



Delft University of Technology

Regularization via mass transportation

Shaéezadeh-Abadeh, Soroosh; Kuhn, Daniel; Esfahani, Peyman Mohajerin

Publication date

2019

Document Version

Final published version

Published in

Journal of Machine Learning Research

Citation (APA)

Shaéezadeh-Abadeh, S., Kuhn, D., & Esfahani, P. M. (2019). Regularization via mass transportation. *Journal of Machine Learning Research*, 20, Article 103.

Important note

To cite this publication, please use the final published version (if applicable).
Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights.
We will remove access to the work immediately and investigate your claim.

Regularization via Mass Transportation

Soroosh Shafieezadeh-Abadeh

Daniel Kuhn

*Risk Analytics and Optimization Chair,
EPFL, Switzerland*

SOROOSH.SHAFIEE@EPFL.CH

DANIEL.KUHN@EPFL.CH

Peyman Mohajerin Esfahani

*Delft Center for Systems and Control
TU Delft, The Netherlands*

P.MOHAJERINESFAHANI@TUDELFT.NL

Editor: Koji Tsuda

Abstract

The goal of regression and classification methods in supervised learning is to minimize the empirical risk, that is, the expectation of some loss function quantifying the prediction error under the empirical distribution. When facing scarce training data, overfitting is typically mitigated by adding regularization terms to the objective that penalize hypothesis complexity. In this paper we introduce new regularization techniques using ideas from distributionally robust optimization, and we give new probabilistic interpretations to existing techniques. Specifically, we propose to minimize the worst-case expected loss, where the worst case is taken over the ball of all (continuous or discrete) distributions that have a bounded transportation distance from the (discrete) empirical distribution. By choosing the radius of this ball judiciously, we can guarantee that the worst-case expected loss provides an upper confidence bound on the loss on test data, thus offering new generalization bounds. We prove that the resulting regularized learning problems are tractable and can be tractably kernelized for many popular loss functions. The proposed approach to regularization is also extended to neural networks. We validate our theoretical out-of-sample guarantees through simulated and empirical experiments.

Keywords: Distributionally robust optimization, optimal transport, Wasserstein distance, robust optimization, regularization

1. Introduction

The fields of machine learning and optimization are closely intertwined. On the one hand, optimization algorithms are routinely used for the solution of classical machine learning problems. Conversely, recent advances in optimization under uncertainty have inspired many new machine learning models.

From a conceptual point of view, many statistical learning tasks give naturally rise to stochastic optimization problems. Indeed, they aim to find an estimator from within a prescribed hypothesis space that minimizes the expected value of some loss function. The loss function quantifies the estimator's ability to correctly predict random outputs (*i.e.*, dependent variables or labels) from random inputs (*i.e.*, independent variables or features). Unfortunately, such stochastic optimization problems cannot be solved exactly because the input-output distribution, which is needed to evaluate the expected loss in the objective

function, is not accessible and only indirectly observable through finitely many training samples. Approximating the expected loss with the empirical loss, that is, the average loss across all training samples, yields fragile estimators that are sensitive to perturbations in the data and suffer from overfitting.

Regularization is the standard remedy to combat overfitting. Regularized learning models minimize the sum of the empirical loss and a penalty for hypothesis complexity, which is typically chosen as a norm of the hypothesis. There is ample empirical evidence that regularization reduces a model’s generalization error. Statistical learning theory reasons that regularization implicitly restricts the hypothesis space, thereby controlling the gap between the training error and the testing error, see, *e.g.*, Bartlett and Mendelson (2002). However, alternative explanations for the practical success of regularization are possible. In particular, ideas from modern robust optimization (Ben-Tal et al. (2009)) recently led to a fresh perspective on regularization.

Robust regression and classification models seek estimators that are immunized against adversarial perturbations in the training data. They have received considerable attention since the seminal treatise on robust least-squares regression by El Ghaoui and Lebret (1997), who seem to be the first authors to discover an intimate connection between robustification and regularization. Specifically, they show that minimizing the worst-case residual error with respect to all perturbations in a Frobenius norm-uncertainty set is equivalent to a Tikhonov regularization procedure. Xu et al. (2010) disclose a similar equivalence between robust least-squares regression with a feature-wise independent uncertainty set and the celebrated Lasso (least absolute shrinkage and selection operator) algorithm. Leveraging this new robustness interpretation, they extend Lasso to a wider class of regularization schemes tailored to regression problems with disturbances that are coupled across features. In the context of classification, Xu et al. (2009) provide a linkage between robustification over non-box-typed uncertainty sets and the standard regularization scheme of support vector machines. A comprehensive characterization of the conditions under which robustification and regularization are equivalent has recently been compiled by Bertsimas and Copenhaver (2017).

New learning models have also been inspired by recent advances in the emerging field of *distributionally* robust optimization, which bridges the gap between the conservatism of robust optimization and the specificity of stochastic programming. Distributionally robust optimization seeks to minimize a worst-case expected loss, where the worst case is taken with respect to all distributions in an ambiguity set, that is, a family of distributions consistent with the given prior information on the uncertainty, see, *e.g.*, Calafiore and El Ghaoui (2006), Delage and Ye (2010), Goh and Sim (2010), Wiesemann et al. (2014) and the references therein. Ambiguity sets are often characterized through generalized moment conditions. For instance, Lanckriet et al. (2002a) propose a distributionally robust minimax probability machine for binary classification, where both classes are encoded by the first and second moments of their features, and the goal is to find a linear classifier that minimizes the worst-case misclassification error in view of all possible input distributions consistent with the given moment information. By construction, this approach forces the worst-case accuracies of both classes to be equal. Huang et al. (2004) propose a generalization of the minimax probability machine that allows for uneven worst-case classification accuracies. Lanckriet et al. (2002b) extend the minimax probability machine to account for estimation errors in the mean vectors and covariance matrices. Strohmann and Grudic (2003) and Bhattacharyya (2004) develop

minimax probability machines for regression and feature selection, respectively. Shivaswamy et al. (2006) study linear classification problems trained with incomplete and noisy features, where each training sample is modeled by an ambiguous distribution with known first and second-order moments. The authors propose to address such classification problems with a distributionally robust soft margin support vector machine and then prove that it is equivalent to a classical robust support vector machine with a feature-wise uncertainty set. Farnia and Tse (2016) investigate distributionally robust learning models with moment ambiguity sets that restrict the marginal of the features to the empirical marginal. The authors highlight similarities and differences to classical regression models.

Ambiguity sets containing all distributions that share certain low-order moments are computationally attractive but fail to converge to a singleton when the number N of training samples tends to infinity. Thus, they preclude any asymptotic consistency results. A possible remedy is to design spherical ambiguity sets with respect to some probability distance functions and to drive their radii to zero as N grows. Examples include the ϕ -divergence ambiguity sets proposed by Ben-Tal et al. (2013) or the Wasserstein ambiguity sets studied by Mohajerin Esfahani and Kuhn (2018) and Zhao and Guan (2018). Blanchet and Murthy (2019) and Gao and Kleywegt (2016) consider generalized Wasserstein ambiguity sets defined over Polish spaces.

In this paper we investigate distributionally robust learning models with Wasserstein ambiguity sets. The Wasserstein distance between two distributions is defined as the minimum cost of transporting one distribution to the other, where the cost of moving a unit point mass is determined by the ground metric on the space of uncertainty realizations. In computer science the Wasserstein distance is therefore sometimes aptly termed the ‘earth mover’s distance’ (Rubner et al. (2000)). Following Mohajerin Esfahani and Kuhn (2018), we define Wasserstein ambiguity sets as balls with respect to the Wasserstein distance that are centered at the empirical distribution on the training samples. These ambiguity sets contain all (continuous or discrete) distributions that can be converted to the (discrete) empirical distribution at bounded transportation cost.

Wasserstein distances are widely used in machine learning to compare histograms. For example, Rubner et al. (2000) use the Wasserstein distance as a metric for image retrieval with a focus on applications to color and texture. Cuturi (2013) and Benamou et al. (2015) propose fast iterative algorithms to compute a regularized Wasserstein distance between two high-dimensional discrete distributions for image classification tasks. Moreover, Cuturi and Doucet (2014) develop first-order algorithms to compute the Wasserstein barycenter between several empirical probability distributions, which has applications in clustering. Arjovsky et al. (2017) utilize the Wasserstein distance to measure the distance between the data distribution and the model distribution in generative adversarial networks. Furthermore, Frogner et al. (2015) propose a learning algorithm based on the Wasserstein distance to predict multi-label outputs.

Distributionally robust optimization models with Wasserstein ambiguity sets were introduced to the realm of supervised learning by Shafieezadeh-Abadeh et al. (2015), who show that distributionally robust logistic regression problems admit a tractable reformulation and encapsulate the classical as well as the popular regularized logistic regression problems as special cases. When the Wasserstein ball is restricted to distributions on a compact set, the problem becomes intractable but can still be addressed with an efficient decomposition

algorithm due to Luo and Mehrotra (2019). Support vector machine models with distributionally robust chance constraints over Wasserstein ambiguity sets are studied by Lee and Mehrotra (2015). These models are equivalent to hard semi-infinite programs and can be solved approximately with a cutting plane algorithm.

Wasserstein ambiguity sets are popular for their attractive statistical properties. For example, Fournier and Guillin (2015) prove that the empirical distribution on N training samples converges in Wasserstein distance to the true distribution at rate $\mathcal{O}(N^{-1/(n+1)})$, where n denotes the feature dimension. This implies that properly scaled Wasserstein balls constitute natural confidence regions for the data-generating distribution. The worst-case expected prediction loss over all distributions in a Wasserstein ball thus provides an upper confidence bound on the expected loss under the unknown true distribution; see Mohajerin Esfahani and Kuhn (2018). Blanchet et al. (2016) show, however, that radii of the order $\mathcal{O}(N^{-1/2})$ are asymptotically optimal even though the corresponding Wasserstein balls are too small to contain the true distribution with constant confidence. For Wasserstein distances of type two (where the transportation cost equals the squared ground metric) Blanchet et al. (2017) develop a systematic methodology for selecting the ground metric. Generalization bounds for the *worst-case* prediction loss with respect to a Wasserstein ball centered at the *true* distribution are derived by Lee and Raginsky (2018) in order to address emerging challenges in domain adaptation problems, where the distributions of the training and test samples can differ.

This paper extends the results by Shafieezadeh-Abadeh et al. (2015) on distributionally robust logistic regression along several dimensions. Our main contributions can be summarized as follows:

- **Tractability:** We propose data-driven distributionally robust regression and classification models that hedge against all input-output distributions in a Wasserstein ball. We demonstrate that the emerging semi-infinite optimization problems admit equivalent reformulations as tractable convex programs for many commonly used loss functions and for spaces of linear hypotheses. We also show that lifted variants of these new learning models are kernelizable and thus offer an efficient procedure for optimizing over all nonlinear hypotheses in a reproducible kernel Hilbert space. Finally, we study distributionally robust learning models over families of feed-forward neural networks. We show that these models can be approximated by regularized empirical loss minimization problems with a convex regularization term and can be addressed with a stochastic proximal gradient descent algorithm.
- **Probabilistic Interpretation of Existing Regularization Techniques:** We show that the classical regularized learning models emerge as special cases of our framework when the cost of moving probability mass along the output space tends to infinity. In this case, the regularization function and its regularization weight are determined by the transportation cost on the input space and the radius of the Wasserstein ball underlying the distributionally robust optimization model, respectively.
- **Generalization Bounds:** We demonstrate that the proposed distributionally robust learning models enjoy new generalization bounds that can be obtained under minimal assumptions. In particular, they do not rely on any notions of hypothesis complexity

and may therefore even extend to hypothesis spaces with infinite VC-dimensions. A naïve generalization bound is obtained by leveraging modern measure concentration results, which imply that Wasserstein balls constitute confidence sets for the unknown data-generating distribution. Unfortunately, this generalization bound suffers from a curse of dimensionality and converges slowly for high input dimensions. By imposing bounds on the hypothesis space, however, we can derive an improved generalization bound, which essentially follows a dimension-independent square root law reminiscent of the central limit theorem.

- **Relation to Robust Optimization:** In classical robust regression and classification the training samples are viewed as uncertain variables that range over a joint uncertainty set, and the best hypothesis is found by minimizing the worst-case loss over this set. We prove that the classical robust and new distributionally robust learning models are equivalent if the data satisfies a dispersion condition (for regression) or a separability condition (for classification). While there is no efficient algorithm for solving the robust learning models in the absence of this condition, the distributionally robust models are efficiently solvable irrespective of the underlying training datasets.
- **Confidence Intervals for Error and Risk:** Using distributionally robust optimization techniques based on the Wasserstein ball, we develop two tractable linear programs whose optimal values provide a confidence interval for the absolute prediction error of any fixed regressor or the misclassification risk of any fixed classifier.
- **Worst-Case Distributions:** We formulate tractable convex programs that enable us to efficiently compute a worst-case distribution in the Wasserstein ball for any fixed hypothesis. This worst-case distribution can be useful for stress tests or contamination experiments.

The rest of the paper develops as follows. Section 2 introduces our new distributionally robust learning models. Section 3 provides finite convex reformulations for learning problems over linear and nonlinear hypothesis spaces and describes efficient procedures for constructing worst-case distributions. Moreover, it compares the new distributionally robust method against existing robust optimization and regularization approaches. Section 4 develops new generalization bounds, while Section 5 addresses error and risk estimation. Numerical experiments are reported in Section 6. All proofs are relegated to the appendix.

1.1. Notation

Throughout this paper, we adopt the conventions of extended arithmetics, whereby $\infty \cdot 0 = 0 \cdot \infty = 0/0 = 0$ and $\infty - \infty = -\infty + \infty = 1/0 = \infty$. The inner product of two vectors $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^n$ is denoted by $\langle \mathbf{x}, \mathbf{x}' \rangle$, and for any norm $\|\cdot\|$ on \mathbb{R}^n , we use $\|\cdot\|_*$ to denote its dual norm defined through $\|\mathbf{x}\|_* = \sup \{ \langle \mathbf{x}, \mathbf{x}' \rangle : \|\mathbf{x}'\| \leq 1 \}$. The conjugate of an extended real-valued function $f(\mathbf{x})$ on \mathbb{R}^n is defined as $f^*(\mathbf{x}) = \sup_{\mathbf{x}'} \langle \mathbf{x}, \mathbf{x}' \rangle - f(\mathbf{x}')$. The indicator function of a set $\mathbb{X} \subseteq \mathbb{R}^n$ is defined as $\delta_{\mathbb{X}}(\mathbf{x}) = 0$ if $\mathbf{x} \in \mathbb{X}$; $= \infty$ otherwise. Its conjugate $S_{\mathbb{X}}(\mathbf{x}) = \sup \{ \langle \mathbf{x}', \mathbf{x} \rangle : \mathbf{x}' \in \mathbb{X} \}$ is termed the support function of \mathbb{X} . The characteristic function of \mathbb{X} is defined through $\mathbf{1}_{\mathbb{X}}(\mathbf{x}) = 1$ if $\mathbf{x} \in \mathbb{X}$; $= 0$ otherwise. For a proper cone $\mathcal{C} \subseteq \mathbb{R}^n$ the relation $\mathbf{x} \succeq_{\mathcal{C}} \mathbf{x}'$ indicates that $\mathbf{x} - \mathbf{x}' \in \mathcal{C}$. The cone dual to \mathcal{C} is defined as

$\mathcal{C}^* = \{\mathbf{x}' : \langle \mathbf{x}', \mathbf{x} \rangle \geq 0 \forall \mathbf{x} \in \mathcal{C}\}$. The Lipschitz modulus of a function $L : \mathbb{X} \rightarrow \mathbb{R}$ is denoted by $\text{lip}(L) = \sup_{\mathbf{x}, \mathbf{x}' \in \mathbb{X}} \{|L(\mathbf{x}) - L(\mathbf{x}')| / \|\mathbf{x} - \mathbf{x}'\| : \mathbf{x} \neq \mathbf{x}'\}$. If \mathbb{P} is a distribution on a set Ξ , then \mathbb{P}^N denotes the N -fold product of \mathbb{P} on the Cartesian product Ξ^N . For $N \in \mathbb{N}$ we define $[N] = \{1, \dots, N\}$. A list of commonly used notations is provided in the following table.

\mathbb{X}	input space	\mathbb{Y}	output space
ℓ	loss function	L	univariate loss function
\mathbb{H}	hypothesis space	\mathcal{K}	kernel matrix
f^*	conjugate of f	$\text{lip}(f)$	Lipschitz modulus of f
\mathcal{C}^*	dual cone of \mathcal{C}	$\ \cdot\ _*$	dual norm of $\ \cdot\ $
$S_{\mathbb{X}}$	support function of \mathbb{X}	$\delta_{\mathbb{X}}$	indicator function of \mathbb{X}
$1_{\mathbb{X}}$	characteristic function of \mathbb{X}	$[N]$	$\{1, \dots, N\}$

2. Problem Statement

We first introduce the basic terminology and then describe our new perspective on regularization.

2.1. Classical Statistical Learning

The goal of supervised learning is to infer an unknown target function $f : \mathbb{X} \rightarrow \mathbb{Y}$ from limited data. The target function maps any input $\mathbf{x} \in \mathbb{X}$ (*e.g.*, information on the frequency of certain keywords in an email) to some output $y \in \mathbb{Y}$ (*e.g.*, a label $+1$ (-1) if the email is likely (unlikely) to be a spam message). If the true target function was accessible, it could be used as a means to reliably predict outputs from inputs (*e.g.*, it could be used to recognize spam messages in an automated fashion). In a supervised learning framework, however, one has only access to finitely many input-output examples $(\hat{\mathbf{x}}_i, \hat{y}_i)$ for $i = 1, \dots, N$ (*e.g.*, a database of emails which have been classified by a human as legitimate or as spam messages). We will henceforth refer to these examples as the *training data* or the *in-sample data*. It is assumed that the training samples are mutually independent and follow an unknown distribution \mathbb{P} on $\mathbb{X} \times \mathbb{Y}$.

The supervised learning problems are commonly subdivided into *regression problems*, where the output y is continuous and $\mathbb{Y} = \mathbb{R}$, and *classification problems*, where y is categorical and $\mathbb{Y} = \{+1, -1\}$. As the space of all functions from \mathbb{X} to \mathbb{Y} is typically vast, it may be very difficult to learn the target function from finitely many training samples. Thus, it is convenient to restrict the search space to a structured family of candidate functions $\mathbb{H} \subseteq \mathbb{R}^{\mathbb{X}}$ such as the space of all linear functions, some reproducible kernel Hilbert space or the family of all feed-forward neural networks with a fixed number of layers. We henceforth refer to each candidate function $h \in \mathbb{H}$ as a *hypothesis* and to \mathbb{H} as the *hypothesis space*.

A *learning algorithm* is a method for finding a hypothesis $h \in \mathbb{H}$ that faithfully replicates the unknown target function f . Specifically, in regression we seek to approximate f with a hypothesis h , and in classification we seek to approximate f with a *thresholded* hypothesis $\text{sign}(h)$. Many learning algorithms achieve this goal by minimizing the in-sample error, that is, the empirical average of a loss function $\ell : \mathbb{R} \times \mathbb{Y} \rightarrow \mathbb{R}_+$ that estimates the mismatch

between the output predicted by $h(\mathbf{x})$ and the actual output y for a particular input-output pair (\mathbf{x}, y) . Any such algorithm solves a minimization problem of the form

$$\inf_{h \in \mathbb{H}} \left\{ \frac{1}{N} \sum_{i=1}^N \ell(h(\hat{\mathbf{x}}_i), \hat{y}_i) = \mathbb{E}^{\hat{\mathbb{P}}_N} [\ell(h(\mathbf{x}), y)] \right\}, \quad (1)$$

where $\hat{\mathbb{P}}_N = \frac{1}{N} \sum_{i=1}^N \delta_{(\hat{\mathbf{x}}_i, \hat{y}_i)}$ denotes the *empirical distribution*, that is, the uniform distribution on the training data. For different choices of the loss function ℓ , the generic supervised learning problem (1) reduces to different popular regression and classification problems from the literature.

EXAMPLES OF REGRESSION MODELS

For ease of exposition, we focus here on learning models with $\mathbb{X} \subseteq \mathbb{R}^n$ and $\mathbb{Y} \subseteq \mathbb{R}$, where \mathbb{H} is set to the space of all linear hypotheses $h(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle$ with $\mathbf{w} \in \mathbb{R}^n$. Thus, there is a one-to-one correspondence between hypotheses and weight vectors \mathbf{w} . Moreover, we focus on loss functions of the form $\ell(h(\mathbf{x}), y) = L(h(\mathbf{x}) - y) = L(\langle \mathbf{w}, \mathbf{x} \rangle - y)$ that are generated by a univariate loss function L .

1. A rich class of *robust regression* problems is obtained from (1) if ℓ is generated by the *Huber loss function* with robustness parameter $\delta > 0$, which is defined as $L(z) = \frac{1}{2}z^2$ if $|z| \leq \delta$; $= \delta(|z| - \frac{1}{2}\delta)$ otherwise. Note that the Huber loss function is both convex and smooth and reduces to the squared loss $L(z) = \frac{1}{2}z^2$ for $\delta \uparrow \infty$, which is routinely used in ordinary least squares regression. Problem (1) with squared loss seeks a hypothesis \mathbf{w} under which $\langle \mathbf{w}, \mathbf{x} \rangle$ approximates the mean of y conditional on \mathbf{x} . The Huber loss function for finite δ favors similar hypotheses but is less sensitive to outliers.
2. The *support vector regression* problem (Smola and Schölkopf, 2004) emerges as a special case of (1) if ℓ is generated by the ϵ -insensitive loss function $L(z) = \max\{0, |z| - \epsilon\}$ with $\epsilon \geq 0$. In this setting, a training sample $(\hat{\mathbf{x}}_i, \hat{y}_i)$ is penalized in (1) only if the output $\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle$ predicted by hypothesis \mathbf{w} differs from the true output \hat{y}_i by more than ϵ . Support vector regression thus seeks hypotheses \mathbf{w} under which all training samples reside within a slab of width 2ϵ centered around the hyperplane $\{(\mathbf{x}, y) : \langle \mathbf{w}, \mathbf{x} \rangle = y\}$.
3. The *quantile regression* problem (Koenker, 2005) is obtained from (1) if ℓ is generated by the *pinball loss function* $L(z) = \max\{-\tau z, (1 - \tau)z\}$ defined for $\tau \in [0, 1]$. Quantile regression seeks hypotheses that approximate the $\tau \times 100\%$ -quantile of the output conditional on the input. More precisely, it seeks hypotheses \mathbf{w} for which $\tau \times 100\%$ of all training samples lie in the halfspace $\{(\mathbf{x}, y) : \langle \mathbf{w}, \mathbf{x} \rangle \geq y\}$.

EXAMPLES OF CLASSIFICATION MODELS

We focus here on linear learning models with $\mathbb{X} \subseteq \mathbb{R}^n$ and $\mathbb{Y} = \{+1, -1\}$, where \mathbb{H} is again identified with the space of all linear hypotheses $h(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle$ with $\mathbf{w} \in \mathbb{R}^n$. Moreover, we focus on loss functions of the form $\ell(\mathbf{x}, y) = L(yh(\mathbf{x})) = L(y\langle \mathbf{w}, \mathbf{x} \rangle)$ generated by a univariate loss function L .

1. The *support vector machine* problem (Cortes and Vapnik, 1995) is obtained from (1) if ℓ is generated by the *hinge loss function* $L(z) = \max\{0, 1 - z\}$, which is large if z is small. Thus, a training sample $(\hat{\mathbf{x}}_i, \hat{y}_i)$ is penalized in (1) if the output $\text{sign}(\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle)$ predicted by hypothesis \mathbf{w} and the true output \hat{y}_i have opposite signs. More precisely, support vector machines seek hypotheses \mathbf{w} under which the inputs of all training samples with output $+1$ reside in the halfspace $\{\mathbf{x} : \langle \mathbf{w}, \mathbf{x} \rangle \geq 1\}$, while the inputs of training samples with output -1 are confined to $\{\mathbf{x} : \langle \mathbf{w}, \mathbf{x} \rangle \leq -1\}$.
2. An alternative support vector machine problem is obtained from (1) if ℓ is generated by the *smooth hinge loss function*, which is defined as $L(z) = \frac{1}{2} - z$ if $z \leq 0$; $= \frac{1}{2}(1 - z)^2$ if $0 < z < 1$; $= 0$ otherwise. The smooth hinge loss inherits many properties of the ordinary hinge loss but has a continuous derivative. Thus, it may be amenable to faster optimization algorithms.
3. The *logistic regression* problem (Hosmer et al., 2013) emerges as a special case of (1) if ℓ is generated by the *logloss function* $L(z) = \log(1 + e^{-z})$, which is large if z is small—similar to the hinge loss function. In this case the objective function of (1) can be viewed as the log-likelihood function corresponding to the logistic model $\mathbb{P}(y = 1|\mathbf{x}) = [1 + \exp(-\langle \mathbf{w}, \mathbf{x} \rangle)]^{-1}$ for the conditional probability of $y = 1$ given \mathbf{x} . Thus, logistic regression allows us to learn the conditional distribution of y given \mathbf{x} .

Remark 1 (Convex approximation) *Note that the hinge loss and the logloss functions represent convex approximations for the (discontinuous) one-zero loss defined through $L(z) = 1$ if $z \leq 0$; $= 0$ otherwise.*

In practice there may be many hypotheses that are compatible with the given training data and thus achieve a small empirical loss in (1). Any such hypothesis would accurately predict outputs from inputs *within the training dataset* (Defourny, 2010). However, due to overfitting, these hypotheses might constitute poor predictors *beyond the training dataset*, that is, on inputs that have not yet been recorded in the database. Mathematically, even if the in-sample error $\mathbb{E}^{\hat{\mathbb{P}}^N}[\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)]$ of a given hypothesis \mathbf{w} is small, the out-of-sample error $\mathbb{E}^{\mathbb{P}}[\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)]$ with respect to the unknown true input-output distribution \mathbb{P} may be large.

Regularization is the standard remedy to combat overfitting. Instead of naïvely minimizing the in-sample error as is done in (1), it may thus be advisable to solve the regularized learning problem

$$\inf_{\mathbf{w}} \mathbb{E}^{\hat{\mathbb{P}}^N} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] + c \Omega(\mathbf{w}), \quad (2)$$

which minimizes the sum of the empirical average loss and a penalty for hypothesis complexity, which consists of a regularization function $\Omega(\mathbf{w})$ and its associated regularization weight c . Tikhonov regularization (Tikhonov et al., 1977), for example, corresponds to the choice $\Omega(\mathbf{w}) = \|\mathbf{\Gamma}\mathbf{w}\|_2^2$ for some Tikhonov matrix $\mathbf{\Gamma} \in \mathbb{R}^{n \times n}$. Setting $\mathbf{\Gamma}$ to the identity matrix gives rise to standard L_2 -regularization. Similarly, Lasso (least absolute shrinkage and selection operator) regularization or L_1 -regularization (Tibshirani, 1996) is obtained by setting $\Omega(\mathbf{w}) = \|\mathbf{w}\|_1$. Lasso regularization has gained popularity because it favors parsimonious interpretable hypotheses.

Most popular regularization methods admit probabilistic interpretations. However, these interpretations typically rely on prior distributional assumptions that remain to some extent arbitrary (*e.g.*, L_2 - and L_1 -regularization can be justified if \mathbf{w} is governed by a Gaussian or Laplacian prior distribution, respectively (Tibshirani, 1996)). Thus, in spite of their many desirable theoretical properties, there is a consensus that “*most of the (regularization) methods used successfully in practice are heuristic methods*” (Abu-Mostafa et al., 2012).

2.2. A New Perspective on Regularization

When *linear* hypotheses are used, problem (1) minimizes the in-sample error $\mathbb{E}^{\hat{\mathbb{P}}_N}[\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)]$. However, a hypothesis \mathbf{w} enjoying a low in-sample error may still suffer from a high out-of-sample error $\mathbb{E}^{\mathbb{P}}[\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)]$ due to overfitting. This is unfortunate as we seek hypotheses that offer high prediction accuracy on *future* data, meaning that the out-of-sample error is the actual quantity of interest. An ideal learning model would therefore minimize the out-of-sample error. This is impossible, however, for the following reasons:

- The true input-output distribution \mathbb{P} is unknown and only indirectly observable through the N training samples. Thus, we lack essential information to compute the out-of-sample error.
- Even if the distribution \mathbb{P} was known, computing the out-of-sample error would typically be hard due to the intractability of high-dimensional integration; see, *e.g.*, (Hanasusanto et al., 2016, Corollary 1).

The regularized loss $\mathbb{E}^{\hat{\mathbb{P}}_N}[\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] + c\Omega(\mathbf{w})$ used in (2), which consists of the in-sample error and an overfitting penalty, can be viewed as an in-sample estimate of the out-of-sample error. Yet, problem (2) remains difficult to justify rigorously. Therefore, we advocate here a more principled approach to regularization. Specifically, we propose to take into account the expected loss of hypothesis \mathbf{w} under *every* distribution \mathbb{Q} that is close to the empirical distribution $\hat{\mathbb{P}}_N$, that is, every \mathbb{Q} that could have generated the training data with high confidence. To this end, we first introduce a distance measure for distributions. For ease of notation, we henceforth denote the input-output pair (\mathbf{x}, y) by $\boldsymbol{\xi}$, and we set $\Xi = \mathbb{X} \times \mathbb{Y}$.

Definition 2 (Wasserstein metric) *The Wasserstein distance between two distributions \mathbb{Q} and \mathbb{Q}' supported on Ξ is defined as*

$$W(\mathbb{Q}, \mathbb{Q}') := \inf_{\Pi} \left\{ \int_{\Xi^2} d(\boldsymbol{\xi}, \boldsymbol{\xi}') \Pi(d\boldsymbol{\xi}, d\boldsymbol{\xi}') : \begin{array}{l} \Pi \text{ is a joint distribution of } \boldsymbol{\xi} \text{ and } \boldsymbol{\xi}' \\ \text{with marginals } \mathbb{Q} \text{ and } \mathbb{Q}', \text{ respectively} \end{array} \right\},$$

where d is a metric on Ξ .

By definition, $W(\mathbb{Q}, \mathbb{Q}')$ represents the solution of an infinite-dimensional transportation problem, that is, it corresponds to the minimal cost for moving the distribution \mathbb{Q} to \mathbb{Q}' , where the cost for moving a unit probability mass from $\boldsymbol{\xi}$ to $\boldsymbol{\xi}'$ is given by the transportation distance $d(\boldsymbol{\xi}, \boldsymbol{\xi}')$. Due to this interpretation, the metric d is often referred to as the transportation cost (Villani, 2008) or ground metric (Cuturi and Avis, 2014), while the Wasserstein metric is sometimes termed the mass transportation distance or earth mover’s distance (Rubner et al., 2000).

Consider now the Wasserstein ball of radius $\rho \geq 0$ around the empirical distribution $\hat{\mathbb{P}}_N$,

$$\mathbb{B}_\rho(\hat{\mathbb{P}}_N) = \{\mathbb{Q} : \mathbb{Q}(\Xi) = 1, W(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \rho\}, \quad (3)$$

which contains all input-output distributions \mathbb{Q} supported on Ξ whose Wasserstein distance from $\hat{\mathbb{P}}_N$ does not exceed ρ . This means that \mathbb{Q} can be transported to $\hat{\mathbb{P}}_N$ (or vice versa) at a cost of at most ρ . The hope is that a large enough Wasserstein ball will contain distributions that are representative of the unknown true input-output distribution \mathbb{P} , such that the worst-case expectation $\sup_{\mathbb{Q} \in \mathbb{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^\mathbb{Q}[\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)]$ can serve as an upper confidence bound on the out-of-sample error $\mathbb{E}^\mathbb{P}[\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)]$. This motivates us to introduce a new regularized learning model, which minimizes precisely this worst-case expectation.

$$\inf_{\mathbf{w}} \sup_{\mathbb{Q} \in \mathbb{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^\mathbb{Q}[\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] \quad (4)$$

Problem (4) represents a distributionally robust convex program of the type considered in (Mohajerin Esfahani and Kuhn, 2018). Note that if $\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)$ is convex in \mathbf{w} for every fixed (\mathbf{x}, y) , *i.e.*, if ℓ is convex in its first argument, then the objective function of (4) is convex because convexity is preserved under integration and maximization. Note also that if ρ is set to zero, then (4) collapses to the unregularized in-sample error minimization problem (1).

