
Optimization for Machine Learning

Editors:

Suvrit Sra

suvrit@gmail.com

Max Planck Institute for Biological Cybernetics

72076 Tübingen, Germany

Sebastian Nowozin

nowozin@gmail.com

Microsoft Research

Cambridge, CB3 0FB, United Kingdom

Stephen J. Wright

swright@cs.uwisc.edu

University of Wisconsin

Madison, WI 53706

This is a draft containing only `sra_chapter.tex` and an abbreviated front matter. Please check that the formatting and small changes have been performed correctly. Please verify the affiliation. Please use this version for sending us future modifications.

The MIT Press
Cambridge, Massachusetts
London, England

Contents

1	Robust Optimization in Machine Learning	1
1.1	Introduction	2
1.2	Background on Robust Optimization	3
1.3	Robust Optimization and Adversary Resistant Learning	5
1.4	Robust Optimization and Regularization	8
1.5	Robustness and Consistency	22
1.6	Robustness and Generalization	26
1.7	Conclusion	30

Constantine Caramanis

caramanis@mail.utexas.edu

*The University of Texas at Austin
Austin, Texas*

Shie Mannor

shie@ee.technion.ac.il

*Technion, the Israel Institute of Technology
Haifa, Israel*

Huan Xu

huan.xu@mail.utexas.edu

*The University of Texas at Austin
Austin, Texas*

Robust optimization is a paradigm that uses ideas from convexity and duality, to immunize solutions of convex problems to bounded uncertainty in the parameters of the problem. Machine learning is fundamentally about making decisions under uncertainty, and optimization has long been a central tool; thus at a high level there is no surprise that robust optimization should have a role to play. Indeed, the first part of the story told in this chapter is about specializing robust optimization to specific optimization problems in machine learning. Yet, beyond this, there have been several surprising and deep developments in the use of robust optimization and machine learning, connecting consistency, generalization ability, and other properties (such as sparsity and stability) to robust optimization.

In addition to surveying the direct applications of robust optimization to machine learning, important in their own right, this chapter explores some of these deeper connections, and points the way towards opportunities for applications, and challenges for further research.

1.1 Introduction

Learning, optimization, and decision-making from data must cope with uncertainty introduced implicitly and explicitly. Uncertainty can be explicitly introduced when the data collection process is noisy, or some data are corrupted. It may be introduced when the model specification is wrong, assumptions are missing, or factors overlooked. Uncertainty is also present in pristine data, implicitly, insofar as a finite sample empirical distribution, or function thereof, cannot exactly describe the true distribution in most cases. In the optimization community, it has long been known that the effect of even small uncertainty can be devastating in terms of the quality or feasibility of a solution. In machine learning, overfitting has long been recognized as a central challenge, and a plethora of techniques, many of them regularization-based, have been developed to combat this problem. The theoretical justification for many of these techniques lies in controlling notions of complexity, such as metric entropy or VC-dimension.

This chapter considers uncertainty in optimization, and overfitting, from a unified perspective: robust optimization. In addition to introducing a novel technique for designing algorithms that are immune to noise and do not overfit data, robust optimization also provides a theoretical justification for the success of these algorithms: algorithms have certain properties, like consistency, good generalization, or sparsity, *because they are robust*.

Robust optimization (e.g., Soyster, 1973; El Ghaoui and Lebret, 1997; Ben-Tal and Nemirovski, 2000; Bertsimas and Sim, 2004; Bertsimas et al., 2010; Ben-Tal et al., 2009, and many others) is designed to deal with parameter uncertainty in convex optimization problems. For example, one can imagine a linear program, $\min : \{\mathbf{c}^\top \mathbf{x} \mid \mathbf{A}\mathbf{x} \leq \mathbf{b}\}$, where there is uncertainty in the constraint matrix A , the objective function, \mathbf{c} , or the right hand side vector, \mathbf{b} . Robust optimization develops immunity to a deterministic or set-based notion of uncertainty. Thus, in the face of uncertainty in A , instead of solving $\min : \{\mathbf{c}^\top \mathbf{x} \mid \mathbf{A}\mathbf{x} \leq \mathbf{b}\}$ one solves $\min : \{\mathbf{c}^\top \mathbf{x} \mid \mathbf{A}\mathbf{x} \leq \mathbf{b}, \forall A \in \mathcal{U}\}$, for some suitably defined *uncertainty set* \mathcal{U} . We give a brief introduction to robust optimization in Section 1.2 below.

The remainder of this chapter is organized as follows. In Section 1.2 we provide a brief review of robust optimization. In Section 1.3 we discuss direct applications of robust optimization to constructing algorithms that are resistant to data corruption. This is a direct application of not only the methodology of robust optimization, but also the motivation behind the development of robust optimization. The focus is developing on computationally efficient algorithms, resistant to bounded but otherwise arbitrary (even

adversarial) noise. In Sections 1.4 - 1.6, we show that robust optimization's impact in machine learning extends far outside the originally envisioned scope as developed in the optimization literature. In Section ??, we show that many existing machine learning algorithms that are based on regularization, including support vector machines (SVMs), ridge-regression, and Lasso, are special cases of robust optimization. Using this re-interpretation, their success can be understood from a unified perspective. We also show how the flexibility of robust optimization paves the way for the design of new regularization-like algorithms. Moreover, we show that robustness can be used directly to prove properties like regularity and sparsity. In Section 1.5, we show that robustness can be used to prove statistical consistency. Then, in Section 1.6, we extend the results of Section 1.5, showing that an algorithm's generalization ability and its robustness are related in a fundamental way.

In summary, we show that robust optimization has deep connections to machine learning. In particular it yields a unified paradigm that (a) explains the success of many existing algorithms; (b) provides a prescriptive algorithmic approach to creating new algorithms with desired properties; and (c) allows us to prove general properties of an algorithm.

