon}{}^M \subseteq \mathcal{P}'_{\epsilon, M}(\eta_M).$$

The argument mirrors the bounded-support case. The only additional ingredient is the technical generalization bound in Proposition 12 in Appendix C, which controls the contribution of the tail event after truncation at  $t_M$ . With this tail correction, the sampled ambiguity set retains the same conservative outer-approximation interpretation. Since  $R(t_M) = \mathcal{O}((\log M)^{1/\vartheta})$ , the relaxation level  $\eta_M$  vanishes at the order  $\tilde{\mathcal{O}}(M^{-1/2})$ .

To compare the sampled problem with the limiting exact problem, we use Lemma 4 in Appendix C, which shows that the moving tail constraint vanishes asymptotically. Recall that  $\hat{v}$  and  $\mathcal{S}$  denote the optimal value and the set of minimizers of the exact problem with ambiguity set  $\mathcal{P}_\epsilon$ .

**THEOREM 7.** *Fix  $\epsilon > 0$ , and let  $\hat{v}'_M$  and  $\hat{\mathbf{x}}'_M$  be the optimal value and an optimal solution of the sampled problem (26) with  $t = t_M$ . Under the assumptions of Theorem 6, as  $M \rightarrow \infty$ ,*

$$\hat{v}'_M \xrightarrow{P} \hat{v} \quad \text{and} \quad \text{dist}(\hat{\mathbf{x}}'_M, \mathcal{S}) \xrightarrow{P} 0.$$

Furthermore, let  $\mathbf{x}_M^*$  be an exact optimizer of the threshold- $t_M$  problem, and assume that there exists an associated optimal dual solution  $(\mathbf{x}_M^*, \alpha_M^*, \beta_M^*, \kappa_M^*, \nu_M^*)$  in Proposition 10 such that  $\gamma\mathbb{Z} - \nu_M^* \in \mathcal{M}_+(\mathcal{X})$  for some  $\gamma \in \mathbb{R}_+$ . Define  $\gamma_M^* := \inf\{\gamma \in \mathbb{R}_+ : \gamma\mathbb{Z} - \nu_M^* \in \mathcal{M}_+(\mathcal{X})\}$ . Then, the following suboptimality bound holds with probability at least  $1 - \tau$ :

$$\begin{aligned} \sup_{\mathbb{Q} \in \mathcal{P}'_{\epsilon, M}} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}'_M, \tilde{\boldsymbol{\xi}})] &\leq \min_{\mathbf{x} \in \mathcal{X}} \sup_{\mathbb{Q} \in \mathcal{P}'_{\epsilon, M}} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] \\ &\quad + \gamma_M^* \left( \mathcal{O} \left( R(t_M) \sqrt{\frac{J \log J (\log M)^3}{M}} \right) + 2(c_1 + c_2 \sqrt{\Omega}) \sqrt{\frac{1}{2M} \log \left( \frac{1}{\tau} \right)} \right. \\ &\quad \left. + \frac{4c_1}{M} + 2c_2 \sqrt{\frac{2\Omega}{M}} \right). \end{aligned}$$

## 4.2. Conic Programming Reformulations

We now show that the approximate problem admits a tractable conic programming reformulation.

THEOREM 8. *Suppose that  $\mathcal{X}$  is convex. The problem (26) is equivalent to the following finite convex conic reformulation:*

$$\begin{aligned}
& \min \alpha + \beta\Omega - \kappa \left( 1 - 2 \exp \left( - (t/K_{\text{sw}})^{\vartheta} \right) \right) + \frac{1}{M} \sum_{m \in [M]} \nu_m (\mathbb{E}_{\hat{\mathbb{P}}_N} [\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] + \epsilon) \\
& \text{s.t. } \mathbf{x} \in \mathcal{X}, \alpha \in \mathbb{R}, \beta, \kappa \in \mathbb{R}_+, \boldsymbol{\nu} \in \mathbb{R}_+^M \\
& \quad \boldsymbol{\lambda}_{jm}, \boldsymbol{\lambda}'_{jm} \in \mathbb{R}_+^J \quad \forall j \in [J] \quad \forall m \in [M] \\
& \quad \boldsymbol{\theta}_j \in \mathbb{R}^{D_\xi}, \zeta_j, \zeta'_j \in \mathbb{R}_+ \quad \forall j \in [J] \\
& \quad \mathbf{e}^\top \boldsymbol{\lambda}_{jm} = \nu_m \quad \forall m \in [M] \quad \forall j \in [J] \\
& \quad \mathbf{e}^\top \boldsymbol{\lambda}'_{jm} = \nu_m \quad \forall m \in [M] \quad \forall j \in [J] \\
& \quad \zeta_j + b_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k b_k(\mathbf{z}_m) \leq \alpha \quad \forall j \in [J] \\
& \quad \left\| \left[ \mathbf{a}_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k \mathbf{a}_k(\mathbf{z}_m) \right] \right\| \leq \zeta_j + \beta \quad \forall j \in [J] \\
& \quad \zeta'_j + t \|\boldsymbol{\theta}_j\| + b_j(\mathbf{x}) + \kappa - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda'_{jm}^k b_k(\mathbf{z}_m) \leq \alpha \quad \forall j \in [J] \\
& \quad \left\| \left[ \mathbf{a}_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda'_{jm}^k \mathbf{a}_k(\mathbf{z}_m) - \boldsymbol{\theta}_j \right] \right\| \leq \zeta'_j + \beta \quad \forall j \in [J].
\end{aligned} \tag{28}$$

In particular, if  $\mathcal{X}$  is second-order cone representable, then the reformulation above is a second-order cone program.

## 5. Numerical Experiments

In this section, we conduct numerical experiments to validate the practical performance of our TIPM-DRO framework across two problem classes: a newsvendor problem and linear regression under adversarial outlier corruption. In each experiment, we describe the problem setting, the competing methods, and the out-of-sample (OOS) evaluation results.

### 5.1. Multi-Item Newsvendor Problem

We consider a multi-item newsvendor problem with order vector  $\mathbf{x} \in \mathbb{R}_+^K$  and random demand vector  $\tilde{\boldsymbol{\xi}} \in \mathbb{R}_+^K$ . The loss is

$$\ell(\mathbf{x}, \boldsymbol{\xi}) = \sum_{j=1}^K [h_j(x_j - \xi_j)_+ + b_j(\xi_j - x_j)_+].$$

We optimize the  $\mathbb{P}$ -CVaR $_\rho$  objective of this loss and keep the formulation generic in  $K$ ; in the numerical study, we set  $K = 3$ . This design serves two purposes. First, we want to evaluate the method under both lighter-tailed and heavier-tailed demand models, since heavy tails are operationally important in inventory systems and can materially change stocking and pooling decisions (Bimpikis

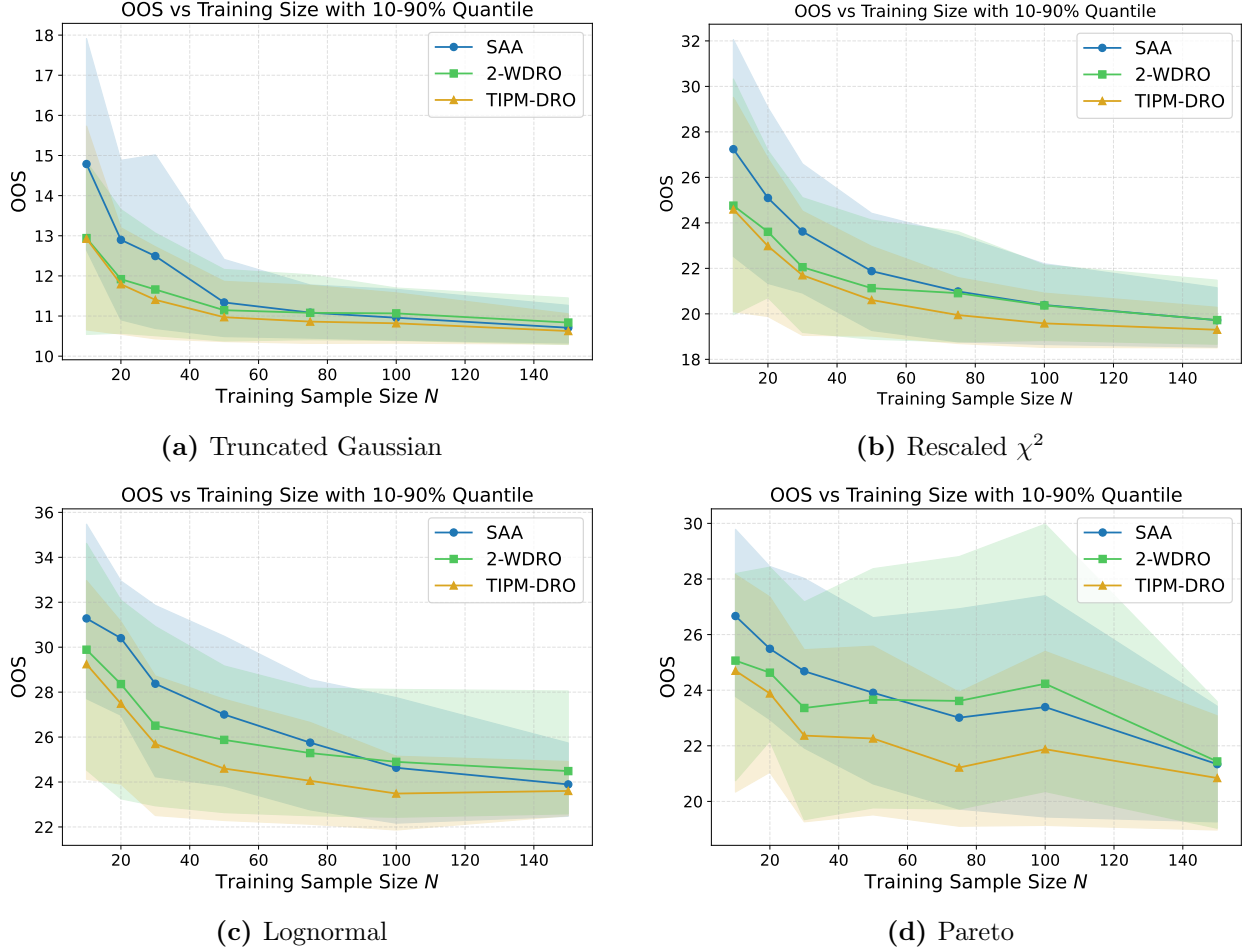
and Markakis 2016). We therefore compare four coordinatewise demand families spanning this range: truncated Gaussian, rescaled  $\chi^2$ , Lognormal, and Pareto. Second, since the paper emphasizes the dimension-robust motivation behind TIPM, the experiment involves a vector-valued decision and uncertainty. In the experiment, we set  $h = (0.1, 0.2, 0.3)$ ,  $b = (1, 1, 1)$ , and  $\rho = 0.05$ , and normalize the four demand models to share the same mean vector  $\boldsymbol{\mu} = (5, 6, 6)$  and standard-deviation vector  $\boldsymbol{\sigma} = (5, 6, 8)$ . Full parameter settings and reformulations are deferred to Appendix E.2.

We compare SAA, TIPM-DRO, and a 2-Wasserstein DRO baseline (2-WDRO) under the same CVaR objective. We consider the 2-Wasserstein DRO, which has recently been shown to generate superior performance in out-of-sample tests (Byeon 2025). We do not separately report 1-WDRO because, for newsvendor losses of this form, 1-Wasserstein DRO shares the same optimal minimizers as SAA, see Mohajerin Esfahani and Kuhn (2018, Remark 6.7). In each trial, the in-sample data are split into training and validation subsets, the ambiguity radius of each robust method is selected from the same grid by validation, and the chosen decision is evaluated on an independent OOS sample. The TIPM-DRO model is solved through the explicit second-moment-constrained reformulation reported in Appendix E.2. Thus, the same experimental pipeline tests both stories at once: how performance changes with tail heaviness and whether the method remains effective in a non-scalar decision problem.

Figure 2 reports mean OOS cost with 10th–90th percentile bands over 50 trials, where lower values are better. TIPM-DRO achieves the best mean OOS performance in all four demand settings. Under truncated Gaussian demand, both robust methods improve on SAA and TIPM-DRO still performs slightly better than 2-WDRO, although the gap narrows as  $N$  increases. Under rescaled  $\chi^2$ , Lognormal, and Pareto demand, the advantage of TIPM-DRO is more persistent, with the clearest separation appearing in the Lognormal and Pareto panels. These patterns support both design objectives. The Gaussian panel shows that TIPM-DRO is already competitive in a lighter-tailed regime, while the heavier-tailed panels show that its advantage over both SAA and 2-WDRO becomes more pronounced. At the same time, these gains are observed in a three-item vector problem rather than only in a scalar benchmark.

## 5.2. Outlier-Corrupted Regression

We consider linear regression with the mean absolute deviation (MAD) loss  $\ell(\mathbf{x}, \boldsymbol{\xi}) = |\mathbf{x}^\top \mathbf{u} - v|$ , where  $\boldsymbol{\xi} = (\mathbf{u}, v) \in \mathbb{R}^{D\xi+1}$  denotes a data point with feature vector  $\mathbf{u} \in \mathbb{R}^{D\xi}$  and response  $v \in \mathbb{R}$ , following the setup of Nietert et al. (2023). The true regression coefficient  $\mathbf{x}^*$  is drawn uniformly from  $\mathbb{S}^{D\xi-1}$ , clean features satisfy  $\mathbf{u}_i \sim \mathcal{N}(\mathbf{0}, \mathbb{I}_{D\xi})$ , and clean responses are  $v_i = \mathbf{x}^{*\top} \mathbf{u}_i$ . An adversary corrupts a fraction  $\omega = 0.2$  of the training samples by replacing each corrupted pair  $(\mathbf{u}_i, v_i)$  with  $(C\mathbf{u}_i, -C^2v_i + \rho)$ , where  $\rho = 0.1$ . This creates leverage-type outliers whose severity is controlled by

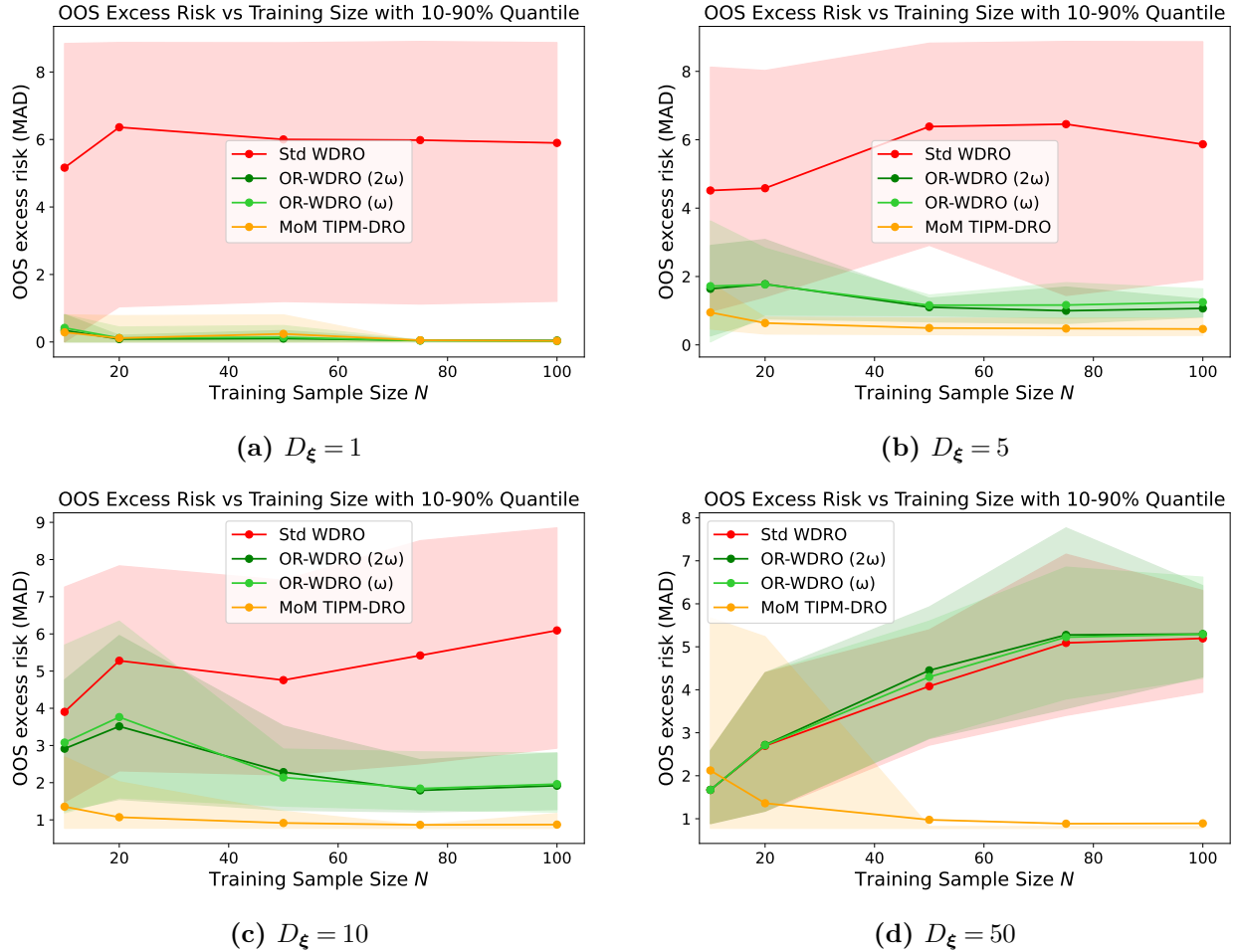


**Figure 2** Mean OOS cost versus the training sample size  $N$  with 10th–90th percentile bands in the multi-item newsvendor experiment with  $K = 3$ .

the scale parameter  $C$ . The experiment is designed to probe two questions. First, what happens as the ambient dimension  $D_\xi$  grows and Wasserstein geometry becomes less informative? Second, does robust performance depend on the outliers being geometrically well separated from the clean sample? Full implementation details and tuning grids are summarized in Appendix E.3.

We benchmark standard Wasserstein DRO (Mohajerin Esfahani and Kuhn 2018) (Std WDRO), the two outlier-robust WDRO variants of Nietert et al. (2023) with TV contamination radii  $\omega$  and  $2\omega$ , and our MoM TIPM-DRO estimator based on the Median-of-Means from Section 2.2. All robust models are tuned by a trimmed validation loss on a held-out split. For the three Wasserstein-based models, we cross-validate the Wasserstein radius over a common grid. For MoM TIPM-DRO, we cross-validate the number of blocks  $K$  and the TIPM radius  $\epsilon$ .

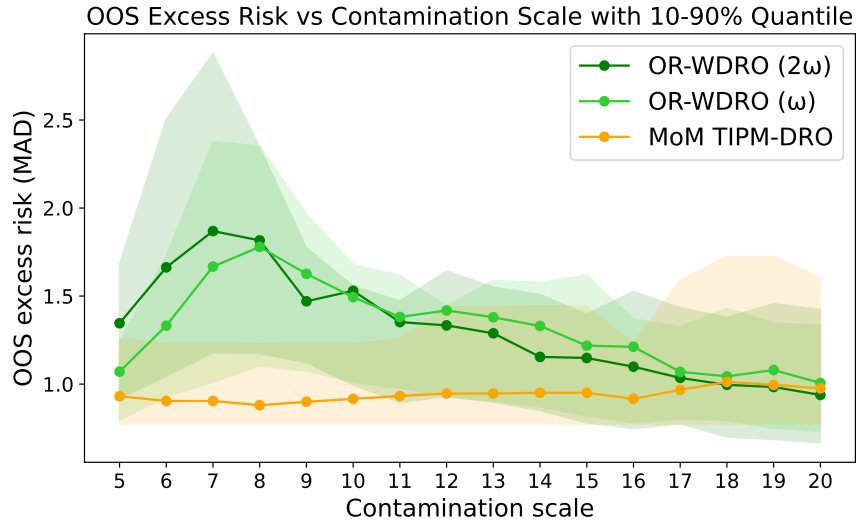
Figure 3 reports the OOS excess risk under the MAD loss as a function of the training sample size  $N$  for four representative dimensions  $D_\xi$ , based on 50 trials and a clean test set of size 1000. Standard WDRO barely improves with  $N$ , because its ambiguity set remains centered at the



**Figure 3** OOS excess risk under the MAD loss versus the training sample size  $N$ , with panels ordered by increasing dimension  $d$  and 10th–90th percentile bands.

contaminated empirical distribution and therefore keeps fitting the outliers as more corrupted data are observed. OR-WDRO works reasonably well in low dimension, where transport costs can still separate inliers from outliers, but this advantage fades as  $D_\xi$  grows and Wasserstein distances concentrate. MoM TIPM-DRO, in contrast, stays stable across dimensions and becomes clearly best at  $D_\xi = 50$ . The message of this figure is that transport-based outlier handling still depends on having informative geometry, whereas the MoM reference is much less sensitive to dimension. This is also consistent with our theory: the MoM contamination tolerance only requires  $\omega < 1/2$ , so it does not worsen with ambient dimension.