Remark 3 (Support information) *The uncertainty set Ξ captures prior information on the range of the inputs and outputs. In image processing, for example, pixel intensities range over a known interval. Similarly, in diagnostic medicine, physiological parameters such as blood glucose or cholesterol concentrations are restricted to be non-negative. Sometimes it is also useful to construct Ξ as a confidence set that covers the support of \mathbb{P} with a prescribed probability. Such confidence sets are often constructed as ellipsoids, as intersections of different norm balls (Ben-Tal et al., 2009; Delage and Ye, 2010) or as sublevel sets of kernel expansions (Schölkopf et al., 2001).*

In the remainder we establish that the distributionally robust learning problem (4) has several desirable properties. (i) Problem (4) is computationally tractable under standard assumptions about the loss function ℓ , the input-output space Ξ and the transportation metric d . For specific choices of d it even reduces to a regularized learning problem of the form (2). (ii) For all univariate loss functions reviewed in Section 2.1, a tight conservative approximation of (4) is kernelizable, that is, it can be solved implicitly over high-dimensional spaces of nonlinear hypotheses at the same computational cost required for linear hypothesis spaces. (iii) Leveraging modern measure concentration results, the optimal value of (4) can be shown to provide an upper confidence bound on the out-of-sample error. This obviates the need to mobilize the full machinery of VC theory and, in particular, to estimate the VC dimension of the hypothesis space in order to establish generalization bounds. (iv) If the number of training samples tends to infinity while the Wasserstein ball shrinks at an appropriate rate, then problem (4) asymptotically recovers the *ex post* optimal hypothesis that attains the minimal out-of-sample error.

3. Tractable Reformulations

In this section we demonstrate that the distributionally robust learning problem (4) over linear hypotheses is amenable to efficient computational solution procedures. We also discuss generalizations to nonlinear hypothesis classes such as reproducing kernel Hilbert spaces and families of feed-forward neural networks.

3.1. Distributionally Robust Linear Regression

Throughout this section we focus on linear regression problems, where $\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y) = L(\langle \mathbf{w}, \mathbf{x} \rangle - y)$ for some convex univariate loss function L . We also assume that \mathbb{X} and \mathbb{Y} are both convex and closed and that the transportation metric d is induced by a norm $\|\cdot\|$ on the input-output space \mathbb{R}^{n+1} . In this setting, the distributionally robust regression problem (4) admits an equivalent reformulation as a finite convex optimization problem if either (i) the univariate loss function L is piecewise affine or (ii) $\Xi = \mathbb{R}^{n+1}$ and L is Lipschitz continuous (but not necessarily piecewise affine).

Theorem 4 (Distributionally robust linear regression) *The following statements hold.*

(i) *If $L(z) = \max_{j \leq J} \{a_j z + b_j\}$, then (4) is equivalent to*

$$\left\{ \begin{array}{ll} \inf_{\substack{\mathbf{w}, \lambda, s_i \\ \mathbf{p}_{ij}, u_{ij}}} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & S_{\Xi}(a_j \mathbf{w} - \mathbf{p}_{ij}, -a_j - u_{ij}) + b_j + \langle \mathbf{p}_{ij}, \hat{\mathbf{x}}_i \rangle + u_{ij} \hat{y}_i \leq s_i \quad i \in [N], j \in [J] \\ & \|(\mathbf{p}_{ij}, u_{ij})\|_* \leq \lambda \quad i \in [N], j \in [J], \end{array} \right. \quad (5)$$

where S_{Ξ} denotes the support function of Ξ .

(ii) *If $\Xi = \mathbb{R}^{n+1}$ and $L(z)$ is Lipschitz continuous, then (4) is equivalent to*

$$\inf_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N L(\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \hat{y}_i) + \rho \text{lip}(L) \|(\mathbf{w}, -1)\|_*. \quad (6)$$

In the following, we exemplify Theorem 4 for the Huber, ϵ -insensitive and pinball loss functions under the assumption that the uncertainty set Ξ admits the conic representation

$$\Xi = \{(\mathbf{x}, y) \in \mathbb{R}^{n+1} : \mathbf{C}_1 \mathbf{x} + \mathbf{c}_2 y \preceq_{\mathcal{C}} \mathbf{d}\} \quad (7)$$

for some matrix \mathbf{C}_1 , vectors \mathbf{c}_2 and \mathbf{d} and proper convex cone \mathcal{C} of appropriate dimensions. We also assume that Ξ admits a Slater point $(\mathbf{x}_S, y_S) \in \mathbb{R}^{n+1}$ with $\mathbf{C}_1 \mathbf{x}_S + \mathbf{c}_2 y_S \prec_{\mathcal{C}} \mathbf{d}$.

Corollary 5 (Robust regression) *If L represents the Huber loss function with threshold $\delta \geq 0$ and $\Xi = \mathbb{R}^{n+1}$, then (4) is equivalent to*

$$\inf_{\mathbf{w}, z_i} \frac{1}{N} \sum_{i=1}^N \frac{1}{2} z_i^2 + \delta |\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \hat{y}_i - z_i| + \rho \delta \|(\mathbf{w}, -1)\|_*. \quad (8)$$

Corollary 6 (Support vector regression) *If L represents the ϵ -insensitive loss function for some $\epsilon \geq 0$ and Ξ is of the form (7), then (4) is equivalent to*

$$\left\{ \begin{array}{ll} \inf_{\substack{\mathbf{w}, \lambda, s_i \\ \mathbf{p}_i^+, \mathbf{p}_i^-}} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & \hat{y}_i - \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \epsilon + \langle \mathbf{p}_i^+, \mathbf{d} - \mathbf{C}_1 \hat{\mathbf{x}}_i - \mathbf{c}_2 \hat{y}_i \rangle \leq s_i & i \in [N] \\ & \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \hat{y}_i - \epsilon + \langle \mathbf{p}_i^-, \mathbf{d} - \mathbf{C}_1 \hat{\mathbf{x}}_i - \mathbf{c}_2 \hat{y}_i \rangle \leq s_i & i \in [N] \\ & \|(\mathbf{C}_1^\top \mathbf{p}_i^+ + \mathbf{w}, \mathbf{c}_2^\top \mathbf{p}_i^+ - 1)\|_* \leq \lambda & i \in [N] \\ & \|(\mathbf{C}_1^\top \mathbf{p}_i^- - \mathbf{w}, \mathbf{c}_2^\top \mathbf{p}_i^- + 1)\|_* \leq \lambda & i \in [N] \\ & \mathbf{p}_i^+, \mathbf{p}_i^- \in \mathcal{C}^* & i \in [N] \\ & s_i \geq 0 & i \in [N]. \end{array} \right. \quad (9)$$

Corollary 7 (Quantile regression) *If L represents the pinball loss function for some $\tau \in [0, 1]$ and Ξ is of the form (7), then (4) is equivalent to*

$$\left\{ \begin{array}{ll} \inf_{\substack{\mathbf{w}, \lambda, s_i \\ \mathbf{p}_i^+, \mathbf{p}_i^-}} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & \tau (\hat{y}_i - \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle) + \langle \mathbf{p}_i^+, \mathbf{d} - \mathbf{C}_1 \hat{\mathbf{x}}_i - \mathbf{c}_2 \hat{y}_i \rangle \leq s_i & i \in [N] \\ & (1 - \tau) (\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \hat{y}_i) + \langle \mathbf{p}_i^-, \mathbf{d} - \mathbf{C}_1 \hat{\mathbf{x}}_i - \mathbf{c}_2 \hat{y}_i \rangle \leq s_i & i \in [N] \\ & \|(\mathbf{C}_1^\top \mathbf{p}_i^+ + \tau \mathbf{w}, \mathbf{c}_2^\top \mathbf{p}_i^+ - \tau)\|_* \leq \lambda & i \in [N] \\ & \|(\mathbf{C}_1^\top \mathbf{p}_i^- - (1 - \tau) \mathbf{w}, \mathbf{c}_2^\top \mathbf{p}_i^- + 1 - \tau)\|_* \leq \lambda & i \in [N] \\ & \mathbf{p}_i^+, \mathbf{p}_i^- \in \mathcal{C}^* & i \in [N] \\ & s_i \geq 0 & i \in [N]. \end{array} \right. \quad (10)$$

Remark 8 (Relation to classical regularization) *Assume now that the mass transportation costs are additively separable with respect to inputs and outputs, that is,*

$$d((\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2)) = \|\mathbf{x}_1 - \mathbf{x}_2\| + \kappa |y_1 - y_2| \quad (11)$$

for some $\kappa > 0$.¹ Note that κ captures the costs of moving probability mass along the output space. For $\kappa = \infty$ all distributions in the Wasserstein ball $\mathbb{B}_\rho(\hat{\mathbb{P}}_N)$ are thus obtained by reshaping $\hat{\mathbb{P}}_N$ only along the input space. It is easy to verify that for $\kappa = \infty$ and $\Xi = \mathbb{R}^{n+1}$ the learning models portrayed in Corollaries 5-7 all simplify to

$$\inf_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N L(\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \hat{y}_i) + c \|\mathbf{w}\|_*, \quad (12)$$

where $c = \rho \delta$ for robust regression with Huber loss, $c = \rho$ for support vector regression with ϵ -insensitive loss and $c = \max\{\tau, 1 - \tau\} \rho$ for quantile regression with pinball loss. Thus, (12) is easily identified as an instance of the classical regularized learning problem (2), where the dual norm term $\|\mathbf{w}\|_*$ plays the role of the regularization function, while c represents

1. By slight abuse of notation, the symbol $\|\cdot\|$ now denotes a norm on \mathbb{R}^n .

the usual regularization weight. By definition of the dual norm, the penalty $\|\mathbf{w}\|_*$ assigned to a hypothesis \mathbf{w} is maximal (minimal) if the cost of moving probability mass along \mathbf{w} is minimal (maximal). We emphasize that if $\kappa = \infty$, then the marginal distribution of y corresponding to every $\mathbf{Q} \in \mathbb{B}_\rho(\widehat{\mathbb{P}}_N)$ coincides with the empirical distribution $\frac{1}{N} \sum_{i=1}^N \delta_{\widehat{y}_i}$. Thus, classical regularization methods, which correspond to $\kappa = \infty$, are explained by a counterintuitive probabilistic model, which pretends that any training sample must have an output that has already been recorded in the training dataset. In other words, classical regularization implicitly assumes that there is no uncertainty in the outputs. More intuitively appealing regularization schemes are obtained for finite values of κ .

To establish a connection between distributionally robust and classical robust regression as discussed in (El Ghaoui and Lebre, 1997; Xu et al., 2010), we further investigate the worst-case expected loss of a fixed linear hypothesis \mathbf{w} .

$$\sup_{\mathbf{Q} \in \mathbb{B}_\rho(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}}[L(\langle \mathbf{w}, \mathbf{x} \rangle - y)] \quad (13)$$

Theorem 9 (Extremal distributions in linear regression) *The following statements hold.*

(i) *If $L(z) = \max_{j \leq J} \{a_j z + b_j\}$, then the worst-case expectation (13) coincides with*

$$\left\{ \begin{array}{ll} \sup_{\alpha_{ij}, \mathbf{q}_{ij}, v_{ij}} & \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij} (a_j (\langle \mathbf{w}, \widehat{\mathbf{x}}_i \rangle - \widehat{y}_i) + b_j) + a_j (\langle \mathbf{w}, \mathbf{q}_{ij} \rangle - v_{ij}) \\ \text{s.t.} & \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \|\mathbf{q}_{ij}, v_{ij}\| \leq \rho \\ & \sum_{j=1}^J \alpha_{ij} = 1 \quad i \in [N] \\ & (\widehat{\mathbf{x}}_i + \mathbf{q}_{ij}/\alpha_{ij}, \widehat{y}_i + v_{ij}/\alpha_{ij}) \in \Xi \quad i \in [N], j \in [J] \\ & \alpha_{ij} \geq 0 \quad i \in [N], j \in [J] \end{array} \right. \quad (14)$$

for any fixed hypothesis \mathbf{w} . Moreover, if $(\alpha_{ij}^*, \mathbf{q}_{ij}^*, v_{ij}^*)$ maximizes (14), then the discrete distribution

$$\mathbf{Q}^* = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij}^* \delta_{(\widehat{\mathbf{x}}_i + \mathbf{q}_{ij}^*/\alpha_{ij}^*, \widehat{y}_i + v_{ij}^*/\alpha_{ij}^*)},$$

represents a maximizer for (13).

(ii) *If $\Xi = \mathbb{R}^{n+1}$ and $L(z)$ is Lipschitz continuous, then the discrete distributions*

$$\mathbf{Q}_\gamma = \frac{1}{N} \sum_{i=2}^N \delta_{(\widehat{\mathbf{x}}_i, \widehat{y}_i)} + \frac{1-\gamma}{N} \delta_{(\widehat{\mathbf{x}}_1, \widehat{y}_1)} + \frac{\gamma}{N} \delta_{(\widehat{\mathbf{x}}_1 + \frac{\rho N}{\gamma} \mathbf{x}^*, \widehat{y}_1 + \frac{\rho N}{\gamma} y^*)} \quad \text{for } \gamma \in (0, 1],$$

where (\mathbf{x}^*, y^*) solves $\max_{\|(\mathbf{x}, y)\| \leq 1} \langle \mathbf{w}, \mathbf{x} \rangle - y$, are feasible and asymptotically optimal in (13) for $\gamma \downarrow 0$.

Recall that $0/0 = 0$ and $1/0 = \infty$ by our conventions of extended arithmetic. Thus, any solution feasible in (14) with $\alpha_{ij} = 0$ must satisfy $\mathbf{q}_{ij} = \mathbf{0}$ and $v_{ij} = 0$ because otherwise $(\hat{\mathbf{x}}_i + \mathbf{q}_{ij}/\alpha_{ij}, \hat{y}_i + v_{ij}/\alpha_{ij}) \notin \Xi$.

Theorem 9 shows how one can use convex optimization to construct a sequence of distributions that are asymptotically optimal in (13). Next, we argue that the worst-case expected cost (13) is equivalent to a (robust) worst-case cost over a suitably defined uncertainty set if the following assumption holds.

Assumption 10 (Minimal dispersion) *For every $\mathbf{w} \in \mathbb{R}^n$ there is a training sample $(\hat{\mathbf{x}}_k, \hat{y}_k)$ for some $k \leq N$ such that the derivative L' exists at $\langle \mathbf{w}, \hat{\mathbf{x}}_k \rangle - \hat{y}_k$ and satisfies $|L'(\langle \mathbf{w}, \hat{\mathbf{x}}_k \rangle - \hat{y}_k)| = \text{lip}(L)$.*

Remark 11 (Minimal dispersion) *Assumption 10 is reminiscent of the non-separability condition in (Xu et al., 2009, Theorem 3), which is necessary to prove the equivalence of robust and regularized support vector machines. In the regression context studied here, Assumption 10 ensures that, for every \mathbf{w} , there exists a training sample that activates the largest absolute slope of L .*

For instance, in support vector regression, it means that for every \mathbf{w} there exists a data point outside of the slab of width $2\epsilon/\|(\mathbf{w}, -1)\|_2$ centered around the hyperplane $H_{\mathbf{w}} = \{(\mathbf{x}, y) \in \mathbb{R}^n \times \mathbb{R} : \langle \mathbf{w}, \mathbf{x} \rangle - y = 0\}$ (i.e., the empirical ϵ -insensitive loss is not zero). Similarly, in robust regression with the Huber loss function, Assumption 10 stipulates that for every \mathbf{w} there exists a data point outside of the slab of width $2\delta/\|(\mathbf{w}, -1)\|_2$ centered around $H_{\mathbf{w}}$. However, quantile regression with $\tau \neq 0.5$ fails to satisfy Assumption 10. Indeed, for any training dataset there always exists some \mathbf{w} such that all data points reside on the side of $H_{\mathbf{w}}$ where the pinball loss function is less steep.

Theorem 12 (Robust regression) *If $\Xi = \mathbb{R}^{n+1}$ and the loss function $L(z)$ is Lipschitz continuous, then the worst-case expected loss (13) provides an upper bound on the (robust) worst-case loss*

$$\left\{ \begin{array}{ll} \sup_{\Delta \mathbf{x}_i, \Delta y_i} & \frac{1}{N} \sum_{i=1}^N [L(\langle \mathbf{w}, \hat{\mathbf{x}}_i + \Delta \mathbf{x}_i \rangle - \hat{y}_i - \Delta y_i)] \\ \text{s.t.} & \frac{1}{N} \sum_{i=1}^N \|(\Delta \mathbf{x}_i, \Delta y_i)\| \leq \rho. \end{array} \right. \quad (15)$$

Moreover, if Assumption 10 holds, then (13) and (15) are equal.

Remark 13 (Tractability of robust regression) *Assume that $\Xi = \mathbb{R}^{n+1}$, while L and $\|\cdot\|$ both admit a tractable conic representation. By Theorem 4, the worst-case expected loss (13) can then be computed in polynomial time by solving a tractable convex program. Theorem 12 thus implies that the worst-case loss (28) can also be computed in polynomial time if Assumption 10 holds. To our best knowledge, there exists no generic efficient method for computing (28) if Assumption 10 fails to hold and L is not piecewise affine. This reinforces our belief that a distributionally robust approach to regression is more natural.*

3.2. Distributionally Robust Linear Classification

Throughout this section we focus on linear classification problems, where $\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y) = L(y\langle \mathbf{w}, \mathbf{x} \rangle)$ for some convex univariate loss function L . We also assume that \mathbb{X} is both convex and closed and that $\mathbb{Y} = \{+1, -1\}$. Moreover, we assume that the transportation metric d is defined via

$$d((\mathbf{x}, y), (\mathbf{x}', y')) = \|\mathbf{x} - \mathbf{x}'\| + \kappa \mathbb{1}_{\{y \neq y'\}}, \quad (16)$$

where $\|\cdot\|$ represents a norm on the input space \mathbb{R}^n , and $\kappa > 0$ quantifies the cost of switching a label. In this setting, the distributionally robust classification problem (4) admits an equivalent reformulation as a finite convex optimization problem if either (i) the univariate loss function L is piecewise affine or (ii) $\mathbb{X} = \mathbb{R}^n$ and L is Lipschitz continuous (but not necessarily piecewise affine).

Theorem 14 (Distributionally robust linear classification) *The following statements hold.*

(i) *If $L(z) = \max_{j \in J} \{a_j z + b_j\}$, then (4) is equivalent to*

$$\left\{ \begin{array}{ll} \inf_{\mathbf{w}, \lambda, s_i} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & S_{\mathbb{X}}(a_j \hat{y}_i \mathbf{w} - \mathbf{p}_{ij}^+) + b_j + \langle \mathbf{p}_{ij}^+, \hat{\mathbf{x}}_i \rangle \leq s_i & i \in [N], j \in [J] \\ & S_{\mathbb{X}}(-a_j \hat{y}_i \mathbf{w} - \mathbf{p}_{ij}^-) + b_j + \langle \mathbf{p}_{ij}^-, \hat{\mathbf{x}}_i \rangle - \kappa \lambda \leq s_i & i \in [N], j \in [J] \\ & \|\mathbf{p}_{ij}^+\|_* \leq \lambda, \|\mathbf{p}_{ij}^-\|_* \leq \lambda & i \in [N], j \in [J], \end{array} \right. \quad (17)$$

where $S_{\mathbb{X}}$ denotes the support function of \mathbb{X} .

(ii) *If $\mathbb{X} = \mathbb{R}^n$ and L is Lipschitz continuous, then (4) is equivalent to*

$$\left\{ \begin{array}{ll} \inf_{\mathbf{w}, \lambda, s_i} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & L(\hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle) \leq s_i & i \in [N] \\ & L(-\hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle) - \kappa \lambda \leq s_i & i \in [N] \\ & \text{lip}(L) \|\mathbf{w}\|_* \leq \lambda. \end{array} \right. \quad (18)$$

In the following, we exemplify Theorem 14 for the hinge loss, logloss and smoothed hinge loss functions under the assumption that the input space \mathbb{X} admits the conic representation

$$\mathbb{X} = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{C}\mathbf{x} \preceq_{\mathcal{C}} \mathbf{d}\} \quad (19)$$

for some matrix \mathbf{C} , vector \mathbf{d} and proper convex cone \mathcal{C} of appropriate dimensions. We also assume that \mathbb{X} admits a Slater point $\mathbf{x}_S \in \mathbb{R}^n$ with $\mathbf{C}\mathbf{x}_S \prec_{\mathcal{C}} \mathbf{d}$.

Corollary 15 (Support vector machine) *If L represents the hinge loss function and \mathbb{X} is of the form (19), then (4) is equivalent to*

$$\left\{ \begin{array}{ll} \inf_{\substack{\mathbf{w}, \lambda \\ s_i, \mathbf{p}_i^+, \mathbf{p}_i^-}} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & 1 - \hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle + \langle \mathbf{p}_i^+, \mathbf{d} - \mathbf{C} \hat{\mathbf{x}}_i \rangle \leq s_i \quad i \in [N] \\ & 1 + \hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle + \langle \mathbf{p}_i^-, \mathbf{d} - \mathbf{C} \hat{\mathbf{x}}_i \rangle - \kappa \lambda \leq s_i \quad i \in [N] \\ & \|\mathbf{C}^\top \mathbf{p}_i^+ + \hat{y}_i \mathbf{w}\|_* \leq \lambda, \quad \|\mathbf{C}^\top \mathbf{p}_i^- - \hat{y}_i \mathbf{w}\|_* \leq \lambda \quad i \in [N] \\ & s_i \geq 0, \quad \mathbf{p}_i^+, \mathbf{p}_i^- \in \mathcal{C}^* \quad i \in [N]. \end{array} \right. \quad (20)$$

Corollary 16 (Support vector machine with smooth hinge loss) *If L represents the smooth hinge loss function and $\mathbb{X} = \mathbb{R}^n$, then (4) is equivalent to*

$$\left\{ \begin{array}{ll} \min_{\substack{\mathbf{w}, \lambda, s_i \\ z_i^+, z_i^-, t_i^+, t_i^-}} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & \frac{1}{2} (z_i^+ - \hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle)^2 + t_i^+ \leq s_i \quad i \in [N] \\ & \frac{1}{2} (z_i^- + \hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle)^2 + t_i^- - \kappa \lambda \leq s_i \quad i \in [N] \\ & 1 - z_i^+ \leq t_i^+, \quad 1 - z_i^- \leq t_i^- \quad i \in [N] \\ & t_i^+, t_i^- \geq 0 \quad i \in [N] \\ & \|\mathbf{w}\|_* \leq \lambda. \end{array} \right. \quad (21)$$

Corollary 17 (Logistic regression) *If L represents the logloss function and $\mathbb{X} = \mathbb{R}^n$, then (4) is equivalent to*

$$\left\{ \begin{array}{ll} \min_{\mathbf{w}, \lambda, s_i} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & \log \left(1 + \exp \left(-\hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle \right) \right) \leq s_i \quad i \in [N] \\ & \log \left(1 + \exp \left(\hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle \right) \right) - \kappa \lambda \leq s_i \quad i \in [N] \\ & \|\mathbf{w}\|_* \leq \lambda. \end{array} \right. \quad (22)$$

Remark 18 (Relation to classical regularization) *If $\mathbb{X} = \mathbb{R}^n$ and the weight parameter κ in the transportation metric (16) is set to infinity, then the learning problems portrayed in Corollaries 15–17 all simplify to*

$$\inf_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N L(\hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle) + \rho \|\mathbf{w}\|_*. \quad (23)$$

Thus, in analogy to the case of regression, (23) reduces to an instance of the classical regularized learning problem (2), where the dual norm term $\|\mathbf{w}\|_*$ plays the role of the regularization function, while the Wasserstein radius ρ represents the usual regularization weight. Note that if $\kappa = \infty$, then mass transportation along the output space is infinitely expensive, that is, any distribution $\mathbf{Q} \in \mathbb{B}_\rho(\hat{\mathbb{P}}_N)$ can smear out the training samples along

\mathbb{X} , but it cannot flip outputs from $+1$ to -1 or vice versa. Thus, classical regularization schemes, which are recovered for $\kappa = \infty$, implicitly assume that output measurements are exact. As this belief is not tenable in most applications, an approach with $\kappa < \infty$ may be more satisfying. We remark that alternative approaches for learning with noisy labels have previously been studied by Lawrence and Schölkopf (2001), Natarajan et al. (2013), and Yang et al. (2012).

Remark 19 (Relation to Tikhonov regularization) *The learning problem*

$$\inf_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N L(\hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle) + c \|\mathbf{w}\|_2^2 \quad (24)$$

with Tikhonov regularizer enjoys wide popularity. If L represents the hinge loss, for example, then (24) reduces to the celebrated soft margin support vector machine problem. However, the Tikhonov regularizer appearing in (24) is not explained by a distributionally robust learning problem of the form (4). It is known, however, that (23) with $\|\cdot\|_* = \|\cdot\|_2$ and (24) are equivalent in the sense that for every $\rho \geq 0$ there exists $c \geq 0$ such that the solution of (23) also solves (24) and vice versa (Xu et al., 2009, Corollary 6).

To establish a connection between distributionally robust and classical robust classification as discussed in (Xu et al., 2009), we further investigate the worst-case expected loss of a fixed linear hypothesis \mathbf{w} .

$$\sup_{\mathbf{Q} \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}}[L(y \langle \mathbf{w}, \mathbf{x} \rangle)] \quad (25)$$

Theorem 20 (Extremal distributions in linear classification) *The following statements hold.*

(i) *If $L(z) = \max_{j \in J} \{a_j z + b_j\}$, then the worst-case expectation (25) coincides with*

$$\left\{ \begin{array}{l} \sup_{\substack{\alpha_{ij}^+, \alpha_{ij}^- \\ \mathbf{q}_{ij}^+, \mathbf{q}_{ij}^-}} \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J (\alpha_{ij}^+ - \alpha_{ij}^-) a_j \hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle + a_j \hat{y}_i \langle \mathbf{w}, \mathbf{q}_{ij}^+ - \mathbf{q}_{ij}^- \rangle + \sum_{j=1}^J b_j \\ \text{s.t.} \quad \sum_{i=1}^N \sum_{j=1}^J \|\mathbf{q}_{ij}^+\| + \|\mathbf{q}_{ij}^-\| + \kappa \alpha_{ij}^- \leq N\rho \\ \sum_{j=1}^J \alpha_{ij}^+ + \alpha_{ij}^- = 1 \quad i \in [N] \\ \hat{\mathbf{x}}_i + \mathbf{q}_{ij}^+ / \alpha_{ij}^+ \in \mathbb{X}, \quad \hat{\mathbf{x}}_i + \mathbf{q}_{ij}^- / \alpha_{ij}^- \in \mathbb{X} \quad i \in [N], j \in [J] \\ \alpha_{ij}^+, \alpha_{ij}^- \geq 0 \quad i \in [N], j \in [J] \end{array} \right. \quad (26)$$

for any fixed \mathbf{w} . Also, if $(\alpha_{ij}^{+*}, \alpha_{ij}^{-*}, \mathbf{q}_{ij}^{+*}, \mathbf{q}_{ij}^{-*})$ maximizes (26), then the discrete distribution

$$\mathbf{Q} = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij}^{+*} \delta_{(\hat{\mathbf{x}}_i - \mathbf{q}_{ij}^{+*} / \alpha_{ij}^{+*}, \hat{y}_i)} + \alpha_{ij}^{-*} \delta_{(\hat{\mathbf{x}}_i - \mathbf{q}_{ij}^{-*} / \alpha_{ij}^{-*}, -\hat{y}_i)}$$

represents a maximizer for (25).

(ii) If $\mathbb{X} = \mathbb{R}^n$, then the worst-case expectation (25) coincides with the optimal value of

$$\left\{ \begin{array}{ll} \sup_{\alpha_i, \theta} & \text{lip}(L) \|\mathbf{w}\|_* \theta + \frac{1}{N} \sum_{i=1}^N (1 - \alpha_i) L(\hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle) + \alpha_i L(-\hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle) \\ \text{s.t.} & \theta + \frac{\kappa}{N} \sum_{i=1}^N \alpha_i = \rho - \gamma \\ & 0 \leq \alpha_i \leq 1 \\ & \theta \geq 0 \end{array} \right. \quad i \in [N] \quad (27)$$

for $\gamma = 0$. Moreover, if $(\alpha_i^*(\gamma), \theta^*(\gamma))$ maximizes (27) for $\gamma > 0$, $\eta(\gamma) = \gamma / (\theta^*(\gamma) + \kappa - \rho + \gamma + 1)$ and \mathbf{x}^* solves $\max_{\mathbf{x}} \{\langle \mathbf{w}, \mathbf{x} \rangle : \|\mathbf{x}\| \leq 1\}$, then the discrete distributions

$$\begin{aligned} \mathbb{Q}_\gamma = & \frac{1}{N} \sum_{i=2}^N (1 - \alpha_i^*(\gamma)) \delta_{(\hat{\mathbf{x}}_i, \hat{y}_i)} + \alpha_i^*(\gamma) \delta_{(\hat{\mathbf{x}}_i, -\hat{y}_i)} + \frac{\eta(\gamma)}{N} \delta_{(\hat{\mathbf{x}}_1 + \frac{\theta^*(\gamma)N}{\eta(\gamma)} \mathbf{x}^*, \hat{y}_1)} \\ & + \frac{1 - \eta(\gamma)}{N} \left[(1 - \alpha_1^*(\gamma)) \delta_{(\hat{\mathbf{x}}_1, \hat{y}_1)} + \alpha_1^*(\gamma) \delta_{(\hat{\mathbf{x}}_1, -\hat{y}_1)} \right] \end{aligned}$$

for $\gamma \in [0, \min\{\rho, 1\}]$ are feasible and asymptotically optimal in (25) for $\gamma \downarrow 0$.

Theorem 20 shows how one can use convex optimization to construct a sequence of distributions that are asymptotically optimal in (25). Next, we show that the worst-case expected cost (25) is equivalent to a (robust) worst-case cost over a suitably defined uncertainty set if the following assumption holds.

Assumption 21 (Non-separability) For every $\mathbf{w} \in \mathbb{R}^n$ there is a training sample $(\hat{\mathbf{x}}_k, \hat{y}_k)$ for some $k \leq N$ such that the derivative L' exists at $\hat{y}_k \langle \mathbf{w}, \hat{\mathbf{x}}_k \rangle$ and satisfies $|L'(\hat{y}_k \langle \mathbf{w}, \hat{\mathbf{x}}_k \rangle)| = \text{lip}(L)$.

Remark 22 (Non-separability) Assumption 21 generalizes the non-separability condition in (Xu et al., 2009, Theorem 3) for the classical and smooth hinge loss functions to more general Lipschitz continuous losses. Note that, in the case of the hinge loss, Assumption 21 effectively stipulates that for any \mathbf{w} there exists a training sample $(\hat{\mathbf{x}}_k, \hat{y}_k)$ with $\hat{y}_k \langle \mathbf{w}, \hat{\mathbf{x}}_k \rangle < 1$, implying that the dataset cannot be perfectly separated by any linear hypothesis \mathbf{w} . An equivalent requirement is that the empirical hinge loss is nonzero for every \mathbf{w} . Similarly, in the case of the smooth hinge loss, Assumption 21 ensures that for any \mathbf{w} there is a training sample with $\hat{y}_k \langle \mathbf{w}, \hat{\mathbf{x}}_k \rangle < 0$, which implies again that the dataset admits no perfect linear separation. Note, however, that the logloss fails to satisfy Assumption 21 as its steepest slope is attained at infinity.

Theorem 23 (Robust classification) Suppose that $\mathbb{X} = \mathbb{R}^n$, the loss function L is Lipschitz continuous and the cost of flipping a label in the transportation metric (16) is set to $\kappa = \infty$. Then, the worst-case expected loss (25) provides an upper bound on the (robust) worst-case loss

$$\left\{ \begin{array}{ll} \sup_{\Delta \mathbf{x}_i} & \frac{1}{N} \sum_{i=1}^N L(\hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i + \Delta \mathbf{x}_i \rangle) \\ \text{s.t.} & \frac{1}{N} \sum_{i=1}^N \|\Delta \mathbf{x}_i\| \leq \rho. \end{array} \right. \quad (28)$$

Moreover, if Assumption 21 holds, then (25) and (28) are equal.

Remark 24 (Tractability of robust classification) Assume that $\mathbb{X} = \mathbb{R}^n$, while L and $\|\cdot\|$ both admit a tractable conic representation. By Theorem 14, the worst-case expected loss (25) can then be computed in polynomial time by solving a tractable convex program. Theorem 23 thus implies that the worst-case loss (28) can also be computed in polynomial time if Assumption 21 holds. This confirms Proposition 4 in (Xu et al., 2009). No efficient method for computing (28) is known if Assumption 21 fails to hold.

3.3. Nonlinear Hypotheses: Reproducing Kernel Hilbert Spaces

We now generalize the learning models from Sections 3.1 and 3.2 to nonlinear hypotheses that range over a *reproducing kernel Hilbert space* (RKHS) $\mathbb{H} \subseteq \mathbb{R}^{\mathbb{X}}$ with inner product $\langle \cdot, \cdot \rangle_{\mathbb{H}}$. By definition, \mathbb{H} thus constitutes a complete metric space with respect to the norm $\|\cdot\|_{\mathbb{H}}$ induced by the inner product, and the point evaluation $h \mapsto h(\mathbf{x})$ of the functions $h \in \mathbb{H}$ represents a continuous linear functional on \mathbb{H} for any fixed $\mathbf{x} \in \mathbb{X}$. The Riesz representation theorem then implies that for every $\mathbf{x} \in \mathbb{X}$ there exists a unique function $\Phi(\mathbf{x}) \in \mathbb{H}$ such that $h(\mathbf{x}) = \langle h, \Phi(\mathbf{x}) \rangle_{\mathbb{H}}$ for all $h \in \mathbb{H}$. We henceforth refer to $\Phi : \mathbb{X} \rightarrow \mathbb{H}$ as the *feature map* and to $k : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}_+$ with $k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle_{\mathbb{H}}$ as the *kernel function*. By construction, the kernel function is symmetric and positive definite, that is, the *kernel matrix* $\mathbf{K} \in \mathbb{R}^{N \times N}$ defined through $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ is positive definite for all $N \in \mathbb{N}$ and $\{\mathbf{x}_i\}_{i \leq N} \subseteq \mathbb{X}$.