1.2 Background on Robust Optimization

In this section we provide a brief background on robust optimization, and refer the reader to the survey (Bertsimas et al., 2010), the textbook (Ben-Tal et al., 2009), and references to the original papers therein, for more details.

Optimization affected by parameter uncertainty has long been a focus of the mathematical programming community. As has been demonstrated in compelling fashion (Ben-Tal and Nemirovski, 2000), solutions to optimization problems can exhibit remarkable sensitivity to perturbations in the problem parameters, thus often rendering a computed solution highly infeasible, suboptimal, or both. This parallels developments in related fields, particularly robust control (we refer to the textbooks Zhou et al., 1996; Dullerud and Paganini, 2000, and the references therein).

Stochastic programming (e.g., Prékopa, 1995; Kall and Wallace, 1994), assumes the uncertainty has a probabilistic description. In contrast, robust optimization is built on the premise that the parameters vary arbitrarily in some *a priori* known bounded set, called the *uncertainty set*. Suppose we are optimizing a function $f_0(\mathbf{x})$, subject to the m constraints $f_i(\mathbf{x}, \mathbf{u}_i) \leq 0$, $i = 1, \dots, m$, where \mathbf{u}_i denotes the parameters of function i . Then where as the nominal optimization problem solves $\min\{f_0(\mathbf{x}) : f_i(\mathbf{x}, \mathbf{u}_i) \leq 0, i =$

$1, \dots, m\}$, assuming that the \mathbf{u}_i are known, robust optimization solves:

$$\begin{aligned} \min_{\mathbf{x}} : & \quad f_0(\mathbf{x}) \\ \text{s.t.} : & \quad f_i(\mathbf{x}, \mathbf{u}_i) \leq 0, \quad \forall \mathbf{u}_i \in \mathcal{U}_i, \quad i = 1, \dots, m. \end{aligned} \tag{1.1}$$

Computational Tractability. The tractability of robust optimization, subject to standard and mild Slater-like regularity conditions, amounts to separation for the convex set: $\mathcal{X}(\mathcal{U}) \triangleq \{\mathbf{x} : f_i(\mathbf{x}, \mathbf{u}_i) \leq 0, \forall \mathbf{u}_i \in \mathcal{U}_i, i = 1, \dots, m\}$. If there is an efficient algorithm that asserts $\mathbf{x} \in \mathcal{X}(\mathcal{U})$ or otherwise provides a separating hyperplane, then problem (1.1) can be solved in polynomial time. While the set $\mathcal{X}(\mathcal{U})$ is a convex set as long as each function f_i is convex in \mathbf{x} , it is not in general true that there is an efficient separation algorithm for the set $\mathcal{X}(\mathcal{U})$. However, in many cases of broad interest and application, solving the robust problem can be done efficiently – the robustified problem may be of complexity comparable to that of the nominal one. We outline some of the main complexity results below.

An Example: Linear Programs with Polyhedral Uncertainty. When the uncertainty set, \mathcal{U} , is polyhedral, the separation problem is not only efficiently solvable, it is in fact linear, thus the robust counterpart is equivalent to a linear optimization problem. To illustrate this, consider the problem with uncertainty in the constraint matrix:

$$\begin{aligned} \min_{\mathbf{x}} : & \quad \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} : & \quad \max_{\{\mathbf{a}_i : \mathbf{D}_i \mathbf{a}_i \leq \mathbf{d}_i\}} [\mathbf{a}_i^\top \mathbf{x}] \leq b_i, \quad i = 1, \dots, m. \end{aligned}$$

The dual of the subproblem (recall that \mathbf{x} is not a variable of optimization in the inner max) becomes again a linear program:

$$\left[\begin{array}{ll} \max_{\mathbf{a}_i} : & \mathbf{a}_i^\top \mathbf{x} \\ \text{s.t.} : & \mathbf{D}_i \mathbf{a}_i \leq \mathbf{d}_i \end{array} \right] \longleftrightarrow \left[\begin{array}{ll} \min_{\mathbf{p}_i} : & \mathbf{p}_i^\top \mathbf{d}_i \\ \text{s.t.} : & \mathbf{p}_i^\top \mathbf{D}_i = \mathbf{x} \\ & \mathbf{p}_i \geq 0 \end{array} \right],$$

and therefore the robust linear optimization now becomes:

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{p}_1, \dots, \mathbf{p}_m} : & \quad \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} : & \quad \mathbf{p}_i^\top \mathbf{d}_i \leq b_i, \quad i = 1, \dots, m \\ & \quad \mathbf{p}_i^\top \mathbf{D}_i = \mathbf{x}, \quad i = 1, \dots, m \\ & \quad \mathbf{p}_i \geq 0, \quad i = 1, \dots, m. \end{aligned}$$

Thus the size of such problems grows polynomially in the size of the nominal problem and the dimensions of the uncertainty set.

Some General Complexity Results. We now list a few of the complexity results that are relevant in the sequel. We refer to Bertsimas et al. (2010); Ben-Tal et al. (2009) and references therein for further details. The robust counterpart for a linear program (LP) with polyhedral uncertainty is again an LP. For an LP with ellipsoidal uncertainty, the counterpart is a second order cone (SOCP). A convex quadratic program with ellipsoidal uncertainty has a robust counterpart that is a semidefinite program (SDP). An SDP with ellipsoidal uncertainty has an NP-hard robust counterpart.

Probabilistic Interpretations and Results. The computational advantage of robust optimization is largely due to the fact that the formulation is deterministic, thus dealing with uncertainty sets rather than probability distributions. While the paradigm makes sense when the disturbances are not stochastic, or the distribution is not known, tractability advantages have made robust optimization an appealing computational framework even when the uncertainty is stochastic, and the distribution is fully or partially known. A major success of robust optimization has been the ability to derive *a priori* probability guarantees – e.g., probability of feasibility – that the solution to a robust optimization will satisfy, under a variety of probabilistic assumptions. Thus robust optimization is a tractable framework one can use to build solutions with probabilistic guarantees such as minimum probability of feasibility, or maximum probability of hinge-loss beyond some threshold level, etc. This probabilistic interpretation of robust optimization is used throughout this chapter.