To isolate the effect of outlier severity, Figure 4 fixes  $N = 50$  and  $D_\xi = 10$  and sweeps  $C$  from 5 to 20. For visual clarity, this figure displays only OR-WDRO( $2\omega$ ), OR-WDRO( $\omega$ ), and MoM TIPM-DRO. OR-WDRO shows a clear hump at moderate contamination scales: the outliers are already harmful enough to distort the empirical reference, but not yet far enough away in transport cost to be trimmed reliably. Only when  $C$  becomes very large does OR-WDRO start to recover. MoM



**Figure 4** OOS excess risk under the MAD loss versus the contamination scale  $C$ , with 10th–90th percentile bands over the trials.

TIPM-DRO stays nearly flat throughout, showing that its performance is much less sensitive to how severe the leverage contamination is. Taken together, the two figures show that MoM TIPM-DRO does not just improve over standard WDRO. It is also more stable than outlier-aware WDRO baselines exactly when transport geometry is least reliable: at higher ambient dimension and at moderate contamination scales.

## 6. Concluding Remarks

We introduced a distributionally robust optimization framework based on targeted integral probability metrics. By defining the ambiguity set directly through the decision-induced loss class, the proposed expected hinge formulation yields a task-aware ambiguity set that is equivalent to a semi-infinite family of loss-discrepancy constraints. This construction enables finite-sample guarantees that bypass the curse of dimensionality whenever a scalar pointwise concentration inequality is available, thereby extending dimension-free DRO guarantees to a broad range of nonstandard data regimes. In many such settings, the framework also simplifies implementation, as it avoids the need to handcraft geometry-driven ambiguity sets. For example, one need not specify a ground cost as in Wasserstein DRO when the uncertainty contains mixed discrete and continuous features, nor combine different discrepancy measures, such as Wasserstein and total variation, to handle outlier contamination. Instead, the framework only requires a loss function that captures how uncertainty affects the decision problem and an estimator that admits an appropriate pointwise guarantee.

Our numerical results highlight the practical value of aligning the ambiguity set with the downstream loss. In both inventory management and outlier-corrupted regression, the proposed frame-

work delivers strong out-of-sample performance and robustness relative to transport-based baselines. Moreover, the experiments suggest that the framework may remain effective even beyond the sub-Weibull regime, thereby motivating further theoretical investigation. Several directions remain for future research. These include extending the framework to heavier-tailed models such as Pareto-type distributions, developing adaptive sampling schemes that may improve the Monte Carlo approximation rate, and designing scalable algorithms for richer multistage decision problems.

## References

- Idan Attias and Aryeh Kontorovich. Fat-shattering dimension of k-fold aggregations. *Journal of Machine Learning Research*, 25(144):1–29, 2024.
- Waïss Azizian, Franck Iutzeler, and Jérôme Malick. Exact generalization guarantees for (regularized) Wasserstein distributionally robust models. *Advances in Neural Information Processing Systems*, 36:14584–14596, 2023.
- A. Ben-Tal, D. den Hertog, A. De Waegenaere, B. Melenberg, and G. Rennen. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 2012. Forthcoming.
- Aharon Ben-Tal, Dick Den Hertog, Anja De Waegenaere, Bertrand Melenberg, and Gijs Rennen. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 59(2):341–357, 2013.
- Aharon Ben-Tal, Omar El Housni, and Vineet Goyal. A tractable approach for designing piecewise affine policies in two-stage adjustable robust optimization. *Mathematical Programming*, 182(1):57–102, 2020.
- Dimitris Bertsimas, Vishal Gupta, and Nathan Kallus. Data-driven robust optimization. *Mathematical Programming*, 167(2):235–292, 2018.
- Patrick Billingsley. *Convergence of probability measures*. John Wiley & Sons, 2013.
- Kostas Bimpikis and Mihalis G. Markakis. Inventory pooling under heavy-tailed demand. *Management Science*, 62(6):1800–1813, 2016. doi: 10.1287/mnsc.2015.2204.
- Jeremiah Birrell, Paul Dupuis, Markos A Katsoulakis, Yannis Pantazis, and Luc Rey-Bellet. (f, gamma)-divergences: Interpolating between f-divergences and integral probability metrics. *Journal of machine learning research*, 23(39):1–70, 2022.
- Jose Blanchet and Yang Kang. Sample out-of-sample inference based on wasserstein distance. *Operations Research*, 69(3):985–1013, 2021.
- Jose Blanchet and Karthyek Murthy. Quantifying distributional model risk via optimal transport. *Mathematics of Operations Research*, 44(2):565–600, 2019.
- Jose Blanchet, Yang Kang, and Karthyek Murthy. Robust wasserstein profile inference and applications to machine learning. *Journal of Applied Probability*, 56(3):830–857, 2019.

- Jose Blanchet, Karthyek Murthy, and Nian Si. Confidence regions in wasserstein distributionally robust estimation. *Biometrika*, 109(2):295–315, 2022.
- Jose Blanchet, Jiajin Li, Markus Pelger, and Greg Zanotti. Automatic outlier rectification via optimal transport. *Advances in Neural Information Processing Systems*, 37:35313–35357, 2024.
- Geunyeong Byeon. Comparative analysis of two-stage distributionally robust optimization over 1-Wasserstein and 2-Wasserstein balls. *arXiv preprint arXiv:2501.05619*, 2025.
- Ruidi Chen and Ioannis Ch Paschalidis. A robust learning approach for regression models based on distributionally robust optimization. *Journal of Machine Learning Research*, 19(13):1–48, 2018.
- Zhi Chen, Melvyn Sim, and Huan Xu. Distributionally robust optimization with infinitely constrained ambiguity sets. *Operations Research*, 67(5):1328–1344, 2019.
- Erick Delage and Yinyu Ye. Distributionally robust optimization under moment uncertainty with application to data-driven problems. *Operations research*, 58(3):595–612, 2010.
- John Duchi and Hongseok Namkoong. Variance-based regularization with convex objectives. *The Journal of Machine Learning Research*, 20(1):2450–2504, 2019.
- John C Duchi, Peter W Glynn, and Hongseok Namkoong. Statistics of robust optimization: A generalized empirical likelihood approach. *Mathematics of Operations Research*, 46(3):946–969, 2021.
- Jianqing Fan, Bai Jiang, and Qiang Sun. Hoeffding’s inequality for general markov chains and its applications to statistical learning. *Journal of Machine Learning Research*, 22(139):1–35, 2021.
- Xiangyi Fan and Grani A Hanasusanto. A decision rule approach for two-stage data-driven distributionally robust optimization problems with random recourse. *INFORMS Journal on Computing*, 36(2):526–542, 2024.
- Nicolas Fournier and Arnaud Guillin. On the rate of convergence in Wasserstein distance of the empirical measure. *Probability Theory and Related Fields*, 162(3):707–738, 2015.
- Rui Gao. Finite-sample guarantees for wasserstein distributionally robust optimization: Breaking the curse of dimensionality. *Operations Research*, 71(6):2291–2306, 2023.
- Rui Gao and Anton Kleywegt. Distributionally robust stochastic optimization with Wasserstein distance. *Mathematics of Operations Research*, 48(2):603–655, 2023.
- Rui Gao, Xi Chen, and Anton J Kleywegt. Wasserstein distributionally robust optimization and variation regularization. *Operations Research*, 72(3):1177–1191, 2024.
- Angelos Georghiou, Wolfram Wiesemann, and Daniel Kuhn. Generalized decision rule approximations for stochastic programming via liftings. *Mathematical Programming*, 152(1):301–338, 2015.
- Joel Goh and Melvyn Sim. Distributionally robust optimization and its tractable approximations. *Operations research*, 58(4-part-1):902–917, 2010.

- Lee-Ad Gottlieb, Aryeh Kontorovich, and Robert Krauthgamer. Adaptive metric dimensionality reduction. *Theoretical Computer Science*, 620:105–118, 2016.
- Nathael Gozlan and Christian Léonard. Transport inequalities. a survey. *arXiv preprint arXiv:1003.3852*, 2010.
- László Györfi, Michael Kohler, Adam Krzyżak, Harro Walk, et al. *A Distribution-Free Theory of Nonparametric Regression*, volume 1. Springer, 2002.
- Grani A Hanasusanto and Daniel Kuhn. Conic programming reformulations of two-stage distributionally robust linear programs over wasserstein balls. *Operations Research*, 66(3):849–869, 2018.
- Haiyun He and Ziv Goldfeld. Information-theoretic generalization bounds for deep neural networks. *IEEE Transactions on Information Theory*, 2025.
- Hisham Husain. Distributional robustness with IPMs and links to regularization and GANs. *Advances in Neural Information Processing Systems*, 33:11816–11827, 2020.
- Garud Iyengar, Henry Lam, and Tianyu Wang. Hedging against complexity: Distributionally robust optimization with parametric approximation. In *International Conference on Artificial Intelligence and Statistics*, pages 9976–10011. PMLR, 2023.
- Zhuangzhuang Jia, Grani A Hanasusanto, Phebe Vayanos, and Weijun Xie. Learning fair policies for multi-stage selection problems from observational data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 21188–21196, 2024.
- Nan Jiang and Weijun Xie. Distributionally favorable optimization: A framework for data-driven decision-making with endogenous outliers. *SIAM Journal on Optimization*, 34(1):419–458, 2024.
- Arun K. Kuchibhotla and Abhishek Chakraborty. Moving beyond sub-gaussianity in high-dimensional statistics: applications in covariance estimation and linear regression. *Information and Inference: A Journal of the IMA*, 11(4):1389–1456, 2022. doi: 10.1093/imaia/iaac012.
- Daniel Kuhn, Soroosh Shafiee, and Wolfram Wiesemann. Distributionally robust optimization. *Acta Numerica*, 34:579–804, 2025.
- Pierre Laforgue, Guillaume Staerman, and Stephan Cléménçon. Generalization bounds in the presence of outliers: a median-of-means study. In *International conference on machine learning*, pages 5937–5947. PMLR, 2021.
- Henry Lam. Recovering best statistical guarantees via the empirical divergence-based distributionally robust optimization. *Operations Research*, 67(4):1090–1105, 2019.
- Tam Le and Jérôme Malick. Universal generalization guarantees for wasserstein distributionally robust models. *arXiv preprint arXiv:2402.11981*, 2024.
- Jaeho Lee and Maxim Raginsky. Minimax statistical learning with Wasserstein distances. *Advances in Neural Information Processing Systems*, 31, 2018.

- Mengmeng Li, Tobias Sutter, and Daniel Kuhn. Distributionally robust optimization with markovian data. In *International Conference on Machine Learning*, pages 6493–6503. PMLR, 2021.
- Adrian Tovar Lopez and Varun Jog. Generalization error bounds using wasserstein distances. In *2018 IEEE Information Theory Workshop (ITW)*, pages 1–5. IEEE, 2018.
- Peyman Mohajerin Esfahani and Daniel Kuhn. Data-driven distributionally robust optimization using the wasserstein metric: Performance guarantees and tractable reformulations. *Mathematical Programming*, 171(1):115–166, 2018.
- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- Dheeraj Nagaraj, Xian Wu, Guy Bresler, Prateek Jain, and Praneeth Netrapalli. Least squares regression with markovian data: Fundamental limits and algorithms. *Advances in neural information processing systems*, 33:16666–16676, 2020.
- Amir Najafi, Shin-ichi Maeda, Masanori Koyama, and Takeru Miyato. Robustness to adversarial perturbations in learning from incomplete data. *Advances in Neural Information Processing Systems*, 32, 2019.
- Sloan Nietert, Ziv Goldfeld, and Soroosh Shafiee. Outlier-robust wasserstein dro. *Advances in Neural Information Processing Systems*, 36:62792–62820, 2023.
- Sloan Nietert, Ziv Goldfeld, and Soroosh Shafiee. Robust distribution learning with local and global adversarial corruptions. In *The Thirty Seventh Annual Conference on Learning Theory*, pages 4007–4008. PMLR, 2024.
- R. T. Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- Daniel Russo and James Zou. How much does your data exploration overfit? controlling bias via information usage. *IEEE Transactions on Information Theory*, 66(1):302–323, 2019.
- Utsav Sadana, Abhilash Chenreddy, Erick Delage, Alexandre Forel, Emma Frejinger, and Thibaut Vidal. A survey of contextual optimization methods for decision-making under uncertainty. *European Journal of Operational Research*, 320(2):271–289, 2025.
- Soroosh Shafieezadeh-Abadeh, Daniel Kuhn, and Peyman Mohajerin Esfahani. Regularization via mass transportation. *Journal of Machine Learning Research*, 20(103):1–68, 2019.
- S. Shalev-Shwartz and S. Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014.
- A. Shapiro. On duality theory of conic linear problems. In M.A. Goberna and M.A. Lopez, editors, *Semi-Infinite Programming: Recent Advances*, pages 135–165. Kluwer Academic Publishers, 2001.
- Aman Sinha, Hongseok Namkoong, Riccardo Volpi, and John Duchi. Certifying some distributional robustness with principled adversarial training. *arXiv preprint arXiv:1710.10571*, 2017.

- Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*, volume 47. Cambridge university press, 2018.
- Roman Vershynin. High-dimensional probability, 2025.
- Cédric Villani et al. *Optimal Transport: Old and New*, volume 338. Springer, 2008.
- Hao Wang, Mario Diaz, José Cândido S Santos Filho, and Flavio P Calmon. An information-theoretic view of generalization via wasserstein distance. In *2019 IEEE International Symposium on Information Theory (ISIT)*, pages 577–581. IEEE, 2019.
- Hao Wang, Yizhe Huang, Rui Gao, and Flavio Calmon. Analyzing the generalization capability of SGLD using properties of Gaussian channels. *Advances in Neural Information Processing Systems*, 34:24222–24234, 2021.
- Yijie Wang, Grani A Hanasusanto, and Chin Pang Ho. Generalization bounds for contextual stochastic optimization using kernel regression. *arXiv preprint arXiv:2407.10764*, 2024.
- Wolfram Wiesemann, Daniel Kuhn, and Melvyn Sim. Distributionally robust convex optimization. *Operations research*, 62(6):1358–1376, 2014.
- Qinyu Wu, Jonathan Yu-Meng Li, and Tiantian Mao. On generalization and regularization via Wasserstein distributionally robust optimization. *Management Science*, 2025.
- Weijun Xie, Jie Zhang, and Shabbir Ahmed. Distributionally robust bottleneck combinatorial problems: uncertainty quantification and robust decision making. *Mathematical Programming*, 196(1):597–640, 2022.
- Aolin Xu and Maxim Raginsky. Information-theoretic analysis of generalization capability of learning algorithms. *Advances in neural information processing systems*, 30, 2017.
- Guanglin Xu and Samuel Burer. A copositive approach for two-stage adjustable robust optimization with uncertain right-hand sides. *Computational Optimization and Applications*, 70(1):33–59, 2018.
- Guanglin Xu and Grani A Hanasusanto. Improved decision rule approximations for multistage robust optimization via copositive programming. *Operations Research*, 73(2):842–861, 2025.
- Yibo Zeng and Henry Lam. Generalization bounds with minimal dependency on hypothesis class via distributionally robust optimization. *Advances in Neural Information Processing Systems*, 35:27576–27590, 2022.
- Jia-Jie Zhu, Wittawat Jitkrittum, Moritz Diehl, and Bernhard Schölkopf. Kernel distributionally robust optimization: Generalized duality theorem and stochastic approximation. In *International Conference on Artificial Intelligence and Statistics*, pages 280–288. PMLR, 2021.

## Appendix A: Proofs of Section 2

### A.1. Ambiguity Set Characterization

*Proof of Proposition 2.* It suffices to show that the semi-infinite constraint

$$\mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq \epsilon \quad \forall \mathbf{z} \in \mathcal{X} \quad (29)$$

is equivalent to the expectation constraint

$$\mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \right] \leq 0. \quad (30)$$

If  $\mathbb{Q}$  is feasible to (29), then  $[\mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \epsilon]_+ \leq 0$  for all  $\mathbf{z} \in \mathcal{X}$ . Taking expectation with respect to  $\mathbb{Z}$  yields (30).

Conversely, suppose  $\mathbb{Q}$  is infeasible to (29); that is, there exists  $\mathbf{z}' \in \mathcal{X}$  such that  $[\mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}', \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}', \tilde{\boldsymbol{\xi}})] - \epsilon]_+ > 0$ . By Assumption (A), for any  $\mathbb{Q}$  feasible to (4),

$$\mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq c_1 + c_2 \mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|] \leq c_1 + c_2 \sqrt{\mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2]} \leq c_1 + c_2 \sqrt{\Omega} < +\infty \quad \forall \mathbf{z} \in \mathcal{X},$$

where the second inequality follows from Jensen's inequality. Since  $\hat{\mathbb{P}}_N$  is also feasible to (4) by construction and by Assumption (B), the dominated convergence theorem implies that  $[\mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \epsilon]_+$  is continuous in  $\mathbf{z}$ . Hence there exists a neighborhood  $\mathcal{Z}_\tau = \{\mathbf{z} \in \mathcal{X} : \|\mathbf{z} - \mathbf{z}'\| < \tau\}$  in which  $[\mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \epsilon]_+ > 0$  for all  $\mathbf{z} \in \mathcal{Z}_\tau$ . We therefore have

$$\mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \right] \geq \mathbb{Z}(\mathbf{z} \in \mathcal{Z}_\tau) \mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \mid \mathbf{z} \in \mathcal{Z}_\tau \right] > 0,$$

because  $\mathbb{Z}$  has full support on  $\mathcal{X}$ . Hence  $\mathbb{Q}$  is infeasible to (30). Replacing the semi-infinite constraint in (4) with the equivalent expectation constraint (30) yields the claim.  $\square$

### A.2. Performance Guarantees

LEMMA 1. For any  $\mathbb{Q}$  such that  $\mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega$ , the following Lipschitz condition holds:

$$\mathbb{E}_{\mathbb{Q}} \left[ \left| \ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) - \ell(\mathbf{z}', \tilde{\boldsymbol{\xi}}) \right| \right] \leq L \|\mathbf{z} - \mathbf{z}'\| \quad \forall \mathbf{z}, \mathbf{z}' \in \mathcal{X},$$

where  $L := \max_{j \in [J]} \|\mathbf{A}_j\|_{\text{op}} \sqrt{\Omega} + \max_{j \in [J]} \|\mathbf{b}_j\|$ .