By the Moore-Aronszajn theorem, any symmetric and positive definite kernel function k on \mathbb{X} induces a unique RKHS $\mathbb{H} \subseteq \mathbb{R}^{\mathbb{X}}$, which can be represented as

$$\mathbb{H} = \left\{ h \in \mathbb{R}^{\mathbb{X}} : \exists \beta_i \in \mathbb{R}, \mathbf{x}_i \in \mathbb{X} \forall i \in \mathbb{N} \text{ with } h(\mathbf{x}) = \sum_{i=1}^{\infty} \beta_i k(\mathbf{x}_i, \mathbf{x}) \text{ and } \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} \beta_i k(\mathbf{x}_i, \mathbf{x}_j) \beta_j < \infty \right\},$$

where the inner product of two arbitrary functions $h_1, h_2 \in \mathbb{H}$ with $h_1(\mathbf{x}) = \sum_{i=1}^{\infty} \beta_i k(\mathbf{x}_i, \mathbf{x})$ and $h_2(\mathbf{x}) = \sum_{j=1}^{\infty} \beta'_j k(\mathbf{x}'_j, \mathbf{x})$ is defined as $\langle h_1, h_2 \rangle_{\mathbb{H}} = \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} \beta_i k(\mathbf{x}_i, \mathbf{x}'_j) \beta'_j$. One may now use the kernel function to define the feature map Φ through $[\Phi(\mathbf{x}')](\mathbf{x}) = k(\mathbf{x}', \mathbf{x})$ for all $\mathbf{x}, \mathbf{x}' \in \mathbb{X}$. This choice is admissible because it respects the consistency condition $\langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle_{\mathbb{H}} = k(\mathbf{x}, \mathbf{x}')$ for all $\mathbf{x}, \mathbf{x}' \in \mathbb{H}$, and because it implies the desired *reproducing property* $\langle f, \Phi(\mathbf{x}') \rangle_{\mathbb{H}} = \sum_{i=1}^{\infty} \beta_i k(\mathbf{x}_i, \mathbf{x}') = f(\mathbf{x}')$ for all $f \in \mathbb{H}$ and $\mathbf{x}' \in \mathbb{X}$.

In summary, given a symmetric and positive definite kernel function k , there exists an associated RKHS \mathbb{H} and a feature map Φ with the reproducing property. As we will see below, however, to optimize over nonlinear hypotheses in \mathbb{H} , knowledge of k is sufficient, and there is no need to construct \mathbb{H} and Φ explicitly.

Assume now that we are given any symmetric and positive definite kernel function k , and construct a distributionally robust learning problem over all *nonlinear* hypotheses in the corresponding RKHS \mathbb{H} via

$$\hat{J}(\rho) = \inf_{h \in \mathbb{H}} \sup_{Q \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^Q[\ell(h(\mathbf{x}), y)], \quad (29)$$

where the transportation metric is given by the Euclidean norm on $\mathbb{X} \times \mathbb{Y}$ (for regression problems) or the separable metric (16) with the Euclidean norm on \mathbb{X} (for classification problems). While problem (29) is hard to solve in general due to the nonlinearity of the hypotheses $h \in \mathbb{H}$, it is easy to solve a lifted learning problem where the inputs $\mathbf{x} \in \mathbb{X}$ are replaced with features $\mathbf{x}_{\mathbb{H}} \in \mathbb{H}$, while each *nonlinear* hypothesis $h \in \mathbb{H}$ over the input space \mathbb{X} is identified with a *linear* hypothesis $h_{\mathbb{H}} \in \mathbb{H}$ over the feature space \mathbb{H} through the identity $h_{\mathbb{H}}(\mathbf{x}_{\mathbb{H}}) = \langle h, \mathbf{x}_{\mathbb{H}} \rangle_{\mathbb{H}}$. Thus, the lifted learning problem can be represented as

$$\widehat{\mathcal{J}}_{\mathbb{H}}(\rho) = \inf_{h \in \mathbb{H}} \sup_{\mathbf{Q} \in \mathcal{B}_{\rho}(\widehat{\mathbb{P}}_{\mathbb{H}}^{\mathbb{H}})} \mathbb{E}^{\mathbf{Q}} [\ell(\langle h, \mathbf{x}_{\mathbb{H}} \rangle_{\mathbb{H}}, y)], \quad (30)$$

where $\widehat{\mathbb{P}}_{\mathbb{H}}^{\mathbb{H}} = 1/N \sum_{i=1}^N \delta_{(\Phi(\widehat{\mathbf{x}}_i), \widehat{y}_i)}$ on $\mathbb{H} \times \mathbb{Y}$ denotes the pushforward measure of the empirical distribution $\widehat{\mathbb{P}}_N$ under the feature map Φ induced by k , while $\mathcal{B}_{\rho}(\widehat{\mathbb{P}}_{\mathbb{H}}^{\mathbb{H}})$ constitutes the Wasserstein ball of radius ρ around $\widehat{\mathbb{P}}_{\mathbb{H}}^{\mathbb{H}}$ corresponding to the transportation metric

$$d_{\mathbb{H}}((\mathbf{x}_{\mathbb{H}}, y), (\mathbf{x}'_{\mathbb{H}}, y')) = \begin{cases} \sqrt{\|\mathbf{x}_{\mathbb{H}} - \mathbf{x}'_{\mathbb{H}}\|_{\mathbb{H}}^2 + (y - y')^2} & \text{for regression problems,} \\ \|\mathbf{x}_{\mathbb{H}} - \mathbf{x}'_{\mathbb{H}}\|_{\mathbb{H}} + \kappa \mathbf{1}_{\{y \neq y'\}} & \text{for classification problems.} \end{cases}$$

Even though $\widehat{\mathbb{P}}_{\mathbb{H}}^{\mathbb{H}}$ constitutes the pushforward measure of $\widehat{\mathbb{P}}_N$ under Φ , not every distribution $\mathbf{Q}^{\mathbb{H}} \in \mathcal{B}_{\rho}(\widehat{\mathbb{P}}_{\mathbb{H}}^{\mathbb{H}})$ can be obtained as the pushforward measure of some $\mathbf{Q} \in \mathcal{B}_{\rho}(\widehat{\mathbb{P}}_N)$. Thus, we should not expect (29) to be equivalent to (30). Instead, one can show that under a judicious transformation of the Wasserstein radius, (30) provides an upper bound on (29) whenever the kernel function satisfies a calmness condition.

Assumption 25 (Calmness of the kernel) *The kernel function k is calm from above, that is, there exist a concave smooth growth function $g : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ with $g(0) = 0$ and $g'(z) \geq 1$ for all $z \in \mathbb{R}_+$ such that*

$$\sqrt{k(\mathbf{x}_1, \mathbf{x}_1) - 2k(\mathbf{x}_1, \mathbf{x}_2) + k(\mathbf{x}_2, \mathbf{x}_2)} \leq g(\|\mathbf{x}_1 - \mathbf{x}_2\|_2) \quad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{X}.$$

The calmness condition is non-restrictive. In fact, it is satisfied by most commonly used kernels.

Example 1 (Growth Functions for Popular Kernels) *For most commonly used kernels k on $\mathbb{X} \subseteq \mathbb{R}^n$, we can construct an explicit growth function g that certifies the calmness of k in the sense of Assumption 25. This construction typically relies on elementary estimates. Derivations are omitted for brevity.*

1. **Linear kernel:** For $k(\mathbf{x}_1, \mathbf{x}_2) = \langle \mathbf{x}_1, \mathbf{x}_2 \rangle$, we may set $g(z) = z$.
2. **Gaussian kernel:** For $k(\mathbf{x}_1, \mathbf{x}_2) = e^{-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2}$ with $\gamma > 0$, we may set $g(z) = \max\{\sqrt{2\gamma}, 1\}z$.
3. **Laplacian kernel:** For $k(\mathbf{x}_1, \mathbf{x}_2) = e^{-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|_1}$ with $\gamma > 0$, we may set $g(z) = \sqrt{2\gamma z \sqrt{n}}$ if $0 \leq z \leq \gamma \sqrt{n}/2$ and $g(z) = z + \gamma \sqrt{n}/2$ otherwise.

4. **Polynomial kernel:** The kernel $k(\mathbf{x}_1, \mathbf{x}_2) = (\gamma \langle \mathbf{x}_1, \mathbf{x}_2 \rangle + 1)^d$ with $\gamma > 0$ and $d \in \mathbb{N}$ fails to satisfy the calmness condition if \mathbb{X} is unbounded and $d > 1$, in which case $\sqrt{k(\mathbf{x}_1, \mathbf{x}_1) - 2k(\mathbf{x}_1, \mathbf{x}_2) + k(\mathbf{x}_2, \mathbf{x}_2)}$ grows superlinearly. If $\mathbb{X} \subseteq \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_2 \leq R\}$ for some $R > 0$, however, the polynomial kernel is calm with respect to the growth function

$$g(z) = \begin{cases} \max\{\frac{1}{2R} \sqrt{2(\gamma R^2 + 1)^d}, 1\}z & \text{if } d \text{ is even,} \\ \max\{\frac{1}{2R} \sqrt{2(\gamma R^2 + 1)^d - 2(1 - \gamma R^2)^d}, 1\}z & \text{if } d \text{ is odd.} \end{cases}$$

Theorem 26 (Lifted learning problems) *If Assumption 25 holds for some growth function g , then the following statements hold for all Wasserstein radii $\rho \geq 0$.*

- (i) *For regression problems we have $\hat{J}(\rho) \leq \hat{J}_{\mathbb{H}}(\sqrt{2}g(\rho))$.*
- (ii) *For classification problems we have $\hat{J}(\rho) \leq \hat{J}_{\mathbb{H}}(g(\rho))$.*

We now argue that the lifted learning problem (30) can be solved efficiently by leveraging the following representer theorem, which generalizes (Schölkopf and Smola, 2001, Theorem 4.2) to non-separable loss functions.

Theorem 27 (Representer theorem) *Assume that we are given a symmetric positive definite kernel k on \mathbb{X} with corresponding RKHS \mathbb{H} , a set of training samples $(\hat{\mathbf{x}}_i, \hat{y}_i) \in \mathbb{X} \times \mathbb{Y}$, $i \leq N$, and an arbitrary loss function $f : (\mathbb{X} \times \mathbb{Y} \times \mathbb{R})^N \times \mathbb{R}_+ \rightarrow \mathbb{R}$ that is non-decreasing in its last argument. Then, there exist $\beta_i \in \mathbb{R}$, $i \leq N$, such that the learning problem*

$$\min_{h \in \mathbb{H}} f((\hat{\mathbf{x}}_1, \hat{y}_1, h(\hat{\mathbf{x}}_1)), \dots, (\hat{\mathbf{x}}_N, \hat{y}_N, h(\hat{\mathbf{x}}_N)), \|h\|_{\mathbb{H}}) \quad (31)$$

is solved by a hypothesis $h^ \in \mathbb{H}$ representable as $h^*(\mathbf{x}) = \sum_{i=1}^N \beta_i k(\mathbf{x}, \hat{\mathbf{x}}_i)$.*

The subsequent results involve the Kernel matrix $\mathcal{K} = [\mathcal{K}_{ij}]$ defined through $\mathcal{K}_{ij} = k(\hat{\mathbf{x}}_i, \hat{\mathbf{x}}_j)$, $i, j \leq N$. The following theorems demonstrate that the lifted learning problem (30) admits a kernel representation.

Theorem 28 (Kernelized distributionally robust regression) *Suppose that $\mathbb{X} = \mathbb{R}^n$, $\mathbb{Y} = \mathbb{R}$ and k is a symmetric positive definite kernel on \mathbb{X} with associated RKHS \mathbb{H} . If ℓ is generated by a convex and Lipschitz continuous loss function L , that is, $\ell(h(\mathbf{x}), y) = L(h(\mathbf{x}) - y)$, then (30) is equivalent to*

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^N} \frac{1}{N} \sum_{i=1}^N L\left(\sum_{j=1}^N \mathcal{K}_{ij} \beta_j - \hat{y}_i\right) + \rho \operatorname{lip}(L) \|(\mathcal{K}^{\frac{1}{2}} \boldsymbol{\beta}, 1)\|_2,$$

and for any of its minimizers $\boldsymbol{\beta}^$ the hypothesis $h^*(\mathbf{x}) = \sum_{i=1}^N \beta_i^* k(\mathbf{x}, \hat{\mathbf{x}}_i)$ is optimal in (30).*

Theorem 29 (Kernelized distributionally robust classification) *Suppose that $\mathbb{X} = \mathbb{R}^n$, $\mathbb{Y} = \{+1, -1\}$ and k is a symmetric positive definite kernel on $\mathbb{X} = \mathbb{R}^n$ with associated*

RKHS \mathbb{H} . If ℓ is generated by a convex and Lipschitz continuous loss function L , that is, $\ell(h(\mathbf{x}), y) = L(yh(\mathbf{x}))$, then (30) is equivalent to

$$\left\{ \begin{array}{ll} \min_{\beta_i, \lambda, s_i} & \lambda\rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & L\left(\sum_{j=1}^N \hat{y}_i \mathbf{K}_{ij} \beta_j\right) \leq s_i \quad i \in [N] \\ & L\left(-\sum_{j=1}^N \hat{y}_i \mathbf{K}_{ij} \beta_j\right) - \kappa\lambda \leq s_i \quad i \in [N] \\ & \text{lip}(L) \|\mathbf{K}^{\frac{1}{2}} \boldsymbol{\beta}\|_2 \leq \lambda, \end{array} \right. \quad (32)$$

and for any of its minimizers $\boldsymbol{\beta}^*$ the hypothesis $h^*(\mathbf{x}) = \sum_{i=1}^N \beta_i^* k(\mathbf{x}, \hat{\mathbf{x}}_i)$ is optimal in (30).

Theorems 28 and 29 show that the lifted learning problem (30) can be solved with similar computational effort as problem (4), that is, optimizing over a possibly infinite-dimensional RKHS of nonlinear hypotheses is not substantially harder than optimizing over the space of linear hypotheses.

Remark 30 (Kernelization in robust regression and classification) Recall from Theorem 12 that distributionally robust and classical robust linear regression are equivalent if $\Xi = \mathbb{R}^{n+1}$ and the training samples are sufficiently dispersed in the sense of Assumption 10. Similarly, Theorem 23 implies that distributionally robust and classical robust linear classification are equivalent if $\kappa = \infty$ and the training samples are non-separable in the sense of Assumption 21. One can show that Theorems 12 and 23 naturally extend to nonlinear regression and classification models over an RKHS induced by some symmetric and positive definite kernel. Specifically, one can show that some lifted robust learning problem is equivalent to the lifted distributionally robust learning problem (30) whenever the lifted training samples $(\Phi(\hat{\mathbf{x}}_1), \hat{y}_1), \dots, (\Phi(\hat{\mathbf{x}}_N), \hat{y}_N)$ satisfy Assumption 10 (for regression) or 21 (for classification). Theorems 28 and 29 thus imply that the lifted robust regression and classification problems can be solved efficiently under mild regularity conditions whenever Assumptions 10 and 21 hold, respectively. Unfortunately, these conditions are often violated for popular kernels. For example, the lifted samples are always linearly separable under the Gaussian kernel (Xu et al., 2009, p. 1496). In this case, the lifted robust classification problem can never be reduced to an efficiently solvable lifted distributionally robust classification problem of the form (30). In fact, no efficient method for solving the lifted robust classification problem seems to be known. In contrast, the lifted distributionally robust learning problems are always efficiently solvable under standard regularity conditions.

3.4. Nonlinear Hypotheses: Neural Networks²

Families of *neural networks* represent particularly expressive classes of nonlinear hypotheses. In the following, we characterize a family \mathbb{H} of neural networks with $M \in \mathbb{N}$ layers through M continuous activation functions $\sigma_m : \mathbb{R}^{n_{m+1}} \rightarrow \mathbb{R}^{n_{m+1}}$ and M weight matrices $\mathbf{W}_m \in$

2. We are grateful to an anonymous referee for encouraging us to write this section.

$\mathbb{R}^{n_{m+1} \times n_m}$, $m \in [M]$. The weight matrices can encode *fully connected* or *convolutional* layers, for example. If $n_1 = n$ and $n_{M+1} = 1$, then we may set

$$\mathbb{H} = \left\{ h \in \mathbb{R}^{\mathbb{X}} : \exists \mathbf{W}_m \in \mathbb{R}^{n_{m+1} \times n_m}, m \in [M], h(\mathbf{x}) = \sigma_M \left(\mathbf{W}_M \cdots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x})) \cdots \right) \right\}.$$

Each hypothesis $h \in \mathbb{H}$ constitutes a neural network and is uniquely determined by the collection of all weight matrices $\mathbf{W}_{[M]} := (\mathbf{W}_1, \dots, \mathbf{W}_M)$. In order to emphasize the dependence on $\mathbf{W}_{[M]}$, we will sometimes use $h(\mathbf{x}; \mathbf{W}_{[M]})$ to denote the hypotheses in \mathbb{H} . Setting $\mathbf{x}_1 = \mathbf{x}$, the features of the neural network are defined recursively through $\mathbf{x}_{m+1} = \sigma_m(\mathbf{z}_m)$, where $\mathbf{z}_m = \mathbf{W}_m \mathbf{x}_m$, $m \in [M]$. The features \mathbf{x}_m , $m = 2, \dots, M$, correspond to the *hidden* layers of the neural network, while \mathbf{x}_{M+1} determines its output.

Example 2 (Activation functions) *The following activation functions are most widely used.*

1. **Hyperbolic tangent:** $[\sigma_m(\mathbf{z}_m)]_i = (\exp(2[\mathbf{z}_m]_i) - 1) / (\exp(2[\mathbf{z}_m]_i) + 1)$
2. **Sigmoid:** $[\sigma_m(\mathbf{z}_m)]_i = 1 / (1 + \exp(-[\mathbf{z}_m]_i))$
3. **Softmax:** $[\sigma_m(\mathbf{z}_m)]_i = \exp([\mathbf{z}_m]_i) / \sum_{j=1}^{n_{m+1}} \exp([\mathbf{z}_m]_j)$
4. **Rectified linear unit (ReLU):** $[\sigma_m(\mathbf{z}_m)]_i = \max\{0, [\mathbf{z}_m]_i\}$
5. **Exponential linear unit (ELU):** $[\sigma_m(\mathbf{z}_m)]_i = \max\{0, [\mathbf{z}_m]_i\} + \min\{0, \alpha(\exp([\mathbf{z}_m]_i) - 1)\}$

The distributionally robust learning model over the hypothesis class \mathbb{H} can now be represented as

$$\inf_{h \in \mathbb{H}} \sup_{\mathbb{Q} \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbb{Q}} [\ell(h(\mathbf{x}), y)] = \inf_{\mathbf{W}_{[M]}} \sup_{\mathbb{Q} \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbb{Q}} [\ell(h(\mathbf{x}; \mathbf{W}_{[M]}), y)], \quad (33)$$

where we use the transportation metrics (11) and (16) for regression and classification problems, respectively. Moreover, we adopt the standard convention that $\ell(h(\mathbf{x}), y) = L(h(\mathbf{x}) - y)$ for regression problems and $\ell(h(\mathbf{x}), y) = L(yh(\mathbf{x}))$ for classification problems, where L is a convex and Lipschitz continuous univariate loss function. In the following we equip each feature space \mathbb{R}^{n_m} with a norm $\|\cdot\|$, $m \in [M+1]$. By slight abuse of notation, we use the same symbol for all norms even though the norms on different feature spaces may differ. Using the norm on $\mathbb{R}^{n_{m+1}}$, we define the Lipschitz modulus of σ_m as

$$\text{lip}(\sigma_m) := \sup_{\mathbf{z}, \mathbf{z}' \in \mathbb{R}^{n_{m+1}}} \left\{ \frac{\|\sigma(\mathbf{z}) - \sigma(\mathbf{z}')\|}{\|\mathbf{z} - \mathbf{z}'\|} : \mathbf{z} \neq \mathbf{z}' \right\}.$$

We are now ready to state the main result of this section, which provides a conservative upper bound on the distributionally robust learning model (33).

Theorem 31 (Distributionally robust learning with neural networks) *The distributionally robust learning model (33) is bounded above by the regularized empirical loss minimization problem*

$$\inf_{\mathbf{W}_{[M]}} \frac{1}{N} \sum_{i=1}^N \ell(h(\hat{\mathbf{x}}_i; \mathbf{W}_{[M]}), \hat{y}_i) + \rho \operatorname{lip}(L) \max \left\{ \prod_{m=1}^M \operatorname{lip}(\sigma_m) \|\mathbf{W}_m\|, \frac{c}{\kappa} \right\}, \quad (34)$$

where $c = 1$ for regression problems and $c = \max\{1, 2 \sup_{h \in \mathbb{H}, \mathbf{x} \in \mathbb{X}} |h(\mathbf{x})|\}$ for classification problems. Moreover, $\|\mathbf{W}_m\| = \sup_{\|\mathbf{x}_m\|=1} \|\mathbf{W}_m \mathbf{x}_m\|$ is the operator norm induced by the norms on \mathbb{R}^{n_m} and $\mathbb{R}^{n_{m+1}}$.

Remark 32 (Uniform upper bound on all neural networks) *For classification problems the constant c in (34) represents a uniform upper bound on all neural networks and may be difficult to evaluate in general. It is easy to estimate c , however, if the last activation function is itself bounded such as the softmax function, which yields a probability distribution over the output space. In this case one may simply set $c = 2$.*

The product term $\prod_{m=1}^M \operatorname{lip}(\sigma_m) \|\mathbf{W}_m\|$ in (34) represents an upper bound on the Lipschitz modulus of $h(\mathbf{x}; \mathbf{W}_{[M]})$. We emphasize that computing the exact Lipschitz modulus of a neural network is NP-hard even if there are only two layers and all activation functions are of the ReLU type (Scaman and Virmaux, 2018, Theorem 2). In contrast, the upper bound at hand is easy to compute as all activation functions listed in Example 2 have Lipschitz modulus 1 with respect to the Euclidean norms on the domain and range spaces (Gouk et al., 2018; Wiatowski et al., 2016). For more details on how to estimate the Lipschitz moduli of neural networks we refer to (Gouk et al., 2018; Miyato et al., 2018; Neyshabur et al., 2018; Szegedy et al., 2013).

Note that even though (34) represents a finite-dimensional optimization problem over the weight matrices of the neural network, both the empirical prediction loss as well as the regularization term are non-convex in $\mathbf{W}_{[M]}$, which complicates numerical solution. If $\kappa = \infty$, however, one can derive an alternative upper bound on the distributionally robust learning model (33) with a *convex* regularization term.

Corollary 33 (Convex regularization term) *If $\kappa = \infty$, then there is $\bar{\rho} \geq 0$ such that the distributionally robust learning model (33) is bounded above by the regularized empirical loss minimization problem*

$$\inf_{\mathbf{W}_{[M]}} \frac{1}{N} \sum_{i=1}^N \ell(h(\hat{\mathbf{x}}_i; \mathbf{W}_{[M]}), \hat{y}_i) + \bar{\rho} \sum_{m=1}^M \|\mathbf{W}_m\|. \quad (35)$$

As the empirical prediction loss remains non-convex, it is expedient to address problem (35) with local optimization methods such as stochastic gradient descent algorithms. For a comprehensive review of first- and the second-order stochastic gradient algorithms we refer to (Agarwal et al., 2017) and the references therein. In the numerical experiments we will use a stochastic proximal gradient descent algorithm that exploits the convexity of the regularization term and generates iterates $\mathbf{W}_{[M]}^k$ for $k \in \mathbb{N}$ according to the update rule

$$\mathbf{W}_m^{k+1} = \operatorname{prox}_{\eta_k \bar{\rho} \|\mathbf{W}_m\|} \left(\mathbf{W}_m^k - \eta_k \nabla_{\mathbf{W}_m} \ell(h(\hat{\mathbf{x}}_{i_k}; \mathbf{W}_{[M]}^k), \hat{y}_{i_k}) \right) \quad \forall m \in [M],$$

where $\eta_k > 0$ is a given step size and i_k is drawn randomly from the index set $[N]$, see, *e.g.*, Nitanda (2014). Here, the proximal operator associated with a convex function $\varphi : \mathbb{R}^{n_{m+1} \times n_m} \rightarrow \mathbb{R}$ is defined through

$$\text{prox}_\varphi(\mathbf{W}_m) := \arg \min_{\mathbf{W}'_m} \varphi(\mathbf{W}'_m) + \frac{1}{2} \|\mathbf{W}'_m - \mathbf{W}_m\|_F^2,$$

where $\|\cdot\|_F$ stands for the Frobenius norm. The algorithm is stopped as soon as the improvement of the objective value falls below a prescribed threshold. As the empirical prediction loss is non-convex and potentially non-smooth, the algorithm fails to offer any strong performance guarantees. For the scalability of the algorithm, however, it is essential that the proximal operator can be evaluated efficiently.

Example 3 (Proximal operator) Suppose that all feature spaces \mathbb{R}^{n_m} are equipped with the p -norm for some $p \in \{1, 2, \infty\}$, which implies that all parameter spaces $\mathbb{R}^{n_{m+1} \times n_m}$ are equipped with the corresponding matrix p -norm. In this case the proximal operator of $\varphi(\mathbf{W}_m) = \eta \|\mathbf{W}_m\|_p$ for some fixed $\eta > 0$ can be evaluated highly efficiently.

1. **MACS** ($p = 1$): The matrix 1-norm returns the maximum absolute column sum (MACS). Evaluating the proximal operator of $\varphi(\mathbf{W}_m) = \eta \|\mathbf{W}_m\|_1$ amounts to solving the minimization problem

$$\text{prox}_\varphi(\mathbf{W}_m) = \begin{cases} \min_{\mathbf{W}'_m, u} & \eta u + \sum_{i=1}^{n_m} \|[\mathbf{W}'_m]_{:,i} - [\mathbf{W}_m]_{:,i}\|_2^2 \\ \text{s.t.} & \|[\mathbf{W}'_m]_{:,i}\|_1 \leq u \quad i \in [n_m], \end{cases}$$

where $[\mathbf{W}_m]_{:,i}$ and $[\mathbf{W}'_m]_{:,i}$ represent the i -th columns of \mathbf{W}_m and \mathbf{W}'_m , respectively. For any fixed u , the above problem decomposes into n_m projections of the vectors $[\mathbf{W}_m]_{:,i}$, $i \in [n_m]$, to the ℓ_1 -ball of radius u centered at the origin. Each of these projections can be computed via an efficient sorting algorithm proposed in (Duchi et al., 2008). Next, we can use any line search method such as the golden-section search algorithm to optimize over u , thereby solving the full proximal problem.

2. **Spectral** ($p = 2$): The matrix 2-norm coincides with the spectral norm, which returns the maximum singular value. In this case, the proximal problem for $\varphi(\mathbf{W}_m) = \eta \|\mathbf{W}_m\|_2$ can be solved analytically via singular value thresholding (Cai et al., 2010, Theorem 2.1), that is, given the singular value decomposition $\mathbf{W}_m = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ with $\mathbf{U} \in \mathbb{R}^{n_{m+1} \times n_{m+1}}$ orthogonal, $\mathbf{S} \in \mathbb{R}_+^{n_{m+1} \times n_m}$ diagonal and $\mathbf{V} \in \mathbb{R}^{n_m \times n_m}$ orthogonal, the proximal operator satisfies

$$\text{prox}_\varphi(\mathbf{W}_m) = \text{prox}_\varphi(\mathbf{U}\mathbf{S}\mathbf{V}^\top) = \mathbf{U}\tilde{\mathbf{S}}\mathbf{V}^\top, \quad \text{where} \quad \tilde{\mathbf{S}}_{ij} = \max\{\mathbf{S}_{ij} - \eta, 0\}.$$

The singular value decomposition can be accelerated using a randomized algorithm proposed in (Halko et al., 2011).

3. **MARS** ($p = \infty$): The matrix ∞ -norm returns the maximum absolute row sum (MARS) and thus satisfies $\|\mathbf{W}_m\|_\infty = \|\mathbf{W}_m^\top\|_1$. Therefore, one can use the iterative scheme developed for MACS to compute the proximal operator of $\varphi(\mathbf{W}_m) = \eta \|\mathbf{W}_m\|_\infty$ by simply transposing the weight matrix \mathbf{W}_m .

The convergence behavior of the stochastic proximal gradient descent algorithm can be further improved by including a momentum term inside the proximal operator, see, *e.g.*, Loizou and Richtárik (2017).

4. Generalization Bounds

Generalization bounds constitute upper confidence bounds on the out-of-sample error. Traditionally, generalization bounds are derived by controlling the complexity of the hypothesis space, which is typically quantified in terms of its VC-dimension or via covering numbers or Rademacher averages (Shalev-Shwartz and Ben-David, 2014). Strengthened generalization bounds for large margin classifiers can be obtained by improving the estimates of the VC-dimension and the Rademacher average (Shivaswamy and Jebara, 2007, 2010). We will now demonstrate that distributionally robust learning models of the type (4) or (30) enjoy simple new generalization bounds that can be obtained under minimal assumptions. In particular, they do not rely on any notions of hypothesis complexity and may therefore even extend to hypothesis spaces with infinite VC-dimensions. Our approach is reminiscent of the generalization theory for robust support vector machines portrayed in (Xu et al., 2009), which also replaces measures of hypothesis complexity with robustness properties. However, we derive explicit finite sample guarantees, while (Xu et al., 2009) establishes asymptotic consistency results. Moreover, we relax some technical conditions used in (Xu et al., 2009) such as the compactness of the input space \mathbb{X} .

The key enabling mechanism of our analysis is a measure concentration property of the Wasserstein metric, which holds whenever the unknown data-generating distribution has exponentially decaying tails.

Assumption 34 (Light-tailed distribution) *There exist constants $a > 1$ and $A > 0$ and a reference point $\xi' \in \mathbb{R}^{n+1}$ such that $\mathbb{E}^{\mathbb{P}}[\exp(d(\xi, \xi')^a)] \leq A$, where d denotes the usual mass transportation cost.*

Theorem 35 (Measure concentration (Fournier and Guillin, 2015, Theorem 2))
If Assumption 34 holds, then we have

$$\mathbb{P}^N \left\{ W(\mathbb{P}, \hat{\mathbb{P}}_N) \geq \rho \right\} \leq \begin{cases} c_1 \exp(-c_2 N \rho^{\max\{n+1, 2\}}) & \text{if } \rho \leq 1, \\ c_1 \exp(-c_2 N \rho^a) & \text{if } \rho > 1, \end{cases} \quad (36)$$

for all $N \geq 1$, $n \neq 1$, and $\rho > 0$, where the constants $c_1, c_2 > 0$ depend only on a, A, d and n .³

Theorem 35 asserts that the empirical distribution $\hat{\mathbb{P}}_N$ converges exponentially fast to the unknown data-generating distribution \mathbb{P} , in probability with respect to the Wasserstein metric, as the sample size N tends to infinity. We can now derive simple generalization bounds by increasing the Wasserstein radius ρ until the violation probability on the right hand side of (36) drops below a prescribed significance level $\eta \in (0, 1]$. Specifically,

3. A similar but slightly more complicated inequality also holds for the special case $n = 1$; see (Fournier and Guillin, 2015, Theorem 2) for details.

Theorem 36 borrowed from (Mohajerin Esfahani and Kuhn, 2018, Theorem 3.5) implies that $\mathbb{P}^N\{\mathbb{P} \in \mathbb{B}_\rho(\hat{\mathbb{P}}_N)\} \geq 1 - \eta$ for any $\rho \geq \rho_N(\eta)$, where

$$\rho_N(\eta) = \begin{cases} \left(\frac{\log(c_1/\eta)}{c_2 N} \right)^{\frac{1}{\max\{n+1, 2\}}} & \text{if } N \geq \frac{\log(c_1/\eta)}{c_2}, \\ \left(\frac{\log(c_1/\eta)}{c_2 N} \right)^{\frac{1}{a}} & \text{if } N < \frac{\log(c_1/\eta)}{c_2}. \end{cases} \quad (37)$$

Theorem 36 (Basic generalization bound) *If Assumption 34 holds, then*

$$\mathbb{P}^N \left\{ \mathbb{E}^{\mathbb{P}} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] \leq \sup_{\mathbb{Q} \in \mathbb{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbb{Q}} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] \quad \forall \mathbf{w} \in \mathbb{R}^n \right\} \geq 1 - \eta \quad (38)$$

for any $N \geq 1$, $n \neq 1$, $\eta \in (0, 1]$ and $\rho \geq \rho_N(\eta)$.