1.3 Robust Optimization and Adversary Resistant Learning

In this section we overview some of the direct applications of robust optimization to coping with uncertainty (adversarial or stochastic) in machine learning problems. The main themes are (a) the formulations one obtains when using different uncertainty sets, (b) the probabilistic interpretation and results one can derive by using robust optimization. Using ellipsoidal uncertainty, we show that the resulting robust problem is tractable. Moreover, we show that this robust formulation has interesting probabilistic interpretations. Then, using a polyhedral uncertainty set, we show that sometimes it is possible to tractably model combinatorial uncertainty, such as missing data.

Robust optimization-based learning algorithms have been proposed for various learning tasks, e.g., learning and planning (Nilim and El Ghaoui, 2005), Fisher linear discriminant analysis (Kim et al., 2005), PCA (d’Aspremont

et al., 2007), and many others. Instead of providing a comprehensive survey, we use support vector machines (SVMs, e.g., Vapnik and Lerner, 1963; Boser et al., 1992; Cortes and Vapnik, 1995) to illustrate the methodology of robust optimization.

Standard SVMs consider the standard binary classification problem, where we are given a finite number of training samples $\{\mathbf{x}_i, y_i\}_{i=1}^m \subseteq \mathbb{R}^n \times \{-1, +1\}$, and must find a linear classifier, specified by the function $h^{\mathbf{w}, b}(\mathbf{x}) = \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle + b)$, where $\langle \cdot, \cdot \rangle$ denotes the standard inner product. The parameters (\mathbf{w}, b) are obtained by solving the following convex optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} : \quad & r(\mathbf{w}, b) + C \sum_{i=1}^m \xi_i \\ \text{s.t.} : \quad & \xi_i \geq [1 - y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b)], \quad i = 1, \dots, m; \\ & \xi_i \geq 0, \quad i = 1, \dots, m; \end{aligned} \quad (1.2)$$

where $r(\mathbf{w}, b)$ is a regularization term, e.g., $r(\mathbf{w}, b) = \frac{1}{2} \|\mathbf{w}\|_2^2$. There are a number of related formulations, some focusing on controlling VC-dimension, promoting sparsity, or some other property; see the textbooks Schölkopf and Smola (2002); Steinwart and Christmann (2008) and references therein.

There are three natural ways uncertainty affects the input data: corruption in the location, \mathbf{x}_i , corruption in the label, y_i , and corruption via altogether missing data. We outline some applications of robust optimization to these three settings.

Corrupted Location. Given observed points $\{\mathbf{x}_i\}$, the additive uncertainty model assumes that $\mathbf{x}_i^{\text{true}} = \mathbf{x}_i + \mathbf{u}_i$. Robust optimization protects against the uncertainty \mathbf{u}_i by minimizing the regularized training loss on all possible locations of the \mathbf{u}_i in some uncertainty set, \mathcal{U}_i .

In Trafalis and Gilbert (2007), the authors consider the ellipsoidal uncertainty set given by:

$$\mathcal{U}_i = \{\mathbf{u}_i : \mathbf{u}_i^\top \Sigma_i \mathbf{u}_i \leq 1\}, \quad i = 1, \dots, m;$$

so that each constraint becomes: $\xi_i \geq [1 - y_i(\langle \mathbf{w}, \mathbf{x}_i + \mathbf{u}_i \rangle + b)]$, $\forall \mathbf{u}_i \in \mathcal{U}_i$. By duality, this is equivalent to $y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 + \|\Sigma_i^{1/2} \mathbf{w}\|_2 - \xi_i$, and hence

their version of robust SVM reduces to

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} : \quad & r(\mathbf{w}, b) + C \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i + \|\Sigma_i^{1/2} \mathbf{w}\|_2, \quad i = 1, \dots, m; \\ & \xi_i \geq 0; \quad i = 1, \dots, m. \end{aligned} \quad (1.3)$$

In Trafalis and Gilbert (2007), $r(\mathbf{w}, b) = \frac{1}{2} \|\mathbf{w}\|_2$, while in Bhattacharyya et al. (2004), the authors use the sparsity-inducing regularizer $r(\mathbf{w}, b) = \|\mathbf{w}\|_1$. In both settings, the robust problem is an instance of a second order cone program (SOCP). Available solvers can solve SOCPs with hundreds of thousands of variables and more.

If the uncertainty \mathbf{u}_i is stochastic, one can use this robust formulation to find a classifier that satisfies constraints on the probability (w.r.t. the distribution of \mathbf{u}_i) that each constraint is violated. In (Shivaswamy et al., 2006), the authors consider two varieties of such chance constraints for $i = 1, \dots, m$:

$$\begin{aligned} (a) \quad & \Pr_{\mathbf{u}_i \sim \mathcal{N}(\bar{\mathbf{0}}, \Sigma_i)} (y_i(\mathbf{w}^\top (\mathbf{x}_i + \mathbf{u}_i) + b) \geq 1 - \xi_i) \geq 1 - \kappa_i; \\ (b) \quad & \inf_{\mathbf{u}_i \sim (\bar{\mathbf{0}}, \Sigma_i)} \Pr_{\mathbf{u}_i} (y_i(\mathbf{w}^\top (\mathbf{x}_i + \mathbf{u}_i) + b) \geq 1 - \xi_i) \geq 1 - \kappa_i; \end{aligned} \quad (1.4)$$

Constraint (a) controls the probability of constraint violation, when the uncertainty follows a known Gaussian distribution. Constraint (b) is more conservative: it controls the worst-case probability of constraint violation, over all centered distributions with variance Σ_i . The next theorem says that the robust formulation with ellipsoidal uncertainty set as above, can be used to control both of these quantities.