*Proof of Lemma 1.* For any  $\mathbf{z}, \mathbf{z}' \in \mathcal{X}$ , since  $\mathbf{a}_j(\mathbf{z}) - \mathbf{a}_j(\mathbf{z}') = \mathbf{A}_j(\mathbf{z} - \mathbf{z}')$  and  $b_j(\mathbf{z}) - b_j(\mathbf{z}') = \mathbf{b}_j^\top(\mathbf{z} - \mathbf{z}')$ , we have

$$\left| \ell(\mathbf{z}, \boldsymbol{\xi}) - \ell(\mathbf{z}', \boldsymbol{\xi}) \right| \leq \max_{j \in [J]} \left| (\mathbf{A}_j(\mathbf{z} - \mathbf{z}'))^\top \boldsymbol{\xi} + \mathbf{b}_j^\top(\mathbf{z} - \mathbf{z}') \right|$$

$$\begin{aligned}
&\leq \max_{j \in [J]} (\|\mathbf{A}_j\|_{\text{op}} \|\boldsymbol{\xi}\| + \|\mathbf{b}_j\|) \cdot \|\mathbf{z} - \mathbf{z}'\| \\
&\leq \left( \max_{j \in [J]} \|\mathbf{A}_j\|_{\text{op}} \|\boldsymbol{\xi}\| + \max_{j \in [J]} \|\mathbf{b}_j\| \right) \cdot \|\mathbf{z} - \mathbf{z}'\|.
\end{aligned}$$

Taking expectation under  $\mathbb{Q}$  yields

$$\begin{aligned}
\mathbb{E}_{\mathbb{Q}} \left[ |\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) - \ell(\mathbf{z}', \tilde{\boldsymbol{\xi}})| \right] &\leq \left( \max_{j \in [J]} \|\mathbf{A}_j\|_{\text{op}} \mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|] + \max_{j \in [J]} \|\mathbf{b}_j\| \right) \|\mathbf{z} - \mathbf{z}'\| \\
&\leq \left( \max_{j \in [J]} \|\mathbf{A}_j\|_{\text{op}} \sqrt{\Omega} + \max_{j \in [J]} \|\mathbf{b}_j\| \right) \|\mathbf{z} - \mathbf{z}'\|.
\end{aligned}$$

This is the desired bound.  $\square$

*Proof of Proposition 3.* Let

$$\eta_N := \frac{1}{L\sqrt{N}},$$

and let  $\mathcal{X}_{\eta_N}$  be an  $\eta_N$ -net of  $\mathcal{X}$  under  $\|\cdot\|$ . Since  $\mathcal{X} \subseteq \mathbb{B}(0, R_{\mathbf{x}})$  is compact, the standard volumetric bound gives

$$|\mathcal{X}_{\eta_N}| \leq \left(1 + \frac{2R_{\mathbf{x}}}{\eta_N}\right)^{D_{\mathbf{x}}} = \left(1 + 2R_{\mathbf{x}}L\sqrt{N}\right)^{D_{\mathbf{x}}}.$$

For each  $\mathbf{z}_0 \in \mathcal{X}_{\eta_N}$ , applying (9) with confidence level

$$\delta' := \frac{\delta}{|\mathcal{X}_{\eta_N}|}$$

and then taking a union bound over  $\mathcal{X}_{\eta_N}$ , we obtain with probability at least  $1 - \delta$ ,

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}_0, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}_0, \tilde{\boldsymbol{\xi}})] \leq \frac{1}{\sqrt{N}} \varepsilon \left( \log \frac{1}{\delta'} \right) \quad \forall \mathbf{z}_0 \in \mathcal{X}_{\eta_N}.$$

Because  $\varepsilon$  is nondecreasing and

$$\log \frac{1}{\delta'} = \log |\mathcal{X}_{\eta_N}| + \log \frac{1}{\delta} \leq D_{\mathbf{x}} \log \left(1 + 2R_{\mathbf{x}}L\sqrt{N}\right) + \log \frac{1}{\delta},$$

the same event implies

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}_0, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}_0, \tilde{\boldsymbol{\xi}})] \leq \frac{1}{\sqrt{N}} \varepsilon \left( D_{\mathbf{x}} \log \left(1 + 2R_{\mathbf{x}}L\sqrt{N}\right) + \log \frac{1}{\delta} \right) \quad \forall \mathbf{z}_0 \in \mathcal{X}_{\eta_N}. \quad (31)$$

Fix any  $\mathbf{z} \in \mathcal{X}$  and choose  $\mathbf{z}_0 \in \mathcal{X}_{\eta_N}$  with  $\|\mathbf{z} - \mathbf{z}_0\| \leq \eta_N$ . By Lemma 1, applied once under  $\mathbb{P}^*$  and once under  $\hat{\mathbb{P}}_N$  (the latter is valid by Assumption (B)), we have

$$\begin{aligned}
\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] &\leq \mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}_0, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}_0, \tilde{\boldsymbol{\xi}})] \\
&\quad + \mathbb{E}_{\mathbb{P}^*}[|\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) - \ell(\mathbf{z}_0, \tilde{\boldsymbol{\xi}})|] + \mathbb{E}_{\hat{\mathbb{P}}_N}[|\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) - \ell(\mathbf{z}_0, \tilde{\boldsymbol{\xi}})|] \\
&\leq \mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}_0, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}_0, \tilde{\boldsymbol{\xi}})] + 2L\eta_N.
\end{aligned}$$

Combining this bound with (31) and using  $2L\eta_N = 2/\sqrt{N}$  yields

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq \epsilon \quad \forall \mathbf{z} \in \mathcal{X}$$

with probability at least  $1 - \delta$ , where  $\epsilon$  is defined in (10). Proposition 2 then implies  $\mathbb{P}^* \in \mathcal{P}_{\epsilon}$ .  $\square$

*Proof of Corollary 1.* On the event  $\{\mathbb{P}^* \in \mathcal{P}_\epsilon\}$ , which occurs with probability at least  $1 - \delta$  by Proposition 3, we have

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})] \leq \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})] = \hat{J}.$$

This is precisely the claimed out-of-sample guarantee.  $\square$

*Proof of Theorem 1.* Let

$$\mathcal{E}_1 := \{\mathbb{P}^* \in \mathcal{P}_\epsilon\} \quad \text{and} \quad \mathcal{E}_2 := \left\{ \left| \mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{x}^*, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{x}^*, \tilde{\boldsymbol{\xi}})] \right| \leq \frac{1}{\sqrt{N}} \epsilon \left( \log \frac{1}{\delta} \right) \right\}.$$

By Proposition 3,  $\text{Prob}(\mathcal{E}_1) \geq 1 - \delta$ , and by (12) applied at the fixed decision  $\mathbf{x}^*$ ,  $\text{Prob}(\mathcal{E}_2) \geq 1 - \delta$ . Hence

$$\text{Prob}(\mathcal{E}_1 \cap \mathcal{E}_2) \geq 1 - 2\delta.$$

On  $\mathcal{E}_1 \cap \mathcal{E}_2$ , we have

$$\begin{aligned} \mathbb{E}_{\mathbb{P}^*}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})] &\leq \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})] \\ &\leq \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}^*, \tilde{\boldsymbol{\xi}})] \\ &\leq \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{x}^*, \tilde{\boldsymbol{\xi}})] + \epsilon \\ &\leq \mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{x}^*, \tilde{\boldsymbol{\xi}})] + \frac{1}{\sqrt{N}} \epsilon \left( \log \frac{1}{\delta} \right) + \epsilon. \end{aligned}$$

Subtracting  $\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{x}^*, \tilde{\boldsymbol{\xi}})]$  from both sides yields the claim.  $\square$

### A.3. Applications

*Proof of Proposition 4.* Fix any  $\mathbf{z} \in \mathcal{X}$  and define the centered random variables

$$X_i := \mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \ell(\mathbf{z}, \hat{\boldsymbol{\xi}}_i), \quad i \in [N].$$

Then  $X_1, \dots, X_N$  are independent, mean-zero, and satisfy  $\|X_i\|_{\psi_\vartheta} \leq K$ . Applying Kuchibhotla and Chakraborty (2022, Theorem 3.1) with weights  $a_i = 1/N$ , we obtain

$$\text{Prob}\left( \left| \mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \right| \geq 2eC(\vartheta) \|b\| \sqrt{t} + 2eL_N^*(\vartheta) \|b\|_{\beta(\vartheta)} t^{1/\vartheta} \right) \leq 2e^{-t} \quad \forall t \geq 0,$$

where

$$b := \left( \frac{\|X_1\|_{\psi_\vartheta}}{N}, \dots, \frac{\|X_N\|_{\psi_\vartheta}}{N} \right), \quad \beta(\vartheta) := \begin{cases} \infty, & \vartheta \leq 1, \\ \frac{\vartheta}{\vartheta-1}, & \vartheta > 1. \end{cases}$$

Since  $\|X_i\|_{\psi_\vartheta} \leq K$ , we have

$$\|b\| \leq \frac{K}{\sqrt{N}}.$$

Moreover,

$$\|b\|_{\beta(\vartheta)} \leq \begin{cases} \|b\|_\infty \leq \frac{K}{N}, & \vartheta \leq 1, \\ \left( N \left( \frac{K}{N} \right)^{\beta(\vartheta)} \right)^{1/\beta(\vartheta)} = K N^{-1/\vartheta}, & \vartheta > 1, \end{cases}$$

and therefore

$$\|b\|_{\beta(\vartheta)} \leq \frac{K}{N^{\min\{1, 1/\vartheta\}}}.$$

Finally, Theorem 3.1 gives

$$L_N^*(\vartheta) = L_N(\vartheta) \frac{C(\vartheta)\|b\|}{\|b\|_{\beta(\vartheta)}} = \begin{cases} \frac{4^{1/\vartheta}}{\sqrt{2}} C(\vartheta), & \vartheta < 1, \\ \frac{4^{1/\vartheta}}{\sqrt{2}} 4e, & \vartheta \geq 1, \end{cases}$$

which depends only on  $\vartheta$ . Hence there exist constants  $c_{1,\vartheta}, c_{2,\vartheta} > 0$ , depending only on  $\vartheta$ , such that

$$\text{Prob}\left(\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \geq c_{1,\vartheta} K \sqrt{\frac{t}{N}} + c_{2,\vartheta} K \frac{t^{1/\vartheta}}{N^{\min\{1, 1/\vartheta\}}}\right) \leq 2e^{-t} \quad \forall t \geq 0.$$

We now absorb the heavy-tail correction term into the leading root- $N$  term.

If  $\vartheta = 2$ , both terms are already of the same order, so there exists  $C_\vartheta > 0$  such that

$$c_{1,\vartheta} K \sqrt{\frac{t}{N}} + c_{2,\vartheta} K \frac{t^{1/\vartheta}}{N^{1/\vartheta}} \leq C_\vartheta K \sqrt{\frac{t}{N}} \quad \forall t \geq 0.$$

If  $\vartheta \in [1, 2)$ , then  $1/\vartheta - 1/2 > 0$ . Therefore, there exists  $c_\vartheta > 0$ , depending only on  $\vartheta$ , such that whenever  $t \leq c_\vartheta N$ ,

$$c_{2,\vartheta} K \frac{t^{1/\vartheta}}{N^{1/\vartheta}} \leq c_{1,\vartheta} K \sqrt{\frac{t}{N}}.$$

Hence, for all  $t \leq c_\vartheta N$ ,

$$c_{1,\vartheta} K \sqrt{\frac{t}{N}} + c_{2,\vartheta} K \frac{t^{1/\vartheta}}{N^{1/\vartheta}} \leq 2c_{1,\vartheta} K \sqrt{\frac{t}{N}}.$$

If  $\vartheta \in (0, 1)$ , then there exists  $c_\vartheta > 0$ , depending only on  $\vartheta$ , such that whenever  $t \leq c_\vartheta N^{\vartheta/(2-\vartheta)}$ ,

$$c_{2,\vartheta} K \frac{t^{1/\vartheta}}{N} \leq c_{1,\vartheta} K \sqrt{\frac{t}{N}}.$$

Hence, for all  $t \leq c_\vartheta N^{\vartheta/(2-\vartheta)}$ ,

$$c_{1,\vartheta} K \sqrt{\frac{t}{N}} + c_{2,\vartheta} K \frac{t^{1/\vartheta}}{N} \leq 2c_{1,\vartheta} K \sqrt{\frac{t}{N}}.$$

Setting  $t = \log(2/\delta)$ , the preceding three cases imply that there exist constants  $C_\vartheta, c_\vartheta > 0$ , depending only on  $\vartheta$ , such that

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq C_\vartheta K \sqrt{\frac{\log(2/\delta)}{N}}$$

with probability at least  $1 - \delta$ , provided that

$$\delta \in \begin{cases} [2 \exp(-c_\vartheta N^{\vartheta/(2-\vartheta)}), 1), & \vartheta \in (0, 1), \\ [2 \exp(-c_\vartheta N), 1), & \vartheta \in [1, 2), \\ (0, 1), & \vartheta = 2. \end{cases}$$

Since  $\log(2/\delta) \leq 1 + \log(1/\delta)$  for all  $\delta \in (0, 1)$ , enlarging  $C_\vartheta$  if necessary yields (13).  $\square$

*Proof of Proposition 5.* By Fan et al. (2021, Theorem 1), for any  $s > 0$ ,

$$\text{Prob}\left(\sum_{i=1}^N \ell(\mathbf{z}, \hat{\boldsymbol{\xi}}_i) - N\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] < -s\right) \leq \exp\left(-\frac{1-\lambda}{1+\lambda} \cdot \frac{s^2}{N(\bar{\ell}-\underline{\ell})^2/2}\right).$$

Dividing by  $N$  and substituting  $s = Nt$  gives, for any  $t \geq 0$ ,

$$\text{Prob}\left(\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\tilde{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \geq t\right) \leq \exp\left(-\frac{2(1-\lambda)Nt^2}{(1+\lambda)(\bar{\ell}-\underline{\ell})^2}\right).$$

Equating the right-hand side to  $\delta$  and solving for  $t$  yields (14). The extension to the time-inhomogeneous setting follows by applying Theorem 5 in Fan et al. (2021).  $\square$

Following Jia et al. (2024), we make the following standard assumptions regarding the logging policy:

- (i) *Conditional exchangeability:* Conditional on observed covariates  $\boldsymbol{\chi}$ , the selection decision is independent of the outcomes  $\boldsymbol{\omega}$ :

$$\tilde{\boldsymbol{\omega}} \perp\!\!\!\perp \tilde{S} \mid \boldsymbol{\chi}.$$

- (ii) *Positivity:* The probability of a candidate being selected is strictly positive for any given covariate values:

$$\pi(\boldsymbol{\chi}) \geq \underline{\pi} \quad \forall \boldsymbol{\chi} \in \mathcal{X},$$

for some positive constant  $\underline{\pi} > 0$ .

**LEMMA 2.** *Suppose Assumptions (i) and (ii) hold. Then the inverse probability weighting estimator is equivalent to the true expectation:*

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] = \mathbb{E}_{\tilde{\mathbb{P}}} \left[ \ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) \frac{\mathbb{I}[\tilde{S} = 1]}{\tilde{\mathbb{P}}(\tilde{S} = 1 \mid \tilde{\boldsymbol{\chi}} = \boldsymbol{\chi})} \right].$$

*Proof of Lemma 2.* We have

$$\begin{aligned} \mathbb{E}_{\tilde{\mathbb{P}}} \left[ \ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) \frac{\mathbb{I}[\tilde{S} = 1]}{\pi(\tilde{\boldsymbol{\chi}})} \right] &= \mathbb{E}_{\tilde{\mathbb{P}}} \left[ \mathbb{E}_{\tilde{\mathbb{P}}} \left[ \ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) \frac{\mathbb{I}[\tilde{S} = 1]}{\pi(\tilde{\boldsymbol{\chi}})} \mid \tilde{\boldsymbol{\chi}}, \tilde{\boldsymbol{\omega}} \right] \right] \\ &= \mathbb{E}_{\tilde{\mathbb{P}}} \left[ \ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) \frac{\tilde{\mathbb{P}}(\tilde{S} = 1 \mid \tilde{\boldsymbol{\chi}}, \tilde{\boldsymbol{\omega}})}{\pi(\tilde{\boldsymbol{\chi}})} \right] \\ &= \mathbb{E}_{\tilde{\mathbb{P}}} \left[ \ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) \frac{\tilde{\mathbb{P}}(\tilde{S} = 1 \mid \tilde{\boldsymbol{\chi}})}{\pi(\tilde{\boldsymbol{\chi}})} \right] \\ &= \mathbb{E}_{\tilde{\mathbb{P}}} \left[ \ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) \right] \\ &= \mathbb{E}_{\mathbb{P}^*} \left[ \ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) \right]. \end{aligned}$$

The first equality is the tower property. The third equality follows from Assumption (i), the fourth uses the definition of  $\pi(\boldsymbol{\chi})$ , and the last holds because  $\mathbb{P}^*$  is the marginal law of  $(\tilde{\boldsymbol{\chi}}, \tilde{\boldsymbol{\omega}})$  under  $\tilde{\mathbb{P}}$ .

$\square$

The notation  $\mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})]$  in (15) denotes an IPW reference functional rather than an expectation under a normalized empirical probability measure. This distinction does not affect the dual reformulations: the generalized moment problem only uses these reference values on the right-hand side of the loss constraints, and the same duality argument applies under the corresponding Slater condition.

*Proof of Proposition 7.* By positivity,

$$\left| \frac{\mathbb{I}[\tilde{S} = 1]}{\overline{\mathbb{P}}(\tilde{S} = 1 | \tilde{\boldsymbol{\chi}} = \boldsymbol{\chi})} \right| \leq \frac{1}{\pi}.$$

Hence, under the assumed bound  $\|\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})\|_{\psi_2} \leq \sigma$ ,

$$\left\| \ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) \frac{\mathbb{I}[\tilde{S} = 1]}{\overline{\mathbb{P}}(\tilde{S} = 1 | \tilde{\boldsymbol{\chi}} = \boldsymbol{\chi})} \right\|_{\psi_2} \leq \frac{\sigma}{\pi}.$$

Centering changes the sub-Gaussian norm by at most a factor of two, so

$$\left\| \ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) \frac{\mathbb{I}[\tilde{S} = 1]}{\overline{\mathbb{P}}(\tilde{S} = 1 | \tilde{\boldsymbol{\chi}} = \boldsymbol{\chi})} - \mathbb{E}_{\overline{\mathbb{P}}} \left[ \ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) \frac{\mathbb{I}[\tilde{S} = 1]}{\overline{\mathbb{P}}(\tilde{S} = 1 | \tilde{\boldsymbol{\chi}} = \boldsymbol{\chi})} \right] \right\|_{\psi_2} \leq \frac{2\sigma}{\pi}.$$

Let

$$Y_z := \ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) \frac{\mathbb{I}[\tilde{S} = 1]}{\overline{\mathbb{P}}(\tilde{S} = 1 | \tilde{\boldsymbol{\chi}} = \boldsymbol{\chi})}, \quad Y_{z,n} := \ell(\mathbf{z}, \boldsymbol{\xi}_n) \frac{\mathbb{I}[S_n = 1]}{\pi(\boldsymbol{\chi}_n)}.$$

Then (15) equals  $N^{-1} \sum_{n=1}^N Y_{z,n}$ , and the preceding bound applies to the centered variable  $Y_z - \mathbb{E}_{\overline{\mathbb{P}}}[Y_z]$ . By Lemma 2,  $\mathbb{E}_{\overline{\mathbb{P}}}[Y_z] = \mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})]$ . Therefore, the estimation error can be written as

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] = -\frac{1}{N} \sum_{n=1}^N (Y_{z,n} - \mathbb{E}_{\overline{\mathbb{P}}}[Y_z]).$$

The summands in this display are independent, centered, and sub-Gaussian with  $\psi_2$  norm bounded by  $2\sigma/\pi$ . The negative average satisfies the same one-sided concentration bound, so taking the deviation level proportional to  $\sqrt{\log(1/\delta)/N}$  and applying the standard concentration inequality for averages of independent centered sub-Gaussian variables (Vershynin 2018, Theorem 2.6.3) yields the stated bound

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq C_{\text{IPW}} \frac{\sigma}{\pi} \sqrt{\frac{\log(1/\delta)}{N}}$$

with the numerical constants absorbed into the universal constant  $C_{\text{IPW}}$ .  $\square$

*Proof of Proposition 8.* Based on (Wang et al. 2024, Proposition 1), for any  $h > 0$  and  $t > 0$  we have