Remark 37 (Discussion of basic generalization bound) *The following comments are in order.*

I. Performance guarantees for optimal hypotheses: If $\hat{J}(\rho)$ denotes the minimum and $\hat{\mathbf{w}}$ a minimizer of the distributionally robust learning problem (4), then Theorem 36 implies that

$$\mathbb{P}^N \left\{ \mathbb{E}^{\mathbb{P}} [\ell(\langle \hat{\mathbf{w}}, \mathbf{x} \rangle, y)] \leq \hat{J}(\rho) \right\} \geq 1 - \eta$$

for any $N \geq 1$, $n > 1$, $\eta \in (0, 1]$ and $\rho \geq \rho_N(\eta)$.

II. Light-tail assumption: Assumption 34 is restrictive but unavoidable for any measure concentration result of the type described in Theorem 35. It is automatically satisfied if the input-output pair has bounded support or is known to follow a Gaussian or exponential distribution, for instance.

III. Asymptotic consistency: It is clear from (37) that for any fixed $\eta \in (0, 1]$, the radius $\rho_N(\eta)$ tends to 0 as N increases. Moreover, Theorem 3.6 in (Mohajerin Esfahani and Kuhn, 2018) implies that if η_N converges to 0 at a carefully chosen rate (e.g., $\eta_N = \exp(-\sqrt{N})$), then the solution of the distributionally robust learning problem (4) with Wasserstein radius $\rho = \rho_N(\eta_N)$ converges almost surely to the solution of the ideal learning problem that minimizes the out-of-sample error under the unknown true distribution \mathbb{P} .

IV. Curse of dimensionality: The Wasserstein radius (37) has two decay regimes. For small N , $\rho_N(\eta)$ decays as $N^{-\frac{1}{a}}$, and for large N it is proportional to $N^{-\frac{1}{n+1}}$. We thus face a curse of dimensionality for large sample sizes. In order to half the Wasserstein radius, one has to increase N by a factor of 2^n . This curse of dimensionality is fundamental, i.e., the dependence of the measure concentration result in Theorem 35 on the input dimension n cannot be improved for generic distributions \mathbb{P} ; see (Weed and Bach, 2019) or (Fournier and Guillin, 2015, Section 1.3). Improvements are only possible in special cases, e.g., if \mathbb{P} is finitely supported.

V. **Extension to nonlinear hypotheses:** Theorem 36 directly extends to any distributionally robust learning problem over an RKHS \mathbb{H} induced by some symmetric and positive definite kernel function k . Specifically, if k is calm in the sense of Assumption 25 with growth function g , then we have

$$\mathbb{P}^N \left\{ \mathbb{E}^{\mathbb{P}}[\ell(h(\mathbf{x}), y)] \leq \sup_{\mathbb{Q} \in \mathbb{B}_\rho(\hat{\mathbb{P}}_N^{\mathbb{H}})} \mathbb{E}^{\mathbb{Q}}[\ell(\langle h, \mathbf{x}_{\mathbb{H}} \rangle_{\mathbb{H}}, y)] \quad \forall h \in \mathbb{H} \right\} \geq 1 - \eta \quad (39)$$

for any $N \geq 1$, $n \neq 1$, $\eta \in (0, 1]$ and $\rho \geq cg(\rho_N(\eta))$, where $c = \sqrt{2}$ for regression problems and $c = 1$ for classification problems. To see this, note that the inclusion $\mathbb{P} \in \mathbb{B}_\rho(\hat{\mathbb{P}}_N)$ implies

$$\mathbb{E}^{\mathbb{P}}[\ell(h(\mathbf{x}), y)] \leq \sup_{\mathbb{Q} \in \mathbb{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbb{Q}}[\ell(h(\mathbf{x}), y)] \leq \sup_{\mathbb{Q} \in \mathbb{B}_\rho(\hat{\mathbb{P}}_N^{\mathbb{H}})} \mathbb{E}^{\mathbb{Q}}[\ell(\langle h, \mathbf{x}_{\mathbb{H}} \rangle_{\mathbb{H}}, y)] \quad \forall h \in \mathbb{H}, \quad (40)$$

where the second inequality follows from the proof of Theorem 26. The generalization bound (39) thus holds because $\mathbb{P}^N\{\mathbb{P} \in \mathbb{B}_\rho(\hat{\mathbb{P}}_N)\} \geq 1 - \eta$ for any $\rho \geq \rho_N(\eta)$. Note that the rightmost term in (40) can be computed for any finitely generated hypothesis $h \in \mathbb{H}$ representable as $h(\mathbf{x}) = \sum_{i=1}^N \beta_i k(\mathbf{x}, \hat{\mathbf{x}}_i)$, which follows from Theorems 28 and 29, while the middle term is hard to compute. We emphasize that the generalization bound (39) does not rely on any notion of hypothesis complexity and remains valid even if \mathbb{H} has infinite VC-dimension (e.g., if \mathbb{H} is generated by the Gaussian kernel).

Theorem 35 provides a confidence set for the unknown probability distribution \mathbb{P} , and Theorem 36 uses this confidence set to construct a uniform generalization bound on the prediction error under \mathbb{P} . The radius of the confidence set for \mathbb{P} decreases slowly due to a curse of dimensionality, but the decay rate is essentially optimal. This does not imply that the decay rate of the generalization bound (38) is optimal, too. In fact, the worst-case expected error over a Wasserstein ball of radius ρ can be a $(1 - \eta)$ -confidence bound on the expected error under \mathbb{P} even if the Wasserstein ball fails to contain \mathbb{P} with confidence $1 - \eta$. Thus, the measure concentration result of Theorem 35 is too powerful for our purposes and leads to an over-conservative generalization bound. Below we will show that the curse of dimensionality in the generalization bound (38) can be broken if we impose the following restriction on the hypothesis space.

Assumption 38 (Hypothesis space) *The space of admissible hypotheses in (4) is restricted to $\mathbb{W} \subseteq \mathbb{R}^n$. There exists $\underline{\Omega} > 0$ with $\inf_{\mathbf{w} \in \mathbb{W}} \|(\mathbf{w}, -1)\|_* \geq \underline{\Omega}$ if (4) is a regression problem and $\inf_{\mathbf{w} \in \mathbb{W}} \|\mathbf{w}\|_* \geq \underline{\Omega}$ if (4) is a classification problem. Similarly, there exists $\bar{\Omega} \geq 0$ with $\sup_{\mathbf{w}, \mathbf{w}' \in \mathbb{W}} \|\mathbf{w} - \mathbf{w}'\|_\infty \leq \bar{\Omega}$.*

Theorem 39 (Improved generalization bound) *Suppose that Assumptions 34 and 38 hold, and the function L is Lipschitz continuous. Moreover, assume that $\Xi = \mathbb{R}^{n+1}$ and $M_n = \max_{i \leq n} \|\mathbf{e}_i^{n+1}\|_*$ if (4) is a regression problem, while $\Xi = \mathbb{R}^n \times \{-1, 1\}$ and $M_n = \max_{i \leq n} \|\mathbf{e}_i^n\|_*$ if (4) is a classification problem, where \mathbf{e}_i^n is the i -th standard basis vector in \mathbb{R}^n . Then, there exist constants $c_3 \geq 1$, $c_4 > 0$ depending only on*

the light tail constants a and A such that the generalization bound (38) holds for any $N \geq \max \{ (16n/c_4)^2, 16 \log(c_3/\eta)/c_4 \}$, $\eta \in (0, 1]$ and $\rho \geq \rho'_N(\eta)$, where

$$\rho'_N(\eta) = \frac{2\bar{\Omega}}{\sqrt{N}\underline{\Omega}} \left[M_n n A + \sqrt{\frac{n \log(\sqrt{N}) + \log(c_3/\eta)}{c_4}} \right].$$

The improved generalization bound from Theorem 38 does not suffer from a curse of dimensionality. In fact, in order to half the Wasserstein radius $\rho'_N(\eta)$, it suffices to increase the sample size N by a factor of 4, irrespective of the input dimension n .

Remark 40 (Discussion of improved generalization bound) *The following comments are in order.*

- I. **Bounds on hypothesis space:** *Assumption 38 imposes upper and lower bounds on \mathbb{W} . The upper bound enables us to control the difference between the empirical and the true expected loss uniformly across all admissible hypotheses. This bound is less restrictive than the uniform bound on the loss function used to derive Rademacher generalization bounds (see, e.g., (Shalev-Shwartz and Ben-David, 2014, Theorem 26.4)), which essentially imposes upper bounds both on the hypotheses and the input-output pairs. The lower bound in Assumption 38 is restrictive for classification problems but trivially holds for regression problems because $\|(\mathbf{w}, -1)\|_*$ is uniformly bounded away from zero for any (dual) norm on \mathbb{R}^{n+1} .*
- II. **Breaking the curse of dimensionality:** *By leveraging Assumption 38, Theorem 39 reduces the critical Wasserstein radius in the generalization bound (38) from $\rho_N(\eta) \propto \mathcal{O}([\log(\eta^{-1})/N]^{1/(n+1)})$, which suffers from a curse of dimensionality, to $\rho'_N(\eta) \propto \mathcal{O}([\log(\eta^{-1}) + n \log(N)]/N]^{1/2})$, which essentially follows a square root law reminiscent of the central limit theorem.*

5. Error and Risk Estimation

Once a hypothesis $h(\mathbf{x})$ has been chosen, it is instructive to derive pessimistic *and* optimistic estimates of its out-of-sample prediction error (in the case of regression) or its out-of-sample risk (in the case of classification). We will argue below that the distributionally robust optimization techniques developed in this paper also offer new perspectives on error and risk estimation. For ease of exposition, we ignore any support constraints, that is, we set $\mathbb{X} = \mathbb{R}^n$ and $\mathbb{Y} = \mathbb{R}$ (for regression) or $\mathbb{X} = \mathbb{R}^n$ and $\mathbb{Y} = \{+1, -1\}$ (for classification). Moreover, we focus on linear hypotheses of the form $h(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle$. Note, however, that all results extend directly to conic representable support sets and to nonlinear hypotheses.

In the context of regression, we aim to estimate the prediction error defined as $\mathcal{E}(\mathbf{w}) = \mathbb{E}^{\mathbb{P}} [|y - \langle \mathbf{w}, \mathbf{x} \rangle|]$ or, more precisely, the mean absolute prediction error under the unknown data-generating distribution \mathbb{P} . As usual, we assume that the transportation metric d is induced by a norm $\|\cdot\|$ on the input-output space \mathbb{R}^{n+1} .

Theorem 41 (Error bounds in linear regression) *The prediction error admits the following estimates.*

(i) The worst-case error $\mathcal{E}_{\max}(\mathbf{w}) = \sup_{\mathbf{Q} \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}}[|y - \langle \mathbf{w}, \mathbf{x} \rangle|]$ is given by

$$\mathcal{E}_{\max}(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N |\hat{y}_i - \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle| + \rho \|(\mathbf{w}, -1)\|_*. \quad (41a)$$

[(ii)] The best-case error $\mathcal{E}_{\min}(\mathbf{w}) = \inf_{\mathbf{Q} \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}}[|y - \langle \mathbf{w}, \mathbf{x} \rangle|]$ is given by

$$\mathcal{E}_{\min}(\mathbf{w}) = \max \left\{ \frac{1}{N} \sum_{i=1}^N |\hat{y}_i - \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle| - \rho \|(\mathbf{w}, -1)\|_*, 0 \right\} \quad (41b)$$

In the context of classification, we aim to quantify the risk $\mathcal{R}(\mathbf{w}) = \mathbb{P}[y \neq \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle)]$, that is, the misclassification probability under the unknown true distribution \mathbb{P} . Note that the risk can equivalently be defined as the expectation of a characteristic function, that is, $\mathcal{R}(\mathbf{w}) = \mathbb{E}^{\mathbb{P}}[\mathbb{1}_{\{y \neq \langle \mathbf{w}, \mathbf{x} \rangle\}}]$. As usual, we assume that the transportation metric d is of the form (16), where $\kappa \geq 0$ is the cost of flipping a label.

Theorem 42 (Risk bounds in linear classification) *The risk admits the following estimates.*

(i) The worst-case risk $\mathcal{R}_{\max}(\mathbf{w}) = \sup_{\mathbf{Q} \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}}[\mathbb{1}_{\{y \neq \langle \mathbf{w}, \mathbf{x} \rangle\}}]$ is given by

$$\mathcal{R}_{\max}(\mathbf{w}) = \begin{cases} \min_{\substack{\lambda, s_i \\ r_i, t_i}} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & 1 - r_i \hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle \leq s_i & i \in [N] \\ & 1 + t_i \hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \lambda \kappa \leq s_i & i \in [N] \\ & r_i \|\mathbf{w}\|_* \leq \lambda, \quad t_i \|\mathbf{w}\|_* \leq \lambda & i \in [N] \\ & r_i, t_i, s_i \geq 0 & i \in [N]. \end{cases} \quad (42a)$$

(ii) The best-case risk $\mathcal{R}_{\min}(\mathbf{w}) = \inf_{\mathbf{Q} \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}}[\mathbb{1}_{\{y \neq \langle \mathbf{w}, \mathbf{x} \rangle\}}]$ is given by

$$\mathcal{R}_{\min}(\mathbf{w}) = 1 - \begin{cases} \min_{\substack{\lambda, s_i \\ r_i, t_i}} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & 1 + r_i \hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle \leq s_i & i \in [N] \\ & 1 - t_i \hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \lambda \kappa \leq s_i & i \in [N] \\ & r_i \|\mathbf{w}\|_* \leq \lambda, \quad t_i \|\mathbf{w}\|_* \leq \lambda & i \in [N] \\ & r_i, t_i, s_i \geq 0 & i \in [N]. \end{cases} \quad (42b)$$

We emphasize that, as the hypothesis \mathbf{w} is fixed, the error and risk estimation problems (41) and (42) constitute tractable linear programs that can be solved highly efficiently.

Remark 43 (Confidence intervals for error and risk) *If the Wasserstein radius is set to $\rho_N(\eta/2)$ defined in (37), where $\eta \in (0, 1]$ is a prescribed significance level, then Theorem 36 implies that $\mathcal{E}(\hat{\mathbf{w}}) \in [\mathcal{E}_{\min}(\hat{\mathbf{w}}), \mathcal{E}_{\max}(\hat{\mathbf{w}})]$ and $\mathcal{R}(\hat{\mathbf{w}}) \in [\mathcal{R}_{\min}(\hat{\mathbf{w}}), \mathcal{R}_{\max}(\hat{\mathbf{w}})]$ with confidence $1 - \eta$ for any $\hat{\mathbf{w}} \in \mathbb{R}^n$ that may even depend on the training data. Theorem 41 implies that the confidence interval for the true error $\mathcal{E}(\hat{\mathbf{w}})$ can be calculated analytically from (41), while Theorem 42 implies that the confidence interval for the true risk $\mathcal{R}(\hat{\mathbf{w}})$ can be computed efficiently by solving the tractable linear programs (42).*

Remark 44 (Extension to nonlinear hypotheses) *By using the tools of Section 3.3, Theorems 41 and 42 generalize immediately to nonlinear hypotheses that range over a RKHS. Specifically, we can formulate lifted error and risk estimation problems where the inputs $\mathbf{x} \in \mathbb{X}$ are replaced with features $\mathbf{x}_{\mathbb{H}} \in \mathbb{H}$, while each nonlinear hypothesis $h \in \mathbb{H}$ over the input space \mathbb{X} is identified with a linear hypothesis $h_{\mathbb{H}} \in \mathbb{H}$ over the feature space \mathbb{H} through the identity $h_{\mathbb{H}}(\mathbf{x}_{\mathbb{H}}) = \langle h, \mathbf{x}_{\mathbb{H}} \rangle_{\mathbb{H}}$. Tractability is again facilitated by Theorem 27, which allows us to focus on finitely parameterized hypotheses of the form $h(\mathbf{x}) = \sum_{i=1}^N \beta_i k(\mathbf{x}, \hat{\mathbf{x}}_i)$.*

6. Numerical Results

We showcase the power of regularization via mass transportation in various applications based on standard datasets from the literature. All optimization problems are implemented in Python and solved with Gurobi 7.5.1. All experiments are run on an Intel XEON CPU (3.40GHz), and the corresponding codes are made publicly available at <https://github.com/sorooshafiee/Regularization-via-Transportation>.

6.1. Regularization with Pre-selected Parameters

We first assess how the out-of-sample performance of a distributionally robust support vector machine (DRSVM) is impacted by the choice of the Wasserstein radius ρ , the cost κ of flipping a label, and the kernel function k . To this end, we solve three binary classification problems from the MNIST database (LeCun et al., 1998) targeted at distinguishing pairs of similar handwritten digits (1-vs-7, 3-vs-8, 4-vs-9). In the first experiment we optimize over linear hypotheses and use the separable transportation metric (16) involving the ∞ -norm on the input space. All results are averaged over 100 independent trials. In each trial, we randomly select 500 images to train the DRSVM model (20) and use the remaining images for testing. The correct classification rate (CCR) on the test data, averaged across all 100 trials, is visualized in Figure 1 as a function of the Wasserstein radius ρ for each $\kappa \in \{0.1, 0.25, 0.5, 0.75, \infty\}$. The best out-of-sample CCR is obtained for $\kappa = 0.25$ uniformly across all Wasserstein radii, and performance deteriorates significantly when κ is reduced or increased. Recall from Remark 18 that, as κ tends to infinity, the DRSVM reduces to the classical regularized support vector machine (RSVM) with 1-norm regularizer. Thus, the results of Figure 1 indicate that regularization via mass transportation may be preferable to classical regularization in terms of the maximum achievable out-of-sample CCR. More specifically, we observe that the out-of-sample CCR of the best DRSVM ($\kappa = 0.25$) displays a slightly higher and significantly wider plateau around the optimal regularization parameter ρ than the classical RSVM ($\kappa = \infty$). This suggests that the regularization parameter in DRSVMs may be easier to calibrate from data than in RSVMs, a conjecture that will be put

Table 1: Average out-of-sample CCR scores of the DRSVM with learned parameters.

	Polynomial		Laplacian		Gaussian	
	RSVM	DRSVM	RSVM	DRSVM	RSVM	DRSVM
1-vs-7	98.9 ± 0.2	99.1 ± 0.2	98.3 ± 0.5	98.5 ± 0.4	99.1 ± 0.2	99.2 ± 0.2
3-vs-8	95.2 ± 0.4	97.0 ± 0.4	96.5 ± 0.4	96.8 ± 0.4	97.0 ± 0.3	97.2 ± 0.3
4-vs-9	95.0 ± 0.4	96.5 ± 0.4	95.8 ± 0.6	96.0 ± 0.6	96.8 ± 0.4	96.9 ± 0.4

to scrutiny in Section 6.2. Finally, Figure 1 reveals that the standard (unregularized) support vector machine (SVM), which can be viewed as a special case of the DRSVM with $\rho = 0$, is dominated by the RSVMs and DRSVMs across a wide range of regularization parameters.⁴ Note that the SVM problem (20) with $\rho = 0$ reduces to a linear program and may thus suffer from multiple optimal solutions. This explains why the limiting out-of-sample CCR for $\rho \downarrow 0$ changes with κ .

6.2. Regularization with Learned Parameters

It is easy to read off the best regularization parameters ρ and κ from the charts in Figure 1. As these charts are constructed from more than 12,000 test samples, however, they are not accessible in the training phase. In practice, ρ and κ must be calibrated from the training data alone. This motivates us to revisit the three classification problems from Section 6.1 using a fully data-driven procedure, where all free model parameters are calibrated via 5-fold cross validation; see, *e.g.*, (Abu-Mostafa et al., 2012, § 4.3.3). Moreover, to evaluate the benefits of kernelization, we now solve a generalized DRSVM model of the form (32), which implicitly optimizes over all nonlinear hypotheses in some RKHS. As explained in Section 3.3, kernelization necessitates the use of the separable transportation metric (16) with the Euclidean norm on the input space.

All free parameters of the resulting DRSVM model are restricted to finite search grids in order to ease the computational burden of cross validation. Specifically, we select the Wasserstein radius ρ from within $\{b \cdot 10^e : b \in \{1, 5\}, e \in \{1, 2, 3, 4\}\}$ and the label flipping cost κ from within $\{0.1, 0.25, 0.5, 0.75, \infty\}$. Moreover, we select the degree d of the polynomial kernel from within $\{1, 2, 3, 4, 5\}$ and the peakedness parameter γ of the Laplacian and Gaussian kernels from within $\{\frac{1}{100}, \frac{1}{81}, \frac{1}{64}, \frac{1}{49}, \frac{1}{36}, \frac{1}{25}\}$. Otherwise, we use the same experimental setup as in Section 6.1. Table 1 reports the averages and standard deviations of the CCR scores on the test data based on 100 independent trials. We observe that the DRSVM (ρ , κ , d , and γ learned by cross validation) outperforms the RSVM (ρ , d and γ learned by cross validation, $\kappa = \infty$) consistently across all tested kernel functions (Polynomial, Laplacian, Gaussian). Note that the DRSVM with polynomial kernel subsumes the non-kernelized DRSVM (20) as a special case because the polynomial kernel with $d = 1$ coincides with the linear kernel.

4. By slight abuse of notation, we use the acronym ‘SVM’ to refer to the unregularized empirical hinge loss minimization problem even though the traditional formulations of the support vector machine involve a Tikhonov regularization term.

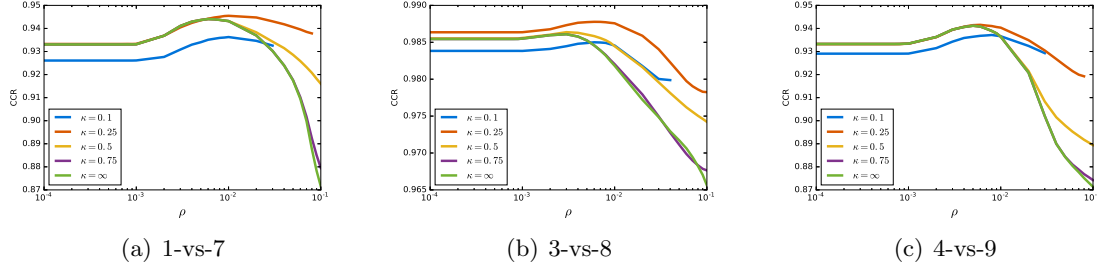


Figure 1: Average out-of-sample CCR scores of the DRSVM with pre-selected parameters.

In the third experiment, we assess the out-of-sample performance of the DRSVM (20) for different transportation metrics on 10 standard datasets from the UCI repository (Bache and Lichman, 2013). Specifically, we use different variants of the separable transportation metric (16), where distances in the input space are measured via a p -norm with $p \in \{1, 2, \infty\}$. We focus exclusively on linear hypotheses because the kernelization techniques described in Section 3.3 are only available for $p = 2$. The DRSVM is compared against the standard (unregularized) SVM and the RSVM with q -norm regularizer ($\frac{1}{p} + \frac{1}{q} = 1$). All results are averaged across 100 independent trials. In each trial, we randomly select 75% of the data for training and the remaining 25% for testing. The training dataset is first standardized to zero mean and unit variance along each coordinate axis. The Wasserstein radius ρ and the label flipping cost κ in the DRSVM as well as the regularization weight ρ in the RSVM are estimated via stratified 5-fold cross validation.

Classifier performance is now quantified in terms of the *receiver operating characteristic* (ROC) curve, which plots the true positive rate (percentage of correctly classified test samples with true label $y = 1$) against the false positive rate (percentage of *incorrectly* classified test samples with true label $y = -1$) by sweeping the discrimination threshold. Specifically, we use the *area under the ROC curve* (AUC) as a measure of classifier performance. AUC does not bias on the size of the test data and is a more appropriate performance measure than CCR in the presence of an unbalanced label distribution in the training data. We emphasize that most of the considered datasets are indeed imbalanced, and thus a high CCR score would not necessarily provide evidence of superior classifier performance. The averages and standard deviations of the AUC scores based on 100 trials are reported in Table 2. The results suggest that the DRSVM outperforms the RSVM in terms of AUC for all norms by about the same amount by which the RSVM outperforms the classical hinge loss minimization, consistently across all datasets.

6.3. Multi-Label Classification

The aim of object recognition is to discover instances of particular object classes in digital images. We now describe an object recognition experiment based on the PASCAL VOC 2007 dataset Everingham et al. (2010) consisting of 9,963 images, which are pre-partitioned into 25% for training, 25% for validation and 50% for testing. Each image is annotated with 20 binary labels corresponding to 20 given object categories (the n -th label is set to

Table 2: Average out-of-sample AUC scores of the SVM, RSVM and DRSVM.

	SVM	$p = \infty / q = 1$		$p = 2 / q = 2$		$p = 1 / q = \infty$	
		RSVM	DRSVM	RSVM	DRSVM	RSVM	DRSVM
Australian	91.6 \pm 3.0	91.5 \pm 3.2	92.0 \pm 2.5	92.0 \pm 2.2	92.3 \pm 2.0	91.9 \pm 2.8	92.2 \pm 2.4
Blood transfusion	73.7 \pm 3.8	73.8 \pm 3.8	75.5 \pm 3.8	74.9 \pm 3.5	75.5 \pm 3.7	75.4 \pm 3.4	75.4 \pm 3.7
Climate model	93.8 \pm 3.9	94.4 \pm 4.0	94.3 \pm 4.0	94.3 \pm 3.8	94.0 \pm 4.0	93.6 \pm 3.9	93.9 \pm 4.0
Cylinder	72.0 \pm 3.7	71.2 \pm 4.0	72.1 \pm 4.1	71.3 \pm 4.0	71.8 \pm 3.8	71.5 \pm 3.8	72.2 \pm 3.7
Heart	90.4 \pm 2.7	90.1 \pm 2.7	90.3 \pm 2.7	90.6 \pm 2.6	90.9 \pm 2.5	90.5 \pm 2.6	90.7 \pm 2.6
Ionosphere	85.0 \pm 4.9	89.7 \pm 4.5	89.2 \pm 4.3	90.4 \pm 3.7	89.9 \pm 3.9	86.0 \pm 4.9	87.2 \pm 4.8
Liver disorders	60.5 \pm 0.0	61.1 \pm 0.7	61.7 \pm 0.7	61.2 \pm 0.3	61.7 \pm 0.5	61.1 \pm 0.4	61.8 \pm 0.5
QSAR	90.5 \pm 1.5	90.5 \pm 1.6	91.0 \pm 1.6	90.5 \pm 1.5	91.2 \pm 1.5	90.6 \pm 1.5	91.1 \pm 1.6
Splice	92.1 \pm 0.0	93.0 \pm 0.4	93.1 \pm 0.1	92.5 \pm 0.1	92.6 \pm 0.1	92.0 \pm 0.1	92.5 \pm 0.1
Thoracic surgery	61.7 \pm 7.1	61.5 \pm 6.5	64.6 \pm 6.6	64.4 \pm 6.4	64.3 \pm 7.0	64.0 \pm 6.3	64.6 \pm 6.3

+1 if the image contains the n -th object and to -1 otherwise). A multi-label classifier is a function that predicts all labels of an unlabelled input image. The ability of a classifier to detect objects belonging to any fixed category is measured by the *average precision* (AP), which is defined in Everingham et al. (2010) as (a proxy for) the area under the classifier’s precision-recall curve. The overall performance of a classifier is quantified by the *mean average precision* (mAP), that is, the arithmetic mean of the AP scores across all object categories.

In the first scenario, we train a separate binary RSVM and DRSVM classifier for each of the 20 object categories. This classifier predicts whether an object of the respective category appears in the input image. At the beginning we preprocess the entire dataset by resizing each image to 256×256 pixels and extracting the central patch of 244×244 pixels. As shown in (Chatfield et al., 2014; Donahue et al., 2014; Zeiler and Fergus, 2014), the features generated by the penultimate layer of a deep convolutional neural network trained on a large image dataset provide a powerful image descriptor. Using the ALEXNET neural network trained on the ImageNet dataset (Krizhevsky et al., 2012), we can thus compress each (preprocessed) image of the PASCAL VOC 2007 dataset into 1,000 meaningful features. We normalize these feature vectors to lie on the unit sphere. When training the RSVM and DRSVM classifiers, we can thus work with these feature vectors instead of the corresponding images. Moreover, we restrict attention to linear hypotheses and assume that transportation distances in the input-output space are measured by the separable metric (16) with the Euclidean norm on the input space. We tune the Wasserstein radius $\rho \in \{b \cdot 10^e : b \in \{1, \dots, 9\}, e \in \{-2, -3, -4\}\}$ and the label flipping cost $\kappa \in \{0.1, 0.2, \dots, 1, \infty\}$ via the holdout method using the validation data. As usual, we fix $\kappa = \infty$ for RSVM. Table 3 reports the AP scores of the RSVM and DRSVM models for each object category. The ensemble of all 20 binary RSVM or DRSVM classifiers, respectively, can be viewed as a naïve multi-label classifier that predicts all labels of an image. As DRSVM outperforms RSVM on an object-by-object basis, it also wins in terms of mAP.

In the second scenario, we construct a proper multi-label classifier by fine-tuning the last layer of the pre-trained ALEXNET network. To this end, we replace the original M -th layer of the network with a new fully connected layer characterized by a parameter matrix $\mathbf{W}_M \in \mathbb{R}^{20 \times 1000}$, and we set σ_M to the Sigmoid activation function. The resulting classifier outputs for each of the 20 object categories a probability that an object from the respective

category appears in the input image. The quality of a classifier (which is encoded by \mathbf{W}_M) is measured by the cross-entropy loss function, which naturally generalizes the logloss to multiple labels. The resulting empirical loss minimization problem is enhanced with a regularization term proportional to $\|\mathbf{W}_M\|_{1,1}$ (Lasso), $\|\mathbf{W}_M\|_F^2$ (Tikhonov), $\|\mathbf{W}_M\|_1$ (MACS), $\|\mathbf{W}_M\|_2$ (Spectral) or $\|\mathbf{W}_M\|_\infty$ (MARS). By using similar arguments as in Section 3.4, one can show that the empirical cross-entropy with MACS, Spectral or MARS regularization term overestimates the worst-case expected cross-entropy over all distributions of $(\mathbf{x}_M, \mathbf{x}_{M+1})$ in a Wasserstein ball provided that the transportation cost is given by

$$d((\mathbf{x}_M, \mathbf{x}_{M+1}), (\mathbf{x}'_M, \mathbf{x}'_{M+1})) = \|\mathbf{x}_M - \mathbf{x}'_M\|_p + \kappa \mathbb{1}_{\{\mathbf{x}_{M+1} \neq \mathbf{x}'_{M+1}\}}$$

for $\kappa = \infty$, whenever $p = 1$, $p = 2$ or $p = \infty$, respectively. Thus, the MACS, Spectral and MARS regularization terms admit a distributionally robust interpretation.





















We use the stochastic proximal gradient descent algorithm of Section 3.4 to tune \mathbf{W}_M , including an additional momentum term with weight 0.9. As in (Krizhevsky et al., 2012), we split the training phase into 100 epochs, each corresponding to a complete pass through the training dataset in a random order. As the ALEXNET requires input images of size 244×244 , in each iteration we extract a random patch of 244×244 pixels from the current image and flip it horizontally at random. This procedure effectively augments the training dataset. The initial step size is set to 10^{-3} and then reduced by a factor of 10 after every 7 epochs. The algorithm terminates after 100 epochs. We preprocess the images in the validation and test datasets as in Scenario 1 and tune the regularization weights via the holdout method using the validation data. Table 3 reports the AP and mAP scores of the different classifiers that were tested. These results suggest that fine-tuning the last layer of a pre-trained neural network may improve classifier performance. We observe that the spectral norm regularizer, which has a distributionally robust interpretation, consistently outperforms almost all other methods. For further details on the experimental setup (such as the exact search grids for all hyperparameters) we refer to the code publicized on Github.

6.4. Generalization Bounds

The next experiment estimates the scaling behavior of the smallest Wasserstein radius that verifies the generalization bound (38) for the synthetic **threenorm** classification problem (Breiman, 1996). The experiment involves 1,000 simulation trials. In each trial we generate N training samples for some $N \in \{10, \dots, 90\} \cup \{100, \dots, 1,000\}$ as well as 10^5 test samples. Each sample $(\mathbf{x}, y) \in \mathbb{R}^{20} \times \{-1, 1\}$ is constructed as follows. The label y is drawn uniformly from $\{-1, 1\}$. If $y = -1$, then \mathbf{x} is drawn from a standard multivariate normal distribution shifted by (c, \dots, c) or $(-c, \dots, -c)$ with equal probabilities, where $c = 2/\sqrt{20}$. If $\hat{y} = 1$, on the other hand, then \mathbf{x} is drawn from a standard multivariate normal distribution shifted by $(c, -c, +c, \dots, -c)$.