Theorem 1.1. *For $i = 1, \dots, m$, consider the robust constraint as given above:*

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i + \gamma_i \|\Sigma_i^{1/2} \mathbf{w}\|_2.$$

If we take $\gamma_i = \Phi^{-1}(\kappa_i)$, for Φ the Gaussian c.d.f., this constraint is equivalent to constraint (a) of (1.4), while taking $\gamma_i = \sqrt{\kappa_i / (1 - \kappa_i)}$ yields constraint (b).

Missing Data. In Globerson and Roweis (2006) the authors use robust optimization with polyhedral uncertainty set to address the problem where some of the features of the testing samples may be deleted (possibly in an adversarial fashion). Using a dummy feature to remove the bias term b if

necessary, we can rewrite the nominal problem as

$$\min_{\mathbf{w}} : \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^m [1 - y_i \mathbf{w}^\top \mathbf{x}_i]_+.$$

For a given choice of \mathbf{w} , the value of the term $[1 - y_i \mathbf{w}^\top \mathbf{x}_i]_+$ in the objective, under an adversarial deletion of K features becomes

$$\begin{aligned} \max_{\alpha_i} \quad & [1 - y_i \mathbf{w}^\top (\mathbf{x}_i \circ (1 - \alpha_i))]_+ \\ \text{s.t.} \quad & \alpha_{ij} \in \{0, 1\}; \quad j = 1, \dots, n; \\ & \sum_{j=1}^n \alpha_{ij} = K. \end{aligned}$$

where \circ denotes pointwise vector multiplication. While this optimization problem is combinatorial, relaxing the integer constraint $\alpha_{ij} \in \{0, 1\}$ to be $0 \leq \alpha_{ij} \leq 1$, does not change the objective value. Thus taking the dual of the maximization, and substituting into the original problem, one obtains the classifier that is maximally resistant to up to K missing features:

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{v}_i, z_i, t_i, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & y_i \mathbf{w}^\top \mathbf{x}_i - t_i \geq 1 - \xi_i; \quad i = 1, \dots, m; \\ & \xi_i \geq 0; \quad i = 1, \dots, m; \\ & t_i \geq K z_i + \sum_{j=1}^n v_{ij}; \quad i = 1, \dots, m; \\ & \mathbf{v}_i \geq \mathbf{0}; \quad i = 1, \dots, m; \\ & z_i + v_{ij} \geq y_i x_{ij} w_{ij}; \quad i = 1, \dots, m; \quad j = 1, \dots, n. \end{aligned}$$

This is again an SOCP, and hence fairly large instances can be solved with specialized software.

Corrupted Labels. When the labels are corrupted, the problem becomes more difficult to address due to its combinatorial nature. However, this too has been recently addressed using robust optimization (Caramanis and Mannor, 2008). While there is still a combinatorial price to pay in the complexity of the classifier class, robust optimization can be used to find the optimal classifier; see Caramanis and Mannor (2008) for the details.

1.4 Robust Optimization and Regularization

In this section and the subsequent two, we demonstrate that robustness can provide a unified explanation for many desirable properties of a learning algorithm, from regularity and sparsity, to consistency and generalization. A main message of this chapter is that many regularized problems exhibit a “hidden robustness” — they are in fact equivalent to a robust optimization problem — which can then be used to directly prove properties like consistency and sparsity, and also to design new algorithms. The main problems that highlight this equivalence are regularized support vector machines, ℓ_2 -regularized regression, and ℓ_1 -regularized regression, also known as Lasso.

1.4.1 Support Vector Machines

We consider regularized SVMs, and show that they are algebraically equivalent to a robust optimization problem. We use this equivalence to provide a probabilistic interpretation of SVMs, which allows us to propose new probabilistic SVM-type formulations. This section is based on Xu et al. (2009).

At a high-level it is known that regularization and robust optimization are related; see, e.g., (El Ghaoui and Le Bret, 1997; Anthony and Bartlett, 1999), and Section 1.3. Yet, the precise connection between robustness and regularized SVMs only first appeared in Xu et al. (2009). One of the mottos of robust optimization is to harness the consequences of probability theory, without paying the computational cost of having to use its axioms. Consider the additive uncertainty model from the previous section: $\mathbf{x}_i + \mathbf{u}_i$. If the uncertainties \mathbf{u}_i are stochastic, various limit results (LLN, CLT, etc.) promise that even independent variables will exhibit strong aggregate coupling behavior. For instance, the set $\{(\mathbf{u}_1, \dots, \mathbf{u}_m) : \sum_{i=1}^m \|\mathbf{u}_i\| \leq c\}$ will have increasing probability as m grows. This motivates designing uncertainty sets with this kind of coupling across uncertainty parameters. We leave it to the reader to check that the *constraint-wise* robustness formulations of the previous section cannot be made to capture such coupling constraints across the disturbances $\{\mathbf{u}_i\}$.

We rewrite SVM without slack variables, as an unconstrained optimization. The natural robust formulation now becomes:

$$\min_{\mathbf{w}, b} \max_{\vec{\mathbf{u}} \in \mathcal{U}} \left\{ r(\mathbf{w}, b) + \sum_{i=1}^m \max [1 - y_i (\langle \mathbf{w}, \mathbf{x}_i - \mathbf{u}_i \rangle + b), 0] \right\}, \quad (1.5)$$

where $\vec{\mathbf{u}}$ denotes the collection of uncertainty vectors, $\{\mathbf{u}_i\}$. Describing our coupled uncertainty set requires a few definitions. The first definition

below characterizes the effect of different uncertainty sets, and captures the coupling that they exhibit. As an immediate consequence we obtain an equivalent robust optimization formulation for regularized SVMs.

Definition 1.2. A set $\mathcal{U}_0 \subseteq \mathbb{R}^n$ is called an Atomic Uncertainty Set if

- (I) $\mathbf{0} \in \mathcal{U}_0$;
- (II) For any $\mathbf{w}_0 \in \mathbb{R}^n$: $\sup_{\mathbf{u} \in \mathcal{U}_0} [\mathbf{w}_0^\top \mathbf{u}] = \sup_{\mathbf{u}' \in \mathcal{U}_0} [-\mathbf{w}_0^\top \mathbf{u}'] < +\infty$.