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) | \tilde{\mathbf{c}} = \mathbf{c}] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}}) | \tilde{\mathbf{c}} = \mathbf{c}] \leq L_c h + t$$

with probability at least  $1 - \exp(-Ncf\underline{h}^{D_c}t^2/2)$ . Setting  $\delta := \exp(-Ncf\underline{h}^{D_c}t^2/2)$  gives

$$\mathbb{E}_{\mathbb{P}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})|\tilde{\mathbf{c}} = \mathbf{c}] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})|\tilde{\mathbf{c}} = \mathbf{c}] \leq L_c h + \sqrt{\frac{2}{Ncf\underline{h}^{D_c}} \log\left(\frac{1}{\delta}\right)}$$

with probability at least  $1 - \delta$ . The choice

$$h := \left(\frac{2 \log(1/\delta)}{Ncf}\right)^{\frac{1}{D_c+2}} \left(\frac{D_c}{2L_c}\right)^{\frac{2}{D_c+2}}$$

minimizes the right-hand side. At this bandwidth, the two terms satisfy

$$\sqrt{\frac{2}{Ncf\underline{h}^{D_c}} \log\left(\frac{1}{\delta}\right)} = \frac{2L_c h}{D_c},$$

and hence the right-hand side equals

$$\frac{D_c + 2}{2(D_c/2)^{\frac{D_c}{D_c+2}}} L_c^{\frac{D_c}{D_c+2}} \left(\frac{2 \log(1/\delta)}{Ncf}\right)^{\frac{1}{D_c+2}},$$

which is the stated bound.  $\square$

#### A.4. Dual Formulations

*Proof of Theorem 2.* Fix  $\mathbf{x} \in \mathcal{X}$ . By Proposition 2, the inner maximization problem in (3) can be written as

$$\begin{aligned} & \sup_{\mathbb{Q} \in \mathcal{M}_+(\Xi)} \int_{\Xi} \ell(\mathbf{x}, \boldsymbol{\xi}) \mathbb{Q}(d\boldsymbol{\xi}) \\ & \text{s.t.} \quad \int_{\Xi} \ell(\mathbf{z}, \boldsymbol{\xi}) \mathbb{Q}(d\boldsymbol{\xi}) \leq \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] + \epsilon \quad \forall \mathbf{z} \in \mathcal{X}, \\ & \quad \int_{\Xi} \|\boldsymbol{\xi}\|^2 \mathbb{Q}(d\boldsymbol{\xi}) \leq \Omega, \\ & \quad \int_{\Xi} 1 \mathbb{Q}(d\boldsymbol{\xi}) = 1. \end{aligned}$$

This is a generalized moment problem over nonnegative measures. Its conic dual introduces a nonnegative measure  $\nu \in \mathcal{M}_+(\mathcal{X})$  for the family of loss constraints, a scalar multiplier  $\beta \in \mathbb{R}_+$  for the second-moment constraint, and a free scalar  $\alpha \in \mathbb{R}$  for the normalization constraint, yielding

$$\begin{aligned} & \min \alpha + \beta \Omega + \int_{\mathcal{X}} \left(\mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] + \epsilon\right) \nu(d\mathbf{z}) \\ & \text{s.t.} \quad \ell(\mathbf{x}, \boldsymbol{\xi}) \leq \alpha + \beta \|\boldsymbol{\xi}\|^2 + \int_{\mathcal{X}} \ell(\mathbf{z}, \boldsymbol{\xi}) \nu(d\mathbf{z}) \quad \forall \boldsymbol{\xi} \in \Xi. \end{aligned}$$

Because  $\epsilon > 0$  and Assumption (B) gives  $\mathbb{E}_{\hat{\mathbb{P}}_N}[\|\tilde{\boldsymbol{\xi}}\|^2] < \Omega$ , the measure  $\hat{\mathbb{P}}_N$  is strictly feasible for the primal inner problem. Standard strong duality for generalized moment problems therefore applies, and dual attainment holds for each fixed  $\mathbf{x}$ ; see Shapiro (2001, Proposition 5.2).

Now let

$$F(\mathbf{x}) := \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})].$$

The mapping  $\mathbf{x} \mapsto F(\mathbf{x})$  is lower semicontinuous as a pointwise supremum of continuous functions, and  $\mathcal{X}$  is compact, so the outer minimization attains its optimum. Combining this minimizer with the attained inner dual solution yields the stated infinite linear programming reformulation and the existence of an optimal tuple  $(\mathbf{x}, \alpha, \beta, \nu)$ .  $\square$

*Proof of Proposition 9.* Fix  $\mathbf{x} \in \mathcal{X}$ , and define

$$f_\epsilon(\mathbf{x}) := \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})].$$

Because  $\mathcal{P}_0 \subseteq \mathcal{P}_\epsilon$  for every  $\epsilon \geq 0$ , we have  $f_0(\mathbf{x}) \leq f_\epsilon(\mathbf{x})$ . Moreover, by Proposition 2,

$$f_0(\mathbf{x}) = \sup_{\mathbb{Q} \in \mathcal{P}_0} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] = \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})],$$

since  $\hat{\mathbb{P}}_N \in \mathcal{P}_0$  and every  $\mathbb{Q} \in \mathcal{P}_0$  satisfies  $\mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] \leq \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})]$ .

It remains to prove  $\limsup_{\epsilon \downarrow 0} f_\epsilon(\mathbf{x}) \leq f_0(\mathbf{x})$ . Let  $\epsilon_k \downarrow 0$  and choose  $\mathbb{Q}_k \in \mathcal{P}_{\epsilon_k}$  such that

$$f_{\epsilon_k}(\mathbf{x}) - \frac{1}{k} \leq \mathbb{E}_{\mathbb{Q}_k}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})].$$

Because every  $\mathbb{Q}_k$  satisfies the common second-moment bound, Markov's inequality implies that the family  $\{\mathbb{Q}_k\}_{k \in \mathbb{N}}$  is tight; see (Billingsley 2013, Section 5). After passing to a subsequence if necessary, Prohorov's theorem (Billingsley 2013, Theorem 5.1) yields  $\mathbb{Q}_k \Rightarrow \mathbb{Q}^*$  for some probability measure  $\mathbb{Q}^*$ . By Assumption (A) and the uniform second-moment bound, the family  $\{\ell(\mathbf{x}, \cdot)\}$  is uniformly integrable under  $\{\mathbb{Q}_k\}_{k \in \mathbb{N}}$ . Since  $\ell(\mathbf{x}, \cdot)$  is continuous, the mapping theorem (Billingsley 2013, Theorem 2.7) implies weak convergence of the pushforward laws, and (Billingsley 2013, Theorem 3.5) therefore gives

$$\lim_{k \rightarrow \infty} \mathbb{E}_{\mathbb{Q}_k}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] = \mathbb{E}_{\mathbb{Q}^*}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})].$$

The same argument holds for every fixed  $\mathbf{z} \in \mathcal{X}$ . Since  $\mathbb{Q}_k \in \mathcal{P}_{\epsilon_k}$ ,

$$\mathbb{E}_{\mathbb{Q}_k}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] + \epsilon_k \quad \forall \mathbf{z} \in \mathcal{X}.$$

Letting  $k \rightarrow \infty$  yields

$$\mathbb{E}_{\mathbb{Q}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \leq \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \quad \forall \mathbf{z} \in \mathcal{X},$$

and the second-moment constraint also passes to the limit by the Portmanteau theorem (Billingsley 2013, Theorem 2.1). Hence  $\mathbb{Q}^* \in \mathcal{P}_0$ , so

$$\limsup_{k \rightarrow \infty} f_{\epsilon_k}(\mathbf{x}) \leq \lim_{k \rightarrow \infty} \mathbb{E}_{\mathbb{Q}_k}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] = \mathbb{E}_{\mathbb{Q}^*}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] \leq f_0(\mathbf{x}).$$

Therefore  $f_\epsilon(\mathbf{x}) \downarrow f_0(\mathbf{x})$  pointwise on  $\mathcal{X}$ .

Let  $v_\epsilon := \min_{\mathbf{x} \in \mathcal{X}} f_\epsilon(\mathbf{x})$  and  $v_0 := \min_{\mathbf{x} \in \mathcal{X}} f_0(\mathbf{x})$ . Since  $f_\epsilon \geq f_0$  pointwise, we have  $v_\epsilon \geq v_0$ . On the other hand, for any minimizer  $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathcal{X}} f_0(\mathbf{x})$ ,

$$\limsup_{\epsilon \downarrow 0} v_\epsilon \leq \limsup_{\epsilon \downarrow 0} f_\epsilon(\mathbf{x}^*) = f_0(\mathbf{x}^*) = v_0.$$

Hence  $v_\epsilon \rightarrow v_0$ .

Finally, let  $\hat{\mathbf{x}}_\epsilon$  be a minimizer of (16) and let  $\hat{\mathbf{x}}^*$  be any cluster point of a sequence  $\{\hat{\mathbf{x}}_{\epsilon_k}\}_{k \in \mathbb{N}}$  with  $\epsilon_k \downarrow 0$ . By compactness of  $\mathcal{X}$ , after passing to a subsequence we may assume  $\hat{\mathbf{x}}_{\epsilon_k} \rightarrow \hat{\mathbf{x}}^*$ . Since  $f_0(\mathbf{x}) = \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})]$  is continuous in  $\mathbf{x}$ ,

$$f_0(\hat{\mathbf{x}}^*) = \lim_{k \rightarrow \infty} f_0(\hat{\mathbf{x}}_{\epsilon_k}) \leq \lim_{k \rightarrow \infty} f_{\epsilon_k}(\hat{\mathbf{x}}_{\epsilon_k}) = \lim_{k \rightarrow \infty} v_{\epsilon_k} = v_0.$$

Thus  $\hat{\mathbf{x}}^*$  is a minimizer of the empirical risk minimization problem (17).  $\square$

## Appendix B: Proofs of Section 3

### B.1. Theoretical Guarantees

PROPOSITION 11. *With probability at least  $1 - \tau$ , we have:*

$$\begin{aligned} & \mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \right] \\ & \leq \frac{1}{M} \sum_{m \in [M]} \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ + \mathcal{O} \left( R \sqrt{\frac{J \log J (\log M)^3}{M}} \right) \\ & \quad + 2(c_1 + c_2 \sqrt{\Omega}) \sqrt{\frac{1}{2M} \log \left( \frac{1}{\tau} \right)} \quad \forall \mathbb{Q} \in \mathcal{P}(\Xi) : \mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega. \end{aligned}$$

*Proof of Proposition 11.* We define  $h_{\mathbb{Q}}(\mathbf{z}) := \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+$  and let

$$\mathcal{H} := \{h_{\mathbb{Q}} : \mathbb{Q} \in \mathcal{P}(\Xi), \mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega\}.$$

Our Assumption (A) implies that

$$\bar{h} := 2(c_1 + c_2 \sqrt{\Omega}) \geq h(\mathbf{z}) \geq 0 \quad \forall h \in \mathcal{H} \forall \mathbf{z} \in \mathcal{X}.$$

For a given sample set  $\{\tilde{\mathbf{z}}_m\}_{m \in [M]}$ , we define the uniform approximation error as:

$$e(\tilde{\mathbf{z}}_1, \dots, \tilde{\mathbf{z}}_M) := \sup_{h \in \mathcal{H}} \left( \mathbb{E}_{\mathbb{Z}}[h(\tilde{\mathbf{z}})] - \frac{1}{M} \sum_{m \in [M]} h(\tilde{\mathbf{z}}_m) \right).$$

Note that changing one data point  $\tilde{z}_m$  changes  $e$  by at most  $\bar{h}/M$ . McDiarmid's inequality (Mohri et al. 2018, Appendix D) therefore yields

$$e(\tilde{z}_1, \dots, \tilde{z}_M) \leq \mathbb{E}[e(\tilde{z}_1, \dots, \tilde{z}_M)] + \bar{h} \sqrt{\frac{1}{2M} \log \left( \frac{1}{\tau} \right)}$$

with probability at least  $1 - \tau$ .

We next upper bound the expectation using Rademacher complexity. Introduce  $M$  i.i.d. ghost samples  $\tilde{z}'_1, \dots, \tilde{z}'_M$  drawn from  $\mathbb{Z}$  and let  $\tilde{s}_1, \dots, \tilde{s}_M \in \{-1, 1\}$  be i.i.d. Rademacher random variables. A standard symmetrization inequality (Mohri et al. 2018, Theorem 3.3) yields

$$\begin{aligned} \mathbb{E} \left[ \sup_{h \in \mathcal{H}} \left( \mathbb{E}_{\mathbb{Z}}[h(\tilde{z})] - \frac{1}{M} \sum_{m \in [M]} h(\tilde{z}_m) \right) \right] &\leq \mathbb{E} \left[ \sup_{h \in \mathcal{H}} \left( \frac{1}{M} \sum_{m \in [M]} h(\tilde{z}'_m) - \frac{1}{M} \sum_{m \in [M]} h(\tilde{z}_m) \right) \right] \\ &= \mathbb{E} \left[ \sup_{h \in \mathcal{H}} \left( \frac{1}{M} \sum_{m \in [M]} \tilde{s}_m (h(\tilde{z}'_m) - h(\tilde{z}_m)) \right) \right] \\ &\leq 2 \mathbb{E} \left[ \sup_{h \in \mathcal{H}} \left( \frac{1}{M} \sum_{m \in [M]} \tilde{s}_m h(\tilde{z}_m) \right) \right] = 2\mathcal{R}_M(\mathcal{H}). \end{aligned}$$

Since the hinge map  $u \mapsto [u]_+$  is 1-Lipschitz, the contraction principle (Shalev-Shwartz and Ben-David 2014, Lemma 26.9) gives

$$\begin{aligned} \mathcal{R}_M(\mathcal{H}) &\leq \mathbb{E} \left[ \sup_{\substack{\mathbb{Q} \in \mathcal{P}(\Xi) \\ \mathbb{E}_{\mathbb{Q}}[\|\tilde{\xi}\|^2] \leq \Omega}} \left( \frac{1}{M} \sum_{m \in [M]} \tilde{s}_m \left( \mathbb{E}_{\mathbb{Q}}[\ell(\tilde{z}_m, \tilde{\xi})] - \mathbb{E}_{\mathbb{P}_N}[\ell(\tilde{z}_m, \tilde{\xi})] - \epsilon \right) \right) \right] \\ &= \mathbb{E} \left[ \sup_{\substack{\mathbb{Q} \in \mathcal{P}(\Xi) \\ \mathbb{E}_{\mathbb{Q}}[\|\tilde{\xi}\|^2] \leq \Omega}} \mathbb{E}_{\mathbb{Q}} \left[ \frac{1}{M} \sum_{m \in [M]} \tilde{s}_m \ell(\tilde{z}_m, \tilde{\xi}) \right] \right] - \mathbb{E} \left[ \frac{1}{M} \sum_{m \in [M]} \tilde{s}_m \left( \mathbb{E}_{\mathbb{P}_N}[\ell(\tilde{z}_m, \tilde{\xi})] + \epsilon \right) \right] \quad (32) \\ &\leq \mathbb{E} \left[ \sup_{\mathbb{Q} \in \mathcal{P}(\Xi)} \mathbb{E}_{\mathbb{Q}} \left[ \frac{1}{M} \sum_{m \in [M]} \tilde{s}_m \ell(\tilde{z}_m, \tilde{\xi}) \right] \right] \\ &= \mathbb{E} \left[ \sup_{\xi \in \Xi} \frac{1}{M} \sum_{m \in [M]} \tilde{s}_m \ell(\tilde{z}_m, \xi) \right]. \end{aligned}$$

In the second line, we moved the term  $\frac{1}{M} \sum_{m \in [M]} \tilde{s}_m \left( \mathbb{E}_{\mathbb{P}_N}[\ell(\tilde{z}_m, \tilde{\xi})] + \epsilon \right)$  outside the supremum as it is independent of  $\mathbb{Q}$ . Since each  $\tilde{s}_m$  is centered and independent of  $\tilde{z}_m$ , this term has expectation zero.

Each affine piece in the loss function can be rewritten as

$$\mathbf{a}_j(\mathbf{x})^\top \boldsymbol{\xi} + b_j(\mathbf{x}) = \left( \begin{bmatrix} \mathbf{A}_j^\top & \mathbf{b}_j \\ \bar{\mathbf{a}}_j & \bar{b}_j \end{bmatrix} \begin{bmatrix} \boldsymbol{\xi} \\ 1 \end{bmatrix} \right)^\top \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}.$$

Let  $\mathcal{B}^{D_x+1}$  be the unit ball in  $\mathbb{R}^{D_x+1}$ . For each  $\boldsymbol{\xi} \in \Xi$  and  $j \in [J]$ , define

$$\mathbf{w}_j(\boldsymbol{\xi}) := \left( \max_{i \in [J]} \left\| \begin{bmatrix} \mathbf{A}_i^\top & \mathbf{b}_i \\ \bar{\mathbf{a}}_i^\top & \bar{b}_i \end{bmatrix} \right\|_{\text{op}} \sqrt{R_\xi^2 + 1} \right)^{-1} \begin{bmatrix} \mathbf{A}_j^\top & \mathbf{b}_j \\ \bar{\mathbf{a}}_j^\top & \bar{b}_j \end{bmatrix} \begin{bmatrix} \boldsymbol{\xi} \\ 1 \end{bmatrix},$$

so that  $\mathbf{w}_j(\boldsymbol{\xi}) \in \mathcal{B}^{D_x+1}$  for all  $j \in [J]$ , and also

$$\left( \sqrt{R_x^2 + 1} \right)^{-1} \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} \in \mathcal{B}^{D_x+1}.$$

Hence

$$\begin{aligned} \mathcal{R}_M(\mathcal{H}) &\leq \mathbb{E} \left[ \sup_{\boldsymbol{\xi} \in \Xi} \frac{1}{M} \sum_{m \in [M]} \tilde{s}_m \ell(\tilde{\mathbf{z}}_m, \boldsymbol{\xi}) \right] \\ &= \mathbb{E} \left[ \sup_{\boldsymbol{\xi} \in \Xi} \frac{1}{M} \sum_{m \in [M]} \tilde{s}_m \max_{j \in [J]} (\mathbf{a}_j(\tilde{\mathbf{z}}_m)^\top \boldsymbol{\xi} + b_j(\tilde{\mathbf{z}}_m)) \right] \\ &\leq \mathbb{E} \left[ \sup_{\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_J \in \Xi} \frac{1}{M} \sum_{m \in [M]} \tilde{s}_m \max_{j \in [J]} (\mathbf{a}_j(\tilde{\mathbf{z}}_m)^\top \boldsymbol{\xi}_j + b_j(\tilde{\mathbf{z}}_m)) \right] \\ &\leq \max_{j \in [J]} \left\| \begin{bmatrix} \mathbf{A}_j^\top & \mathbf{b}_j \\ \bar{\mathbf{a}}_j^\top & \bar{b}_j \end{bmatrix} \right\|_{\text{op}} \sqrt{R_\xi^2 + 1} \sqrt{R_x^2 + 1} \\ &\quad \times \mathbb{E} \left[ \sup_{\mathbf{w}_1, \dots, \mathbf{w}_J \in \mathcal{B}^{D_x+1}} \frac{1}{M} \sum_{m \in [M]} \tilde{s}_m \max_{j \in [J]} \mathbf{w}_j^\top \left( \sqrt{R_x^2 + 1} \right)^{-1} \begin{bmatrix} \tilde{\mathbf{z}}_m \\ 1 \end{bmatrix} \right] \\ &\leq \mathcal{O} \left( \max_{j \in [J]} \left\| \begin{bmatrix} \mathbf{A}_j^\top & \mathbf{b}_j \\ \bar{\mathbf{a}}_j^\top & \bar{b}_j \end{bmatrix} \right\|_{\text{op}} \sqrt{R_\xi^2 + 1} \sqrt{R_x^2 + 1} \sqrt{\frac{J \log J (\log M)^3}{M}} \right), \end{aligned} \tag{33}$$

where the last step uses the Rademacher complexity of  $J$ -fold maxima of hyperplanes (Attias and Kontorovich 2024, Corollary 5). Combining the McDiarmid and Rademacher bounds yields the claim.  $\square$

*Proof of Theorem 3.* Consider the high-probability event of Proposition 11, which occurs with probability at least  $1 - \tau$ . Fix any  $\mathbb{Q} \in \mathcal{P}_\epsilon^M$ . By definition of  $\mathcal{P}_\epsilon^M$ ,

$$\frac{1}{M} \sum_{m \in [M]} \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \leq 0 \quad \text{and} \quad \mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega.$$

Applying Proposition 11 to this  $\mathbb{Q}$  shows that

$$\mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \right] \leq \eta,$$

where  $\eta$  is the bound in (21). Together with the shared second-moment constraint, this means  $\mathbb{Q} \in \mathcal{P}_\epsilon(\eta)$ . Since  $\mathbb{Q} \in \mathcal{P}_\epsilon^M$  was arbitrary, we obtain  $\mathcal{P}_\epsilon^M \subseteq \mathcal{P}_\epsilon(\eta)$ .  $\square$

**LEMMA 3.** *As  $\eta \downarrow 0$ , the optimal value of the relaxed DRO problem*

$$\min_{\mathbf{x} \in \mathcal{X}} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon(\eta)} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})]$$

*converges to that of the exact DRO problem.*

*Proof of Lemma 3.* We follow the development of the proof of Proposition 9. Fix  $\mathbf{x} \in \mathcal{X}$ , and define

$$f_\eta(\mathbf{x}) := \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon(\eta)} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})].$$

Since  $\mathcal{P}_\epsilon = \mathcal{P}_\epsilon(0) \subseteq \mathcal{P}_\epsilon(\eta)$  for every  $\eta \geq 0$ , we have  $f_0(\mathbf{x}) \leq f_\eta(\mathbf{x})$ . Thus, it remains to prove

$$\lim_{\eta \downarrow 0} f_\eta(\mathbf{x}) \leq f_0(\mathbf{x}). \quad (34)$$

Let  $\eta_k \downarrow 0$  be any decreasing sequence converging to zero. For each  $k \in \mathbb{N}$ , choose  $\mathbb{Q}_k \in \mathcal{P}_\epsilon(\eta_k)$  such that

$$f_{\eta_k}(\mathbf{x}) - \frac{1}{k} \leq \mathbb{E}_{\mathbb{Q}_k}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})]. \quad (35)$$

Because each  $\mathbb{Q}_k$  satisfies the second-moment bound in  $\mathcal{P}_\epsilon(\eta_k)$ , the family  $\{\mathbb{Q}_k\}_{k \in \mathbb{N}}$  is tight. Hence, after passing to a subsequence if necessary,  $\mathbb{Q}_k \Rightarrow \mathbb{Q}^*$  for some probability measure  $\mathbb{Q}^*$ .