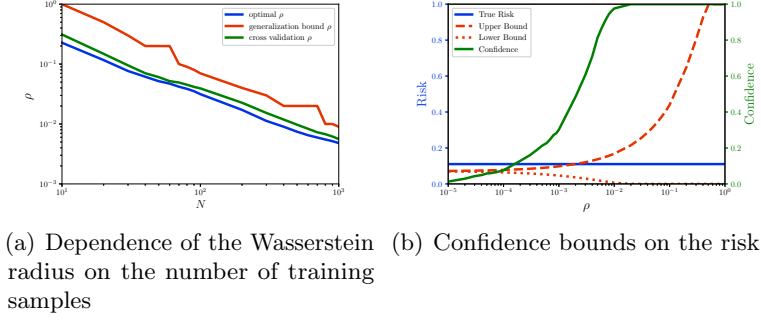
We now describe three different approaches to choose the Wasserstein radius ρ in the DRSVM (20) with transportation cost (16), where $\kappa = \infty$ and $\|\cdot\|$ represents the ∞ -norm on the input space. Throughout the experiment we use $P = \{b \cdot 10^{-e} : b \in \{1, \dots, 10\}, e \in \{1, \dots, 5\}\}$ as the search space for ρ . Approach 1 (‘cross validation’) calibrates the Wasserstein radius as before via 5-fold cross validation based solely on the N training samples. This approach reflects what would typically be done in practice.

Table 3: AP scores of different multi-label classifiers.

	Scenario 1		Scenario 2				
	RSVM	DRSVM	Lasso	Tikhonov	MACS	Spectral	MARS
	84.80	84.85	85.50	84.26	84.35	83.89	85.53
	78.54	78.49	76.18	76.55	76.37	75.67	76.22
	82.19	82.19	83.37	83.11	83.52	84.19	83.08
	79.55	79.56	77.65	78.04	77.70	78.92	77.53
	37.52	37.98	38.10	39.73	39.53	38.75	38.10
	72.05	72.03	70.05	69.43	69.17	70.28	69.77
	83.08	83.12	83.10	83.67	83.50	82.92	83.17
	80.60	80.56	79.87	79.86	80.07	79.85	79.82
	54.64	54.64	54.40	54.76	53.94	54.72	54.55
	47.82	53.13	52.06	52.10	51.62	55.54	51.19
	54.26	58.88	63.41	65.23	65.15	66.79	62.95
	75.81	75.81	76.92	77.39	77.26	76.54	76.97
	82.74	82.72	82.17	81.89	81.6	80.81	81.9
	72.69	72.88	74.70	75.41	74.76	76.48	74.31
	90.36	90.36	90.07	90.30	90.22	90.35	90.09
	50.22	51.90	50.20	50.18	50.27	51.39	50.20
	60.75	63.57	71.78	71.3	70.39	71.64	71.40
	56.85	56.98	52.15	54.36	54.65	55.12	51.94
	85.09	85.03	84.89	84.55	84.41	85.43	84.96
	69.02	69.08	64.73	65.85	65.63	64.26	64.48
mAP	69.92	70.69	70.56	70.90	70.71	71.20	70.40

Approaches 2 and 3 both solve (20) based on the empirical distribution induced by the N training samples and select the Wasserstein radius using the 10^5 test samples. Specifically, approach 2 (‘optimal’) chooses the Wasserstein radius that leads to the lowest test error, while approach 3 (‘generalization bound’) selects the smallest Wasserstein radius for which the optimal value of (20) exceeds the expected loss on the test samples in at least 95% of all trials, that is, it approximates the smallest Wasserstein radius that verifies the generalization bound (38) for $\eta = 5\%$. As the test samples are not available in the training phase, the last two approaches are not implementable in practice, and we merely study them to gain insights. Figure 2(a) visualizes all resulting Wasserstein radii as a function of N . Note that the radii obtained with the first two approaches are uncertain as they depend on a particular choice of the training samples. Figure 2(a) thus only shows their averages across all simulation trials. In contrast, the radii obtained with the third approach depend on the training sample sets of all 1,000 trials and are thus essentially deterministic.

We observe that the Wasserstein radii of all three approaches decay approximately as $1/\sqrt{N}$, which is in line with the theoretical generalization bound of Theorem 39. We expect


 Figure 2: Results of the `threenorm` classification problem.

this decay rate to be optimal because any faster decay would be in conflict with the central limit theorem. Note also that our results empirically confirm Theorem 39 even though we did not impose any restrictions on \mathbb{W} as dictated by Assumption 38. This suggests that Theorem 39 might remain valid under weaker conditions.

In the experiment underlying Figure 2(b), we first fix $\hat{\mathbf{w}}$ to an optimal solution of (20) for $\rho = 0.1$ and $N = 100$. Figure 2(b) shows the true risk $\mathcal{R}(\hat{\mathbf{w}})$ and its confidence bounds given by Theorem 42. As expected, for $\rho = 0$ the upper and lower bounds coincide with the empirical risk on the training data, which is a lower bound for the true risk on the test data due to over-fitting effects. As ρ increases, the confidence interval between the bounds widens and eventually covers the true risk. For instance, at $\rho \approx 0.009$ the confidence interval is given by $[0.008, 0.162]$ and contains the true risk with probability $1 - \eta = 95\%$.

6.5. Worst-Case Distributions

Consider again the 3-vs-8 classification problems from the MNIST database (LeCun et al., 1998) and fix \mathbf{w}^* to an optimal solution of the empirical hinge loss minimization problem. The goal of the last experiment is to evaluate the *worst-case* hinge loss of \mathbf{w}^* for different Wasserstein radii $\rho \in \{0, 0.01, 0.05, 0.1, 0.5, 1\}$ and label flipping costs $\kappa \in \{0, \infty\}$ and to investigate the corresponding worst-case distributions, which are computable by virtue of Theorem 20(i). As each input constitutes a vector of pixels intensities between zero and one, we impose support constraints of the form $\mathbf{C}\mathbf{x} \leq \mathbf{d}$ with $\mathbf{C} = [\mathbf{I}, -\mathbf{I}]^\top$ and $\mathbf{d} = [\mathbf{1}^\top, \mathbf{0}^\top]^\top$.

For illustrative purposes we only use the $N = 10$ first datapoints in the MNIST dataset as training samples. Each training sample $\hat{\mathbf{x}}_i$ corresponds to four discretization points ($\hat{\mathbf{x}}_i + \mathbf{q}_{ij}^+/\alpha_{ij}^+$ and $\hat{\mathbf{x}}_i + \mathbf{q}_{ij}^-/\alpha_{ij}^-$ for $j = 1, 2$) in the worst-case distribution obtained from (26). We observe that for every i exactly one out of these four points has probability $\frac{1}{N}$, while all others have probability 0. Figure 3 depicts only those 10 discretization points that have nonzero probability for a fixed ρ and κ . As expected, the perturbations of the training samples are more severe for larger Wasserstein radii. For $\kappa = \infty$ these scenarios must have the same labels as the corresponding training samples. For $\kappa = 0$, on the other hand, the labels can be flipped at no cost (flipped labels are indicated by red frames). Each scenario group shown in Figure 3 can thus be viewed as a worst-case training dataset for the corresponding Wasserstein radius and label flipping cost.

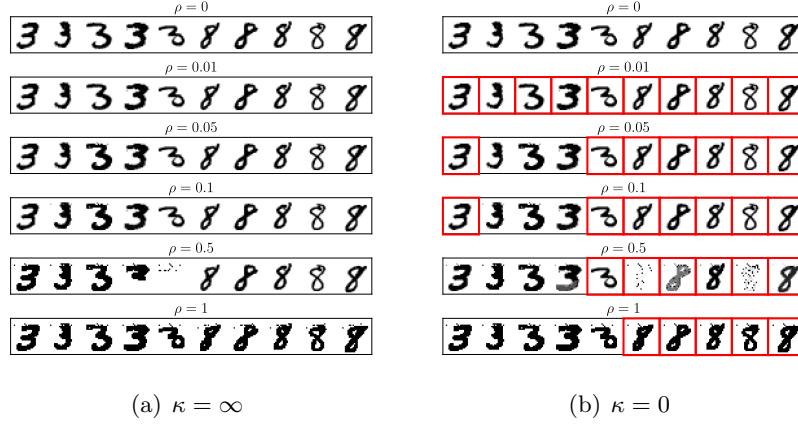


Figure 3: Discretization points (input images) of the worst-case distribution for different ρ and κ . Red frames indicate that the corresponding labels are flipped under the worst-case distribution.

Acknowledgments

We gratefully acknowledge financial support from the Swiss National Science Foundation under grants BSCGI0_157733 and P2EZP2_165264.

Appendix: Proofs

A.1. Proofs of Section 3

The proof of Theorem 4 requires three preparatory Lemmas. The first lemma is adapted from (Mohajerin Esfahani and Kuhn, 2018) and asserts that the worst-case expectation over a Wasserstein ball can be re-expressed as a classical robust optimization problem.

Lemma 45 (Robust reformulation) *Set $\hat{\xi}_i = (\hat{x}_i, \hat{y}_i)$ for all $i \leq N$. For any measurable integrand $I(\xi)$ that is bounded above by a Lipschitz continuous function we have*

$$\sup_{Q \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^Q[I(\xi)] = \inf_{\lambda \geq 0} \lambda \rho + \frac{1}{N} \sum_{i=1}^N \sup_{\xi \in \Xi} I(\xi) - \lambda d(\xi, \hat{\xi}_i). \quad (\text{A.1})$$

Proof By the definition of the Wasserstein ball we have

$$\sup_{Q \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^Q[I(\xi)] = \begin{cases} \sup_{\Pi} \int_{\Xi^2} I(\xi) \Pi(d\xi, d\xi') \\ \text{s.t.} & \Pi \text{ is a joint distribution of } \xi \\ & \text{and } \xi' \text{ with marginals } Q \text{ and } \hat{\mathbb{P}}_N \\ & \int_{\Xi^2} d(\xi, \xi') \Pi(d\xi, d\xi') \leq \rho \end{cases}$$

$$= \begin{cases} \sup_{\mathbf{Q}^i} & \frac{1}{N} \sum_{i=1}^N \int_{\Xi} I(\boldsymbol{\xi}) \mathbf{Q}^i(d\boldsymbol{\xi}) \\ \text{s.t.} & \int_{\Xi} \mathbf{Q}^i(d\boldsymbol{\xi}) = 1 \quad i \in [N] \\ & \frac{1}{N} \sum_{i=1}^N \int_{\Xi} d(\boldsymbol{\xi}, \hat{\boldsymbol{\xi}}_i) \mathbf{Q}^i(d\boldsymbol{\xi}) \leq \rho. \end{cases}$$

Note that the integral of $I(\boldsymbol{\xi})$ exists under every $\mathbf{Q} \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)$ because $I(\boldsymbol{\xi})$ admits a Lipschitz continuous majorant. The last equality in the above expression holds because the marginal distribution of $\boldsymbol{\xi}'$ is the uniform distribution on the training samples, which implies that Π is completely determined by the conditional distributions \mathbf{Q}^i of $\boldsymbol{\xi}$ given $\boldsymbol{\xi}' = \hat{\boldsymbol{\xi}}_i$, that is, $\Pi(d\boldsymbol{\xi}, d\boldsymbol{\xi}') = \frac{1}{N} \sum_{i=1}^N \delta_{\hat{\boldsymbol{\xi}}_i}(d\boldsymbol{\xi}') \mathbf{Q}^i(d\boldsymbol{\xi})$. The resulting generalized moment problem over the normalized measures \mathbf{Q}^i admits the semi-infinite dual

$$\sup_{\mathbf{Q} \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}}[I(\boldsymbol{\xi})] = \begin{cases} \inf_{\lambda, s_i} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & \sup_{\boldsymbol{\xi} \in \Xi} I(\boldsymbol{\xi}) - \lambda d(\boldsymbol{\xi}, \hat{\boldsymbol{\xi}}_i) \leq s_i \quad i \in [N] \\ & \lambda \geq 0. \end{cases}$$

Strong duality holds for any $\rho > 0$ due to (Shapiro, 2001, Proposition 3.4). The claim then follows by eliminating s_i . \blacksquare

Lemma 46 *For any $a \in \mathbb{R}$, $\boldsymbol{\beta}, \hat{\boldsymbol{\zeta}} \in \mathbb{R}^d$, $\gamma \in \mathbb{R}_+$ and $\boldsymbol{\zeta} \in \mathcal{Z}$, where $\mathcal{Z} \subseteq \mathbb{R}^d$ is a closed convex set, we have*

$$\sup_{\boldsymbol{\zeta} \in \mathcal{Z}} a \langle \boldsymbol{\beta}, \boldsymbol{\zeta} \rangle - \gamma \|\boldsymbol{\zeta} - \hat{\boldsymbol{\zeta}}\| = \begin{cases} \inf_{\mathbf{p}} & S_{\mathcal{Z}}(a\boldsymbol{\beta} - \mathbf{p}) + \langle \mathbf{p}, \hat{\boldsymbol{\zeta}} \rangle \\ \text{s.t.} & \|\mathbf{p}\|_* \leq \gamma. \end{cases}$$

Proof We have

$$\begin{aligned} \sup_{\boldsymbol{\zeta} \in \mathcal{Z}} a \langle \boldsymbol{\beta}, \boldsymbol{\zeta} \rangle - \gamma \|\boldsymbol{\zeta} - \hat{\boldsymbol{\zeta}}\| &= \sup_{\boldsymbol{\zeta} \in \mathcal{Z}} \inf_{\|\mathbf{p}\|_* \leq \gamma} a \langle \boldsymbol{\beta}, \boldsymbol{\zeta} \rangle - \langle \mathbf{p}, \boldsymbol{\zeta} - \hat{\boldsymbol{\zeta}} \rangle \\ &= \inf_{\|\mathbf{p}\|_* \leq \gamma} \sup_{\boldsymbol{\zeta} \in \mathcal{Z}} a \langle \boldsymbol{\beta}, \boldsymbol{\zeta} \rangle - \langle \mathbf{p}, \boldsymbol{\zeta} - \hat{\boldsymbol{\zeta}} \rangle \\ &= \inf_{\|\mathbf{p}\|_* \leq \gamma} \sup_{\boldsymbol{\zeta} \in \mathbb{R}^d} \langle a\boldsymbol{\beta} - \mathbf{p}, \boldsymbol{\zeta} \rangle - \delta_{\mathcal{Z}}(\boldsymbol{\zeta}) + \langle \mathbf{p}, \hat{\boldsymbol{\zeta}} \rangle \\ &= \inf_{\|\mathbf{p}\|_* \leq \gamma} S_{\mathcal{Z}}(a\boldsymbol{\beta} - \mathbf{p}) + \langle \mathbf{p}, \hat{\boldsymbol{\zeta}} \rangle, \end{aligned}$$

where the first equality follows from the definition of the dual norm, the second equality holds due to the minimax theorem (Bertsekas, 2009, Proposition 5.5.4), and the last equality holds because the support function $S_{\mathcal{Z}}$ is the conjugate of the indicator function $\delta_{\mathcal{Z}}$. Thus, the claim follows. \blacksquare

Lemma 47 *If $L(z)$ is a convex and Lipschitz continuous loss function, $\beta, \hat{\zeta} \in \mathbb{R}^d$ and $\gamma > 0$, then*

$$\sup_{\zeta \in \mathbb{R}^d} L(\langle \beta, \zeta \rangle) - \gamma \|\zeta - \hat{\zeta}\| = \begin{cases} L(\langle \beta, \hat{\zeta} \rangle) & \text{if } \text{lip}(L) \|\beta\|_* \leq \gamma \\ +\infty & \text{otherwise.} \end{cases}$$

Proof Note that $L(\langle \beta, \zeta \rangle) - \gamma \|\zeta - \hat{\zeta}\|$ constitutes a difference of convex functions and may thus be neither convex nor concave in ζ . In order to maximize this function, we re-write $I(\zeta) = L(\langle \beta, \zeta \rangle)$ as an upper envelope of infinitely many affine functions. To this end, we express the conjugate of $I(\zeta)$ as

$$I^*(z) = \sup_{\zeta} \langle z, \zeta \rangle - L(\langle \beta, \zeta \rangle) = \sup_{t, \zeta} \{ \langle z, \zeta \rangle - L(t) : t = \langle \beta, \zeta \rangle \} = \inf_{\theta} \{ L^*(\theta) : \theta \beta = z \},$$

where the last equality follows from strong Lagrangian duality, which holds because Slater's constraint qualification is trivially satisfied in the absence of inequality constraints (Bertsekas, 2009, Proposition 5.3.1). Defining $\Theta = \{\theta \in \mathbb{R} : L^*(\theta) < \infty\}$ as the effective domain of $L^*(\theta)$, we may then replace $\theta \in \mathbb{R}$ with $\theta \in \Theta$ in the last expression. As $I(\zeta)$ is convex and continuous, it coincides with its bi-conjugate, that is,

$$I(\zeta) = I^{**}(\zeta) = \sup_z \langle z, \zeta \rangle - I^*(z) = \sup_{\theta \in \Theta} \langle \theta \beta, \zeta \rangle - L^*(\theta).$$

In other words, we have represented $I(\zeta)$ as the upper envelope of infinitely many linear functions. Using this representation, we obtain

$$\begin{aligned} \sup_{\zeta} I(\zeta) - \gamma \|\zeta - \hat{\zeta}\| &= \sup_{\zeta} I^{**}(\zeta) - \gamma \|\zeta - \hat{\zeta}\| \\ &= \sup_{\zeta} \sup_{\theta \in \Theta} \langle \theta \beta, \zeta \rangle - L^*(\theta) - \gamma \|\zeta - \hat{\zeta}\| \\ &= \sup_{\theta \in \Theta} \sup_{\zeta} \inf_{\|p\|_* \leq \gamma} \theta \langle \beta, \zeta \rangle - L^*(\theta) - \langle p, \zeta - \hat{\zeta} \rangle \\ &= \sup_{\theta \in \Theta} \inf_{\|p\|_* \leq \gamma} \sup_{\zeta} \langle \theta \beta - p, \zeta \rangle - L^*(\theta) + \langle p, \hat{\zeta} \rangle, \end{aligned}$$

where the last equality holds due to (Bertsekas, 2009, Proposition 5.5.4). Evaluating the maximization over ζ yields

$$\begin{aligned} \sup_{\zeta} I(\zeta) - \gamma \|\zeta - \hat{\zeta}\| &= \sup_{\theta \in \Theta} \inf_{\|p\|_* \leq \gamma} \begin{cases} \langle p, \hat{\zeta} \rangle - L^*(\theta) & \text{if } p = \theta \beta \\ +\infty & \text{otherwise} \end{cases} \\ &= \sup_{\theta \in \Theta} \begin{cases} \langle \theta \beta, \hat{\zeta} \rangle - L^*(\theta) & \text{if } \|\theta \beta\|_* \leq \gamma \\ +\infty & \text{otherwise} \end{cases} \\ &= \begin{cases} \sup_{\theta \in \Theta} \theta \langle \beta, \hat{\zeta} \rangle - L^*(\theta) & \text{if } \sup_{\theta \in \Theta} \|\theta \beta\|_* \leq \gamma \\ +\infty & \text{otherwise} \end{cases} \\ &= \begin{cases} L(\langle \beta, \hat{\zeta} \rangle) & \text{if } \sup_{\theta \in \Theta} |\theta| \cdot \|\beta\|_* \leq \gamma \\ +\infty & \text{otherwise.} \end{cases} \end{aligned}$$

Thus, the claim follows by noting that $\sup_{\theta} \{|\theta| : L^*(\theta) < \infty\}$ represents the Lipschitz modulus of L . \blacksquare

Proof of Theorem 4 To prove assertion (i), we apply Lemma 45 to the integrand $f(\mathbf{x}, y) = L(\langle \mathbf{w}, \mathbf{x} \rangle - y)$ with $L(z) = \max_{j \leq J} \{a_j z + b_j\}$ to obtain

$$\begin{aligned} & \sup_{\mathbf{Q} \in \mathcal{B}_\rho(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] \\ &= \inf_{\lambda \geq 0} \lambda \rho + \frac{1}{N} \sum_{i=1}^N \sup_{(\mathbf{x}, y) \in \Xi} L(\langle \mathbf{w}, \mathbf{x} \rangle - y) - \lambda \|(\mathbf{x}, y) - (\widehat{\mathbf{x}}_i, \widehat{y}_i)\| \\ &= \inf_{\lambda \geq 0} \lambda \rho + \frac{1}{N} \sum_{i=1}^N \max_{j \in [J]} \sup_{(\mathbf{x}, y) \in \Xi} a_j (\langle \mathbf{w}, \mathbf{x} \rangle - y) + b_j - \lambda \|(\mathbf{x}, y) - (\widehat{\mathbf{x}}_i, \widehat{y}_i)\| \\ &= \begin{cases} \inf_{\lambda, \mathbf{p}_{ij}, u_{ij}} \lambda \rho + \sum_{i=1}^N \max_{j \leq J} S_{\Xi}(a_j \mathbf{w} - \mathbf{p}_{ij}, -a_j - u_{ij}) + \langle \mathbf{p}_{ij}, \widehat{\mathbf{x}}_i \rangle + u_{ij} \widehat{y}_i + b_j \\ \text{s.t.} \quad \|(\mathbf{p}_{ij}, u_{ij})\|_* \leq \lambda \quad i \in [N], j \in [J], \end{cases} \end{aligned}$$

where the last equality follows from Lemma 46. The claim now follows by introducing auxiliary epigraphical variables s_i for the max-terms in the objective function and by including \mathbf{w} as a decision variable.

To prove assertion (ii), we apply Lemma 45 to the integrand $f(\mathbf{x}, y) = L(\langle \mathbf{w}, \mathbf{x} \rangle - y)$, where L is a Lipschitz continuous convex loss function. Thus we find

$$\begin{aligned} \sup_{\mathbf{Q} \in \mathcal{B}_\rho(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}} [\ell_{\mathbf{w}}(\mathbf{x}, y)] &= \inf_{\lambda \geq 0} \lambda \rho + \sum_{i=1}^N \sup_{\mathbf{x}, y} L(\langle \mathbf{w}, \mathbf{x} \rangle - y) - \lambda \|(\mathbf{x}, y) - (\widehat{\mathbf{x}}_i, \widehat{y}_i)\| \\ &= \begin{cases} \inf_{\lambda} \lambda \rho + \frac{1}{N} \sum_{i=1}^N L(\langle \mathbf{w}, \widehat{\mathbf{x}}_i \rangle - \widehat{y}_i) \\ \text{s.t.} \quad \text{lip}(L)|\theta| \cdot \|(\mathbf{w}, -1)\|_* \leq \lambda, \end{cases} \end{aligned}$$

where the last equality uses Lemma 47. Next, we eliminate λ and include \mathbf{w} as a decision variable. \blacksquare

Proof of Corollary 5 Note that the Huber loss function $L(z)$ coincides with the inf-convolution of $\frac{1}{2}z^2$ and $\delta|z|$ and can thus be expressed as $L(z) = \min_{z_1} \frac{1}{2}z_1^2 + \delta|z - z_1|$. Moreover, the Lipschitz modulus of the Huber loss function is δ . The rest of the proof follows from Theorem 4(ii). \blacksquare

Proof of Corollary 6 Notice that the ϵ -insensitive loss function is a piecewise linear function with $J = 3$ pieces, see Section 2.1. By strong conic duality, the support function of $\Xi = \{(\mathbf{x}, y) \in \mathbb{R}^{n+1} : \mathbf{C}_1 \mathbf{x} + \mathbf{c}_2 y \preceq_{\mathcal{C}} \mathbf{d}\}$ can be re-expressed as

$$S_{\Xi}(z_1, z_2) = \sup_{\mathbf{x}, y} \{\langle z_1, \mathbf{x} \rangle + z_2 y : \mathbf{C}_1 \mathbf{x} + \mathbf{c}_2 y \preceq_{\mathcal{C}} \mathbf{d}\} = \inf_{\mathbf{q} \in \mathcal{C}^*} \left\{ \langle \mathbf{q}, \mathbf{d} \rangle : \mathbf{C}_1^{\top} \mathbf{q} = z_1, \mathbf{c}_2^{\top} \mathbf{q} = z_2 \right\}.$$

Strong duality holds because Ξ admits a Slater point. The rest of proof follows from Theorem 4(i). \blacksquare

Proof of Corollary 7 The pinball loss function is a piecewise linear function with $J = 2$ pieces, see Section 2.1. The rest of proof follows from the dual representation of the support function $S_{\Xi}(z_1, z_2)$, which is known from the proof of Corollary 6, and from Theorem 4(i). \blacksquare

The proof of Theorem 9 is based on the following preparatory lemma.

Lemma 48 *If $\mathcal{Z} \subseteq \mathbb{R}^d$ is a non-empty convex closed set, $\hat{\zeta} \in \mathcal{Z}$, $\beta \in \mathbb{R}^d$ and $\alpha, \gamma \geq 0$, then we have*

$$\inf_{\mathbf{p}} \alpha S_{\mathcal{Z}}(\beta - \mathbf{p}) + \alpha \langle \mathbf{p}, \hat{\zeta} \rangle + \gamma \|\mathbf{p}\|_* = \begin{cases} \sup_{\|\mathbf{q}\| \leq \gamma} \alpha \langle \beta, \hat{\zeta} \rangle + \langle \beta, \mathbf{q} \rangle \\ \text{s.t. } \hat{\zeta} + \mathbf{q}/\alpha \in \mathcal{Z}. \end{cases}$$

Proof If $\alpha = 0$, then the optimal values of both optimization problems vanish due to our conventions of extended arithmetic, and thus the claim trivially holds. If $\alpha > 0$, however, we have

$$\begin{aligned} \inf_{\mathbf{p}} \alpha S_{\mathcal{Z}}(\beta - \mathbf{p}) + \alpha \langle \mathbf{p}, \hat{\zeta} \rangle + \gamma \|\mathbf{p}\|_* &= \inf_{\mathbf{p}} \sup_{\|\mathbf{q}\| \leq \gamma} \alpha S_{\mathcal{Z}}(\beta - \mathbf{p}) + \alpha \langle \mathbf{p}, \hat{\zeta} \rangle + \langle \mathbf{p}, \mathbf{q} \rangle \\ &= \sup_{\|\mathbf{q}\| \leq \gamma} \inf_{\mathbf{p}} \alpha S_{\mathcal{Z}}(\beta - \mathbf{p}) + \langle \mathbf{p}, \alpha \hat{\zeta} + \mathbf{q} \rangle \\ &= \sup_{\|\mathbf{q}\| \leq \gamma} \inf_{\mathbf{z}} \alpha S_{\mathcal{Z}}(\mathbf{z}) + \langle \beta - \mathbf{z}, \alpha \hat{\zeta} + \mathbf{q} \rangle \\ &= \sup_{\|\mathbf{q}\| \leq \gamma} \alpha \langle \beta, \hat{\zeta} \rangle + \langle \beta, \mathbf{q} \rangle - \alpha \left(\sup_{\mathbf{z}} \langle \mathbf{z}, \hat{\zeta} + \mathbf{q}/\alpha \rangle - S_{\mathcal{Z}}(\mathbf{z}) \right) \\ &= \sup_{\|\mathbf{q}\| \leq \gamma} \alpha \langle \beta, \hat{\zeta} \rangle + \langle \beta, \mathbf{q} \rangle - \alpha \delta_{\mathcal{Z}}(\hat{\zeta} + \mathbf{q}/\alpha), \end{aligned}$$

where the first equality follows from the definition of the dual norm, the second equality exploits (Bertsekas, 2009, Proposition 5.5.4), and the last equality holds because, for any convex closed set, the indicator function is the conjugate of the support function. \blacksquare

Proof of Theorem 9 We first prove assertion (i). By Theorem 4(i), the worst-case expectation problem (13) constitutes a restriction of (5) where \mathbf{w} is fixed, and thus it coincides with the minimax problem

$$\begin{aligned} \inf_{\substack{\lambda, s_i \\ \mathbf{p}_{ij}, u_{ij}}} \sup_{\alpha_{ij} \geq 0, \gamma_{ij} \geq 0} \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i + \sum_{i=1}^N \sum_{j=1}^J \gamma_{ij} (\|(\mathbf{p}_{ij}, u_{ij})\|_* - \lambda) \\ + \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij} (S_{\Xi}(-a_j \mathbf{w} - \mathbf{p}_{ij}, a_j - u_{ij}) + b_j + \langle \mathbf{p}_{ij}, \hat{\mathbf{x}}_i \rangle + u_{ij} \hat{y}_i - s_i). \end{aligned}$$

The minimization and the maximization may be interchanged by strong duality, which holds because the convex program (5) satisfies Slater's constraint qualification for every fixed

\mathbf{w} (Bertsekas, 2009, Proposition 5.3.1). Indeed, note that S_Ξ is proper, convex and lower semi-continuous and appears in constraints that are always satisfiable because they involve a free decision variable. Thus, the above minimax problem is equivalent to

$$\left\{ \begin{array}{l} \sup_{\alpha_{ij}, \gamma_{ij}} \inf_{\mathbf{p}_{ij}, u_{ij}} \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij} (S_\Xi(a_j \mathbf{w} - \mathbf{p}_{ij}, a_j - u_{ij}) + b_j + \langle \mathbf{p}_{ij}, \hat{\mathbf{x}}_i \rangle + u_{ij} \hat{y}_i) \\ \quad + \sum_{i=1}^N \sum_{j=1}^J \gamma_{ij} \|\mathbf{p}_{ij}\|_* \\ \text{s.t.} \quad \sum_{i=1}^N \sum_{j=1}^J \gamma_{ij} = \rho \\ \quad \sum_{j=1}^J \alpha_{ij} = \frac{1}{N} \quad i \in [N] \\ \quad \alpha_{ij}, \gamma_{ij} \geq 0 \quad i \in [N], j \in [J]. \end{array} \right.$$

By Lemma 48, which applies because $(\hat{\mathbf{x}}_i, \hat{y}_i) \in \Xi$ for all $i \leq N$, the above dual problem simplifies to

$$\left\{ \begin{array}{l} \sup_{\substack{\alpha_{ij}, \gamma_{ij} \\ \mathbf{q}_{ij}, v_{ij}}} \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij} (a_j (\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \hat{y}_i) + b_j) + a_j (\langle \mathbf{w}, \mathbf{q}_{ij} \rangle - v_{ij}) \\ \text{s.t.} \quad \sum_{i=1}^N \sum_{j=1}^J \gamma_{ij} = \rho \\ \quad \sum_{j=1}^J \alpha_{ij} = \frac{1}{N} \quad i \in [N] \\ \quad \|(\mathbf{q}_{ij}, v_{ij})\| \leq \gamma_{ij} \quad i \in [N], j \in [J] \\ \quad (\hat{\mathbf{x}}_i - \mathbf{q}_{ij}/\alpha_{ij}, \hat{y}_i - v_{ij}/\alpha_{ij}) \in \Xi \quad i \in [N], j \in [J] \\ \quad \alpha_{ij}, \gamma_{ij} \geq 0 \quad i \in [N], j \in [J]. \end{array} \right.$$

Problem (14) is now obtained by eliminating the variables γ_{ij} and by substituting α_{ij} , \mathbf{q}_{ij} , and v_{ij} with α_{ij}/N , \mathbf{q}_{ij}/N , and v_{ij}/N , respectively.