Definition 1.3. Let \mathcal{U}_0 be an atomic uncertainty set. A set $\mathcal{U} \subseteq \mathbb{R}^{n \times m}$ is called a Sublinear Aggregated Uncertainty Set of \mathcal{U}_0 , if

$$\mathcal{U}^- \subseteq \mathcal{U} \subseteq \mathcal{U}^+,$$

$$\text{where: } \mathcal{U}^- \triangleq \bigcup_{t=1}^m \mathcal{U}_t^-; \quad \mathcal{U}_t^- \triangleq \{(\mathbf{u}_1, \dots, \mathbf{u}_m) \mid \mathbf{u}_t \in \mathcal{U}_0; \mathbf{u}_{i \neq t} = \mathbf{0}\}.$$

$$\mathcal{U}^+ \triangleq \{(\alpha_1 \mathbf{u}_1, \dots, \alpha_m \mathbf{u}_m) \mid \sum_{i=1}^m \alpha_i = 1; \alpha_i \geq 0, \mathbf{u}_i \in \mathcal{U}_0, i = 1, \dots, m\}.$$

Sublinear aggregated uncertainty models the case where the disturbances on each sample are treated identically, but their aggregate behavior across multiple samples is controlled. Some interesting examples include

- (1) $\mathcal{U} = \{(\mathbf{u}_1, \dots, \mathbf{u}_m) \mid \sum_{i=1}^m \|\mathbf{u}_i\| \leq c\}$;
- (2) $\mathcal{U} = \{(\mathbf{u}_1, \dots, \mathbf{u}_m) \mid \exists t \in [1 : m]; \|\mathbf{u}_t\| \leq c; \mathbf{u}_i = \mathbf{0}, \forall i \neq t\}$; and
- (3) $\mathcal{U} = \{(\mathbf{u}_1, \dots, \mathbf{u}_m) \mid \sum_{i=1}^m \sqrt{c \|\mathbf{u}_i\|} \leq c\}$.

All these examples share the same atomic uncertainty set $\mathcal{U}_0 = \{\mathbf{u} \mid \|\mathbf{u}\| \leq c\}$. Figure 1.1 provides an illustration of a sublinear aggregated uncertainty set for $n = 1$ and $m = 2$, i.e., the training set consists of two univariate samples.

Theorem 1.4. Assume $\{\mathbf{x}_i, y_i\}_{i=1}^m$ are non-separable, $r(\cdot, \cdot) : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$ is an arbitrary function, \mathcal{U} is a sublinear aggregated uncertainty set with corresponding atomic uncertainty set \mathcal{U}_0 . Then the min-max problem

$$\min_{\mathbf{w}, b} \sup_{(\mathbf{u}_1, \dots, \mathbf{u}_m) \in \mathcal{U}} \left\{ r(\mathbf{w}, b) + \sum_{i=1}^m \max [1 - y_i(\langle \mathbf{w}, \mathbf{x}_i - \mathbf{u}_i \rangle + b), 0] \right\} \quad (1.6)$$

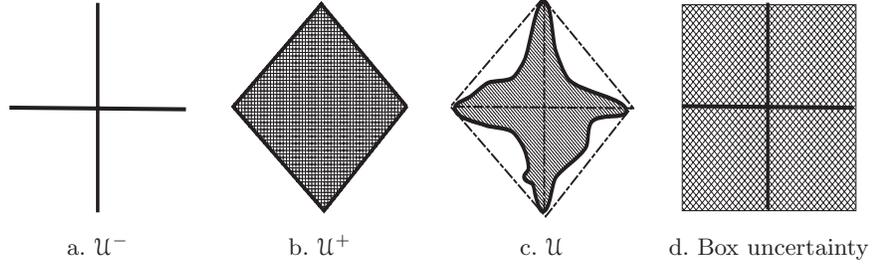


Figure 1.1: Illustration of a sublinear aggregated uncertainty set \mathcal{U} , and the contrast with the box uncertainty set.

is equivalent to the following optimization problem on \mathbf{w}, b, ξ :

$$\begin{aligned}
 \min_{\mathbf{w}, b, \xi} : \quad & r(\mathbf{w}, b) + \sup_{\mathbf{u} \in \mathcal{U}_0} (\mathbf{w}^\top \mathbf{u}) + \sum_{i=1}^m \xi_i, \\
 \text{s.t.} : \quad & \xi_i \geq 1 - [y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b)], \quad i = 1, \dots, m; \\
 & \xi_i \geq 0, \quad i = 1, \dots, m.
 \end{aligned} \tag{1.7}$$

The minimization of Problem (1.7) is attainable when $r(\cdot, \cdot)$ is lower semi-continuous.

Proof. We give only the proof idea. The details can be found in Xu et al. (2009). Define

$$v(\mathbf{w}, b) \triangleq \sup_{\mathbf{u} \in \mathcal{U}_0} (\mathbf{w}^\top \mathbf{u}) + \sum_{i=1}^m \max [1 - y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b), 0].$$

In the first step, we show

$$v(\hat{\mathbf{w}}, \hat{b}) \leq \sup_{(\mathbf{u}_1, \dots, \mathbf{u}_m) \in \mathcal{U}^-} \sum_{i=1}^m \max [1 - y_i(\langle \hat{\mathbf{w}}, \mathbf{x}_i - \mathbf{u}_i \rangle + \hat{b}), 0]. \tag{1.8}$$

This follows because the samples are non-separable. In the second step, we prove the reverse inequality:

$$\sup_{(\mathbf{u}_1, \dots, \mathbf{u}_m) \in \mathcal{U}^+} \sum_{i=1}^m \max [1 - y_i(\langle \hat{\mathbf{w}}, \mathbf{x}_i - \mathbf{u}_i \rangle + \hat{b}), 0] \leq v(\hat{\mathbf{w}}, \hat{b}). \tag{1.9}$$