By the same argument as in Proposition 9, weak convergence together with the uniform second-moment bound implies  $\lim_{k \rightarrow \infty} \mathbb{E}_{\mathbb{Q}_k}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] = \mathbb{E}_{\mathbb{Q}^*}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})]$ , and  $\mathbb{Q}^*$  satisfies the second-moment constraint in  $\mathcal{P}_\epsilon$ . It remains to verify the expected hinge constraint. For each fixed  $\mathbf{z} \in \mathcal{X}$ , the convergence of expectations implies

$$\lim_{k \rightarrow \infty} \left[ \mathbb{E}_{\mathbb{Q}_k}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ = \left[ \mathbb{E}_{\mathbb{Q}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+.$$

By Assumptions (A) and (B), the hinge term is bounded uniformly in  $\mathbf{z} \in \mathcal{X}$  and  $k \in \mathbb{N}$  by a finite constant, so dominated convergence yields

$$\lim_{k \rightarrow \infty} \mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}_k}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \right] = \mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}^*}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \right].$$

Since  $\mathbb{Q}_k \in \mathcal{P}_\epsilon(\eta_k)$ , we have

$$\mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}_k}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \right] \leq \eta_k.$$

Letting  $k \rightarrow \infty$ , we obtain  $\mathbb{Q}^* \in \mathcal{P}_\epsilon$ . Taking the limit in (35) gives

$$\lim_{k \rightarrow \infty} f_{\eta_k}(\mathbf{x}) \leq \lim_{k \rightarrow \infty} \mathbb{E}_{\mathbb{Q}_k}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] = \mathbb{E}_{\mathbb{Q}^*}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] \leq f_0(\mathbf{x}).$$

Since  $\eta_k \downarrow 0$  was arbitrary, (34) holds. Thus, for any  $\mathbf{x} \in \mathcal{X}$ ,

$$\lim_{\eta \downarrow 0} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon(\eta)} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] = \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})]. \quad (36)$$

The optimal value convergence then follows from the monotonicity of  $f_\eta(\mathbf{x})$  in  $\eta$ .  $\square$

*Proof of Theorem 4.* For  $\eta \geq 0$ , define

$$\hat{v}(\eta) := \min_{\mathbf{x} \in \mathcal{X}} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon(\eta)} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})].$$

By Theorem 3, setting  $\eta_M$  to (21), we have for any  $\tau > 0$ ,

$$\text{Prob}\left(\inf_{\mathbf{x} \in \mathcal{X}} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon^M} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] \leq \inf_{\mathbf{x} \in \mathcal{X}} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon(\eta_M)} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})]\right) \geq 1 - \tau.$$

Equivalently,

$$\text{Prob}(\hat{v}_M \leq \hat{v}(\eta_M)) \geq 1 - \tau.$$

Since  $\mathcal{P}_\epsilon \subseteq \mathcal{P}_\epsilon^M$ , we also have  $\hat{v} \leq \hat{v}_M$ . Therefore,

$$\text{Prob}(\hat{v} \leq \hat{v}_M \leq \hat{v}(\eta_M)) \geq 1 - \tau.$$

By Lemma 3,  $\hat{v}(\eta_M) \rightarrow \hat{v}$  as  $M \rightarrow \infty$ . Hence, for any  $\rho > 0$ , there exists  $M' \in \mathbb{N}$  such that

$$\hat{v}(\eta_M) - \hat{v} \leq \rho \quad \forall M \geq M'.$$

Combining this bound with the preceding high-probability event yields

$$\text{Prob}(\hat{v}_M - \hat{v} \leq \rho) \geq 1 - \tau \quad \forall M \geq M'.$$

Since also  $\hat{v}_M - \hat{v} \geq 0$ , we conclude that  $\hat{v}_M \xrightarrow{P} \hat{v}$ .

We next establish convergence of optimal solutions. For  $\delta > 0$ , define

$$\mathcal{A}_\delta := \{\mathbf{x} \in \mathcal{X} : \text{dist}(\mathbf{x}, \mathcal{S}) \geq \delta\}.$$

This set is closed because the distance function is continuous. Since  $\mathcal{X}$  is compact,  $\mathcal{A}_\delta$  is compact.

Moreover, by definition of  $\mathcal{S}$ , no point in  $\mathcal{A}_\delta$  is optimal. Hence,

$$\rho_\delta := \inf_{\mathbf{x} \in \mathcal{A}_\delta} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] - \inf_{\mathbf{x} \in \mathcal{X}} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] > 0.$$

Now, if  $\text{dist}(\hat{\mathbf{x}}_M, \mathcal{S}) \geq \delta$ , then  $\hat{\mathbf{x}}_M \in \mathcal{A}_\delta$ , and therefore

$$\sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}_M, \tilde{\boldsymbol{\xi}})] - \inf_{\mathbf{x} \in \mathcal{X}} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] \geq \rho_\delta.$$

Since  $\mathcal{P}_\epsilon \subseteq \mathcal{P}_\epsilon^M$ , we also have

$$\sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}_M, \tilde{\boldsymbol{\xi}})] \leq \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon^M} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}_M, \tilde{\boldsymbol{\xi}})] = \hat{v}_M.$$

Thus,  $\{\text{dist}(\hat{\mathbf{x}}_M, \mathcal{S}) \geq \delta\} \subseteq \{\hat{v}_M - \hat{v} \geq \rho_\delta\}$ , and consequently

$$\text{Prob}(\text{dist}(\hat{\mathbf{x}}_M, \mathcal{S}) \geq \delta) \leq \text{Prob}(\hat{v}_M - \hat{v} \geq \rho_\delta).$$

Since  $\hat{v}_M \xrightarrow{P} \hat{v}$ , the right-hand side converges to zero. Therefore,  $\text{dist}(\hat{\mathbf{x}}_M, \mathcal{S}) \xrightarrow{P} 0$ .

Lastly, we prove the suboptimality bound. We have

$$\begin{aligned} \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}_M, \tilde{\boldsymbol{\xi}})] - \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})] &= \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}_M, \tilde{\boldsymbol{\xi}})] - \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon^M} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}_M, \tilde{\boldsymbol{\xi}})] \\ &\quad + \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon^M} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}_M, \tilde{\boldsymbol{\xi}})] - \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})] \\ &\leq \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon^M} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})] - \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})]. \end{aligned}$$

Here, the first difference is nonpositive because  $\mathcal{P}_\epsilon \subseteq \mathcal{P}_\epsilon^M$ , and the second difference is bounded above by replacing  $\hat{\mathbf{x}}_M$  with  $\hat{\mathbf{x}}$ , since  $\hat{\mathbf{x}}$  is suboptimal for the approximate problem.

Next, by Theorem 3, setting  $\eta$  to (21), we obtain the high-probability bound

$$\sup_{\mathbb{Q} \in \mathcal{P}_\epsilon^M} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})] - \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})] \leq \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon(\eta)} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})] - \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})]$$

with probability at least  $1 - \tau$ . For any fixed  $\mathbf{x}$  and  $\eta \geq 0$ , the relaxed problem (20) has the dual reformulation

$$\begin{aligned} \min \quad & \alpha + \beta\Omega + \gamma\eta + \int_{\mathcal{X}} \left( \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] + \epsilon \right) \nu(d\mathbf{z}) \\ \text{s.t.} \quad & \mathbf{x} \in \mathcal{X}, \alpha \in \mathbb{R}, \beta, \gamma \in \mathbb{R}_+, \nu \in \mathcal{M}_+(\mathcal{X}), \\ & \ell(\mathbf{x}, \boldsymbol{\xi}) \leq \alpha + \beta\|\boldsymbol{\xi}\|^2 + \int_{\mathcal{X}} \ell(\mathbf{z}, \boldsymbol{\xi}) \nu(d\mathbf{z}) \quad \forall \boldsymbol{\xi} \in \Xi, \\ & \gamma\mathbb{Z} - \nu \in \mathcal{M}_+(\mathcal{X}), \end{aligned}$$

which follows from the same generalized-moment duality argument as in Theorem 2. Let  $(\hat{\mathbf{x}}, \alpha^*, \beta^*, \nu^*)$  be an optimal solution to the exact dual problem, and let  $\gamma^*$  be as defined in the theorem statement. Then  $(\hat{\mathbf{x}}, \alpha^*, \beta^*, \gamma^*, \nu^*)$  is feasible, though generally suboptimal, for the relaxed dual problem. We can therefore further upper bound the preceding display by

$$\begin{aligned} & \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon(\eta)} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})] - \sup_{\mathbb{Q} \in \mathcal{P}_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}, \tilde{\boldsymbol{\xi}})] \\ & \leq \alpha^* + \beta^*\Omega + \int_{\mathcal{X}} \left( \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] + \epsilon \right) \nu^*(d\mathbf{z}) + \gamma^*\eta \\ & \quad - \left( \alpha^* + \beta^*\Omega + \int_{\mathcal{X}} \left( \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] + \epsilon \right) \nu^*(d\mathbf{z}) \right) \\ & = \gamma^*\eta. \end{aligned}$$

Substituting the expression for  $\eta$  in (21) completes the proof.  $\square$

## B.2. Conic Programming Reformulations

*Proof of Theorem 5.* For each sampled point  $\mathbf{z}_m$ , define

$$g_m(\boldsymbol{\xi}) := \max_{k \in [J]} (\mathbf{a}_k(\mathbf{z}_m)^\top \boldsymbol{\xi} + b_k(\mathbf{z}_m)).$$

The semi-infinite constraint in (19) is equivalent to

$$\mathbf{a}_j(\mathbf{x})^\top \boldsymbol{\xi} + b_j(\mathbf{x}) \leq \alpha + \beta \|\boldsymbol{\xi}\|^2 + \frac{1}{M} \sum_{m \in [M]} \nu_m g_m(\boldsymbol{\xi}) \quad \forall \boldsymbol{\xi} \in \Xi, \forall j \in [J].$$

For each pair  $(j, m)$ , the term  $\nu_m g_m(\boldsymbol{\xi})$  admits the simplex representation

$$\nu_m g_m(\boldsymbol{\xi}) = \max_{\substack{\boldsymbol{\lambda}_{jm} \in \mathbb{R}_+^J \\ \mathbf{e}^\top \boldsymbol{\lambda}_{jm} = \nu_m}} \sum_{k \in [J]} \lambda_{jm}^k (\mathbf{a}_k(\mathbf{z}_m)^\top \boldsymbol{\xi} + b_k(\mathbf{z}_m)).$$

Hence the semi-infinite constraint is equivalent to the existence of multipliers  $\boldsymbol{\lambda}_{jm} \in \mathbb{R}_+^J$  with  $\mathbf{e}^\top \boldsymbol{\lambda}_{jm} = \nu_m$  such that, for every  $j \in [J]$ ,

$$\sup_{\boldsymbol{\xi} \in \Xi} \left\{ \left( \mathbf{a}_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k \mathbf{a}_k(\mathbf{z}_m) \right)^\top \boldsymbol{\xi} - \beta \|\boldsymbol{\xi}\|^2 \right\} + b_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k b_k(\mathbf{z}_m) \leq \alpha.$$

Let

$$\mathbf{c}_j := \mathbf{a}_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k \mathbf{a}_k(\mathbf{z}_m).$$

Using Fenchel duality with the support function of  $\Xi$  and the infimal convolution identity for Fenchel conjugates (Rockafellar 1970, Theorem 16.4),

$$\sup_{\boldsymbol{\xi} \in \Xi} \{ \mathbf{c}_j^\top \boldsymbol{\xi} - \beta \|\boldsymbol{\xi}\|^2 \} = \inf_{\boldsymbol{\theta}_j \in \mathbb{R}^{D_\xi}} \left\{ \sigma_\Xi(\boldsymbol{\theta}_j) + \sup_{\boldsymbol{\xi} \in \mathbb{R}^{D_\xi}} [(\mathbf{c}_j - \boldsymbol{\theta}_j)^\top \boldsymbol{\xi} - \beta \|\boldsymbol{\xi}\|^2] \right\}.$$

The remaining quadratic supremum equals  $\|\mathbf{c}_j - \boldsymbol{\theta}_j\|^2 / (4\beta)$  when  $\beta > 0$ , and its epigraph admits the standard second-order-cone representation

$$\frac{\|\mathbf{c}_j - \boldsymbol{\theta}_j\|^2}{4\beta} \leq \zeta_j \iff \left\| \begin{bmatrix} \mathbf{c}_j - \boldsymbol{\theta}_j \\ \zeta_j - \beta \end{bmatrix} \right\| \leq \zeta_j + \beta, \quad \zeta_j \geq 0.$$

When  $\beta = 0$ , the semi-infinite constraint is satisfied if and only if

$$\sigma_\Xi(\mathbf{c}_j) + b_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k b_k(\mathbf{z}_m) \leq \alpha \quad \forall j \in [J].$$

Moreover, when  $\beta = 0$ , the second-order conic constraint above reduces to

$$\left\| \begin{bmatrix} \mathbf{c}_j - \boldsymbol{\theta}_j \\ \zeta_j \end{bmatrix} \right\| \leq \zeta_j,$$

which implies  $\boldsymbol{\theta}_j = \mathbf{c}_j$ . Substituting this identity into the scalar inequality of the conic reformulation gives

$$\zeta_j + \sigma_\Xi(\mathbf{c}_j) + b_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k b_k(\mathbf{z}_m) \leq \alpha.$$

In this scalar inequality, smaller values of  $\zeta_j$  weakly enlarge the feasible region. Since the reduced conic constraint imposes no positive lower bound on  $\zeta_j$  beyond  $\zeta_j \geq 0$ , the least restrictive feasible

choice is  $\zeta_j = 0$ , which recovers exactly the support-function constraint above. Hence the conic reformulation in the theorem remains valid also in the case  $\beta = 0$ . Substituting  $\mathbf{c}_j$  back into the inequality therefore produces exactly the conic reformulation stated in the theorem.

It remains to verify convexity of the displayed reformulation. The objective is affine in the decision variables, and the sign, simplex, and equality constraints are affine. Since  $\mathbf{a}_j(\mathbf{x})$  and  $b_j(\mathbf{x})$  are affine in  $\mathbf{x}$ , the scalar inequality

$$\zeta_j + \sigma_{\Xi}(\boldsymbol{\theta}_j) + b_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k b_k(\mathbf{z}_m) \leq \alpha$$

has a convex left-hand side, because  $\sigma_{\Xi}$  is convex as the support function of a closed convex set. The remaining nonlinear constraint is the second-order-cone membership

$$\left( \mathbf{a}_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k \mathbf{a}_k(\mathbf{z}_m) - \boldsymbol{\theta}_j, \zeta_j - \beta, \zeta_j + \beta \right) \in \mathcal{Q},$$

where  $\mathcal{Q} := \{(\mathbf{u}, t) : \|\mathbf{u}\| \leq t\}$  denotes the standard second-order cone. Its argument is affine in the decision variables, so its inverse image is convex. Together with the convexity of the constraint  $\mathbf{x} \in \mathcal{X}$ , this shows that the finite reformulation is a convex conic program. The final claims follow because the support function of a second-order-cone representable set is second-order-cone representable.  $\square$

## Appendix C: Proofs of Section 4

### C.1. Theoretical Guarantees

*Proof of Proposition 10.* Fix  $\mathbf{x} \in \mathcal{X}$ . The inner maximization problem associated with (23) is again a generalized moment problem, now with the additional constraint

$$\mathbb{Q}(\|\tilde{\boldsymbol{\xi}}\| \leq t) \geq 1 - 2 \exp\left(- (t/K_{\text{sw}})^\vartheta\right).$$

Introducing a multiplier  $\kappa \in \mathbb{R}_+$  for this constraint, together with the multipliers  $\nu \in \mathcal{M}_+(\mathcal{X})$ ,  $\beta \in \mathbb{R}_+$ , and  $\alpha \in \mathbb{R}$  used in the proof of Theorem 2, yields the dual problem

$$\begin{aligned} \min \quad & \alpha + \beta \Omega - \kappa \left(1 - 2 \exp\left(- (t/K_{\text{sw}})^\vartheta\right)\right) + \int_{\mathcal{X}} \left(\mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] + \epsilon\right) \nu(d\mathbf{z}) \\ \text{s.t.} \quad & \ell(\mathbf{x}, \boldsymbol{\xi}) + \kappa \mathbb{I}_{\{\|\boldsymbol{\xi}\| \leq t\}} \leq \alpha + \beta \|\boldsymbol{\xi}\|^2 + \int_{\mathcal{X}} \ell(\mathbf{z}, \boldsymbol{\xi}) \nu(d\mathbf{z}) \quad \forall \boldsymbol{\xi} \in \Xi. \end{aligned}$$

The strict feasibility argument from Theorem 2 still applies, so strong duality and dual attainment hold. Minimizing over  $\mathbf{x} \in \mathcal{X}$  gives (24) and yields the existence of an optimal tuple.  $\square$

PROPOSITION 12. *With probability at least  $1 - \tau$ , the following bound holds for all  $\mathbb{Q} \in \mathcal{P}(\Xi)$  satisfying  $\mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega$  and  $\mathbb{Q}(\|\tilde{\boldsymbol{\xi}}\| \leq t) \geq 1 - 2\exp(- (t/K_{\text{sw}})^\vartheta)$ :*

$$\begin{aligned} & \mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \right] \\ & \leq \frac{1}{M} \sum_{m \in [M]} \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ + \mathcal{O} \left( R(t) \sqrt{\frac{J \log J (\log M)^3}{M}} \right) \\ & \quad + 2(c_1 + c_2 \sqrt{\Omega}) \sqrt{\frac{1}{2M} \log \left( \frac{1}{\tau} \right)} + 4c_1 \exp \left( - (t/K_{\text{sw}})^\vartheta \right) + 2c_2 \sqrt{2\Omega \exp \left( - (t/K_{\text{sw}})^\vartheta \right)}. \end{aligned}$$

Here  $R(t)$  is defined as in Theorem 6.