As for assertion (ii), we first show that the discrete distribution \mathbf{Q}_γ belongs to the Wasserstein ball $\mathbb{B}_\rho(\hat{\mathbb{P}}_N)$ for all $\gamma \in (0, 1]$. Indeed, the Wasserstein distance between \mathbf{Q}_γ and $\hat{\mathbb{P}}_N$ amounts to

$$d(\mathbf{Q}_\gamma, \hat{\mathbb{P}}_N) \leq \frac{\gamma}{N} \left\| \left(\hat{\mathbf{x}}_1 + \frac{\rho N}{\gamma} \mathbf{x}^*, \hat{y}_1 + \frac{\rho N}{\gamma} y^* \right) - (\hat{\mathbf{x}}_1, \hat{y}_1) \right\| = \rho \|(\mathbf{x}^*, y^*)\| \leq \rho,$$

where the first inequality holds because the Wasserstein distance coincides with the optimal mass transportation cost, and the last inequality holds because the norm of (\mathbf{x}^*, y^*) is at most 1 by construction. Thus, $\mathbf{Q}_\gamma \in \mathbb{B}_\rho(\hat{\mathbb{P}}_N)$ for all $\gamma \in (0, 1]$. Denoting the optimal value

of (13) by $J^*(\mathbf{w})$ and using $\ell_{\mathbf{w}}(\mathbf{x}, y)$ as a shorthand for $L(\langle \mathbf{w}, \mathbf{x} \rangle - y)$, we find

$$\begin{aligned}
J^*(\mathbf{w}) &\geq \mathbb{E}^{\mathbb{Q}_\gamma}[\ell_{\mathbf{w}}(\mathbf{x}, y)] \\
&= \frac{1}{N} \sum_{i=1}^N \ell_{\mathbf{w}}(\hat{\mathbf{x}}_i, \hat{y}_i) - \frac{\gamma}{N} \ell_{\mathbf{w}}(\hat{\mathbf{x}}_1, \hat{y}_1) + \frac{\gamma}{N} \ell_{\mathbf{w}}(\hat{\mathbf{x}}_1 + \frac{\rho N}{\gamma} \mathbf{x}^*, \hat{y}_1 + \frac{\rho N}{\gamma} y^*) \\
&\geq \frac{1}{N} \sum_{i=1}^N \ell_{\mathbf{w}}(\hat{\mathbf{x}}_i, \hat{y}_i) - \frac{\gamma}{N} \ell_{\mathbf{w}}(\hat{\mathbf{x}}_1, \hat{y}_1) + \frac{\gamma}{N} (\langle \mathbf{w}, \mathbf{x}^* \rangle, \langle \hat{\mathbf{x}}_1 + \frac{\rho N}{\gamma} \mathbf{x}^*, \hat{y}_1 + \frac{\rho N}{\gamma} y^* \rangle) \\
&\quad - \ell_{\mathbf{w}}^*(\mathbf{x}, y) \quad \forall (\mathbf{x}, y) \in \mathbb{R}^{n+1},
\end{aligned}$$

where the last estimate follows from Fenchel's inequality. Setting $(\mathbf{x}, y) = \text{lip}(L)(\mathbf{w}, -1)$ we thus have

$$\begin{aligned}
J^*(\mathbf{w}) &\geq \lim_{\gamma \rightarrow 0^+} \frac{1}{N} \sum_{i=1}^N \ell_{\mathbf{w}}(\hat{\mathbf{x}}_i, \hat{y}_i) - \frac{\gamma}{N} \ell_{\mathbf{w}}(\hat{\mathbf{x}}_1, \hat{y}_1) + \frac{\gamma}{N} \text{lip}(L)(\langle \mathbf{w}, \hat{\mathbf{x}}_1 \rangle - \hat{y}_1) \\
&\quad + \rho \text{lip}(L)\|(\mathbf{w}, -1)\|_* - \frac{\gamma}{N} \ell_{\mathbf{w}}^*(\text{lip}(L)(\mathbf{w}, -1)) \\
&= \frac{1}{N} \sum_{i=1}^N \ell_{\mathbf{w}}(\hat{\mathbf{x}}_i, \hat{y}_i) + \rho \text{lip}(L)\|(\mathbf{w}, -1)\|_* = J^*(\mathbf{w}),
\end{aligned}$$

where the equality follows from Theorem 4(ii). The above reasoning implies that

$$\lim_{\gamma \rightarrow 0^+} \mathbb{E}^{\mathbb{Q}_\gamma}[\ell_{\mathbf{w}}(\mathbf{x}, y)] = J^*(\mathbf{w}),$$

and thus the claim follows. ■

Proof of Theorem 12 Assume first that the loss function L is convex piecewise linear, that is, $L(z) = \max_{j \in J} \{a_j z + b_j\}$. As $\Xi \in \mathbb{R}^{n+1}$, Theorem 9(i) implies that the worst-case expectation (13) is given by

$$\begin{aligned}
&\left\{ \begin{array}{ll} \sup_{\alpha_{ij}, \mathbf{q}_{ij}, v_{ij}} & \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij} [a_j (\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \hat{y}_i) + b_j] + a_j (\langle \mathbf{w}, \mathbf{q}_{ij} \rangle - v_{ij}) \\ \text{s.t.} & \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \|\mathbf{q}_{ij}, v_{ij}\| \leq \rho \\ & \sum_{j=1}^J \alpha_{ij} = 1 \quad i \in [N] \\ & \alpha_{ij} \geq 0 \quad i \in [N], j \in [J] \end{array} \right\} \\
&\geq \left\{ \begin{array}{ll} \sup_{\alpha_{ij}, \Delta \mathbf{x}_{ij}, \Delta y_{ij}} & \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij} [a_j (\langle \mathbf{w}, \hat{\mathbf{x}}_i + \Delta \mathbf{x}_{ij} \rangle - \hat{y}_i - \Delta y_{ij}) + b_j] \\ \text{s.t.} & \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij} \|\Delta \mathbf{x}_{ij}, \Delta y_{ij}\| \leq \rho \\ & \sum_{j=1}^J \alpha_{ij} = 1 \quad i \in [N] \\ & \alpha_{ij} \geq 0 \quad i \in [N], j \in [J] \end{array} \right\}
\end{aligned}$$

$$\geq \left\{ \begin{array}{ll} \sup_{\alpha_{ij}, \Delta \mathbf{x}_i, \Delta y_i} & \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij} [a_j(\langle \mathbf{w}, \hat{\mathbf{x}}_i + \Delta \mathbf{x}_i \rangle - \hat{y}_i - \Delta y_i) + b_j] \\ \text{s.t.} & \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij} \|(\Delta \mathbf{x}_i, \Delta y_i)\| \leq \rho \\ & \sum_{j=1}^J \alpha_{ij} = 1 \quad i \in [N] \\ & \alpha_{ij} \geq 0 \quad i \in [N], j \in [J]. \end{array} \right.$$

The first inequality holds because for any feasible solution $\{\alpha_{ij}, \Delta \mathbf{x}_{ij}, \Delta y_{ij}\}$ to the second optimization problem, the solution $\{\alpha_{ij}, \mathbf{q}_{ij}, v_{ij}\}$ with $\mathbf{q}_{ij} = \alpha_{ij} \Delta \mathbf{x}_{ij}$ and $v_{ij} = \alpha_{ij} \Delta y_{ij}$ is feasible in the first problem and attains the same objective value (conversely, note that the first problem admits feasible solutions with $\alpha_{ij} = 0$ and $\mathbf{q}_{ij} \neq \mathbf{0}$ that have no counterpart in the second problem). The second inequality follows from the restriction that $\Delta \mathbf{x}_{ij}$ and Δy_{ij} must be independent of j . It is easy to verify that the last optimization problem in the above expression is equivalent to (15) because $(\alpha_{i1}, \dots, \alpha_{iJ})$ ranges over a simplex for every $i \leq N$, and thus (15) provides a lower bound on (13).

Suppose now that Assumption 10 holds, and note that the worst-case loss (15) can be expressed as

$$\geq \left\{ \begin{array}{ll} \sup_{\Delta \mathbf{x}_i, \Delta y_i} & \frac{1}{N} \sum_{i=1}^N \max_{j \leq J} [a_j(\langle \mathbf{w}, \hat{\mathbf{x}}_i + \Delta \mathbf{x}_i \rangle - \hat{y}_i - \Delta y_i) + b_j] \\ \text{s.t.} & \frac{1}{N} \sum_{i=1}^N \|(\Delta \mathbf{x}_i, \Delta y_i)\| \leq \rho \end{array} \right. \\ \geq \left\{ \begin{array}{ll} \sup_{\Delta \mathbf{x}, \Delta y} & \frac{1}{N} \sum_{i \neq k} \max_{j \leq J} [a_j(\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \hat{y}_i) + b_j] + \frac{1}{N} \max_{j \leq J} [a_j(\langle \mathbf{w}, \hat{\mathbf{x}}_k + \Delta \mathbf{x} \rangle - \hat{y}_k - \Delta y) + b_j] \\ \text{s.t.} & \frac{1}{N} \|(\Delta \mathbf{x}, \Delta y)\| \leq \rho. \end{array} \right.$$

The above inequality follows from setting $\Delta \mathbf{x}_i = 0$ and $\Delta y_i = 0$ for all $i \neq k$, where $(\hat{\mathbf{x}}_k, \hat{y}_k)$ is a training sample satisfying $|L'(\langle \mathbf{w}, \hat{\mathbf{x}}_k \rangle - \hat{y}_k)| = \text{lip}(L)$, which exists due to Assumption 10. The last expression equals

$$\begin{aligned} &= \frac{1}{N} \sum_{i \neq k} \max_{j \leq J} [a_j(\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \hat{y}_i) + b_j] + \max_{j \leq J} \left\{ \begin{array}{ll} \sup_{\Delta \mathbf{x}, \Delta y} & \frac{1}{N} [a_j(\langle \mathbf{w}, \hat{\mathbf{x}}_k \rangle - \hat{y}_k) + b_j] \\ & + \frac{1}{N} [a_j(\langle \mathbf{w}, \Delta \mathbf{x} \rangle - \Delta y)] \\ \text{s.t.} & \frac{1}{N} \|(\Delta \mathbf{x}, \Delta y)\| \leq \rho \end{array} \right. \\ &= \frac{1}{N} \sum_{i \neq k} \max_{j \leq J} [a_j(\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \hat{y}_i) + b_j] + \frac{1}{N} \max_{j \leq J} \left[a_j(\langle \mathbf{w}, \hat{\mathbf{x}}_k \rangle - \hat{y}_k) + b_j + \rho N \|(\mathbf{w}, -1)\|_* |a_j| \right] \\ &= \frac{1}{N} \sum_{i \neq k} \max_{j \leq J} [a_j(\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \hat{y}_i) + b_j] + \frac{1}{N} \max_{j \leq J} [a_j(\langle \mathbf{w}, \hat{\mathbf{x}}_k \rangle - \hat{y}_k) + b_j] + \max_{j \leq J} \rho \|(\mathbf{w}, -1)\|_* |a_j| \end{aligned}$$

$$= \frac{1}{N} \sum_{i=1}^N L(\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \hat{y}_i) + \max_{j \leq J} \rho \|(\mathbf{w}, -1)\|_* |a_j|,$$

where the penultimate equality holds because $\langle \mathbf{w}, \hat{\mathbf{x}}_k \rangle - \hat{y}_k$ resides on the steepest linear piece of the loss function L by virtue of Assumption 10. The claim then follows from Theorem 4(ii) because $\text{lip}(L) = \max_{j \leq J} |a_j|$. Note that generic convex Lipschitz continuous loss functions can be uniformly approximated as closely as desired with convex piecewise linear functions. Thus, the above argument extends directly to generic convex Lipschitz continuous loss functions. Details are omitted for brevity. \blacksquare

Proof of Theorem 14 To prove assertion (i), we apply Lemma 45 to problem (4) with the transportation distance (16), where $\boldsymbol{\xi} = (\mathbf{x}, y)$ and $\Xi = \mathbb{X} \times \{-1, +1\}$. Thus, we obtain

$$\begin{aligned} & \sup_{\mathbf{Q} \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] \\ &= \begin{cases} \inf_{\lambda, s_i} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & \sup_{(\mathbf{x}, y) \in \Xi} L(y \langle \mathbf{w}, \mathbf{x} \rangle) - \lambda (\|\mathbf{x} - \hat{\mathbf{x}}_i\| + \frac{\kappa}{2} |y - \hat{y}_i|) \leq s_i \quad i \in [N] \\ & \lambda \geq 0. \end{cases} \\ &= \begin{cases} \inf_{\lambda, s_i} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & \sup_{\mathbf{x} \in \mathbb{X}} L(\hat{y}_i \langle \mathbf{w}, \mathbf{x} \rangle) - \lambda \|\mathbf{x} - \hat{\mathbf{x}}_i\| \leq s_i \quad i \in [N] \\ & \sup_{\mathbf{x} \in \mathbb{X}} L(-\hat{y}_i \langle \mathbf{w}, \mathbf{x} \rangle) - \lambda \|\mathbf{x} - \hat{\mathbf{x}}_i\| - \kappa \lambda \leq s_i \quad i \in [N] \\ & \lambda \geq 0 \end{cases} \\ &= \begin{cases} \inf_{\lambda, s_i} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & \sup_{\mathbf{x} \in \mathbb{X}} a_j \hat{y}_i \langle \mathbf{w}, \mathbf{x} \rangle + b_j - \lambda \|\mathbf{x} - \hat{\mathbf{x}}_i\| \leq s_i \quad i \in [N], j \in [J] \\ & \sup_{\mathbf{x} \in \mathbb{X}} -a_j \hat{y}_i \langle \mathbf{w}, \mathbf{x} \rangle + b_j - \lambda \|\mathbf{x} - \hat{\mathbf{x}}_i\| - \kappa \lambda \leq s_i \quad i \in [N], j \in [J] \\ & \lambda \geq 0, \end{cases} \end{aligned}$$

where the second equality holds because, for every i , y can be either equal to \hat{y}_i or to $-\hat{y}_i$. Reformulating the constraints using Lemma 46 and including \mathbf{w} as a decision variable then yields (17).

When $\mathbb{X} = \mathbb{R}^n$ and L is Lipschitz continuous, we can use similar arguments as above to prove that

$$\sup_{\mathbf{Q} \in \mathcal{B}_\rho(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}}[\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] = \begin{cases} \inf_{\lambda, s_i} & \lambda\rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & \sup_{\mathbf{x} \in \mathbb{R}^n} L(\widehat{y}_i \langle \mathbf{w}, \mathbf{x} \rangle) - \lambda \|\mathbf{x} - \widehat{\mathbf{x}}_i\| \leq s_i & i \in [N] \\ & \sup_{\mathbf{x} \in \mathbb{R}^n} L(-\widehat{y}_i \langle \mathbf{w}, \mathbf{x} \rangle) - \lambda \|\mathbf{x} - \widehat{\mathbf{x}}_i\| - \kappa\lambda \leq s_i & i \in [N] \\ & \lambda \geq 0. \end{cases}$$

Applying Lemma 47 to the subordinate maximization problems in the constraints yields

$$\sup_{\mathbf{Q} \in \mathcal{B}_\rho(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}}[\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] = \begin{cases} \inf_{\lambda, s_i} & \lambda\rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & L(\widehat{y}_i \langle \mathbf{w}, \widehat{\mathbf{x}}_i \rangle) \leq s_i & i \in [N], j \in [J] \\ & L(-\widehat{y}_i \langle \mathbf{w}, \widehat{\mathbf{x}}_i \rangle) - \kappa\lambda \leq s_i & i \in [N], j \in [J] \\ & \sup_{\theta \in \Theta} |\theta| \cdot \|\mathbf{w}\|_* \leq \lambda. \end{cases}$$

Thus, assertion (ii) follows by recalling that $\text{lip}(L) = \sup_{\theta} \{|\theta| : L^*(\theta) < \infty\} = \sup_{\theta \in \Theta} |\theta|$ and by including \mathbf{w} as a decision variable. \blacksquare

Proof of Corollary 15 Notice that the hinge loss function is piecewise linear with $J = 2$ pieces, see Section 2.1. Moreover, by strong conic duality the support function of \mathbb{X} can be re-expressed as

$$S_{\mathbb{X}}(\mathbf{z}) = \sup_{\mathbf{x}} \{\langle \mathbf{z}, \mathbf{x} \rangle : \mathbf{C}\mathbf{x} \preceq_{\mathcal{C}} \mathbf{d}\} = \inf_{\mathbf{q} \in \mathcal{C}^*} \left\{ \langle \mathbf{q}, \mathbf{d} \rangle : \mathbf{C}^\top \mathbf{q} = \mathbf{z} \right\}.$$

Strong duality holds because \mathbb{X} admits a Slater point. The proof thus follows from Theorem 14(i). \blacksquare

Proof of Corollary 16 The smooth hinge loss $L(z)$ coincides with the inf-convolution of $\frac{1}{2}z^2$ and $\max\{0, 1 - z\}$ and can thus be expressed as $L(z) = \min_{z_1} \frac{1}{2}z_1^2 + \max\{0, 1 - z - z_1\}$. Moreover, the Lipschitz modulus of the smooth hinge loss function is 1. The proof thus follows from Theorem 14(ii). \blacksquare

Proof of Corollary 17 Notice that the logloss function is convex and has Lipschitz modulus 1, see Section 2.1. The rest of proof follows from Theorem 14(ii). For more details see (Shafieezadeh-Abadeh et al., 2015). \blacksquare

Proof of Theorem 20 We first prove assertion (i). By Theorem 14(i), the worst-case expectation problem (25) constitutes a restriction of (17) where \mathbf{w} is fixed, and thus it is

equivalent to the minimax problem

$$\begin{aligned}
\inf_{\substack{\lambda, s_i \\ \mathbf{p}_{ij}^+, \mathbf{p}_{ij}^-}} \sup_{\substack{\alpha_{ij}^+ \geq 0, \gamma_{ij}^+ \geq 0 \\ \alpha_{ij}^- \geq 0, \gamma_{ij}^- \geq 0}} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i + \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij}^+ (S_{\mathbb{X}}(a_j \hat{y}_i \mathbf{w} - \mathbf{p}_{ij}^+) + b_j + \langle \mathbf{p}_{ij}^+, \hat{\mathbf{x}}_i \rangle - s_i) \\
& + \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij}^- (S_{\mathbb{X}}(-a_j \hat{y}_i \mathbf{w} - \mathbf{p}_{ij}^-) + b_j + \langle \mathbf{p}_{ij}^-, \hat{\mathbf{x}}_i \rangle - \kappa \lambda - s_i) \\
& + \sum_{i=1}^N \sum_{j=1}^J \gamma_{ij}^+ (\|\mathbf{p}_{ij}^+\|_* - \lambda) + \sum_{i=1}^N \sum_{j=1}^J \gamma_{ij}^- (\|\mathbf{p}_{ij}^-\|_* - \lambda).
\end{aligned}$$

The minimization and the maximization may be interchanged by strong Lagrangian duality, which holds because the convex program (17) satisfies Slater's constraint qualification for any fixed \mathbf{w} (Bertsekas, 2009, Proposition 5.3.1). Indeed, note that $S_{\mathbb{X}}$ is proper, convex and lower semi-continuous and appears in constraints that are always satisfiable because they involve a free decision variable. Thus, problem (25) is equivalent to

$$\left\{ \begin{array}{l} \sup_{\substack{\alpha_{ij}^+, \gamma_{ij}^+ \\ \alpha_{ij}^-, \gamma_{ij}^-}} \inf_{\mathbf{p}_{ij}^+, \mathbf{p}_{ij}^-} \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij}^+ (S_{\mathbb{X}}(a_j \hat{y}_i \mathbf{w} - \mathbf{p}_{ij}^+) + b_j + \langle \mathbf{p}_{ij}^+, \hat{\mathbf{x}}_i \rangle) + \gamma_{ij}^+ \|\mathbf{p}_{ij}^+\|_* \\ \quad + \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij}^- (S_{\mathbb{X}}(-a_j \hat{y}_i \mathbf{w} - \mathbf{p}_{ij}^-) + b_j + \langle \mathbf{p}_{ij}^-, \hat{\mathbf{x}}_i \rangle) + \gamma_{ij}^- \|\mathbf{p}_{ij}^-\|_* \\ \text{s.t.} \quad \sum_{i=1}^N \sum_{j=1}^J \gamma_{ij}^+ + \gamma_{ij}^- + \kappa \alpha_{ij}^- = \rho \\ \quad \sum_{j=1}^J \alpha_{ij}^+ + \alpha_{ij}^- = \frac{1}{N} \quad i \in [N] \\ \quad \alpha_{ij}^+, \alpha_{ij}^-, \gamma_{ij}^+, \gamma_{ij}^- \geq 0 \quad i \in [N], j \in [J]. \end{array} \right.$$

By Lemma 48, the above dual problem simplifies to

$$\left\{ \begin{array}{l} \sup_{\substack{\alpha_{ij}^+, \gamma_{ij}^+, \mathbf{q}_{ij}^+ \\ \alpha_{ij}^-, \gamma_{ij}^-, \mathbf{q}_{ij}^-}} \sum_{i=1}^N \sum_{j=1}^J (\alpha_{ij}^+ - \alpha_{ij}^-) a_j \hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle + a_j \hat{y}_i \langle \mathbf{w}, \mathbf{q}_{ij}^+ - \mathbf{q}_{ij}^- \rangle + \sum_{j=1}^J b_j \\ \text{s.t.} \quad \sum_{i=1}^N \sum_{j=1}^J \gamma_{ij}^+ + \gamma_{ij}^- + \kappa \alpha_{ij}^- = \rho \\ \quad \sum_{j=1}^J \alpha_{ij}^+ + \alpha_{ij}^- = \frac{1}{N} \quad i \in [N] \\ \quad \|\mathbf{q}_{ij}^+\| \leq \gamma_{ij}^+, \quad \|\mathbf{q}_{ij}^-\| \leq \gamma_{ij}^- \quad i \in [N], j \in [J] \\ \quad \hat{\mathbf{x}}_i + \mathbf{q}_{ij}^+ / \alpha_{ij}^+ \in \mathbb{X}, \quad \hat{\mathbf{x}}_i + \mathbf{q}_{ij}^- / \alpha_{ij}^- \in \mathbb{X} \quad i \in [N], j \in [J] \\ \quad \alpha_{ij}^+, \alpha_{ij}^-, \gamma_{ij}^+, \gamma_{ij}^- \geq 0 \quad i \in [N], j \in [J]. \end{array} \right.$$

Problem (26) is now obtained by eliminating γ_{ij}^+ and γ_{ij}^- and by substituting α_{ij}^+ , α_{ij}^- , \mathbf{q}_{ij}^+ , and \mathbf{q}_{ij}^- with α_{ij}^+/N , α_{ij}^-/N , \mathbf{q}_{ij}^+/N , and \mathbf{q}_{ij}^-/N , respectively.

As for assertion (ii), we use L_i^+ and L_i^- to abbreviate $L(\widehat{y}_i\langle \mathbf{w}, \widehat{\mathbf{x}}_i \rangle)$ and $L(-\widehat{y}_i\langle \mathbf{w}, \widehat{\mathbf{x}}_i \rangle)$, respectively. By Theorem 14(ii), we have

$$\begin{aligned} \sup_{\mathbf{Q} \in \mathbb{B}_\rho(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}}[L(y\langle \mathbf{w}, \mathbf{x} \rangle)] &= \begin{cases} \inf_{\mathbf{w}, \lambda, s_i} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & L_i^+ \leq s_i & i \in [N] \\ & L_i^- - \kappa \lambda \leq s_i & i \in [N] \\ & \text{lip}(L) \|\mathbf{w}\|_* \leq \lambda \end{cases} \\ &= \begin{cases} \sup_{\alpha_i^+, \alpha_i^-, \theta} & \text{lip}(L) \|\mathbf{w}\|_* \theta + \sum_{i=1}^N \alpha_i^+ L_i^+ + \alpha_i^- L_i^- \\ \text{s.t.} & \alpha_i^+ + \alpha_i^- = \frac{1}{N} & i \in [N] \\ & \theta + \kappa \sum_{i=1}^N \alpha_i^- = \rho \\ & \alpha_i^+ \geq 0, \alpha_i^- \geq 0 & i \in [N] \\ & \theta \geq 0, \end{cases} \end{aligned}$$

where the second equality follows from strong linear programming duality, which holds because the primal problem is manifestly feasible. Eliminating the first constraint and replacing α_i^- with α_i/N and α_i^+ with $(1 - \alpha_i)/N$ allows us to reformulate the dual linear program as

$$\begin{cases} \sup_{\alpha_i, \theta} & \text{lip}(L) \|\mathbf{w}\|_* \theta + \frac{1}{N} \sum_{i=1}^N (1 - \alpha_i) L_i^+ + \alpha_i L_i^- \\ \text{s.t.} & \theta + \frac{\kappa}{N} \sum_{i=1}^N \alpha_i = \rho \\ & 0 \leq \alpha_i \leq 1 & i \in [N] \\ & \theta \geq 0. \end{cases}$$

Thus, the worst-case expectation (25) coincides with the optimal value of (27) for $\gamma = 0$. Next let $(\alpha_i^*(\gamma), \theta^*(\gamma))$ be an optimal solution of (27) for $\gamma \geq 0$, and define \mathbf{Q}_γ as in the theorem statement. Note that (27) is infeasible for $\gamma > \rho$. Moreover, note that the atoms of \mathbf{Q}_γ have non-negative probabilities if $\eta(\gamma) \in [0, 1]$, which holds whenever $\gamma \in [0, 1]$. We thus focus on parameter values $\gamma \in [0, \min\{\rho, 1\}]$. By construction, the Wasserstein distance between \mathbf{Q}_γ and the empirical distribution satisfies

$$\begin{aligned} d(\mathbf{Q}_\gamma, \widehat{\mathbb{P}}_N) &\leq \frac{\kappa}{N} \sum_{i=1}^N \alpha_i^*(\gamma) - \frac{\eta(\gamma)\kappa}{N} \alpha_1^*(\gamma) + \frac{\eta(\gamma)}{N} d((\widehat{\mathbf{x}}_1 + \frac{\theta^*(\gamma)N}{\eta(\gamma)} \mathbf{x}^*, \widehat{y}_1) - (\widehat{\mathbf{x}}_1, \widehat{y}_1)) \\ &\leq \rho - \gamma - \theta^*(\gamma) + \theta^*(\gamma) \|\mathbf{x}^*\| \leq \rho, \end{aligned}$$

where the first inequality holds because the Wasserstein distance is defined as the minimum cost of moving \mathbf{Q}_γ to $\widehat{\mathbb{P}}_N$, the second inequality follows from the feasibility of $(\alpha_i^*(\gamma), \theta^*(\gamma))$ in (27) and the non-negativity of $\eta(\gamma)$, and the last inequality holds because $\|\mathbf{x}^*\| \leq 1$,

$\theta^*(\gamma) \geq 0$ and $\gamma \geq 0$. Thus, $\mathbb{Q}_\gamma \in \mathbb{B}_\rho(\widehat{\mathbb{P}}_N)$ for all $\gamma \in [0, \min\{\rho, 1\}]$. Denoting the optimal value of (25) by $J^*(\mathbf{w})$, we find

$$\begin{aligned}
J^*(\mathbf{w}) &\geq \mathbb{E}^{\mathbb{Q}_\gamma} [L(y\langle \mathbf{w}, \mathbf{x} \rangle)] \\
&= \frac{1}{N} \left(\sum_{i=1}^N (1 - \alpha_i^*(\gamma)) L_i^+ + \alpha_i^*(\gamma) L_i^- \right) - \frac{\eta(\gamma)}{N} ((1 - \alpha_1^*(\gamma)) L_1^+ + \alpha_1^*(\gamma) L_1^-) \\
&\quad + \frac{\eta(\gamma)}{N} L(\widehat{y}_1 \langle \mathbf{w}, \widehat{\mathbf{x}}_1 + \frac{\theta^*(\gamma)N}{\eta(\gamma)} \mathbf{x}^* \rangle) \\
&\geq \frac{1}{N} \left(\sum_{i=1}^N (1 - \alpha_i^*(\gamma)) L_i^+ + \alpha_i^*(\gamma) L_i^- \right) - \frac{\eta(\gamma)}{N} ((1 - \alpha_1^*(\gamma)) L_1^+ + \alpha_1^*(\gamma) L_1^-) \\
&\quad + \frac{\eta(\gamma)}{N} \left(\langle \mathbf{x}, \widehat{\mathbf{x}}_1 + \frac{\theta^*(\gamma)N}{\eta(\gamma)} \mathbf{x}^* \rangle - L^*(\widehat{y}_1 \langle \mathbf{w}, \mathbf{x} \rangle) \right) \quad \forall \mathbf{x} \in \mathbb{R}^n,
\end{aligned}$$

where the last estimate follows from Fenchel's inequality. Setting $\mathbf{x} = \text{lip}(L)\mathbf{w}$ and driving γ to zero yields

$$\begin{aligned}
J^*(\mathbf{w}) &\geq \lim_{\gamma \rightarrow 0^+} \frac{1}{N} \left(\sum_{i=1}^N (1 - \alpha_i^*(\gamma)) L_i^+ + \alpha_i^*(\gamma) L_i^- \right) - \frac{\eta(\gamma)}{N} ((1 - \alpha_1^*(\gamma)) L_1^+ + \alpha_1^*(\gamma) L_1^-) \\
&\quad + \frac{\eta(\gamma)}{N} (\text{lip}(L) \langle \mathbf{w}, \widehat{\mathbf{x}}_1 \rangle - L^*(\widehat{y}_1 \text{lip}(L) \langle \mathbf{w}, \mathbf{w} \rangle)) + \text{lip}(L) \|\mathbf{w}\|_* \cdot \theta^*(\gamma) \\
&= \lim_{\gamma \rightarrow 0^+} \frac{1}{N} \left(\sum_{i=1}^N (1 - \alpha_i^*(\gamma)) L_i^+ + \alpha_i^*(\gamma) L_i^- \right) + \text{lip}(L) \|\mathbf{w}\|_* \cdot \theta^*(\gamma) = J^*(\mathbf{w}),
\end{aligned}$$

where the first equality follows from the observation that $\eta(\gamma) \in [0, \gamma]$, which implies that $\eta(\gamma)$ converges to zero as γ tends to zero. The second equality holds because the optimal value of (27) is concave and non-increasing and—*a fortiori*—continuous in $\gamma \in [0, \min\{\rho, 1\}]$ and because $J^*(\mathbf{w})$ coincides with the optimal value of (27) when $\gamma = 0$. The above reasoning implies that $\lim_{\gamma \rightarrow 0^+} \mathbb{E}^{\mathbb{Q}_\gamma} [L(y\langle \mathbf{w}, \mathbf{x} \rangle)] = J^*(\mathbf{w})$, and thus the claim follows. \blacksquare

Proof of Theorem 23 Assume first that the loss function L is convex piecewise linear, that is, $L(z) = \max_{j \in J} \{a_j z + b_j\}$. As $\mathbb{X} = \mathbb{R}^n$ and $\kappa = \infty$, Theorem 20(i) implies that (25) can be expressed as

$$\left\{ \begin{array}{ll} \sup_{\alpha_{ij}^+, \mathbf{q}_{ij}^+} & \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij}^+ a_j \widehat{y}_i \langle \mathbf{w}, \widehat{\mathbf{x}}_i \rangle + a_j \widehat{y}_i \langle \mathbf{w}, \mathbf{q}_{ij}^+ \rangle + \sum_{j=1}^J b_j \\ \text{s.t.} & \sum_{i=1}^N \sum_{j=1}^J \|\mathbf{q}_{ij}^+\| \leq N\rho \\ & \sum_{j=1}^J \alpha_{ij}^+ = 1 \quad i \in [N] \\ & \alpha_{ij}^+ \geq 0 \quad i \in [N], j \in [J] \end{array} \right.$$

$$\begin{aligned}
 &\geq \left\{ \begin{array}{l} \sup_{\alpha_{ij}, \Delta \mathbf{x}_{ij}} \quad \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij} a_j \widehat{y}_i \langle \mathbf{w}, \widehat{\mathbf{x}}_i + \Delta \mathbf{x}_{ij} \rangle + \sum_{j=1}^J b_j \\ \text{s.t.} \quad \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij} \|\Delta \mathbf{x}_{ij}\| \leq N\rho \\ \sum_{j=1}^J \alpha_{ij} = 1 \quad i \in [N] \\ \alpha_{ij} \geq 0 \quad i \in [N], j \in [J] \end{array} \right. \\
 &\geq \left\{ \begin{array}{l} \sup_{\alpha_{ij}, \Delta \mathbf{x}_i} \quad \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij} a_j \widehat{y}_i \langle \mathbf{w}, \widehat{\mathbf{x}}_i + \Delta \mathbf{x}_i \rangle + \sum_{j=1}^J b_j \\ \text{s.t.} \quad \sum_{i=1}^N \sum_{j=1}^J \alpha_{ij} \|\Delta \mathbf{x}_i\| \leq N\rho \\ \sum_{j=1}^J \alpha_{ij} = 1 \quad i \in [N] \\ \alpha_{ij} \geq 0 \quad i \in [N], j \in [J]. \end{array} \right.
 \end{aligned}$$

The first optimization problem constitutes a special case of (26). Indeed, as $\kappa = \infty$, the first constraint in (26) implies that $\alpha_{ij}^- = 0$, which in turn implies via the fourth constraint and our conventions of extended arithmetics that $\mathbf{q}_{ij}^- = \mathbf{0}$. The first inequality in the above expression holds because for any feasible solution $\{\alpha_{ij}, \Delta \mathbf{x}_{ij}\}$ to the second problem, the solution $\{\alpha_{ij}^+, \mathbf{q}_{ij}^+\}$ with $\mathbf{q}_{ij}^+ = \alpha_{ij}^+ \Delta \mathbf{x}_{ij}$ and $\alpha_{ij} = \alpha_{ij}^+$ for all $i \leq N$ and $j \leq J$ is feasible in the first problem and attains the same objective value. The second inequality in the above expression follows from the restriction that $\Delta \mathbf{x}_{ij}$ must be independent of j . It is easy to verify that the last optimization problem is equivalent to (28) because $(\alpha_{i1}, \dots, \alpha_{iJ})$ ranges over a simplex for every $i \leq N$, and thus (28) provides a lower bound on (25).