This holds regardless of separability. Combining the two, adding the regularizer, and then infimizing both sides concludes the proof. \square

An immediate corollary is that a special case of our robust formulation is equivalent to the norm-regularized SVM setup:

Corollary 1.5. Let $\mathcal{T} \triangleq \{(\mathbf{u}_1, \dots, \mathbf{u}_m) \mid \sum_{i=1}^m \|\mathbf{u}_i\|^* \leq c\}$, where $\|\cdot\|^*$ stands for the dual norm of $\|\cdot\|$. If the training samples $\{\mathbf{x}_i, y_i\}_{i=1}^m$ are non-separable, then the following two optimization problems on (\mathbf{w}, b) are equivalent.

$$\min_{\mathbf{w}, b} : \quad \max_{(\mathbf{u}_1, \dots, \mathbf{u}_m) \in \mathcal{T}} \sum_{i=1}^m \max [1 - y_i (\langle \mathbf{w}, \mathbf{x}_i - \mathbf{u}_i \rangle + b), 0], \quad (1.10)$$

$$\min_{\mathbf{w}, b} : \quad c \|\mathbf{w}\| + \sum_{i=1}^m \max [1 - y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b), 0]. \quad (1.11)$$

Proof. Let \mathcal{U}_0 be the dual-norm ball $\{\mathbf{u} \mid \|\mathbf{u}\|^* \leq c\}$ and $r(\mathbf{w}, b) \equiv 0$. Then $\sup_{\|\mathbf{u}\|^* \leq c} (\mathbf{w}^\top \mathbf{u}) = c \|\mathbf{w}\|$. The corollary follows from Theorem 1.4. Notice that the equivalence holds for any \mathbf{w} and b . \square

This corollary explains the common belief that regularized classifiers tend to be more robust. Specifically, it explains the observation that when the disturbance is noise-like and neutral rather than adversarial, a norm-regularized classifier (without explicit robustness) has a performance often superior to a *box-type* robust classifier (see Trafalis and Gilbert, 2007). One take-away message is that while robust optimization is in its formulation adversarial, it can be quite flexible, and can be designed to yield solutions, such as the regularized solution above, that are appropriate for a non-adversarial setting.

One interesting research direction is to use this equivalence to find good regularizers without the need for cross validation. This could be done by mapping a measure of the variation in the training data to an appropriate uncertainty set, and then using the above equivalence to map back to a regularizer.

1.4.1.1 Kernelization

The previous results can be easily generalized to the kernelized setting. The kernelized SVM formulation considers a linear classifier in the feature space \mathcal{H} , a Hilbert space containing the range of some feature mapping $\Phi(\cdot)$. The standard formulation is as follows,

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} : \quad & r(\mathbf{w}, b) + \sum_{i=1}^m \xi_i \\ \text{s.t.} : \quad & \xi_i \geq [1 - y_i (\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b)], \quad i = 1, \dots, m; \\ & \xi_i \geq 0, \quad i = 1, \dots, m; \end{aligned}$$

where we use the representer theorem (see Schölkopf and Smola (2002)).

The definitions of an atomic uncertainty set and a sublinear aggregated uncertainty set in the feature space are identical to Definitions 1.2 and 1.3, with \mathbb{R}^n replaced by \mathcal{H} . The following theorem is a feature-space counterpart of Theorem 1.4, and the proof follows from a similar argument.

Theorem 1.6. *Assume $\{\Phi(\mathbf{x}_i), y_i\}_{i=1}^m$ are not linearly separable, $r(\cdot) : \mathcal{H} \times \mathbb{R} \rightarrow \mathbb{R}$ is an arbitrary function, $\mathcal{U} \subseteq \mathcal{H}^m$ is a sublinear aggregated uncertainty set with corresponding atomic uncertainty set $\mathcal{U}_0 \subseteq \mathcal{H}$. Then the following min-max problem*

$$\min_{\mathbf{w}, b} \sup_{(\mathbf{u}_1, \dots, \mathbf{u}_m) \in \mathcal{U}} \left\{ r(\mathbf{w}, b) + \sum_{i=1}^m \max [1 - y_i(\langle \mathbf{w}, \Phi(\mathbf{x}_i) - \mathbf{u}_i \rangle + b), 0] \right\}$$

is equivalent to

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} : \quad & r(\mathbf{w}, b) + \sup_{\mathbf{u} \in \mathcal{U}_0} (\langle \mathbf{w}, \mathbf{u} \rangle) + \sum_{i=1}^m \xi_i, \\ \text{s.t. :} \quad & \xi_i \geq 1 - y_i(\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b), \quad i = 1, \dots, m; \\ & \xi_i \geq 0, \quad i = 1, \dots, m. \end{aligned} \quad (1.12)$$

The minimization of Problem (1.12) is attainable when $r(\cdot, \cdot)$ is lower semi-continuous.

For some widely used feature mappings (e.g., RKHS of a Gaussian kernel), $\{\Phi(\mathbf{x}_i), y_i\}_{i=1}^m$ are always separable. In this case, the equivalence reduces to a bound.

The next corollary is the feature-space counterpart of Corollary 1.5, where $\|\cdot\|_{\mathcal{H}}$ stands for the RKHS norm, i.e., for $\mathbf{z} \in \mathcal{H}$, $\|\mathbf{z}\|_{\mathcal{H}} = \sqrt{\langle \mathbf{z}, \mathbf{z} \rangle}$.