*Proof of Proposition 12.* We define the ambiguity set

$$\mathcal{P}_{\text{sw}} := \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : \begin{array}{l} \mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega \\ \mathbb{Q}(\|\tilde{\boldsymbol{\xi}}\| \leq t) \geq 1 - 2\exp \left( - (t/K_{\text{sw}})^\vartheta \right) \end{array} \right\},$$

which contains both  $\mathcal{P}'_\epsilon$  and  $\mathcal{P}'_\epsilon{}^M$ . Let  $h_{\mathbb{Q}}(\mathbf{z}) := \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+$  and  $\mathcal{H} := \{h_{\mathbb{Q}} : \mathbb{Q} \in \mathcal{P}_{\text{sw}}\}$ . We now repeat the steps of the proof of Proposition 11. Assumption (A) implies that  $\bar{h} := 2(c_1 + c_2 \sqrt{\Omega})$  satisfies  $0 \leq h(\mathbf{z}) \leq \bar{h}$  for all  $h \in \mathcal{H}$  and  $\mathbf{z} \in \mathcal{X}$ .

By McDiarmid's inequality (Mohri et al. 2018, Appendix D), we have

$$\sup_{h \in \mathcal{H}} \left( \mathbb{E}_{\mathbb{Z}}[h(\tilde{\mathbf{z}})] - \frac{1}{M} \sum_{m \in [M]} h(\tilde{\mathbf{z}}_m) \right) \leq \mathbb{E} \left[ \sup_{h \in \mathcal{H}} \left( \mathbb{E}_{\mathbb{Z}}[h(\tilde{\mathbf{z}})] - \frac{1}{M} \sum_{m \in [M]} h(\tilde{\mathbf{z}}_m) \right) \right] + \bar{h} \sqrt{\frac{1}{2M} \log \left( \frac{1}{\tau} \right)}$$

with probability at least  $1 - \tau$ . The expectation term can be bounded by twice the Rademacher complexity of  $\mathcal{H}$ . Let  $F_M(\boldsymbol{\xi}) := M^{-1} \sum_{m \in [M]} \tilde{s}_m \ell(\tilde{\mathbf{z}}_m, \boldsymbol{\xi})$ . Then

$$\begin{aligned} \mathcal{R}_M(\mathcal{H}) & \leq \mathbb{E} \left[ \sup_{\mathbb{Q} \in \mathcal{P}_{\text{sw}}} \mathbb{E}_{\mathbb{Q}}[F_M(\tilde{\boldsymbol{\xi}})] \right] \\ & = \mathbb{E} \left[ \sup_{\mathbb{Q} \in \mathcal{P}_{\text{sw}}} \left\{ \mathbb{E}_{\mathbb{Q}} \left[ F_M(\tilde{\boldsymbol{\xi}}) \mathbb{I}_{\{\|\tilde{\boldsymbol{\xi}}\| \leq t\}} \right] + \mathbb{E}_{\mathbb{Q}} \left[ F_M(\tilde{\boldsymbol{\xi}}) \mathbb{I}_{\{\|\tilde{\boldsymbol{\xi}}\| > t\}} \right] \right\} \right] \\ & \leq \mathbb{E} \left[ \sup_{\mathbb{Q} \in \mathcal{P}_{\text{sw}}} \mathbb{E}_{\mathbb{Q}} \left[ F_M(\tilde{\boldsymbol{\xi}}) \mathbb{I}_{\{\|\tilde{\boldsymbol{\xi}}\| \leq t\}} \right] \right] \\ & \quad + \mathbb{E} \left[ \sup_{\mathbb{Q} \in \mathcal{P}_{\text{sw}}} \mathbb{E}_{\mathbb{Q}} \left[ F_M(\tilde{\boldsymbol{\xi}}) \mathbb{I}_{\{\|\tilde{\boldsymbol{\xi}}\| > t\}} \right] \right]. \end{aligned}$$

For the first expectation, since the indicator restricts the integrand to the ball  $\{\|\boldsymbol{\xi}\| \leq t\}$ ,

$$\sup_{\mathbb{Q} \in \mathcal{P}_{\text{sw}}} \mathbb{E}_{\mathbb{Q}} \left[ F_M(\tilde{\boldsymbol{\xi}}) \mathbb{I}_{\{\|\tilde{\boldsymbol{\xi}}\| \leq t\}} \right] \leq \sup_{\|\boldsymbol{\xi}\| \leq t} |F_M(\boldsymbol{\xi})|.$$

By symmetry of the Rademacher signs, the expectation of the right-hand side is bounded, up to an absolute constant, by the localized version of the complexity bound in (33):

$$\begin{aligned}
& \mathbb{E} \left[ \sup_{\mathbb{Q} \in \mathcal{P}_{\text{sw}}} \mathbb{E}_{\mathbb{Q}} \left[ F_M(\tilde{\boldsymbol{\xi}}) \mathbb{I}_{\{\|\tilde{\boldsymbol{\xi}}\| \leq t\}} \right] \right] \\
& \leq \mathbb{E} \left[ \sup_{\boldsymbol{\xi} \in \mathbb{R}^D: \|\boldsymbol{\xi}\| \leq t} \left| \frac{1}{M} \sum_{m \in [M]} \tilde{s}_m \ell(\tilde{\mathbf{z}}_m, \boldsymbol{\xi}) \right| \right] \\
& \leq 2 \mathbb{E} \left[ \sup_{\boldsymbol{\xi} \in \mathbb{R}^D: \|\boldsymbol{\xi}\| \leq t} \frac{1}{M} \sum_{m \in [M]} \tilde{s}_m \ell(\tilde{\mathbf{z}}_m, \boldsymbol{\xi}) \right] \\
& \leq \mathcal{O} \left( R(t) \sqrt{\frac{J \log J (\log M)^3}{M}} \right).
\end{aligned}$$

For the second expectation, Assumption (A), Cauchy–Schwarz, and the defining constraints of  $\mathcal{P}_{\text{sw}}$  yield

$$\begin{aligned}
& \sup_{\mathbb{Q} \in \mathcal{P}_{\text{sw}}} \mathbb{E}_{\mathbb{Q}} \left[ \left| \frac{1}{M} \sum_{m \in [M]} \tilde{s}_m \ell(\tilde{\mathbf{z}}_m, \tilde{\boldsymbol{\xi}}) \right| \mathbb{I}_{\{\|\tilde{\boldsymbol{\xi}}\| > t\}} \right] \\
& \leq c_1 \mathbb{Q}(\|\tilde{\boldsymbol{\xi}}\| > t) + c_2 \sqrt{\mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \mathbb{Q}(\|\tilde{\boldsymbol{\xi}}\| > t)} \\
& \leq 2c_1 \exp\left(- (t/K_{\text{sw}})^\vartheta\right) + c_2 \sqrt{2\Omega \exp\left(- (t/K_{\text{sw}})^\vartheta\right)},
\end{aligned}$$

where the last step uses the tail-mass constraint  $\mathbb{Q}(\|\tilde{\boldsymbol{\xi}}\| > t) \leq 2 \exp\left(- (t/K_{\text{sw}})^\vartheta\right)$ . The factor 2 in the tail correction of the statement comes from the symmetrization step. Combining these bounds with the same symmetrization argument as in Proposition 11 yields the claim.  $\square$

*Proof of Theorem 6.* On the high-probability event of Proposition 12 applied with  $t = t_M$ , fix any  $\mathbb{Q} \in \mathcal{P}'_\epsilon$ . By definition of (25), the sample-average hinge term is nonpositive,  $\mathbb{E}_{\mathbb{Q}}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega$ , and  $\mathbb{Q}(\|\tilde{\boldsymbol{\xi}}\| \leq t_M) \geq 1 - 2/M$ . Proposition 12 therefore implies

$$\mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\mathbb{P}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \right] \leq \eta,$$

where  $\eta = \eta_M$  is the bound stated in Theorem 6. Thus  $\mathbb{Q} \in \mathcal{P}'_{\epsilon, M}(\eta_M)$ . Since  $\mathbb{Q}$  was arbitrary, we obtain  $\mathcal{P}'_\epsilon \subseteq \mathcal{P}'_{\epsilon, M}(\eta_M)$ . The inclusion  $\mathcal{P}'_{\epsilon, M} \subseteq \mathcal{P}'_\epsilon$  is immediate from (25), because the empirical average of nonnegative terms is zero whenever the expected hinge constraint is satisfied exactly. Finally,  $R(t_M) = \mathcal{O}((\log M)^{1/\vartheta})$ , and hence  $\eta_M \rightarrow 0$ .  $\square$

LEMMA 4. Let  $\mathcal{P}'_{\epsilon, M}$  denote (23) with  $t = t_M$ , and let  $\mathcal{P}'_{\epsilon, M}(\eta)$  denote (27). Fix  $\epsilon > 0$ . For any sequence  $\eta_M \geq 0$  with  $\eta_M \rightarrow 0$  as  $M \rightarrow \infty$ , define

$$\begin{aligned}
f_0(\mathbf{x}) &:= \sup_{\mathbb{Q} \in \mathcal{P}'_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})], \\
f_M(\mathbf{x}) &:= \sup_{\mathbb{Q} \in \mathcal{P}'_{\epsilon, M}} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})], \\
g_M(\mathbf{x}) &:= \sup_{\mathbb{Q} \in \mathcal{P}'_{\epsilon, M}(\eta_M)} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})].
\end{aligned}$$

Then

$$\sup_{\mathbf{x} \in \mathcal{X}} |f_M(\mathbf{x}) - f_0(\mathbf{x})| \rightarrow 0, \quad \sup_{\mathbf{x} \in \mathcal{X}} |g_M(\mathbf{x}) - f_0(\mathbf{x})| \rightarrow 0 \quad \text{as } M \rightarrow \infty.$$

*Proof of Lemma 4.* For a probability measure  $\mathbb{Q}$ , define

$$H(\mathbb{Q}) := \mathbb{E}_{\mathbb{Z}} \left[ \left[ \mathbb{E}_{\mathbb{Q}}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\xi}})] - \epsilon \right]_+ \right].$$

We first show that  $f_M(\mathbf{x}) \rightarrow f_0(\mathbf{x})$  and  $g_M(\mathbf{x}) \rightarrow f_0(\mathbf{x})$  for each fixed  $\mathbf{x} \in \mathcal{X}$ .

Since  $\mathcal{P}'_{\epsilon, M} \subseteq \mathcal{P}_{\epsilon}$ , we have  $f_M(\mathbf{x}) \leq f_0(\mathbf{x})$ . To prove the reverse limit inequality, fix  $\delta > 0$  and choose  $\mathbb{Q}^{\delta} \in \mathcal{P}_{\epsilon}$  such that

$$\mathbb{E}_{\mathbb{Q}^{\delta}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] \geq f_0(\mathbf{x}) - \delta.$$

Let  $\Pi_M(\boldsymbol{\xi}) = \boldsymbol{\xi}$  if  $\|\boldsymbol{\xi}\| \leq t_M$  and  $\Pi_M(\boldsymbol{\xi}) = t_M \boldsymbol{\xi} / \|\boldsymbol{\xi}\|$  otherwise, and let  $\tilde{\mathbb{Q}}_M = \mathbb{Q}^{\delta} \circ \Pi_M^{-1}$  be the push-forward of  $\mathbb{Q}^{\delta}$  under  $\Pi_M$ ; equivalently, if  $\tilde{\boldsymbol{\xi}} \sim \mathbb{Q}^{\delta}$ , then  $\Pi_M(\tilde{\boldsymbol{\xi}}) \sim \tilde{\mathbb{Q}}_M$ . Then  $\tilde{\mathbb{Q}}_M$  is supported on  $\{\boldsymbol{\xi} \in \Xi : \|\boldsymbol{\xi}\| \leq t_M\}$  and satisfies the second-moment constraint. Moreover, by Assumption (A), as  $M \rightarrow \infty$ ,

$$\begin{aligned} \Delta_M &:= \sup_{\mathbf{z} \in \mathcal{X}} \left| \mathbb{E}_{\tilde{\mathbb{Q}}_M}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\mathbb{Q}^{\delta}}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \right| \\ &= \sup_{\mathbf{z} \in \mathcal{X}} \left| \mathbb{E}_{\mathbb{Q}^{\delta}}[\ell(\mathbf{z}, \Pi_M(\tilde{\boldsymbol{\xi}})) - \ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \right| \\ &\leq \mathbb{E}_{\mathbb{Q}^{\delta}} \left[ (2c_1 + 2c_2 \|\tilde{\boldsymbol{\xi}}\|) \mathbb{I}_{\{\|\tilde{\boldsymbol{\xi}}\| > t_M\}} \right] \rightarrow 0. \end{aligned}$$

Since  $\epsilon > 0$ , define

$$\lambda_M := \frac{\Delta_M}{\epsilon + \Delta_M}, \quad \mathbb{Q}_M := (1 - \lambda_M) \tilde{\mathbb{Q}}_M + \lambda_M \hat{\mathbb{P}}_N. \quad (37)$$

Because  $\hat{\mathbb{P}}_N$  is the fixed empirical reference distribution, it has finite support. Hence, for all sufficiently large  $M$ ,  $\hat{\mathbb{P}}_N$  is supported on  $\{\boldsymbol{\xi} \in \Xi : \|\boldsymbol{\xi}\| \leq t_M\}$ , and therefore  $\mathbb{Q}_M(\{\boldsymbol{\xi} \in \Xi : \|\boldsymbol{\xi}\| \leq t_M\}) = 1$ . The second-moment constraint also holds, since

$$\mathbb{E}_{\mathbb{Q}_M}[\|\tilde{\boldsymbol{\xi}}\|^2] = (1 - \lambda_M) \mathbb{E}_{\tilde{\mathbb{Q}}_M}[\|\tilde{\boldsymbol{\xi}}\|^2] + \lambda_M \mathbb{E}_{\hat{\mathbb{P}}_N}[\|\tilde{\boldsymbol{\xi}}\|^2] \leq \Omega,$$

where  $\tilde{\mathbb{Q}}_M$  satisfies the second-moment constraint by construction and  $\hat{\mathbb{P}}_N$  satisfies it by Assumption (B). Since  $\mathbb{Q}^{\delta} \in \mathcal{P}_{\epsilon}$ , feasibility gives  $H(\mathbb{Q}^{\delta}) \leq 0$ . Since  $H(\mathbb{Q}^{\delta})$  is the expectation of a nonnegative hinge term,  $H(\mathbb{Q}^{\delta}) \geq 0$ , and hence  $H(\mathbb{Q}^{\delta}) = 0$ . Therefore, for  $\mathbb{Z}$ -almost every realization  $\mathbf{z} \in \mathcal{X}$ ,

$$\begin{aligned} \mathbb{E}_{\mathbb{Q}_M}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \epsilon &= (1 - \lambda_M) \left( \mathbb{E}_{\tilde{\mathbb{Q}}_M}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \right) - \epsilon \\ &\leq (1 - \lambda_M) \left( \mathbb{E}_{\mathbb{Q}^{\delta}}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] - \mathbb{E}_{\hat{\mathbb{P}}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] + \Delta_M \right) - \epsilon \\ &\leq (1 - \lambda_M)(\epsilon + \Delta_M) - \epsilon \\ &= (1 - \lambda_M)\Delta_M - \lambda_M \epsilon = 0. \end{aligned}$$

Here, the first equality uses the definition of  $\mathbb{Q}_M$  in (37); the first inequality uses the definition of  $\Delta_M$ ; the second inequality uses  $H(\mathbb{Q}^\delta) = 0$ ; and the last equality follows from the definition of  $\lambda_M$  in (37).

Thus  $\mathbb{Q}_M \in \mathcal{P}'_{\epsilon, M}$  for all sufficiently large  $M$ . Because  $\lambda_M \rightarrow 0$  and  $\Delta_M \rightarrow 0$ ,

$$\liminf_{M \rightarrow \infty} f_M(\mathbf{x}) \geq \liminf_{M \rightarrow \infty} \mathbb{E}_{\mathbb{Q}_M}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] = \mathbb{E}_{\mathbb{Q}^\delta}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] \geq f_0(\mathbf{x}) - \delta.$$

Letting  $\delta \downarrow 0$  gives  $f_M(\mathbf{x}) \rightarrow f_0(\mathbf{x})$ .

Next, consider  $g_M$ . The proof follows the same compactness argument as Lemma 3. The only difference is that the sets  $\mathcal{P}'_{\epsilon, M}(\eta_M)$  also impose the moving tail constraint, but this constraint does not need to be inherited by the limit because the limiting set  $\mathcal{P}_\epsilon$  has no tail constraint.

The lower bound  $\liminf_M g_M(\mathbf{x}) \geq f_0(\mathbf{x})$  follows from  $g_M(\mathbf{x}) \geq f_M(\mathbf{x})$  and the convergence of  $f_M(\mathbf{x})$ . For the upper bound, take any subsequence and choose  $\mathbb{Q}_M \in \mathcal{P}'_{\epsilon, M}(\eta_M)$  such that  $g_M(\mathbf{x}) - \frac{1}{M} \leq \mathbb{E}_{\mathbb{Q}_M}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})]$ . The common second-moment bound gives tightness, so along a further subsequence  $\mathbb{Q}_M \Rightarrow \mathbb{Q}^*$ . By uniform integrability and Assumption (A), the loss expectations converge for every fixed  $\mathbf{z} \in \mathcal{X}$ :  $\mathbb{E}_{\mathbb{Q}_M}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] \rightarrow \mathbb{E}_{\mathbb{Q}^*}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})]$ . The hinge terms are uniformly bounded by Assumptions (A) and (B); hence dominated convergence gives  $H(\mathbb{Q}_M) \rightarrow H(\mathbb{Q}^*)$ . Since  $H(\mathbb{Q}_M) \leq \eta_M$  and  $\eta_M \rightarrow 0$ , we have  $H(\mathbb{Q}^*) = 0$ . Lower semicontinuity preserves the second-moment constraint, and hence  $\mathbb{Q}^* \in \mathcal{P}_\epsilon$ . Therefore,  $\limsup_{M \rightarrow \infty} g_M(\mathbf{x}) \leq \mathbb{E}_{\mathbb{Q}^*}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})] \leq f_0(\mathbf{x})$ . Together with the lower bound, this proves  $g_M(\mathbf{x}) \rightarrow f_0(\mathbf{x})$ .

The convergence is uniform over  $\mathcal{X}$ . Indeed, the Lipschitz estimate in Lemma 1 holds uniformly over all measures satisfying the second-moment constraint, and taking suprema over such measures preserves the same Lipschitz constant for  $f_0$ ,  $f_M$ , and  $g_M$ . Since  $\mathcal{X}$  is compact, a finite-net argument upgrades the pointwise convergence above to

$$\sup_{\mathbf{x} \in \mathcal{X}} |f_M(\mathbf{x}) - f_0(\mathbf{x})| \rightarrow 0, \quad \sup_{\mathbf{x} \in \mathcal{X}} |g_M(\mathbf{x}) - f_0(\mathbf{x})| \rightarrow 0.$$

Thus, the claim follows.  $\square$

*Proof of Theorem 7.* Let  $f_0$ ,  $f_M$ , and  $g_M$  be as in Lemma 4, with  $\eta_M$  chosen as in Theorem 6. Let

$$h_M(\mathbf{x}) := \sup_{\mathbb{Q} \in \mathcal{P}'_{\epsilon^M}} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})]$$

be the sampled objective. The value convergence follows from the same sandwich argument used in Theorem 4, with Lemma 4 replacing Lemma 3; the solution convergence then follows from the same separation argument. Indeed, on the high-probability event of Theorem 6, we have

$$f_M(\mathbf{x}) \leq h_M(\mathbf{x}) \leq g_M(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{X}.$$

To obtain convergence in probability, fix  $\rho > 0$  and  $\nu > 0$ , and apply Theorem 6 with confidence parameter  $\nu$ . The corresponding sequence  $\eta_M$  still satisfies  $\eta_M \rightarrow 0$ . Hence Lemma 4 implies that, for all sufficiently large  $M$ ,

$$\sup_{\mathbf{x} \in \mathcal{X}} |f_M(\mathbf{x}) - f_0(\mathbf{x})| \leq \rho \quad \text{and} \quad \sup_{\mathbf{x} \in \mathcal{X}} |g_M(\mathbf{x}) - f_0(\mathbf{x})| \leq \rho.$$

On the containment event, the sandwich then gives  $\sup_{\mathbf{x} \in \mathcal{X}} |h_M(\mathbf{x}) - f_0(\mathbf{x})| \leq \rho$ . Therefore,

$$\text{Prob} \left( \sup_{\mathbf{x} \in \mathcal{X}} |h_M(\mathbf{x}) - f_0(\mathbf{x})| > \rho \right) \leq \nu$$

for all sufficiently large  $M$ . Since  $\nu > 0$  is arbitrary,

$$\sup_{\mathbf{x} \in \mathcal{X}} |h_M(\mathbf{x}) - f_0(\mathbf{x})| \xrightarrow{P} 0.$$

Consequently,

$$|\hat{v}'_M - \hat{v}| \leq \sup_{\mathbf{x} \in \mathcal{X}} |h_M(\mathbf{x}) - f_0(\mathbf{x})| \xrightarrow{P} 0.$$

It remains to prove convergence of minimizers. For  $\delta > 0$ , define

$$\mathcal{A}_\delta := \{\mathbf{x} \in \mathcal{X} : \text{dist}(\mathbf{x}, \mathcal{S}) \geq \delta\}.$$

If  $\mathcal{A}_\delta$  is empty, the claim is trivial. Otherwise, the function  $f_0$  is continuous and  $\mathcal{X}$  is compact, so

$$\rho_\delta := \inf_{\mathbf{x} \in \mathcal{A}_\delta} f_0(\mathbf{x}) - \hat{v} > 0.$$

Whenever  $\sup_{\mathbf{x} \in \mathcal{X}} |h_M(\mathbf{x}) - f_0(\mathbf{x})| < \rho_\delta/3$ , no minimizer of  $h_M$  can belong to  $\mathcal{A}_\delta$ . Therefore

$$\text{Prob}(\text{dist}(\hat{\mathbf{x}}'_M, \mathcal{S}) \geq \delta) \leq \text{Prob} \left( \sup_{\mathbf{x} \in \mathcal{X}} |h_M(\mathbf{x}) - f_0(\mathbf{x})| \geq \rho_\delta/3 \right) \rightarrow 0.$$

This proves  $\text{dist}(\hat{\mathbf{x}}'_M, \mathcal{S}) \xrightarrow{P} 0$ .