Suppose now that Assumption 21 holds, and note that (28) can be expressed as

$$\begin{aligned}
 &\left\{ \begin{array}{l} \sup_{\Delta \mathbf{x}_i, \Delta y_i} \quad \frac{1}{N} \sum_{i=1}^N \max_{j \leq J} [a_j \widehat{y}_i \langle \mathbf{w}, \widehat{\mathbf{x}}_i + \Delta \mathbf{x}_i \rangle + b_j] \\ \text{s.t.} \quad \frac{1}{N} \sum_{i=1}^N \|\Delta \mathbf{x}_i\| \leq \rho \end{array} \right. \\
 &\geq \left\{ \begin{array}{l} \sup_{\Delta \mathbf{x}, \Delta y} \quad \frac{1}{N} \sum_{i \neq k} \max_{j \leq J} [a_j \widehat{y}_i \langle \mathbf{w}, \widehat{\mathbf{x}}_i \rangle + b_j] + \frac{1}{N} \max_{j \leq J} [a_j \widehat{y}_k \langle \mathbf{w}, \widehat{\mathbf{x}}_k + \Delta \mathbf{x} \rangle + b_j] \\ \text{s.t.} \quad \frac{1}{N} \|\Delta \mathbf{x}\| \leq \rho. \end{array} \right.
 \end{aligned}$$

The above inequality follows from setting $\Delta \mathbf{x}_i = 0$ and $\Delta y_i = 0$ for all $i \neq k$, where $(\widehat{\mathbf{x}}_k, \widehat{y}_k)$ is a training sample satisfying $|L'(\widehat{y}_k \langle \mathbf{w}, \widehat{\mathbf{x}}_k \rangle)| = \text{lip}(L)$, which exists due to Assumption 21. The last expression equals

$$\frac{1}{N} \sum_{i \neq k} \max_{j \leq J} [a_j \widehat{y}_i \langle \mathbf{w}, \widehat{\mathbf{x}}_i \rangle + b_j] + \max_{j \leq J} \left\{ \begin{array}{l} \sup_{\Delta \mathbf{x}, \Delta y} \quad \frac{1}{N} [a_j \widehat{y}_k \langle \mathbf{w}, \widehat{\mathbf{x}}_k \rangle + b_j] + \frac{1}{N} [a_j \widehat{y}_k \langle \mathbf{w}, \Delta \mathbf{x} \rangle] \\ \text{s.t.} \quad \frac{1}{N} \|(\Delta \mathbf{x}, \Delta y)\| \leq \rho \end{array} \right.$$

$$\begin{aligned}
&= \frac{1}{N} \sum_{i \neq k} \max_{j \leq J} [a_j \hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle + b_j] + \frac{1}{N} \max_{j \leq J} [a_j \hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_k \rangle + b_j + \rho N \|\mathbf{w}\|_* |a_j|] \\
&= \frac{1}{N} \sum_{i \neq k} \max_{j \leq J} [a_j \hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle + b_j] + \frac{1}{N} \max_{j \leq J} [a_j \hat{y}_k \langle \mathbf{w}, \hat{\mathbf{x}}_k \rangle + b_j] + \max_{j \leq J} \rho \|\mathbf{w}\|_* |a_j| \\
&= \frac{1}{N} \sum_{i=1}^N L(\hat{y}_i \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle) + \max_{j \leq J} \rho \|\mathbf{w}\|_* |a_j|,
\end{aligned}$$

where the penultimate equality holds because $\hat{y}_k \langle \mathbf{w}, \hat{\mathbf{x}}_k \rangle$ resides on the steepest linear piece of the loss function L by virtue of Assumption 21. The claim then follows from Theorem 14(ii) because $\text{lip}(L) = \max_{j \leq J} |a_j|$. Note that generic convex Lipschitz continuous loss functions can be uniformly approximated as closely as desired with convex piecewise linear functions. Thus, the above arguments extend directly to generic convex Lipschitz continuous loss functions. Details are omitted for brevity. \blacksquare

Proof of Theorem 26 By the definition of the feature map Φ corresponding to the kernel k , we have

$$\begin{aligned}
\|\Phi(\mathbf{x}_1) - \Phi(\mathbf{x}_2)\|_{\mathbb{H}} &= \sqrt{\langle \Phi(\mathbf{x}_1), \Phi(\mathbf{x}_1) \rangle_{\mathbb{H}} - 2\langle \Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2) \rangle_{\mathbb{H}} + \langle \Phi(\mathbf{x}_2), \Phi(\mathbf{x}_2) \rangle_{\mathbb{H}}} \\
&= \sqrt{k(\mathbf{x}_1, \mathbf{x}_1) - 2k(\mathbf{x}_1, \mathbf{x}_2) + k(\mathbf{x}_2, \mathbf{x}_2)} \leq g(\|\mathbf{x}_1 - \mathbf{x}_2\|_2)
\end{aligned} \tag{A.3}$$

for all $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{X}$, where the inequality follows from Assumption 25. As for assertion (i), we may use similar argument as in the proof of Lemma 45 to reformulate the worst-case expectation in (29) as

$$\begin{aligned}
&\sup_{\mathbf{Q} \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}} [\ell(h(\mathbf{x}), y)] \tag{A.4} \\
&= \left\{ \begin{array}{l} \sup_{\mathbf{Q}^i} \frac{1}{N} \sum_{i=1}^N \int_{\mathbb{X} \times \mathbb{Y}} \ell(h(\mathbf{x}), y) \mathbf{Q}^i(d\mathbf{x}, dy) \\ \text{s.t.} \quad \frac{1}{N} \sum_{i=1}^N \int_{\mathbb{X} \times \mathbb{Y}} d((\mathbf{x}, y), (\hat{\mathbf{x}}_i, \hat{y}_i)) \mathbf{Q}^i(d\mathbf{x}, dy) \leq \rho \\ \int_{\mathbb{X} \times \mathbb{Y}} \mathbf{Q}^i(d\mathbf{x}, dy) = 1 \quad i \in [N] \end{array} \right. \\
&= \left\{ \begin{array}{l} \sup_{\mathbf{Q}^i} \frac{1}{N} \sum_{i=1}^N \int_{\mathbb{X} \times \mathbb{Y}} \ell(h(\mathbf{x}), y) \mathbf{Q}^i(d\mathbf{x}, dy) \\ \text{s.t.} \quad g \left(\frac{1}{N} \sum_{i=1}^N \int_{\mathbb{X} \times \mathbb{Y}} \sqrt{2} d((\mathbf{x}, y), (\hat{\mathbf{x}}_i, \hat{y}_i)) \mathbf{Q}^i(d\mathbf{x}, dy) \right) \leq g(\sqrt{2}\rho) \\ \int_{\mathbb{X} \times \mathbb{Y}} \mathbf{Q}^i(d\mathbf{x}, dy) = 1 \quad i \in [N] \end{array} \right. \\
&\leq \left\{ \begin{array}{l} \sup_{\mathbf{Q}^i} \frac{1}{N} \sum_{i=1}^N \int_{\mathbb{X} \times \mathbb{Y}} \ell(h(\mathbf{x}), y) \mathbf{Q}^i(d\mathbf{x}, dy) \\ \text{s.t.} \quad \frac{1}{N} \sum_{i=1}^N \int_{\mathbb{X} \times \mathbb{Y}} g \left(\sqrt{2} d((\mathbf{x}, y), (\hat{\mathbf{x}}_i, \hat{y}_i)) \right) \mathbf{Q}^i(d\mathbf{x}, dy) \leq g(\sqrt{2}\rho) \\ \int_{\mathbb{X} \times \mathbb{Y}} \mathbf{Q}^i(d\mathbf{x}, dy) = 1 \quad i \in [N], \end{array} \right. \tag{A.5}
\end{aligned}$$

where the second equality holds because g is strictly monotonically increasing, and the inequality follows from Jensen inequality, which applies because g is concave. By the definition of the transportation metric on $\mathbb{H} \times \mathbb{Y}$ for regression problems, we then find

$$\begin{aligned}
 g\left(\sqrt{2}d((\mathbf{x}, y), (\widehat{\mathbf{x}}_i, \widehat{y}_i))\right) &= g\left(\sqrt{2\|\mathbf{x} - \widehat{\mathbf{x}}_i\|_2^2 + 2(y - \widehat{y}_i)^2}\right) \\
 &\geq g(\|\mathbf{x} - \widehat{\mathbf{x}}_i\|_2 + |y - \widehat{y}_i|) \\
 &\geq g(\|\mathbf{x} - \widehat{\mathbf{x}}_i\|_2) + |y - \widehat{y}_i| \\
 &\geq \|\Phi(\mathbf{x}) - \Phi(\widehat{\mathbf{x}}_i)\|_{\mathbb{H}} + |y - \widehat{y}_i| \\
 &\geq \sqrt{\|\Phi(\mathbf{x}) - \Phi(\widehat{\mathbf{x}}_i)\|_{\mathbb{H}}^2 + (y - \widehat{y}_i)^2} \\
 &= d_{\mathbb{H}}((\Phi(\mathbf{x}), y), (\Phi(\widehat{\mathbf{x}}_i), \widehat{y}_i)),
 \end{aligned} \tag{A.6}$$

where the first inequality holds because $2a^2 + 2b^2 \geq (a + b)^2$ for all $a, b \geq 0$ and because g is strictly monotonically increasing, the second inequality exploits the assumption that $g'(z) \geq 1$ for all $z \geq 0$, the third inequality follows from (A.3), and the last equality holds because $a^2 + b^2 \leq (a + b)^2$ for all $a, b \geq 0$. Substituting the above estimate into (A.5) and using the reproducing property $h(\mathbf{x}) = \langle h, \Phi(\mathbf{x}) \rangle_{\mathbb{H}}$ yields

$$\begin{aligned}
 &\sup_{\mathbf{Q} \in \mathbb{B}_\rho(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}}[\ell(h(\mathbf{x}), y)] \\
 &\leq \begin{cases} \sup_{\mathbf{Q}^i} \frac{1}{N} \sum_{i=1}^N \int_{\mathbb{X} \times \mathbb{Y}} \ell(\langle h, \Phi(\mathbf{x}) \rangle_{\mathbb{H}}, y) \mathbf{Q}^i(d\mathbf{x}, dy) \\ \text{s.t.} \quad \frac{1}{N} \sum_{i=1}^N \int_{\mathbb{X} \times \mathbb{Y}} d_{\mathbb{H}}((\Phi(\mathbf{x}), y), (\Phi(\widehat{\mathbf{x}}_i), \widehat{y}_i)) \mathbf{Q}^i(d\mathbf{x}, dy) \leq g(\sqrt{2}\rho) \\ \int_{\mathbb{X} \times \mathbb{Y}} \mathbf{Q}^i(d\mathbf{x}, dy) = 1 \quad i \in [N] \end{cases} \\
 &\leq \begin{cases} \sup_{\mathbf{Q}^i} \frac{1}{N} \sum_{i=1}^N \int_{\mathbb{H} \times \mathbb{Y}} \ell(\mathbf{x}_{\mathbb{H}}, y) \mathbf{Q}^i(d\mathbf{x}_{\mathbb{H}}, dy) \\ \text{s.t.} \quad \frac{1}{N} \sum_{i=1}^N \int_{\mathbb{H} \times \mathbb{Y}} d_{\mathbb{H}}((\mathbf{x}_{\mathbb{H}}, y), (\Phi(\widehat{\mathbf{x}}_i), \widehat{y}_i)) \mathbf{Q}^i(d\mathbf{x}_{\mathbb{H}}, dy) \leq g(\sqrt{2}\rho) \\ \int_{\mathbb{H} \times \mathbb{Y}} \mathbf{Q}^i(d\mathbf{x}_{\mathbb{H}}, dy) = 1 \quad i \in [N] \end{cases} \\
 &= \sup_{\mathbf{Q} \in \mathbb{B}_{g(\sqrt{2}\rho)}(\widehat{\mathbb{P}}_N^{\mathbb{H}})} \mathbb{E}^{\mathbf{Q}}[\ell(\mathbf{x}_{\mathbb{H}}, y)],
 \end{aligned}$$

where the second inequality follows from relaxing the implicit condition that the random variable $\mathbf{x}_{\mathbb{H}}$ must be supported on $\{\Phi(\mathbf{x}) : \mathbf{x} \in \mathbb{X}\} \subseteq \mathbb{H}$. This completes the proof of assertion (i) for regression problems.

The proof of assertion (ii) parallels that of assertion (i) with obvious modifications. Due to the different transportation metric for classification problems, however, the estimate (A.6) changes to

$$\begin{aligned}
 g(d((\mathbf{x}, y), (\widehat{\mathbf{x}}_i, \widehat{y}_i))) &= g(\|\mathbf{x} - \widehat{\mathbf{x}}_i\|_2 + \kappa|y - \widehat{y}_i|) \\
 &\geq g(\|\mathbf{x} - \widehat{\mathbf{x}}_i\|_2) + \kappa|y - \widehat{y}_i|
 \end{aligned}$$

$$\geq \|\Phi(\mathbf{x}) - \Phi(\hat{\mathbf{x}}_i)\|_{\mathbb{H}} + \kappa|y - \hat{y}_i| = d_{\mathbb{H}}((\Phi(\mathbf{x}), y), (\Phi(\hat{\mathbf{x}}_i), \hat{y}_i)),$$

where the first inequality exploits the assumption that $g'(z) \geq 1$ for all $z \geq 0$, and the second inequality follows from (A.3). Further details are omitted for brevity. \blacksquare

Proof of Theorem 27 The proof follows immediately from (Schölkopf and Smola, 2001, Theorem 4.2), which applies to loss functions representable as a sum of an empirical loss depending on $(\hat{\mathbf{x}}_i, \hat{y}_i, h(\hat{\mathbf{x}}_i))$, $i \leq N$, and a regularization term that is strictly monotonically increasing in $\|h\|_{\mathbb{H}}$. However, the additive separability is not needed for the proof. We remark that the optimal solution of (31) is unique if f is strictly increasing in $\|h\|_{\mathbb{H}}$. If f is only non-decreasing in $\|h\|_{\mathbb{H}}$, on the other hand, uniqueness may be lost. Details are omitted for brevity. \blacksquare

Proof of Theorem 28 Using similar arguments as in the proof of Theorem 4(ii) and observing that any Hilbert norm is self-dual, one can show that

$$\inf_{h \in \mathbb{H}} \sup_{Q \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N^{\mathbb{H}})} \mathbb{E}^Q [L(\langle h, \mathbf{x}_{\mathbb{H}} \rangle_{\mathbb{H}} - y)] = \min_{h \in \mathbb{H}} \frac{1}{N} \sum_{i=1}^N L(h(\hat{\mathbf{x}}_i) - \hat{y}_i) + \rho \operatorname{lip}(L) \sqrt{\|h\|_{\mathbb{H}}^2 + 1}.$$

By the representer theorem, which applies because the objective function of the above optimization problem is non-decreasing in $\|h\|_{\mathbb{H}}$, we may restrict the feasible set from \mathbb{H} to the subset of all linearly parametrized hypotheses of the form $h(\mathbf{x}) = \sum_{j=1}^N \beta_j k(\mathbf{x}, \hat{\mathbf{x}}_j)$ for some $\beta \in \mathbb{R}^N$ without sacrificing optimality. The claim then follows by observing that $h(\hat{\mathbf{x}}_i) = \sum_{j=1}^N \mathcal{K}_{ij} \beta_j$ and $\|h\|_2^2 = \langle \beta, \mathcal{K} \beta \rangle$. \blacksquare

Proof of Theorem 29 Using similar arguments as in the proof of Theorem 14(ii) and observing that any Hilbert norm is self-dual, one can show that

$$\inf_{h \in \mathbb{H}} \sup_{Q \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N^{\mathbb{H}})} \mathbb{E}^Q [L(y \langle h, \mathbf{x}_{\mathbb{H}} \rangle_{\mathbb{H}})] = \begin{cases} \min_{h, \lambda, s_i} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & L(\hat{y}_i h(\hat{\mathbf{x}}_i)) \leq s_i & i \in [N] \\ & L(-\hat{y}_i h(\hat{\mathbf{x}}_i)) - \kappa \lambda \leq s_i & i \in [N] \\ & \operatorname{lip}(L) \|h\|_{\mathbb{H}} \leq \lambda, \end{cases}$$

see (Gao and Kleywegt, 2016, Theorem 1) for a full proof. By the representer theorem, which applies because the loss function

$$\begin{aligned} & f((\hat{\mathbf{x}}_1, \hat{y}_1, h(\hat{\mathbf{x}}_1)), \dots, (\hat{\mathbf{x}}_N, \hat{y}_N, h(\hat{\mathbf{x}}_N)), \|h\|_{\mathbb{H}}) \\ &= \begin{cases} \min_{\lambda} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N \max\{L(\hat{y}_i h(\hat{\mathbf{x}}_i)), L(-\hat{y}_i h(\hat{\mathbf{x}}_i)) - \kappa \lambda\} \\ \text{s.t.} & \lambda \geq \operatorname{lip}(L) \|h\|_{\mathbb{H}} \end{cases} \end{aligned}$$

is non-decreasing in $\|h\|_{\mathbb{H}}$, we may restrict attention to all linearly parametrized hypotheses of the form $h(\mathbf{x}) = \sum_{j=1}^N \beta_j k(\mathbf{x}, \hat{\mathbf{x}}_j)$ for some $\beta \in \mathbb{R}^N$ without sacrificing optimality. Thus,

the claim follows. ■

Proof of Theorem 31 For regression problems, the worst-case expected prediction loss satisfies

$$\begin{aligned}
 & \sup_{\mathbf{Q} \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}} [\ell(h(\mathbf{x}), y)] \\
 &= \inf_{\lambda \geq 0} \lambda \rho + \frac{1}{N} \sum_{i=1}^N \sup_{\substack{\mathbf{x} \in \mathbb{X} \\ y \in \mathbb{Y}}} L(h(\mathbf{x}) - y) - \lambda(\|\mathbf{x} - \hat{\mathbf{x}}_i\| + \kappa|y - \hat{y}_i|) \\
 &\leq \inf_{\lambda \geq 0} \lambda \rho + \frac{1}{N} \sum_{i=1}^N \sup_{\substack{\mathbf{x} \in \mathbb{X} \\ y \in \mathbb{Y}}} L(h(\hat{\mathbf{x}}_i) - \hat{y}_i) + \text{lip}(L)(|h(\mathbf{x}) - y - h(\hat{\mathbf{x}}_i) + \hat{y}_i|) \\
 &\quad - \lambda(\|\mathbf{x} - \hat{\mathbf{x}}_i\| + \kappa|y - \hat{y}_i|) \\
 &\leq \inf_{\lambda \geq 0} \lambda \rho + \frac{1}{N} \sum_{i=1}^N \sup_{\substack{\mathbf{x} \in \mathbb{X} \\ y \in \mathbb{Y}}} L(h(\hat{\mathbf{x}}_i) - \hat{y}_i) + \text{lip}(L) \text{lip}(h) \|\mathbf{x} - \hat{\mathbf{x}}_i\| + \text{lip}(L)|y - \hat{y}_i| \\
 &\quad - \lambda(\|\mathbf{x} - \hat{\mathbf{x}}_i\| + \kappa|y - \hat{y}_i|) \\
 &\leq \frac{1}{N} \sum_{i=1}^N \ell(h(\hat{\mathbf{x}}_i), \hat{y}_i) + \rho \text{lip}(L) \max \{\text{lip}(h), 1/\kappa\},
 \end{aligned}$$

where the equality holds due to Lemma 45, while the first and the second inequalities follow from the Lipschitz continuity of L and h , respectively. The last inequality holds by setting $\lambda = \text{lip}(L) \max \{\text{lip}(h), 1/\kappa\}$. Note that $\text{lip}(\varphi(\psi)) \leq \text{lip}(\varphi) \text{lip}(\psi)$ for any functions φ and ψ defined on appropriate normed spaces; see for example (Rockafellar and Wets, 2009, Exercise 9.8). Thus, the Lipschitz modulus of h can be estimated in terms of the Lipschitz moduli of the activation functions σ_m and the operator norms of the weight matrices \mathbf{W}_M , which coincide with the Lipschitz moduli of the corresponding linear mappings. Formally, we have $\text{lip}(h) \leq \prod_{m=1}^M \text{lip}(\sigma_m) \|\mathbf{W}_m\|$. Substituting this estimate into the last line of the above display equation yields the postulated upper bound on the worst-case expectation for regression problems.

Similarly, for classification problems we have

$$\begin{aligned}
 & \sup_{\mathbf{Q} \in \mathcal{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}} [\ell(h(\mathbf{x}), y)] \\
 &= \inf_{\lambda \geq 0} \lambda \rho + \frac{1}{N} \sum_{i=1}^N \sup_{\substack{\mathbf{x} \in \mathbb{X} \\ y \in \mathbb{Y}}} L(yh(\mathbf{x})) - \lambda(\|\mathbf{x} - \hat{\mathbf{x}}_i\| + \kappa \mathbb{1}_{\{y \neq \hat{y}_i\}}) \\
 &\leq \inf_{\lambda \geq 0} \lambda \rho + \frac{1}{N} \sum_{i=1}^N \sup_{\substack{\mathbf{x} \in \mathbb{X} \\ y \in \mathbb{Y}}} L(\hat{y}_i h(\hat{\mathbf{x}}_i)) + \text{lip}(L)(|yh(\mathbf{x}) - \hat{y}_i h(\hat{\mathbf{x}}_i)|) - \lambda(\|\mathbf{x} - \hat{\mathbf{x}}_i\| + \kappa \mathbb{1}_{\{y \neq \hat{y}_i\}}) \\
 &\leq \inf_{\lambda \geq 0} \lambda \rho + \frac{1}{N} \sum_{i=1}^N \sup_{\substack{\mathbf{x} \in \mathbb{X} \\ y \in \mathbb{Y}}} L(\hat{y}_i h(\hat{\mathbf{x}}_i)) + \text{lip}(L) \text{lip}(h) \|\mathbf{x} - \hat{\mathbf{x}}_i\| \mathbb{1}_{\{y \neq \hat{y}_i\}}
 \end{aligned}$$

$$\begin{aligned}
& + \text{lip}(L) |h(\mathbf{x}) + h(\widehat{\mathbf{x}}_i)| \mathbb{1}_{\{y \neq \widehat{y}_i\}} - \lambda(\|\mathbf{x} - \widehat{\mathbf{x}}_i\| + \kappa \mathbb{1}_{\{y \neq \widehat{y}_i\}}) \\
\leq & \inf_{\lambda \geq 0} \lambda \rho + \frac{1}{N} \sum_{i=1}^N \sup_{\substack{\mathbf{x} \in \mathbb{X} \\ y \in \mathbb{Y}}} L(\widehat{y}_i h(\widehat{\mathbf{x}}_i)) - \lambda(\|\mathbf{x} - \widehat{\mathbf{x}}_i\| + \kappa \mathbb{1}_{\{y \neq \widehat{y}_i\}}) \\
& + \text{lip}(L) \max \{2c/\kappa, \text{lip}(h)\} (\|\mathbf{x} - \widehat{\mathbf{x}}_i\| + \kappa \mathbb{1}_{\{y \neq \widehat{y}_i\}}),
\end{aligned}$$

where $c = \sup_{h \in \mathbb{H}, \mathbf{x} \in \mathbb{X}} |h(\mathbf{x})|$ is defined as in the claim. Setting $\lambda = \text{lip}(L) \max \{1/\kappa, 2c/\kappa, \text{lip}(h)\}$ and using the estimate $\text{lip}(h) \leq \prod_{m=1}^M \text{lip}(\sigma_m) \|\mathbf{W}_m\|$ yields the postulated upper bound on the worst-case expectation for classification problems. \blacksquare

Proof of Corollary 33 Set $\bar{\sigma} = \prod_{m \in [M]} \text{lip}(\sigma_m)$. As $\kappa = \infty$, Theorem 31 implies that the worst-case expected prediction loss, both for classification and regression problems, satisfies

$$\begin{aligned}
\sup_{\mathbf{Q} \in \mathcal{B}_\rho(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}} [\ell(h(\mathbf{x}), y)] & \leq \frac{1}{N} \sum_{i=1}^N \sup_{\substack{\mathbf{x} \in \mathbb{X} \\ y \in \mathbb{Y}}} \ell(h(\widehat{\mathbf{x}}_i), \widehat{y}_i) + \rho \bar{\sigma} \text{lip}(L) \prod_{m=1}^M \|\mathbf{W}_m\| \\
& \leq \frac{1}{N} \sum_{i=1}^N \sup_{\substack{\mathbf{x} \in \mathbb{X} \\ y \in \mathbb{Y}}} \ell(h(\widehat{\mathbf{x}}_i), \widehat{y}_i) + \rho \bar{\sigma} \text{lip}(L) \left(\sum_{m=1}^M \frac{\|\mathbf{W}_m\|}{M} \right)^M,
\end{aligned}$$

where the second inequality follows from the inequality of arithmetic and geometric means. By (Everett III, 1963, Theorem 1), if $\mathbf{W}_{[M]}^*$ is a minimizer of the optimization problem

$$\min_{\mathbf{W}_{[M]}} \frac{1}{N} \sum_{i=1}^N \sup_{\substack{\mathbf{x} \in \mathbb{X} \\ y \in \mathbb{Y}}} \ell(h(\widehat{\mathbf{x}}_i), \widehat{y}_i) + \rho \bar{\sigma} \text{lip}(\ell) \left(\sum_{k=1}^M \frac{\|\mathbf{W}_k\|}{M} \right)^M,$$

then the same $\mathbf{W}_{[M]}^*$ also minimizes the constrained optimization problem

$$\begin{cases} \inf_{\mathbf{W}_{[M]}} \frac{1}{N} \sum_{i=1}^N \ell(h(\widehat{\mathbf{x}}_i), \widehat{y}_i) \\ \text{s.t.} \quad \left(\sum_{m=1}^M \frac{\|\mathbf{W}_m\|}{M} \right)^M \leq \left(\frac{\theta}{M} \right)^M, \end{cases}$$

for $\theta = \sum_{m=1}^M \|\mathbf{W}_m^*\|$. Notice that the constraint in the above optimization problem can be simplified to $\sum_{m=1}^M \|\mathbf{W}_m\| \leq \theta$. Hence, there exists a Lagrange multiplier $\bar{\rho}$ for the simplified constraint such that $\mathbf{W}_{[M]}^*$ is a minimizer of the penalized problem

$$\inf_{\mathbf{W}_{[M]}} \frac{1}{N} \sum_{i=1}^N \ell(h(\widehat{\mathbf{x}}_i), \widehat{y}_i) + \bar{\rho} \sum_{m=1}^M \|\mathbf{W}_m\|.$$

This observation completes the proof. \blacksquare

A.2. Proofs of Section 4

The proof of Theorem 39 relies on the following preparatory lemma, which basically asserts that the sample average of a linearly growing function of ξ has sub-Gaussian tails.

Lemma 49 (Sub-Gaussian tails) *If Assumption 34 holds, then there exist constants $c_3 \geq 1$ and $c_4 > 0$ that depend only on the light tail constants a and A of \mathbb{P} such that*

$$\mathbb{P}^N \left\{ \left| \mathbb{E}^{\mathbb{P}}[f(\xi)] - \mathbb{E}^{\widehat{\mathbb{P}}_N}[f(\xi)] \right| \geq \delta \right\} \leq c_3 \exp(-c_4 N \delta^2)$$

for any $N \in \mathbb{N}$, $\delta \in [0, 1]$ and function $f : \Xi \rightarrow \mathbb{R}$ with $|f(\xi) - f(\xi')| \leq d(\xi, \xi')$ for all $\xi \in \Xi$ and some reference point $\xi' \in \Xi$.

Proof Assume that $f : \Xi \rightarrow \mathbb{R}$ is a linear growth function with $|f(\xi) - f(\xi')| \leq d(\xi, \xi')$ for all $\xi \in \Xi$ and some reference point $\xi' \in \Xi$. Set $\xi_f = f(\xi)$ and $\xi'_f = f(\xi')$. Thus, the distribution of the scalar random variable ξ_f is given by the pushforward measure $f_*(\mathbb{P})$ of \mathbb{P} . By construction, we have

$$\mathbb{E}^{f_*(\mathbb{P})} [\exp(|\xi_f - \xi'_f|^a)] = \mathbb{E}^{\mathbb{P}} [\exp(|f(\xi) - f(\xi')|^a)] \leq \mathbb{E}^{\mathbb{P}} [\exp(d(\xi, \xi')^a)] \leq A,$$

where the first inequality follows from the growth condition of f , while the second inequality holds because \mathbb{P} satisfies Assumption 34. Hence, the distribution $f_*(\mathbb{P})$ satisfies Assumption 34 for $n = 0$ when distances on \mathbb{R} are measured by the absolute value, and it inherits the light-tail constants a and A from \mathbb{P} . By using Theorem 35 for $n = 0$, we may thus conclude that there exist constants $c_3, c_4 > 0$ with

$$\mathbb{P}^N \left\{ W(f_*(\mathbb{P}), f_*(\widehat{\mathbb{P}}_N)) \geq \delta \right\} \leq c_3 \exp(-c_4 N \delta^2) \quad \forall \delta \in [0, 1],$$

where $f_*(\widehat{\mathbb{P}}_N)$ represents the empirical distribution of ξ_f , which coincides with the pushforward measure of $\widehat{\mathbb{P}}_N$ under f . By slight abuse of notation, W stands here for the Wasserstein distance between *univariate* distributions, where the absolute value is used as the ground metric. Note that the above univariate measure concentration result holds for any linear growth function f with asymptotic growth rate ≤ 1 . We emphasize that $c_3 \geq 1$ because otherwise the above estimate would fail to hold for $\delta = 0$.

By construction of the Wasserstein distance in Definition 2, we have $W(f_*(\mathbb{P}), f_*(\widehat{\mathbb{P}}_N)) < \delta$ if and only if the scalar random variables ξ_f and ξ'_f admit a joint distribution Π with $\mathbb{E}^{\Pi}[|\xi_f - \xi'_f|] < \delta$ under which ξ_f and ξ'_f have marginals $f_*(\mathbb{P})$ and $f_*(\widehat{\mathbb{P}}_N)$, respectively. The inequality $W(f_*(\mathbb{P}), f_*(\widehat{\mathbb{P}}_N)) < \delta$ thus implies

$$\left| \mathbb{E}^{\mathbb{P}}[f(\xi)] - \mathbb{E}^{\widehat{\mathbb{P}}_N}[f(\xi')] \right| = \left| \mathbb{E}^{f_*(\mathbb{P})}[\xi_f] - \mathbb{E}^{f_*(\widehat{\mathbb{P}}_N)}[\xi'_f] \right| \leq \mathbb{E}^{\Pi}[|\xi_f - \xi'_f|] = \mathbb{E}^{\Pi}[|\xi_f - \xi'_f|] < \delta.$$

By contraposition, we then obtain the implication

$$\left| \mathbb{E}^{\mathbb{P}}[f(\xi)] - \mathbb{E}^{\widehat{\mathbb{P}}_N}[f(\xi')] \right| \geq \delta \quad \implies \quad W(f_*(\mathbb{P}), f_*(\widehat{\mathbb{P}}_N)) \geq \delta,$$

which leads to the desired inequality

$$\mathbb{P}^N \left\{ \left| \mathbb{E}^{\mathbb{P}}[f(\xi)] - \mathbb{E}^{\widehat{\mathbb{P}}_N}[f(\xi)] \right| \geq \delta \right\} \leq \mathbb{P}^N \left\{ W(f_*(\mathbb{P}), f_*(\widehat{\mathbb{P}}_N)) \geq \delta \right\} \leq c_3 \exp(-c_4 N \delta^2).$$

The last inequality holds for all $\delta \in [0, 1]$ and $N \in \mathbb{N}$, irrespective of the linear growth function f . \blacksquare

Remark 50 (Hoeffding's inequality) *If it is known that $\mathbb{P}\{\underline{f} \leq f(\boldsymbol{\xi}) \leq \bar{f}\} = 1$, then Lemma 49 reduces to Hoeffding's inequality (Boucheron et al., 2013, Theorem 2.8), in which case we may set $c_3 = 2$ and $c_4 = 2/(\bar{f} - \underline{f})^2$.*

Proof of Theorem 39 To avoid cumbersome case distinctions, we prove the theorem only in the case when (4) is a classification problem. Thus, we assume that $\Xi = \mathbb{R}^n \times \{-1, 1\}$ and that the transportation cost is of the form (16), where $\|\cdot\|$ denotes a norm on the input space \mathbb{R}^n . The proof for regression problems is similar and only requires minor modifications. It will be omitted for brevity.

From Theorem 14(ii) we know that

$$\sup_{\mathbf{Q} \in \mathbb{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] = \mathbb{E}^{\hat{\mathbb{P}}_N} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] + \rho \text{lip}(L) \Omega(\mathbf{w}) \quad \forall \mathbf{w} \in \mathbb{W},$$

where $\Omega(\mathbf{w}) = \|\mathbf{w}\|_*$ can be viewed as a regularization function. For every $\rho \geq \rho'_N(\eta)$ we thus have

$$\begin{aligned} & \mathbb{P}^N \left\{ \mathbb{E}^{\mathbb{P}} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] \leq \sup_{\mathbf{Q} \in \mathbb{B}_\rho(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbf{Q}} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] \quad \forall \mathbf{w} \in \mathbb{W} \right\} \\ & \geq \mathbb{P}^N \left\{ 0 \leq \min_{\mathbf{w} \in \mathbb{W}} \mathbb{E}^{\hat{\mathbb{P}}_N} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] + \rho'_N(\eta) \text{lip}(L) \Omega(\mathbf{w}) - \mathbb{E}^{\mathbb{P}} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] \right\} \\ & = 1 - \mathbb{P}^N \left\{ \min_{\mathbf{w} \in \mathbb{W}} \mathbb{E}^{\hat{\mathbb{P}}_N} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] + \rho'_N(\eta) \text{lip}(L) \Omega(\mathbf{w}) - \mathbb{E}^{\mathbb{P}} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] < 0 \right\}. \quad (\text{A.7}) \end{aligned}$$

Observe that $\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)$ is Lipschitz continuous in \mathbf{w} for every fixed \mathbf{x} and y because the underlying univariate loss function L is Lipschitz continuous by assumption. Specifically, we have

$$|\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y) - \ell(\langle \mathbf{w}', \mathbf{x} \rangle, y)| \leq \text{lip}(L) |\langle \mathbf{w} - \mathbf{w}', \mathbf{x} \rangle| \leq \text{lip}(L) \|\mathbf{w} - \mathbf{w}'\|_\infty \|\mathbf{x}\|_1 \quad \forall \mathbf{w}, \mathbf{w}' \in \mathbb{W}.$$

For any $\Delta > 0$ there exists a finite set $\mathbb{W}_\Delta \subseteq \mathbb{W}$ with $\Delta = \sup_{\mathbf{w} \in \mathbb{W}} \inf_{\mathbf{w}' \in \mathbb{W}_\Delta} \|\mathbf{w} - \mathbf{w}'\|_\infty$ and whose cardinality satisfies $|\mathbb{W}_\Delta| < (\bar{\Omega}/\Delta - 1)^n < (\bar{\Omega}/\Delta)^n - 1$ where the second inequality holds because $\Delta \leq \bar{\Omega}$ by construction. In the following we set $\Delta = \bar{\Omega}/\sqrt{N}$.