Corollary 1.7. *Let $\mathcal{T}_{\mathcal{H}} \triangleq \{(\mathbf{u}_1, \dots, \mathbf{u}_m) \mid \sum_{i=1}^m \|\mathbf{u}_i\|_{\mathcal{H}} \leq c\}$. If $\{\Phi(\mathbf{x}_i), y_i\}_{i=1}^m$ are non-separable, then the following two optimization problems on (\mathbf{w}, b) are equivalent*

$$\begin{aligned} \min_{\mathbf{w}, b} : \quad & \max_{(\mathbf{u}_1, \dots, \mathbf{u}_m) \in \mathcal{T}_{\mathcal{H}}} \sum_{i=1}^m \max [1 - y_i(\langle \mathbf{w}, \Phi(\mathbf{x}_i) - \mathbf{u}_i \rangle + b), 0], \\ \min_{\mathbf{w}, b} : \quad & c\|\mathbf{w}\|_{\mathcal{H}} + \sum_{i=1}^m \max [1 - y_i(\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b), 0]. \end{aligned} \quad (1.13)$$

Equation (1.13) is a variant form of the standard SVM that has a squared RKHS norm regularization term, and by convexity arguments the two formulations are equivalent up to a change of tradeoff parameter c . Therefore, Corollary 1.7 essentially means that the standard kernelized SVM is im-

licitly a robust classifier (without regularization) with disturbance in the feature-space, and the sum of the magnitude of the disturbance is bounded.

Disturbance in the feature-space is less intuitive than disturbance in the sample space, and the next lemma relates these two different notions.

Lemma 1.8. *Suppose there exists $\mathcal{X} \subseteq \mathbb{R}^n$, $\rho > 0$, and a continuous non-decreasing function $f : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ satisfying $f(0) = 0$, such that*

$$k(\mathbf{x}, \mathbf{x}) + k(\mathbf{x}', \mathbf{x}') - 2k(\mathbf{x}, \mathbf{x}') \leq f(\|\mathbf{x} - \mathbf{x}'\|_2^2), \quad \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}, \|\mathbf{x} - \mathbf{x}'\|_2 \leq \rho.$$

Then,

$$\|\Phi(\hat{\mathbf{x}} + \mathbf{u}) - \Phi(\hat{\mathbf{x}})\|_{\mathcal{H}} \leq \sqrt{f(\|\mathbf{u}\|_2^2)}, \quad \forall \|\mathbf{u}\|_2 \leq \rho, \hat{\mathbf{x}}, \hat{\mathbf{x}} + \boldsymbol{\delta} \in \mathcal{X}.$$

Lemma 1.8 essentially says that under certain conditions, robustness in the feature space is a stronger requirement than robustness in the sample space. Therefore, a classifier that achieves robustness in the feature space also achieves robustness in the sample space. Notice that the condition of Lemma 1.8 is rather weak. In particular, it holds for any continuous $k(\cdot, \cdot)$ and bounded domain \mathcal{X} .

1.4.1.2 Probabilistic Interpretations

As discussed and demonstrated above, robust optimization can often be used for *probabilistic analysis*. In this section, we show that robust optimization and the equivalence theorem can be used to construct a classifier with *probabilistic margin protection*, i.e., a classifier with probabilistic constraints on the chance of violation beyond a given threshold. Second, we show that in the Bayesian setup, if one has a prior only on the total magnitude of the disturbance vector, robust optimization can be used to tune the regularizer.

Probabilistic Protection. We can use Problem (1.6) to obtain an upper bound for a chance-constrained classifier. Suppose the disturbance is stochastic with known distribution. We denote the disturbance vector by $(\mathbf{u}_1^r, \dots, \mathbf{u}_m^r)$ to emphasize that it is now a random variable. The chance-constrained classifier minimizes the hinge loss that occurs with probability above some given confidence level $\eta \in [0, 1]$. The classifier is given by the optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, b, l} : & \quad l & (1.14) \\ \text{s.t.} : & \quad \mathbb{P}\left\{\sum_{i=1}^m \max[1 - y_i(\langle \mathbf{w}, \mathbf{x}_i - \mathbf{u}_i^r \rangle + b), 0] \leq l\right\} \geq 1 - \eta. \end{aligned}$$

The constraint controls the η -quantile of the average (or equivalently the sum of) empirical error. In Shivaswamy et al. (2006), Lanckriet et al. (2003) and Bhattacharyya et al. (2004), the authors explore a different direction, and starting from the constraint formulation of SVM as in (1.2), they impose probabilistic constraints on each random variable individually. This formulation requires all constraints to be satisfied with high probability *simultaneously*. Thus, instead of controlling the η -quantile of the average loss, they control the η -quantile of the hinge-loss for each sample. For the same reason that box uncertainty in the robust setting may be too conservative, this constraint-wise formulation may also be too conservative.

Problem (1.14) is generally intractable. However, we can approximate it as follows. Let

$$\hat{c} \triangleq \inf \left\{ \alpha \mid \mathbb{P} \left(\sum_{i=1}^m \|\mathbf{u}_i\|^* \leq \alpha \right) \geq 1 - \eta \right\}.$$

Notice that \hat{c} is easily simulated given μ . Then for any (\mathbf{w}, b) , with probability no less than $1 - \eta$, the following holds,

$$\begin{aligned} & \sum_{i=1}^m \max [1 - y_i (\langle \mathbf{w}, \mathbf{x}_i - \mathbf{u}_i^r \rangle + b), 0] \\ & \leq \max_{\sum_i \|\mathbf{u}_i\|^* \leq \hat{c}} \sum_{i=1}^m \max [1 - y_i (\langle \mathbf{w}, \mathbf{x}_i - \mathbf{u}_i \rangle + b), 0]. \end{aligned}$$

Thus (1.14) is upper bounded by (1.11) with $c = \hat{c}$. This gives an additional probabilistic robustness property of the standard regularized classifier. We observe that we can follow a similar approach using the constraint-wise robust setup, i.e., the box uncertainty set. The interested reader can check that this would lead to considerably more pessimistic approximations of the chance constraint.