Lastly, we prove the suboptimality bound. Let  $\mathbf{x}_M^*$  be the exact optimizer from the statement. Since  $\hat{\mathbf{x}}'_M$  minimizes the sampled objective  $h_M$  and  $\mathcal{P}'_{\epsilon, M} \subseteq \mathcal{P}'_\epsilon$ , we have

$$\begin{aligned} & \sup_{\mathbb{Q} \in \mathcal{P}'_{\epsilon, M}} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}'_M, \tilde{\boldsymbol{\xi}})] - \sup_{\mathbb{Q} \in \mathcal{P}'_{\epsilon, M}} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}_M^*, \tilde{\boldsymbol{\xi}})] \\ & \leq \sup_{\mathbb{Q} \in \mathcal{P}'_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\hat{\mathbf{x}}'_M, \tilde{\boldsymbol{\xi}})] - \sup_{\mathbb{Q} \in \mathcal{P}'_{\epsilon, M}} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}_M^*, \tilde{\boldsymbol{\xi}})] \\ & \leq \sup_{\mathbb{Q} \in \mathcal{P}'_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}_M^*, \tilde{\boldsymbol{\xi}})] - \sup_{\mathbb{Q} \in \mathcal{P}'_{\epsilon, M}} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}_M^*, \tilde{\boldsymbol{\xi}})]. \end{aligned}$$

Next, by Theorem 6, setting  $\eta = \eta_M$ , we obtain the high-probability bound

$$\sup_{\mathbb{Q} \in \mathcal{P}'_\epsilon} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}_M^*, \tilde{\boldsymbol{\xi}})] - \sup_{\mathbb{Q} \in \mathcal{P}'_{\epsilon, M}} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}_M^*, \tilde{\boldsymbol{\xi}})] \leq \sup_{\mathbb{Q} \in \mathcal{P}'_{\epsilon, M}(\eta_M)} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}_M^*, \tilde{\boldsymbol{\xi}})] - \sup_{\mathbb{Q} \in \mathcal{P}'_{\epsilon, M}} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}_M^*, \tilde{\boldsymbol{\xi}})]$$

with probability at least  $1 - \tau$ . By the same moment-duality argument as in Proposition 10, the relaxed unbounded-support problem admits the dual obtained by augmenting (24) with the slack multiplier  $\gamma\eta_M$  and the domination constraint  $\gamma\mathbb{Z} - \nu \in \mathcal{M}_+(\mathcal{X})$ . Thus, the associated exact dual optimizer  $(\mathbf{x}_M^*, \alpha_M^*, \beta_M^*, \kappa_M^*, \nu_M^*)$  from the statement yields a feasible, though generally suboptimal, relaxed-dual solution  $(\mathbf{x}_M^*, \alpha_M^*, \beta_M^*, \kappa_M^*, \gamma_M^*, \nu_M^*)$ . We can therefore further upper bound the preceding display by

$$\begin{aligned} & \sup_{\mathbb{Q} \in \mathcal{P}'_{\epsilon, M}(\eta_M)} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}_M^*, \tilde{\boldsymbol{\xi}})] - \sup_{\mathbb{Q} \in \mathcal{P}'_{\epsilon, M}} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{x}_M^*, \tilde{\boldsymbol{\xi}})] \\ & \leq \alpha_M^* + \beta_M^* \Omega - \kappa_M^* (1 - 2/M) + \int_{\mathcal{X}} \left( \mathbb{E}_{\mathbb{P}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] + \epsilon \right) \nu_M(d\mathbf{z}) + \gamma_M^* \eta_M \\ & \quad - \left( \alpha_M^* + \beta_M^* \Omega - \kappa_M^* (1 - 2/M) + \int_{\mathcal{X}} \left( \mathbb{E}_{\mathbb{P}_N}[\ell(\mathbf{z}, \tilde{\boldsymbol{\xi}})] + \epsilon \right) \nu_M(d\mathbf{z}) \right) \\ & = \gamma_M^* \eta_M. \end{aligned}$$

Substituting the explicit expression for  $\eta_M$  proves the theorem.  $\square$

## C.2. Conic Programming Reformulations

*Proof of Theorem 8.* For each sampled point  $\mathbf{z}_m$ , define

$$g_m(\boldsymbol{\xi}) := \max_{k \in [J]} (\mathbf{a}_k(\mathbf{z}_m)^\top \boldsymbol{\xi} + b_k(\mathbf{z}_m)).$$

The semi-infinite constraint in (26) can be equivalently rewritten as

$$\ell(\mathbf{x}, \boldsymbol{\xi}) + \kappa \leq \alpha + \beta \|\boldsymbol{\xi}\|^2 + \frac{1}{M} \sum_{m \in [M]} \nu_m \ell(\mathbf{z}_m, \boldsymbol{\xi}) \quad \forall \boldsymbol{\xi} \in \mathbb{R}^D \boldsymbol{\xi} : \|\boldsymbol{\xi}\| \leq t, \quad (38)$$

$$\ell(\mathbf{x}, \boldsymbol{\xi}) \leq \alpha + \beta \|\boldsymbol{\xi}\|^2 + \frac{1}{M} \sum_{m \in [M]} \nu_m \ell(\mathbf{z}_m, \boldsymbol{\xi}) \quad \forall \boldsymbol{\xi} \in \mathbb{R}^D \boldsymbol{\xi} : \|\boldsymbol{\xi}\| > t. \quad (39)$$

Since  $\kappa \geq 0$ , (38) implies

$$\ell(\mathbf{x}, \boldsymbol{\xi}) \leq \alpha + \beta \|\boldsymbol{\xi}\|^2 + \frac{1}{M} \sum_{m \in [M]} \nu_m \ell(\mathbf{z}_m, \boldsymbol{\xi}) \quad \forall \boldsymbol{\xi} \in \mathbb{R}^D \boldsymbol{\xi} : \|\boldsymbol{\xi}\| \leq t.$$

Hence (39) can be simplified to

$$\ell(\mathbf{x}, \boldsymbol{\xi}) \leq \alpha + \beta \|\boldsymbol{\xi}\|^2 + \frac{1}{M} \sum_{m \in [M]} \nu_m \ell(\mathbf{z}_m, \boldsymbol{\xi}) \quad \forall \boldsymbol{\xi} \in \mathbb{R}^D \boldsymbol{\xi}.$$

We now reformulate these two semi-infinite systems separately. Assume first that  $\beta > 0$ ; the case  $\beta = 0$  will be treated separately after deriving the conic reformulations for the global and ball constraints.

For the global constraint, the same simplex linearization as in Theorem 5 shows that it is equivalent to the existence of multipliers  $\lambda_{jm} \in \mathbb{R}_+^J$  with  $\mathbf{e}^\top \lambda_{jm} = \nu_m$  such that, for every  $j \in [J]$ ,

$$\sup_{\boldsymbol{\xi} \in \mathbb{R}^{D\xi}} \left\{ \left( \mathbf{a}_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k \mathbf{a}_k(\mathbf{z}_m) \right)^\top \boldsymbol{\xi} - \beta \|\boldsymbol{\xi}\|^2 \right\} + b_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k b_k(\mathbf{z}_m) \leq \alpha.$$

Under this assumption, the quadratic supremum above is equal to  $\|\mathbf{c}_j\|^2/(4\beta)$ , where

$$\mathbf{c}_j := \mathbf{a}_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k \mathbf{a}_k(\mathbf{z}_m).$$

Therefore the global constraint is equivalent to the existence of  $\zeta_j \in \mathbb{R}_+$  satisfying

$$\zeta_j + b_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda_{jm}^k b_k(\mathbf{z}_m) \leq \alpha$$

and

$$\left\| \begin{bmatrix} \mathbf{c}_j \\ \zeta_j - \beta \end{bmatrix} \right\| \leq \zeta_j + \beta.$$

For the ball constraint (38), introduce simplex multipliers  $\lambda'_{jm} \in \mathbb{R}_+^J$  with  $\mathbf{e}^\top \lambda'_{jm} = \nu_m$ . Then the constraint is equivalent to

$$\sup_{\|\boldsymbol{\xi}\| \leq t} \left\{ \left( \mathbf{a}_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda'_{jm}^k \mathbf{a}_k(\mathbf{z}_m) \right)^\top \boldsymbol{\xi} - \beta \|\boldsymbol{\xi}\|^2 \right\} + b_j(\mathbf{x}) + \kappa - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda'_{jm}^k b_k(\mathbf{z}_m) \leq \alpha$$

for every  $j \in [J]$ . Let

$$\mathbf{c}'_j := \mathbf{a}_j(\mathbf{x}) - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda'_{jm}^k \mathbf{a}_k(\mathbf{z}_m).$$

Using Fenchel duality (Rockafellar 1970, Theorem 16.4), with the support function of the Euclidean ball  $\{\boldsymbol{\xi} : \|\boldsymbol{\xi}\| \leq t\}$ , whose support function is  $t\|\cdot\|$ , we obtain

$$\sup_{\|\boldsymbol{\xi}\| \leq t} \left\{ \mathbf{c}'_j{}^\top \boldsymbol{\xi} - \beta \|\boldsymbol{\xi}\|^2 \right\} = \inf_{\boldsymbol{\theta}_j \in \mathbb{R}^{D\xi}} \left\{ t\|\boldsymbol{\theta}_j\| + \sup_{\boldsymbol{\xi} \in \mathbb{R}^{D\xi}} [(\mathbf{c}'_j - \boldsymbol{\theta}_j)^\top \boldsymbol{\xi} - \beta \|\boldsymbol{\xi}\|^2] \right\}.$$

Since  $\beta > 0$ , the remaining unconstrained quadratic supremum equals  $\|\mathbf{c}'_j - \boldsymbol{\theta}_j\|^2/(4\beta)$ . Hence the ball constraint is equivalent to the existence of  $(\boldsymbol{\theta}_j, \zeta'_j) \in \mathbb{R}^{D\xi} \times \mathbb{R}_+$  such that

$$\zeta'_j + t\|\boldsymbol{\theta}_j\| + b_j(\mathbf{x}) + \kappa - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda'_{jm}^k b_k(\mathbf{z}_m) \leq \alpha$$

and

$$\left\| \begin{bmatrix} \mathbf{c}'_j - \boldsymbol{\theta}_j \\ \zeta'_j - \beta \end{bmatrix} \right\| \leq \zeta'_j + \beta.$$

The case  $\beta = 0$  is handled exactly as in Theorem 5. For the global constraint, finiteness of the supremum over  $\mathbb{R}^{D\xi}$  requires  $\mathbf{c}_j = \mathbf{0}$ , and the conic constraint reduces to

$$\left\| \begin{bmatrix} \mathbf{c}_j \\ \zeta_j \end{bmatrix} \right\| \leq \zeta_j,$$

so the reduced conic constraint imposes no positive lower bound on  $\zeta_j$  beyond  $\zeta_j \geq 0$ . Since  $\zeta_j$  enters the scalar inequality for the global constraint with coefficient one, the least restrictive feasible choice is  $\zeta_j = 0$ . For the ball constraint, the conic constraint reduces to

$$\left\| \begin{bmatrix} \mathbf{c}'_j - \boldsymbol{\theta}_j \\ \zeta'_j \end{bmatrix} \right\| \leq \zeta'_j,$$

which forces  $\boldsymbol{\theta}_j = \mathbf{c}'_j$  and imposes no positive lower bound on  $\zeta'_j$ . Since  $\zeta'_j$  enters the scalar inequality for the ball constraint with coefficient one, the least restrictive feasible choice is  $\zeta'_j = 0$ , and the scalar inequality becomes precisely the support-function condition for the Euclidean ball. Collecting the global and ball constraints yields exactly the finite conic reformulation stated in the theorem.

The displayed reformulation is convex. The objective is affine, while the sign, simplex, and equality constraints are affine. Since  $\mathbf{a}_j(\mathbf{x})$  and  $b_j(\mathbf{x})$  are affine in  $\mathbf{x}$ , the scalar global constraint is affine. The scalar ball constraint has the form

$$\zeta'_j + t \|\boldsymbol{\theta}_j\| + b_j(\mathbf{x}) + \kappa - \frac{1}{M} \sum_{m \in [M], k \in [J]} \lambda'_{jm^k} b_k(\mathbf{z}_m) \leq \alpha,$$

which is convex because  $t \geq 0$  and the Euclidean norm is convex. The two norm inequalities are standard second-order-cone constraints with affine arguments in the decision variables, and are therefore convex. Together with the convexity of the constraint  $\mathbf{x} \in \mathcal{X}$ , this shows that the finite reformulation is a convex conic program. The final tractability claim follows as in Theorem 5, since the Euclidean norm is directly second-order cone representable.  $\square$

## Appendix D: Tractable approximations for two-stage linear loss

In this section, we develop a tractable approximation to the two-stage DRO problem via linear decision rules.

THEOREM 9. *The problem (19) under the two-stage loss function in (6) can be conservatively approximated by the following conic program:*

$$\begin{aligned}
& \min \alpha + \beta\Omega + \frac{1}{M} \sum_{m \in [M]} \nu_m (\mathbb{E}_{\tilde{\mathbf{p}}_N} [\ell(\mathbf{z}_m, \tilde{\boldsymbol{\xi}})] + \epsilon) \\
& \text{s.t. } \mathbf{x} \in \mathcal{X}, \alpha \in \mathbb{R}, \beta \in \mathbb{R}_+, \boldsymbol{\nu} \in \mathbb{R}_+^M \\
& \quad \boldsymbol{\lambda}_m \in \mathbb{R}_+^L \quad \forall m \in [M], \boldsymbol{\theta} \in \mathbb{R}^{D\xi}, \zeta \in \mathbb{R}_+ \\
& \quad \mathbf{Y} \in \mathbb{R}^{D\mathbf{y} \times D\xi}, \mathbf{y}_0 \in \mathbb{R}^{D\mathbf{y}} \\
& \quad \sigma_{\Xi}([\mathbf{T}(\mathbf{x}) - \mathbf{W}\mathbf{Y}]_{l:}^{\top}) \leq [\mathbf{W}\mathbf{y}_0 - \mathbf{h}(\mathbf{x})]_l \quad \forall l \in [L] \\
& \quad \zeta + \sigma_{\Xi}(\boldsymbol{\theta}) + \mathbf{c}^{\top} \mathbf{x} + \mathbf{q}^{\top} \mathbf{y}_0 - \frac{1}{M} \sum_{m \in [M]} (\nu_m \mathbf{c}^{\top} \mathbf{z}_m + \mathbf{h}(\mathbf{z}_m)^{\top} \boldsymbol{\lambda}_m) \leq \alpha \\
& \quad \left\| \begin{bmatrix} \mathbf{Y}^{\top} \mathbf{q} - \frac{1}{M} \sum_{m \in [M]} \mathbf{T}(\mathbf{z}_m)^{\top} \boldsymbol{\lambda}_m - \boldsymbol{\theta} \\ \zeta - \beta \end{bmatrix} \right\| \leq \zeta + \beta \\
& \quad \mathbf{W}^{\top} \boldsymbol{\lambda}_m = \nu_m \mathbf{q} \quad \forall m \in [M].
\end{aligned}$$

*Proof* Consider the semi-infinite constraint in (19):

$$\sup_{\boldsymbol{\xi} \in \Xi} \left[ \ell(\mathbf{x}, \boldsymbol{\xi}) - \beta \|\boldsymbol{\xi}\|^2 - \frac{1}{M} \sum_{m \in [M]} \nu_m \ell(\mathbf{z}_m, \boldsymbol{\xi}) \right] \leq \alpha. \quad (40)$$

We approximate the recourse decision in the first term by the linear decision rule

$$\mathbf{y}(\boldsymbol{\xi}) = \mathbf{Y}\boldsymbol{\xi} + \mathbf{y}_0,$$

where  $\mathbf{Y} \in \mathbb{R}^{D\mathbf{y} \times D\xi}$  and  $\mathbf{y}_0 \in \mathbb{R}^{D\mathbf{y}}$  satisfy

$$\mathbf{T}(\mathbf{x})\boldsymbol{\xi} + \mathbf{h}(\mathbf{x}) \leq \mathbf{W}(\mathbf{Y}\boldsymbol{\xi} + \mathbf{y}_0) \quad \forall \boldsymbol{\xi} \in \Xi.$$

This robust linear constraint is equivalent, row by row, to

$$\sup_{\boldsymbol{\xi} \in \Xi} [\mathbf{T}(\mathbf{x}) - \mathbf{W}\mathbf{Y}]_{l:} \boldsymbol{\xi} \leq [\mathbf{W}\mathbf{y}_0 - \mathbf{h}(\mathbf{x})]_l \quad \forall l \in [L],$$

or, equivalently,

$$\sigma_{\Xi}([\mathbf{T}(\mathbf{x}) - \mathbf{W}\mathbf{Y}]_{l:}^{\top}) \leq [\mathbf{W}\mathbf{y}_0 - \mathbf{h}(\mathbf{x})]_l \quad \forall l \in [L].$$

Then, the two-stage loss function in (6) can be upper bounded by

$$\ell(\mathbf{x}, \boldsymbol{\xi}) \leq \mathbf{c}^{\top} \mathbf{x} + \mathbf{q}^{\top} (\mathbf{Y}\boldsymbol{\xi} + \mathbf{y}_0) \quad \forall \boldsymbol{\xi} \in \Xi.$$

Replacing only  $\ell(\mathbf{x}, \boldsymbol{\xi})$  by this upper bound, while keeping each sampled term  $\ell(\mathbf{z}_m, \boldsymbol{\xi})$  exact, yields a conservative approximation to (40):

$$\left. \begin{aligned}
& \max \mathbf{c}^{\top} \mathbf{x} + \mathbf{q}^{\top} (\mathbf{Y}\boldsymbol{\xi} + \mathbf{y}_0) - \beta \|\boldsymbol{\xi}\|^2 - \frac{1}{M} \sum_{m \in [M]} \nu_m (\mathbf{c}^{\top} \mathbf{z}_m + \mathbf{q}^{\top} \mathbf{y}_m) \\
& \text{s.t. } \boldsymbol{\xi} \in \Xi, \mathbf{y}_m \in \mathbb{R}^{D\mathbf{y}} \quad \forall m \in [M] \\
& \quad \mathbf{T}(\mathbf{z}_m)\boldsymbol{\xi} + \mathbf{h}(\mathbf{z}_m) \leq \mathbf{W}\mathbf{y}_m \quad \forall m \in [M]
\end{aligned} \right\} \leq \alpha.$$

Next, dualizing the recourse variables  $\{\mathbf{y}_m\}_{m \in [M]}$  and applying the support-function reformulation to the remaining maximization over  $\boldsymbol{\xi}$  yields the equivalent conic conditions: there exist  $\boldsymbol{\lambda}_m \in \mathbb{R}_+^L \forall m \in [M]$ ,  $\boldsymbol{\theta} \in \mathbb{R}^D \boldsymbol{\xi}$ ,  $\zeta \in \mathbb{R}_+$ , such that

$$\begin{aligned} \zeta + \sigma_{\Xi}(\boldsymbol{\theta}) + \mathbf{c}^\top \mathbf{x} + \mathbf{q}^\top \mathbf{y}_0 - \frac{1}{M} \sum_{m \in [M]} (\nu_m \mathbf{c}^\top \mathbf{z}_m + \mathbf{h}(\mathbf{z}_m)^\top \boldsymbol{\lambda}_m) &\leq \alpha \\ \left\| \begin{bmatrix} \mathbf{Y}^\top \mathbf{q} - \frac{1}{M} \sum_{m \in [M]} \mathbf{T}(\mathbf{z}_m)^\top \boldsymbol{\lambda}_m - \boldsymbol{\theta} \\ \zeta - \beta \end{bmatrix} \right\| &\leq \zeta + \beta \\ \mathbf{W}^\top \boldsymbol{\lambda}_m = \nu_m \mathbf{q} \quad \forall m \in [M]. & \end{aligned}$$

This proves the result.  $\square$

**REMARK 11 (TWO-STAGE APPROXIMATION ALTERNATIVES).** Two-stage (distributionally) robust optimization is generically NP-hard, which motivates the approximation in Theorem 9. If desired, tighter approximations via piecewise-linear decision rules (Georghiou et al. 2015, Ben-Tal et al. 2020) or quadratic decision rules (Fan and Hanasusanto 2024, Xu and Hanasusanto 2025) can also be developed. One could also, in principle, formulate a convex copositive program for the two-stage problem via the techniques of Hanasusanto and Kuhn (2018), Xu and Burer (2018). The resulting conic program admits tractable conservative semidefinite programming approximations that can be solved with standard off-the-shelf solvers.