As $\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)$ is Lipschitz continuous in \mathbf{w} , for every $\mathbf{w} \in \mathbb{W}$ there is $\mathbf{w}' \in \mathbb{W}_\Delta$ with

$$|\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)| \geq |\ell(\langle \mathbf{w}', \mathbf{x} \rangle, y)| - \text{lip}(L) \Delta \|\mathbf{x}\|_1.$$

Applying this estimate twice and recalling the assumption that $\Omega(\mathbf{w}) \geq \underline{\Omega}$ for all $\mathbf{w} \in \mathbb{W}$, we may thus conclude that the probability in (A.7) is smaller or equal to

$$\mathbb{P}^N \left\{ \min_{\mathbf{w} \in \mathbb{W}} \mathbb{E}^{\hat{\mathbb{P}}_N} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] - \mathbb{E}^{\mathbb{P}} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] + \rho'_N(\eta) \text{lip}(L) \underline{\Omega} < 0 \right\}$$

$$\begin{aligned}
 &\leq \mathbb{P}^N \left\{ \min_{\mathbf{w} \in \mathbb{W}_\Delta} \mathbb{E}^{\hat{\mathbb{P}}^N} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] - \mathbb{E}^{\mathbb{P}} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] - \text{lip}(L) \Delta \left(\mathbb{E}^{\hat{\mathbb{P}}^N} [\|\mathbf{x}\|_1] + \mathbb{E}^{\mathbb{P}} [\|\mathbf{x}\|_1] \right) \right. \\
 &\quad \left. < -\rho'_N(\eta) \text{lip}(L) \underline{\Omega} \right\} \\
 &\leq \mathbb{P}^N \left\{ \min_{\mathbf{w} \in \mathbb{W}_\Delta} \mathbb{E}^{\hat{\mathbb{P}}^N} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] - \mathbb{E}^{\mathbb{P}} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] - \text{lip}(L) \Delta \left(\mathbb{E}^{\hat{\mathbb{P}}^N} [\|\mathbf{x}\|_1] - \mathbb{E}^{\mathbb{P}} [\|\mathbf{x}\|_1] \right) \right. \\
 &\quad \left. < 2 \text{lip}(L) \Delta M_n n A - \rho'_N(\eta) \text{lip}(L) \underline{\Omega} \right\}. \tag{A.8}
 \end{aligned}$$

The second inequality in the above expression follows from the estimate

$$M_n = \max_{i \leq n} \|\mathbf{e}_i^n\|_* = \max_{i \leq n} \sup_{\|\mathbf{x}\| \leq 1} \langle \mathbf{e}_i^n, \mathbf{x} \rangle = \max_{i \leq n} \sup_{\|\mathbf{x}\| \leq 1} |x_i| = \sup_{\|\mathbf{x}\| \leq 1} \|\mathbf{x}\|_\infty \geq \sup_{\mathbf{x} \in \mathbb{R}^n} \frac{\|\mathbf{x}\|_1}{n \|\mathbf{x}\|},$$

which can be paraphrased as $\|\mathbf{x}\|_1 \leq M_n n \|\mathbf{x}\|$ for every $\mathbf{x} \in \mathbb{R}^n$ and thus implies

$$\mathbb{E}^{\mathbb{P}} [\|\mathbf{x}\|_1] \leq M_n n \mathbb{E}^{\mathbb{P}} [\|\mathbf{x}\|] \leq M_n n \mathbb{E}^{\mathbb{P}} [\exp(\|\mathbf{x}\|^a)] \leq M_n n \mathbb{E}^{\mathbb{P}} [\exp(d(\boldsymbol{\xi}, \boldsymbol{\xi}'))^a] \leq M_n n A.$$

Next, we introduce an auxiliary parameter

$$\delta = \frac{\rho'_N(\eta) \underline{\Omega} - M_n A}{\Delta + \bar{\Omega}} = \frac{2\sqrt{n \log(\sqrt{N})/c_4 + \log(c_3/\eta)/c_4}}{1 + \sqrt{N}},$$

where the second equality follows from the definition of $\rho'_N(\eta)$ in the theorem statement and the convention that $\Delta = \bar{\Omega}/\sqrt{N}$. One can prove that $\delta \in [0, 1]$. Indeed, the nonnegativity of δ is immediate because $c_3 \geq 1$ and $c_4 > 0$. Moreover, we find

$$\delta \leq \frac{2\sqrt{n \log(\sqrt{N})/c_4}}{1 + \sqrt{N}} + \frac{2\sqrt{\log(c_3/\eta)/c_4}}{1 + \sqrt{N}} \leq 1,$$

where the first inequality follows from the observation that $\sqrt{x_1 + x_2} \leq \sqrt{x_1} + \sqrt{x_2}$ for all $x_1, x_2 \geq 0$, while the second inequality holds because $\log(\sqrt{N}) \leq \sqrt{N}$ for all $N \in \mathbb{N}$ and $N \geq \max\{(16n/c_4)^2, 16 \log(c_3/\eta)/c_4\}$, which implies that both fractions in the middle of the above expression are smaller or equal to $\frac{1}{2}$.

Multiplying the definition of δ with $-\text{lip}(L)(\Delta + \bar{\Omega})$ yields the identity

$$-\text{lip}(L) \Delta \delta - \text{lip}(L) \bar{\Omega} \delta = 2 \text{lip}(L) \Delta M_n n A - \text{lip}(L) \rho_N(\eta) \underline{\Omega},$$

and thus the probability (A.8) can be bounded above by

$$\begin{aligned}
 &\mathbb{P}^N \left\{ \min_{\mathbf{w} \in \mathbb{W}_\Delta} \mathbb{E}^{\hat{\mathbb{P}}^N} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] - \mathbb{E}^{\mathbb{P}} [\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] \leq -\text{lip}(L) \bar{\Omega} \delta \quad \text{or} \right. \\
 &\quad \left. -\text{lip}(L) \Delta \left(\mathbb{E}^{\hat{\mathbb{P}}^N} [\|\mathbf{x}\|] - \mathbb{E}^{\mathbb{P}} [\|\mathbf{x}\|] \right) \leq -\text{lip}(L) \Delta \delta \right\} \\
 &\leq \mathbb{P}^N \left\{ \min_{\mathbf{w} \in \mathbb{W}_\Delta} \mathbb{E}^{\hat{\mathbb{P}}^N} \left[\frac{\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)}{\text{lip}(L) \bar{\Omega}} \right] - \mathbb{E}^{\mathbb{P}} \left[\frac{\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)}{\text{lip}(L) \bar{\Omega}} \right] \leq -\delta \right\} \tag{A.9a}
 \end{aligned}$$

$$+ \mathbb{P}^N \left\{ \left(\mathbb{E}^{\hat{\mathbb{P}}_N} [\|\mathbf{x}\|] - \mathbb{E}^{\mathbb{P}} [\|\mathbf{x}\|] \right) \geq \delta \right\}, \quad (\text{A.9b})$$

where the inequality follows from the subadditivity of probability measures.

For any fixed $\mathbf{w} \in \mathbb{W}$ one can show that the function $f(\boldsymbol{\xi}) = \ell(\langle \mathbf{w}, \mathbf{x} \rangle, y) / (\text{lip}(L)\bar{\Omega})$ with $\boldsymbol{\xi} = (\mathbf{x}, y)$ satisfies the linear growth condition $|f(\boldsymbol{\xi}) - f(\boldsymbol{\xi}')| \leq d(\boldsymbol{\xi}, \boldsymbol{\xi}')$ for all $\boldsymbol{\xi} \in \Xi$ if $\boldsymbol{\xi}' = (\mathbf{0}, 1)$. Details are omitted for brevity. By the subadditivity of probability measures, the probability (A.9a) is thus smaller or equal to

$$\sum_{\mathbf{w} \in \mathbb{W}_\Delta} \mathbb{P}^N \left\{ \left| \mathbb{E}^{\hat{\mathbb{P}}_N} \left[\frac{\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)}{\text{lip}(L)\bar{\Omega}} \right] - \mathbb{E}^{\mathbb{P}} \left[\frac{\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)}{\text{lip}(L)\bar{\Omega}} \right] \right| > \delta \right\} \leq |\mathbb{W}_\Delta| c_3 \exp(-c_4 N \delta^2),$$

where the inequality follows from Lemma 49, which applies because $\delta \in [0, 1]$. Moreover, the function $f(\boldsymbol{\xi}) = \|\mathbf{x}\|$ with $\boldsymbol{\xi} = (\mathbf{x}, y)$ trivially satisfies the linear growth condition $|f(\boldsymbol{\xi}) - f(\boldsymbol{\xi}')| \leq d(\boldsymbol{\xi}, \boldsymbol{\xi}')$ for all $\boldsymbol{\xi} \in \Xi$ if $\boldsymbol{\xi}' = (\mathbf{0}, 1)$. By Lemma 49, the probability (A.9b) is thus smaller or equal to

$$\mathbb{P}^N \left\{ \left| \mathbb{E}^{\mathbb{P}} [\|\mathbf{x}\|] - \mathbb{E}^{\hat{\mathbb{P}}_N} [\|\mathbf{x}\|] \right| > \delta \right\} \leq c_3 \exp(-c_4 N \delta^2).$$

By combining the above estimates, we may conclude that the probability in (A.7) does not exceed

$$\begin{aligned} & (|\mathbb{W}_\Delta| + 1) c_3 \exp(-c_4 N \delta^2) \\ & \leq (\bar{\Omega}/\Delta)^n c_3 \exp \left(-c_4 N \left(\frac{\rho'_N(\eta)\bar{\Omega} - 2\Delta M_n n A}{\Delta + \bar{\Omega}} \right)^2 \right) \\ & = N^{\frac{n}{2}} c_3 \exp \left(-c_4 N \left(\frac{\rho'_N(\eta)\bar{\Omega}\sqrt{N} - 2\bar{\Omega} M_n n A}{\bar{\Omega}(\sqrt{N} + 1)} \right)^2 \right) \\ & \leq N^{\frac{n}{2}} c_3 \exp \left(-c_4 N \left(\frac{2\bar{\Omega}\sqrt{n \log(\sqrt{N})/c_4 + \log(c_3/\eta)/c_4}}{2\bar{\Omega}\sqrt{N}} \right)^2 \right) \\ & = N^{\frac{n}{2}} c_3 \exp \left(-n \log(\sqrt{N}) - \log(c_3/\eta) \right) = \eta, \end{aligned}$$

where the first inequality follows from the definition of δ and the assumption that $|\mathbb{W}_\Delta| < (\bar{\Omega}/\Delta)^n - 1$, the first equality holds because $\bar{\Omega}/\Delta = \sqrt{N}$, and the second inequality holds due to the definition of $\rho'_N(\eta)$. In summary, we have shown that the probability in (A.7) is at most η , and thus the claim follows. \blacksquare

A.3. Proofs of Section 5

Proof of Theorem 41 As for assertion (i), note that the absolute value function coincides with the ϵ -insensitive loss for $\epsilon = 0$. Thus, (41a) follows immediately from Corollary 6 by

fixing \mathbf{w} and by setting $\epsilon = 0$ and $\Xi = \mathbb{R}^{n+1}$. As for assertion (ii), similar arguments as in the proof of Lemma 45 show that

$$\mathcal{E}_{\min}(\mathbf{w}) = \sup_{\lambda \geq 0} -\lambda\rho + \frac{1}{N} \sum_{i=1}^N \inf_{\mathbf{x}, y} |y - \langle \mathbf{w}, \mathbf{x} \rangle| + \lambda \|(\mathbf{x}, y) - (\hat{\mathbf{x}}_i, \hat{y}_i)\|. \quad (\text{A.10})$$

The subordinate minimization problem in the first constraint of (A.10) is equivalent to

$$\begin{aligned} & \inf_{\mathbf{x}, y} |y - \langle \mathbf{w}, \mathbf{x} \rangle| + \lambda \|(\mathbf{x}, y) - (\hat{\mathbf{x}}_i, \hat{y}_i)\| \\ &= \inf_{\mathbf{x}, y} \sup_{\|(\mathbf{q}_i, v_i)\|_* \leq \lambda} |y - \langle \mathbf{w}, \mathbf{x} \rangle| + \langle \mathbf{q}_i, \mathbf{x} - \hat{\mathbf{x}}_i \rangle + v_i(y_i - \hat{y}_i) \\ &= \sup_{\|(\mathbf{q}_i, v_i)\|_* \leq \lambda} \begin{cases} \inf_{\mathbf{x}, y, z} & z + \langle \mathbf{q}_i, \mathbf{x} - \hat{\mathbf{x}}_i \rangle + v_i(y_i - \hat{y}_i) \\ \text{s.t.} & z \geq y - \langle \mathbf{w}, \mathbf{x} \rangle, \quad z \geq \langle \mathbf{w}, \mathbf{x} \rangle - y \end{cases} \\ &= \begin{cases} \sup_{\mathbf{q}_i, v_i, r_i, t_i} & -\langle \mathbf{q}_i, \hat{\mathbf{x}}_i \rangle - \hat{y}_i v_i \\ \text{s.t.} & t_i + r_i = 1, \quad t_i - r_i = v_i \\ & (r_i - t_i)\mathbf{w} = \mathbf{q}_i, \quad \|(\mathbf{q}_i, v_i)\|_* \leq \lambda \\ & r_i, t_i \geq 0 \end{cases} \\ &= \begin{cases} \sup_{v_i} & v_i(\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle - \hat{y}_i) \\ \text{s.t.} & -1 \leq v_i \leq 1 \\ & v_i \|(\mathbf{w}, -1)\|_* \leq \lambda, \quad -v_i \|(\mathbf{w}, -1)\|_* \leq \lambda, \end{cases} \end{aligned}$$

where the second equality holds due to Proposition 5.5.4 in (Bertsekas, 2009), and the third equality holds due to strong linear programming duality. By substituting the last optimization problem into (A.10) and replacing v_i with $-v_i$, we have

$$\begin{aligned} \mathcal{E}_{\min}(\mathbf{w}) &= \begin{cases} \sup_{v_i, \lambda} & -\lambda\rho + \frac{1}{N} \sum_{i=1}^N v_i(\hat{y}_i - \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle) \\ \text{s.t.} & -1 \leq v_i \leq 1 & i \in [N] \\ & v_i \|(\mathbf{w}, -1)\|_* \leq \lambda & i \in [N] \\ & -v_i \|(\mathbf{w}, -1)\|_* \leq \lambda & i \in [N] \end{cases} \\ &= \begin{cases} \sup_{v_i, \lambda} & -\lambda\rho + \frac{1}{N} \sum_{i=1}^N v_i |\hat{y}_i - \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle| \\ \text{s.t.} & 0 \leq v_i \leq 1 & i \in [N] \\ & v_i \|(\mathbf{w}, -1)\|_* \leq \lambda & i \in [N] \end{cases} \\ &= \begin{cases} \sup_{v, \lambda} & -\lambda\rho + \frac{1}{N} \sum_{i=1}^N v |\hat{y}_i - \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle| \\ \text{s.t.} & 0 \leq v \leq 1 \\ & v \|(\mathbf{w}, -1)\|_* \leq \lambda \end{cases} \\ &= \begin{cases} \sup_v & v \left(\frac{1}{N} \sum_{i=1}^N |\hat{y}_i - \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle| - \rho \|(\mathbf{w}, -1)\|_* \right) \\ \text{s.t.} & 0 \leq v \leq 1 \end{cases} \end{aligned}$$

$$= \max \left\{ \frac{1}{N} \sum_{i=1}^N |\hat{y}_i - \langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle| - \rho \|(\mathbf{w}, -1)\|_*, 0 \right\}.$$

Thus, the claim follows. \blacksquare

Proof of Theorem 42 As for assertion (i), similar arguments as in the proof of Lemma 45 show that

$$\mathcal{R}_{\max}(\mathbf{w}) = \begin{cases} \inf_{\lambda, s_i} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & \sup_{\mathbf{x}} \mathbb{1}_{\{\hat{y}_i \langle \mathbf{w}, \mathbf{x} \rangle \leq 0\}} - \lambda \|\hat{\mathbf{x}}_i - \mathbf{x}\| \leq s_i & \forall i \leq N \\ & \sup_{\mathbf{x}} \mathbb{1}_{\{-\hat{y}_i \langle \mathbf{w}, \mathbf{x} \rangle \leq 0\}} - \lambda \|\hat{\mathbf{x}}_i - \mathbf{x}\| - \kappa \lambda \leq s_i & \forall i \leq N \\ & \lambda \geq 0. \end{cases} \quad (\text{A.11a})$$

Next, observe that the indicator functions in (A.11a) can be represented as pointwise maxima of extended real-valued concave functions of the form $\mathbb{1}_{\{\hat{y}_i \langle \mathbf{w}, \mathbf{x} \rangle \leq 0\}} = \max\{I_1(\mathbf{x}), 0\}$ and $\mathbb{1}_{\{-\hat{y}_i \langle \mathbf{w}, \mathbf{x} \rangle \leq 0\}} = \max\{I_2(\mathbf{x}), 0\}$, respectively, where

$$I_1(\mathbf{x}) = \begin{cases} 1 & \hat{y}_i \langle \mathbf{w}, \mathbf{x} \rangle \leq 0, \\ -\infty & \text{otherwise,} \end{cases} \quad \text{and} \quad I_2(\mathbf{x}) = \begin{cases} 1 & \hat{y}_i \langle \mathbf{w}, \mathbf{x} \rangle \geq 0, \\ -\infty & \text{otherwise.} \end{cases}$$

This allows us to reformulate (A.11a) as

$$\begin{cases} \inf_{\lambda, s_i} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & \sup_{\mathbf{x} \in \mathbb{R}^n} I_1(\mathbf{x}) - \lambda \|\hat{\mathbf{x}}_i - \mathbf{x}\| \leq s_i & \forall i \leq N \\ & \sup_{\mathbf{x}} 0 - \lambda \|\hat{\mathbf{x}}_i - \mathbf{x}\| \leq s_i & \forall i \leq N \\ & \sup_{\mathbf{x}} I_2(\mathbf{x}) - \lambda \|\hat{\mathbf{x}}_i - \mathbf{x}\| - \kappa \lambda \leq s_i & i \in [N] \\ & \sup_{\mathbf{x}} 0 - \lambda \|\hat{\mathbf{x}}_i - \mathbf{x}\| - \kappa \lambda \leq s_i & \forall i \leq N \\ & \lambda \geq 0. \end{cases}$$

Using the definition of the dual norm and applying the duality theorem (Bertsekas, 2009, Proposition 5.5.4), we find

$$\mathcal{R}_{\max}(\mathbf{w}) = \begin{cases} \inf_{\lambda, s_i, \mathbf{p}_i, \mathbf{q}_i} & \lambda \rho + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} & \sup_{\mathbf{x}} I_1(\mathbf{x}) + \langle \mathbf{p}_i, \mathbf{x} \rangle - \langle \mathbf{p}_i, \hat{\mathbf{x}}_i \rangle \leq s_i & \forall i \leq N \\ & \sup_{\mathbf{x}} I_2(\mathbf{x}) + \langle \mathbf{q}_i, \mathbf{x} \rangle - \langle \mathbf{q}_i, \hat{\mathbf{x}}_i \rangle - \kappa \lambda \leq s_i & \forall i \leq N \\ & s_i \geq 0, \|\mathbf{p}_i\|_* \leq \lambda, \|\mathbf{q}_i\|_* \leq \lambda & i \in [N]. \end{cases} \quad (\text{A.11b})$$

Moreover, by strong linear programming duality we have

$$\sup_{\mathbf{x}} I_1(\mathbf{x}) + \langle \mathbf{p}_i, \mathbf{x} \rangle = \sup_{\mathbf{x}} \{1 + \langle \mathbf{p}_i, \mathbf{x} \rangle : \hat{y}_i \langle \mathbf{w}, \mathbf{x} \rangle \leq 0\} = \inf_{r_i \geq 0} \{1 : \hat{y}_i r_i \mathbf{w} = \mathbf{p}_i\} \quad (\text{A.11c})$$

and

$$\sup_{\mathbf{x}} I_2(\mathbf{x}) + \langle \mathbf{q}_i, \mathbf{x} \rangle = \sup_{\mathbf{x}} \{1 + \langle \mathbf{q}_i, \mathbf{x} \rangle : \hat{y}_i \langle \mathbf{w}, \mathbf{x} \rangle \geq 0\} = \inf_{t_i \geq 0} \{1 : \hat{y}_i t_i \mathbf{w} = \mathbf{q}_i\}. \quad (\text{A.11d})$$

Substituting (A.11c) and (A.11d) into (A.11b) yields (42a). The expression (42b) for the best-case risk can be proved in a similar fashion. Details are omitted for brevity. ■

References

- Y. S. Abu-Mostafa, M. Magdon-Ismail, and H.-T. Lin. *Learning from Data*. AMLBook, 2012.
- N. Agarwal, B. Bullins, and E. Hazan. Second-order stochastic optimization for machine learning in linear time. *Journal of Machine Learning Research*, 18:4148–4187, 2017.
- M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223, 2017.
- K. Bache and M. Lichman. UCI machine learning repository, 2013. URL <http://archive.ics.uci.edu/ml>. Available from <http://archive.ics.uci.edu/ml>.
- P. L. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- A. Ben-Tal, L. El Ghaoui, and A. Nemirovski. *Robust Optimization*. Princeton University Press, 2009.
- A. Ben-Tal, D. Den Hertog, A. De Waegenare, B. Melenberg, and G. Rennen. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 59(2):341–357, 2013.
- J.-D. Benamou, G. Carlier, M. Cuturi, L. Nenna, and G. Peyré. Iterative Bregman projections for regularized transportation problems. *SIAM Journal on Scientific Computing*, 37(2): A1111–A1138, 2015.
- D. Bertsekas. *Convex Optimization Theory*. Athena Scientific, 2009.
- D. Bertsimas and M. S. Copenhaver. Characterization of the equivalence of robustification and regularization in linear and matrix regression. *European Journal of Operational Research*, 2017.
- C. Bhattacharyya. Second order cone programming formulations for feature selection. *Journal of Machine Learning Research*, 5:1417–1433, 2004.
- J. Blanchet and K. Murthy. Quantifying distributional model risk via optimal transport. *Mathematics of Operations Research*, 2019.
- J. Blanchet, Y. Kang, and K. Murthy. Robust Wasserstein profile inference and applications to machine learning. *arXiv preprint arXiv:1610.05627*, 2016.

- J. Blanchet, Y. Kang, F. Zhang, and K. Murthy. Data-driven optimal transport cost selection for distributionally robust optimization. *arXiv preprint arXiv:1705.07152*, 2017.
- S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 2013.
- L. Breiman. Bias, variance, and arcing classifiers. Technical report, University of California, Berkeley, 1996.
- J.-F. Cai, E. J. Candès, and Z. Shen. A singular value thresholding algorithm for matrix completion. *SIAM Journal on Optimization*, 20(4):1956–1982, 2010.
- G. C. Calafiore and L. El Ghaoui. On distributionally robust chance-constrained linear programs. *Journal of Optimization Theory and Applications*, 130(1):1–22, 2006.
- K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman. Return of the devil in the details: Delving deep into convolutional nets. In *British Machine Vision Conference*, 2014.
- C. Cortes and V. Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
- M. Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in Neural Information Processing Systems*, pages 2292–2300, 2013.
- M. Cuturi and D. Avis. Ground metric learning. *Journal of Machine Learning Research*, 15: 533–564, 2014.
- M. Cuturi and A. Doucet. Fast computation of Wasserstein barycenters. In *International Conference on Machine Learning*, pages 685–693, 2014.
- B. Defourny. *Machine Learning Solution Methods for Multistage Stochastic Programming*. PhD thesis, Institut Montefiore, Université de Liege, 2010.
- E. Delage and Y. Ye. Distributionally robust optimization under moment uncertainty with application to data-driven problems. *Operations Research*, 58(3):595–612, 2010.
- J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. DeCAF: A deep convolutional activation feature for generic visual recognition. In *International Conference on Machine Learning*, pages 647–655, 2014.
- J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra. Efficient projections onto the ℓ_1 -ball for learning in high dimensions. In *International Conference on Machine Learning*, pages 272–279, 2008.
- L. El Ghaoui and H. Lebret. Robust solutions to least-squares problems with uncertain data. *SIAM Journal on Matrix Analysis and Applications*, 18(4):1035–1064, 1997.
- H. Everett III. Generalized Lagrange multiplier method for solving problems of optimum allocation of resources. *Operations Research*, 11(3):399–417, 1963.
- M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The Pascal visual object classes (VOC) challenge. *International Journal of Computer Vision*, 88(2): 303–338, 2010.

- F. Farnia and D. Tse. A minimax approach to supervised learning. In *Advances in Neural Information Processing Systems*, pages 4240–4248, 2016.
- N. Fournier and A. Guillin. On the rate of convergence in Wasserstein distance of the empirical measure. *Probability Theory and Related Fields*, 162(3-4):707–738, 2015.
- C. Frogner, C. Zhang, H. Mobahi, M. Araya, and T. A. Poggio. Learning with a Wasserstein loss. In *Advances in Neural Information Processing Systems*, pages 2053–2061, 2015.
- R. Gao and A. Kleywegt. Distributionally robust stochastic optimization with Wasserstein distance. *arXiv preprint arXiv:1604.02199*, 2016.
- J. Goh and M. Sim. Distributionally robust optimization and its tractable approximations. *Operations Research*, 58(4):902–917, 2010.
- H. Gouk, E. Frank, B. Pfahringer, and M. Cree. Regularisation of neural networks by enforcing Lipschitz continuity. *arXiv preprint arXiv:1804.04368*, 2018.
- N. Halko, P.-G. Martinsson, and J. A. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Review*, 53(2):217–288, 2011.
- G. A. Hanasusanto, D. Kuhn, and W. Wiesemann. A comment on computational complexity of stochastic programming problems. *Mathematical Programming A*, 159(1-2):557–569, 2016.
- D. Hosmer, S. Lemeshow, and R. X. Sturdivant. *Applied Logistic Regression*. John Wiley & Sons, 2013.
- K. Huang, H. Yang, I. King, M. R. Lyu, and L. Chan. The minimum error minimax probability machine. *Journal of Machine Learning Research*, 5:1253–1286, 2004.
- R. Koenker. *Quantile Regression*. Cambridge University Press, 2005.
- A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- G. Lanckriet, L. El Ghaoui, C. Bhattacharyya, and M. I. Jordan. Minimax probability machine. In *Advances in Neural Information Processing Systems*, pages 801–807, 2002a.
- G. R. Lanckriet, L. El Ghaoui, C. Bhattacharyya, and M. I. Jordan. A robust minimax approach to classification. *Journal of Machine Learning Research*, 3:555–582, 2002b.
- N. Lawrence and B. Schölkopf. Estimating a kernel Fisher discriminant in the presence of label noise. In *International Conference on Machine Learning*, pages 306–306, 2001.
- Y. LeCun, C. Cortes, and C. J. Burges. The MNIST database of handwritten digits, 1998. Available from <http://yann.lecun.com/exdb/mnist>.

- C. Lee and S. Mehrotra. A distributionally-robust approach for finding support vector machines. *Available from Optimization Online*, 2015.
- J. Lee and M. Raginsky. Minimax statistical learning with Wasserstein distances. In *Advances in Neural Information Processing Systems*, pages 2687–2696, 2018.
- N. Loizou and P. Richtárik. Momentum and stochastic momentum for stochastic gradient, Newton, proximal point and subspace descent methods. *arXiv preprint arXiv:1712.09677*, 2017.
- F. Luo and S. Mehrotra. Decomposition algorithm for distributionally robust optimization using Wasserstein metric with an application to a class of regression models. *European Journal of Operational Research*, 278(1):20–35, 2019.
- T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. Spectral normalization for generative adversarial networks. In *International Conference on Learning Representations*, 2018.
- P. Mohajerin Esfahani and D. Kuhn. Data-driven distributionally robust optimization using the Wasserstein metric: Performance guarantees and tractable reformulations. *Mathematical Programming A*, 171(1-2):115–166, 2018.
- N. Natarajan, I. S. Dhillon, P. K. Ravikumar, and A. Tewari. Learning with noisy labels. In *Advances in Neural Information Processing Systems*, pages 1196–1204, 2013.
- B. Neyshabur, S. Bhojanapalli, and N. Srebro. A PAC-Bayesian approach to spectrally-normalized margin bounds for neural networks. In *International Conference on Learning Representations*, 2018.
- A. Nitanda. Stochastic proximal gradient descent with acceleration techniques. In *Advances in Neural Information Processing Systems*, pages 1574–1582, 2014.
- R. T. Rockafellar and R. J.-B. Wets. *Variational Analysis*. Springer, 2009.
- Y. Rubner, C. Tomasi, and L. J. Guibas. The earth mover’s distance as a metric for image retrieval. *International Journal of Computer Vision*, 40(2):99–121, 2000.
- K. Scaman and A. Virmaux. Lipschitz regularity of deep neural networks: analysis and efficient estimation. In *Advances in Neural Information Processing Systems*, pages 3835–3844, 2018.
- B. Schölkopf and A. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, 2001.
- B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson. Estimating the support of a high-dimensional distribution. *Neural Computation*, 13(7):1443–1471, 2001.
- S. Shafieezadeh-Abadeh, P. M. Esfahani, and D. Kuhn. Distributionally robust logistic regression. In *Advances in Neural Information Processing Systems*, pages 1576–1584, 2015.

- 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. Á. Goberna and M. A. López, editors, *Semi-Infinite Programming*, pages 135–165. Kluwer Academic Publishers, 2001.
- P. K. Shivaswamy and T. Jebara. Ellipsoidal kernel machines. In *Artificial Intelligence and Statistics*, pages 484–491, 2007.
- P. K. Shivaswamy and T. Jebara. Maximum relative margin and data-dependent regularization. *Journal of Machine Learning Research*, 11:747–788, 2010.
- P. K. Shivaswamy, C. Bhattacharyya, and A. J. Smola. Second order cone programming approaches for handling missing and uncertain data. *Journal of Machine Learning Research*, 7:1283–1314, 2006.
- A. Smola and B. Schölkopf. A tutorial on support vector regression. *Statistics and Computing*, 14(3):199–222, 2004.
- T. Strohmann and G. Z. Grudic. A formulation for minimax probability machine regression. In *Advances in Neural Information Processing Systems*, pages 785–792, 2003.
- C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society Series B*, 58(1):267–288, 1996.
- A. N. Tikhonov, V. I. Arsenin, and F. John. *Solutions of Ill-Posed Problems*. Winston, 1977.
- C. Villani. *Optimal Transport: Old and New*. Springer, 2008.
- J. Weed and F. Bach. Sharp asymptotic and finite-sample rates of convergence of empirical measures in Wasserstein distance. *To appear in Bernoulli*, 2019.
- T. Wiatowski, M. Tschannen, A. Stanic, P. Grohs, and H. Bölcskei. Discrete deep feature extraction: A theory and new architectures. In *International Conference on Machine Learning*, pages 2149–2158, 2016.
- W. Wiesemann, D. Kuhn, and M. Sim. Distributionally robust convex optimization. *Operations Research*, 62(6):1358–1376, 2014.
- H. Xu, C. Caramanis, and S. Mannor. Robustness and regularization of support vector machines. *Journal of Machine Learning Research*, 10:1485–1510, 2009.
- H. Xu, C. Caramanis, and S. Mannor. Robust regression and Lasso. *IEEE Transactions on Information Theory*, 56(7):3561–3574, 2010.
- T. Yang, M. Mahdavi, R. Jin, L. Zhang, and Y. Zhou. Multiple kernel learning from noisy labels by stochastic programming. In *International Conference on Machine Learning*, pages 123–130, 2012.

- M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*, pages 818–833, 2014.
- C. Zhao and Y. Guan. Data-driven risk-averse stochastic optimization with Wasserstein metric. *Operations Research Letters*, 46(2):262–267, 2018.