## Appendix E: Details of experiments

### E.1. Numerical illustration of the Monte Carlo approximation

This experiment corresponds to Remark 9. The purpose is to isolate the sampling error in the finite-sample dual approximation, so we set  $\epsilon = 0$  and evaluate the sampled-dual solution against the empirical objective  $E_{\hat{P}_N}[\ell(\mathbf{x}, \tilde{\boldsymbol{\xi}})]$ . The symmetric construction below gives a known optimizer of this benchmark, avoiding the need for a large-reference Monte Carlo solution.

For each dimension  $d \in \{5, 25, 100\}$ , we take  $D_{\mathbf{x}} = D_{\boldsymbol{\xi}} = d$  and set both the decision set and uncertainty support to the Euclidean unit ball,

$$\mathcal{X} = \Xi = \{\mathbf{u} \in \mathbb{R}^d : \|\mathbf{u}\| \leq 1\}.$$

The second-moment cap is set to the empirical training second moment,  $\Omega = N^{-1} \sum_{i=1}^N \|\hat{\boldsymbol{\xi}}_i\|^2$ . The known optimizer is chosen as  $\mathbf{x}_* = 0.1\mathbf{e}$ , which belongs to  $\mathcal{X}$  for all three dimensions. The empirical reference distribution contains  $N = 100$  points and is constructed symmetrically: we draw 50 independent points uniformly from the unit ball and add their negatives. This symmetry makes  $\mathbf{x}_*$  an optimizer of the empirical objective for the shifted loss below.

The loss is a shifted weighted projection max-affine loss,

$$\ell(\mathbf{x}, \boldsymbol{\xi}) = \max_{r=1, \dots, K} c_r |\mathbf{q}_r^\top ((\mathbf{x} - \mathbf{x}_*) - \boldsymbol{\xi})|,$$

with  $K = 4$ . For each dimension, the directions  $\mathbf{q}_r \in \mathbb{R}^d$  are generated as independent dense Gaussian vectors and normalized to unit length. The weights  $c_r$  are linearly spaced between 0.5 and 1.5 and then rescaled to have average one. Equivalently, the loss has  $J = 2K = 8$  affine pieces with slopes  $\{\pm c_r \mathbf{q}_r : r = 1, \dots, K\}$ , so the conic reformulation in Section 3 can be solved directly.

For each dimension, the experiment uses 50 independent replications. In each replication, it draws a pool of 500 auxiliary decisions uniformly from  $\mathcal{X}$  and uses nested prefixes of this pool for  $M \in \{2, 5, 10, 20, 50, 100, 250, 500\}$ . For each  $M$ , the sampled dual (19) is solved using the compact-support conic reformulation with  $\sigma_{\Xi}(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|$ . As in Remark 9, the reported percentage gap is

$$\frac{E_{\hat{P}_N}[\ell(\hat{x}_M, \tilde{\boldsymbol{\xi}})] - E_{\hat{P}_N}[\ell(x_*, \tilde{\boldsymbol{\xi}})]}{E_{\hat{P}_N}[\ell(x_*, \tilde{\boldsymbol{\xi}})]} \times 100\%.$$

The plotted curves report the mean percentage gap over the 10 replications for each dimension, the shaded bands show the interquartile range, and the dashed reference line is a pooled log-scale fit with fixed slope  $-1/2$ , i.e.,  $C/\sqrt{M}$ .

## E.2. Newsvendor

The newsvendor figures in Section 5.1 are generated by the final selected-trial plotting blocks in the four three-dimensional notebooks. Earlier exploratory cells use different parameter values, but those preliminary settings are not used for the figure files included in the manuscript.

The experiment studies a three-product newsvendor with order vector  $\mathbf{x} \in \mathbb{R}_+^3$ , demand vector  $\tilde{\boldsymbol{\xi}} \in \mathbb{R}_+^3$ , and loss

$$\ell(\mathbf{x}, \boldsymbol{\xi}) = \sum_{j=1}^3 [h_j(x_j - \xi_j)_+ + b_j(\xi_j - x_j)_+].$$

Across all four final plotting blocks, the common decision-side parameters are

$$h = (0.1, 0.2, 0.3), \quad b = (1, 1, 1), \quad \rho = 0.05, \quad U = (50, 50, 50).$$

The demand coordinates are generated independently with mean vector  $\mu = (5, 6, 6)$  and standard-deviation vector  $\sigma = (5, 6, 8)$ . The in-sample data are split into training and validation subsets using the ratio 0.8/0.2. We use  $N \in \{10, 20, 30, 50, 75, 100, 150\}$  with  $T = 50$  independent trials for each sample size, and each trial is evaluated on an independent OOS sample of size  $N_{\text{OOS}} = 10000$ .

For both 2-WDRO and TIPM-DRO, the ambiguity radius is tuned over the common grid

$$\varepsilon \in \{a \times 10^b : a \in \{1, 3, 5, 7, 9\}, b \in \{-3, -2, -1, 0, 1, 2\}\}.$$

For IPM-DRO, the expectation over the auxiliary distribution is approximated by Monte Carlo sampling. The final plotting blocks use  $M = 800$  sampled decision points drawn uniformly from  $[0, 50]^3$ , and the sampled set is augmented with the SAA solution from the same trial. The second-moment cap is set to the empirical training second moment,  $\Omega = \frac{1}{N_{\text{tr}}} \sum_{i=1}^{N_{\text{tr}}} \|\hat{\xi}_i\|^2$ .

To state the dual reformulation used in the code, let  $\mathcal{B}_{\text{nv}} := \{0, 1\}^3$  denote the eight affine pieces of the separable newsvendor loss, and let  $\mathcal{R}_{\text{nv}} := \mathcal{B}_{\text{nv}} \cup \{T\}$  include the additional CVaR epigraph piece. For each  $s \in \mathcal{B}_{\text{nv}}$ , define

$$p_j^s = \begin{cases} h_j, & s_j = 0, \\ -b_j, & s_j = 1, \end{cases} \quad g_j^s = \begin{cases} -h_j, & s_j = 0, \\ b_j, & s_j = 1, \end{cases} \quad j = 1, 2, 3.$$

For the epigraph piece, set  $\mathbf{p}^T = \mathbf{0}$ ,  $\mathbf{g}^T = \mathbf{0}$ , and  $q_T = 1$ , while  $q_r = 0$  for  $r \in \mathcal{B}_{\text{nv}}$ . Then

$$\ell(\mathbf{x}, \boldsymbol{\xi}) = \max_{s \in \mathcal{B}_{\text{nv}}} \{(\mathbf{p}^s)^\top \mathbf{x} + (\mathbf{g}^s)^\top \boldsymbol{\xi}\}.$$

If  $\mathbf{z}_1, \dots, \mathbf{z}_M \in [0, 50]^3$  are the sampled decision points and

$$\hat{L}_m := \frac{1}{N_{\text{tr}}} \sum_{i=1}^{N_{\text{tr}}} \ell(\mathbf{z}_m, \hat{\xi}_i), \quad m \in [M],$$

then the sampled second-moment-constrained TIPM-DRO problem solved by the implementation is

$$\begin{aligned} & \min \left( 1 - \frac{1}{\rho} \right) t + \frac{1}{\rho} \left( \alpha + \frac{1}{M} \sum_{m \in [M]} \beta_m (\hat{L}_m + \epsilon) + \gamma \Omega \right) \\ & \text{s.t. } \mathbf{x} \in \mathbb{R}_+^3, \mathbf{x} \leq U, t, \alpha \in \mathbb{R}, \boldsymbol{\beta} \in \mathbb{R}_+^M, \gamma \in \mathbb{R}_+, \\ & \eta_r \in \mathbb{R}_+, \mathbf{u}_r \in \mathbb{R}^3, \mathbf{u}_r \leq \mathbf{0}, \lambda_{rms} \in \mathbb{R}_+ \quad \forall r \in \mathcal{R}_{\text{nv}}, \forall m \in [M], \forall s \in \mathcal{B}_{\text{nv}}, \\ & \eta_r + (\mathbf{p}^r)^\top \mathbf{x} + q_r t - \frac{1}{M} \sum_{m \in [M]} \sum_{s \in \mathcal{B}_{\text{nv}}} \lambda_{rms} (\mathbf{p}^s)^\top \mathbf{z}_m \leq \alpha \quad \forall r \in \mathcal{R}_{\text{nv}}, \\ & \sum_{s \in \mathcal{B}_{\text{nv}}} \lambda_{rms} = \beta_m \quad \forall r \in \mathcal{R}_{\text{nv}}, \forall m \in [M], \\ & \left\| \left[ \mathbf{g}^r - \frac{1}{M} \sum_{m \in [M]} \sum_{s \in \mathcal{B}_{\text{nv}}} \lambda_{rms} \mathbf{g}^s - \mathbf{u}_r \right] \right\| \leq \eta_r + \gamma \quad \forall r \in \mathcal{R}_{\text{nv}}. \end{aligned} \tag{41}$$

This is exactly the explicit three-dimensional dual/SOCP reformulation implemented in the experiment and used to generate the TIPM-DRO curves in Section 5.1.

For the four demand models, the Gaussian benchmark samples each coordinate from a Normal distribution with the target mean and variance and rejects negative draws. The rescaled  $\chi^2$  case

uses  $\xi_j = s_j Y_j$ , where  $Y_j \sim \chi_{\nu_j}^2$ ,  $\nu_j = 2(\mu_j/\sigma_j)^2$ , and  $s_j = \sigma_j^2/(2\mu_j)$ . The Lognormal case uses  $\xi_j \sim \text{LogNormal}(\mu_{\log,j}, \sigma_{\log,j})$ , where  $\sigma_{\log,j}^2 = \log(1 + \sigma_j^2/\mu_j^2)$  and  $\mu_{\log,j} = \log(\mu_j) - \frac{1}{2}\sigma_{\log,j}^2$ . The Pareto case uses  $\xi_j = x_{m,j}(1 + Y_j)$ , where  $Y_j \sim \text{Pareto}(\alpha_j)$ ,  $\alpha_j = 1 + \sqrt{1 + (\mu_j/\sigma_j)^2}$ , and  $x_{m,j} = \mu_j(\alpha_j - 1)/\alpha_j$ .

### E.3. Outlier-Corrupted Regression

The figures in Section 5.2 are generated by the final plotting blocks in our implementation, which call a common worker routine for each trial. In each trial, the code draws  $\mathbf{x}^*$  uniformly from the unit sphere in  $\mathbb{R}^d$ , generates clean covariates  $\mathbf{u}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  and responses  $v_i = \mathbf{x}^{*\top} \mathbf{u}_i$ , and evaluates the fitted estimator on an independent clean test sample of size  $N_{\text{OOS}}$ . It then corrupts  $\lfloor \omega N \rfloor$  training samples via  $\mathbf{u}_i \leftarrow C\mathbf{u}_i$  and  $v_i \leftarrow -C^2 v_i + \rho$ , randomly shuffles the sample, and splits it into training and validation subsets with ratio 0.8/0.2. Hyperparameters are selected by a trimmed validation loss that removes the largest  $\lceil \omega_{\text{tv}} N_{\text{val}} \rceil$  absolute residuals.

For the dimension sweep underlying Figure 3, the final plotting block uses  $N \in \{10, 20, 50, 75, 100\}$  and  $D_{\boldsymbol{\xi}} \in \{1, 5, 10, 50\}$  with  $T = 50$ ,  $N_{\text{OOS}} = 1000$ ,  $\omega = \omega_{\text{tv}} = 0.2$ ,  $C = 10$ ,  $\rho = 0.1$ , and  $\delta = 0.1$ . The OR-WDRO scale parameter is set to  $\omega_{\text{tv}} = 0.2$ , and Std WDRO and OR-WDRO tune the Wasserstein radius over  $\rho \in \{0.01, 0.1, 0.2, 0.5, 1.0\}$ . The worker also computes the plain LAD estimator, but the final plotting cell reports only Std WDRO, OR-WDRO( $2\omega$ ), OR-WDRO( $\omega$ ), and MoM TIPM-DRO.

For MoM TIPM-DRO, the code uses  $M = 800$  Monte Carlo test functions. The decision set radius is  $R = \sqrt{D_{\boldsymbol{\xi}}}$ , and the constraint  $\|\boldsymbol{\theta}\| \leq R$  is enforced with the same value of  $R$ . Let  $N_{\text{tr}} = \lfloor 0.8N \rfloor$ . The block count is chosen from  $K \in \{K_{\text{theory}}, \min(\lfloor N_{\text{tr}}/3 \rfloor, 2K_{\text{theory}}), \lfloor N_{\text{tr}}/2 \rfloor\} \cap \{1, \dots, N_{\text{tr}}\}$ , where  $K_{\text{theory}} = \left\lceil \frac{4(1+2\omega)}{(1-2\omega)^2} \log \frac{1}{\delta} \right\rceil$ . The IPM radius is tuned over  $\epsilon_{\text{IPM}} \in \{\epsilon_{\text{theory}} : s \in \{0.01, 0.05, 0.1, 0.2, 0.5, 1, 5, 10\}\}$  with  $\epsilon_{\text{theory}} = \frac{4\sqrt{e}\sigma_{\text{loss}}\Gamma_{\omega}}{\sqrt{N_{\text{tr}}}}$  and  $\Gamma_{\omega} = (1 - 2\omega_{\text{tv}})^{-1}$ . The second-moment cap is set to  $\Omega = 2\Omega_{\text{tr}}$ , where  $\Omega_{\text{tr}} = N_{\text{tr}}^{-1} \sum_{i=1}^{N_{\text{tr}}} \|(\mathbf{u}_i, v_i)\|^2$ . For each candidate  $R$ , the same set of sampled test functions is reused across all  $(K, \epsilon_{\text{IPM}})$  candidates in that trial to stabilize validation comparisons.

For scalar-response MAD regression, with decision variable  $\boldsymbol{\theta} \in \mathbb{R}^{D_{\boldsymbol{\xi}}}$ , the loss admits the piecewise-affine representation

$$\ell(\boldsymbol{\theta}, \boldsymbol{\xi}) = \max \left\{ \begin{bmatrix} \boldsymbol{\theta} \\ -1 \end{bmatrix}^{\top} \boldsymbol{\xi}, \begin{bmatrix} -\boldsymbol{\theta} \\ 1 \end{bmatrix}^{\top} \boldsymbol{\xi} \right\}, \quad \boldsymbol{\xi} = (\mathbf{u}, v) \in \mathbb{R}^{D_{\boldsymbol{\xi}}+1}.$$

Let  $\{\mathbf{z}_m\}_{m=1}^M \subset \mathbb{R}^{D\xi}$  be the sampled test functions and let  $\mathbf{Z} \in \mathbb{R}^{M \times D\xi}$  collect the row vectors  $\mathbf{z}_m^\top$ . Writing  $\widehat{\mu}_{\text{MoM}}(\mathbf{z}_m)$  for the MoM reference from Section 2.2, the implementation solves the conic program

$$\begin{aligned}
& \min \alpha + \beta\Omega + \frac{1}{M} \sum_{m \in [M]} \nu_m (\widehat{\mu}_{\text{MoM}}(\mathbf{z}_m) + \epsilon) \\
& \text{s.t. } \boldsymbol{\theta} \in \mathbb{R}^{D\xi}, \|\boldsymbol{\theta}\| \leq R, \alpha \in \mathbb{R}, \beta \in \mathbb{R}_+, \boldsymbol{\nu} \in \mathbb{R}_+^M, \\
& \quad \boldsymbol{\lambda}_1^+, \boldsymbol{\lambda}_1^-, \boldsymbol{\lambda}_2^+, \boldsymbol{\lambda}_2^- \in \mathbb{R}_+^M, \zeta_1, \zeta_2 \in \mathbb{R}_+, \\
& \quad \boldsymbol{\lambda}_1^+ + \boldsymbol{\lambda}_1^- = \boldsymbol{\nu}, \quad \boldsymbol{\lambda}_2^+ + \boldsymbol{\lambda}_2^- = \boldsymbol{\nu}, \\
& \quad \zeta_1 \leq \alpha, \quad \zeta_2 \leq \alpha, \\
& \quad \left\| \begin{bmatrix} \boldsymbol{\theta} - \frac{1}{M} \mathbf{Z}^\top (\boldsymbol{\lambda}_1^+ - \boldsymbol{\lambda}_1^-) \\ -1 + \frac{1}{M} \mathbf{e}^\top (\boldsymbol{\lambda}_1^+ - \boldsymbol{\lambda}_1^-) \\ \zeta_1 - \beta \end{bmatrix} \right\| \leq \zeta_1 + \beta, \\
& \quad \left\| \begin{bmatrix} -\boldsymbol{\theta} - \frac{1}{M} \mathbf{Z}^\top (\boldsymbol{\lambda}_2^+ - \boldsymbol{\lambda}_2^-) \\ 1 + \frac{1}{M} \mathbf{e}^\top (\boldsymbol{\lambda}_2^+ - \boldsymbol{\lambda}_2^-) \\ \zeta_2 - \beta \end{bmatrix} \right\| \leq \zeta_2 + \beta.
\end{aligned} \tag{42}$$

This is exactly the formulation implemented by the MoM TIPM-DRO solver used in Section 5.2, written in the notation of Section 2.

For Figure 4, the final plotting block fixes  $N = 50$ ,  $D_\xi = 10$ , and  $C \in \{5, 6, \dots, 20\}$  while keeping  $T = 50$ ,  $N_{\text{OS}} = 1000$ ,  $\omega = \omega_{\text{tv}} = 0.2$ ,  $\rho = 0.1$ ,  $\delta = 0.1$ , and  $M = 800$ . In this block, the OR-WDRO scale parameter is fixed at  $\omega_{\text{tv}} = 0.2$ , and the Wasserstein radius grid is expanded to  $\rho \in \{0.01, 0.05, 0.1, 0.5, 1.0, 2.0\}$ . The plotting cell for Figure 4 reports only OR-WDRO( $2\omega$ ), OR-WDRO( $\omega$ ), and MoM TIPM-DRO. Apart from the fixed dimension and the sweep over  $C$ , these are the only substantive differences relative to the dimension-sweep experiments.