A Bayesian Regularizer. Next, we show how the above can be used in a Bayesian setup, to obtain an appropriate regularization coefficient. Suppose the total disturbance $c^r \triangleq \sum_{i=1}^m \|\mathbf{u}_i^r\|^*$ is a random variable, and follows a prior distribution $\rho(\cdot)$. This can model for example the case that the training sample set is a mixture of several data sets where the disturbance magnitude of each set is known. Such a setup leads to the following classifier which minimizes the Bayesian (robust) error:

$$\min_{\mathbf{w}, b} : \int \left\{ \max_{\sum \|\delta_i\|^* \leq c} \sum_{i=1}^m \max [1 - y_i (\langle \mathbf{w}, \mathbf{x}_i - \mathbf{u}_i \rangle + b), 0] \right\} d\rho(c). \quad (1.15)$$

By Corollary 1.5, the Bayesian classifier (1.15) is equivalent to

$$\min_{\mathbf{w}, b} : \int \left\{ c \|\mathbf{w}\| + \sum_{i=1}^m \max [1 - y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b), 0] \right\} d\rho(c),$$

which can be further simplified as

$$\min_{\mathbf{w}, b} : \bar{c} \|\mathbf{w}\| + \sum_{i=1}^m \max [1 - y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b), 0],$$

where $\bar{c} \triangleq \int c d\rho(c)$. This provides a justifiable parameter tuning method different from cross validation: simply using the expected value of c^r .

1.4.2 Tikhonov regularized ℓ_2 regression

We now move from classification and SVMs, to regression, and show that ℓ_2 -regularized regression, like SVM, is equivalent to a robust optimization problem. This equivalence is then used to define new regularization-like algorithms, and also to prove properties of the regularized solution.

Given input output pairs \mathbf{x}_i, y_i forming the rows of X and the elements of vector \mathbf{y} , respectively, the goal is to find a predictor $\boldsymbol{\beta}$ that minimizes the squared loss $\|\mathbf{y} - X\boldsymbol{\beta}\|_2^2$. As is well-known, this problem is often notoriously ill-conditioned, and may not have a unique solution. The classical, and much-explored remedy has been, as in the SVM case, regularization. Regularizing with an ℓ_2 -norm, known in statistics as ridge regression (Hoerl, 1962), and in analysis as Tikhonov regularization (Tikhonov and Arsenin, 1977), solves the problem¹

$$\min_{\boldsymbol{\beta}} : \|\mathbf{y} - X\boldsymbol{\beta}\|_2 + \lambda \|\boldsymbol{\beta}\|_2. \quad (1.16)$$

The main result of this section states that Tikhonov regularized regression is the solution to a robust optimization, where X is subject to matrix-disturbance U with a bounded Frobenius norm.

Theorem 1.9. *The following robust optimization formulation*

$$\min_{\boldsymbol{\beta}} : \max_{U: \|U\|_F \leq \lambda} \|\mathbf{y} - (X + U)\boldsymbol{\beta}\|_2,$$

is equivalent to Tikhonov-regularized regression (1.16).

Proof. For any perturbation U , we have $\|\mathbf{y} - (X + U)\boldsymbol{\beta}\|_2 = \|\mathbf{y} - X\boldsymbol{\beta} - U\boldsymbol{\beta}\|_2$.

1. This problem is equivalent to one where we square the norm, up to a change in the regularization coefficient, λ .

By the triangle inequality and because $\|U\|_F \leq \lambda$, we thus have $\|\mathbf{y} - (X + U)\boldsymbol{\beta}\|_2 \leq \|\mathbf{y} - X\boldsymbol{\beta}\|_2 + \lambda\|\boldsymbol{\beta}\|_2$. On the other hand, for any given $\boldsymbol{\beta}$, we can choose a rank one U so that $U\boldsymbol{\beta}$ is aligned with $(\mathbf{y} - X\boldsymbol{\beta})$, and thus equality is attained. \square

This connection was first explored in the seminal work of El Ghaoui and Lebret (1997). There, they further show that the solution to the robust counterpart is almost as easily determined as that to the nominal problem: one need only perform a line search, in the case where the SVD of A is available. Thus, the computational cost of the robust regression is comparable to the original formulation.

As with SVMs, the “hidden robustness” has several consequences. By changing the uncertainty set, robust optimization allows for a rich class of regularization-like algorithms. Motivated by problems from robust control, El Ghaoui and Lebret (1997) then consider perturbations that have structure, leading to structured robust least squares problems. They then analyze tractability and approximations to these structured least squares.² Finally, they use the robustness equivalence to prove regularity properties of the solution. We refer to El Ghaoui and Lebret (1997) for further details about structured robustness, tractability, and regularity.

1.4.3 Lasso

In this section, we consider a similar problem: ℓ_1 -regularized regression, also known as Lasso (Tibshirani, 1996). Lasso has been explored extensively for its remarkable sparsity properties (e.g., Tibshirani, 1996; Bickel et al., 2009; Wainwright, 2009) most recently under the banner of compressed sensing (e.g., Chen et al., 1999; Candès et al., 2006; Candès and Tao, 2006; Candès and Tao, 2007; Candès and Tao, 2008; Donoho, 2006, for an incomplete list). Following the theme of this section, we show that the solution to Lasso is the solution to a robust optimization problem. As with Tikhonov regularization, robustness provides a connection of the regularizer to a physical property, namely, protection from noise. This allows a principled selection of the regularizer. Moreover, by considering different uncertainty sets, we obtain generalizations of Lasso. Next, we go on to show that robustness can itself be used as an avenue for exploring different properties of the solution. In particular, we show that robustness explains why the solution is sparse – that

2. Note that arbitrary uncertainty sets may lead to intractable problems. This is because the inner maximization in the robust formulation is of a convex function, and hence is nonconvex.

is, *Lasso is sparse because it is robust*. The analysis as well as the specific results obtained differ from standard sparsity results, providing different geometric intuition. This section is based on results reported in Xu et al. (2010a), where full proofs to all stated results can be found.

Lasso, or ℓ_1 -regularized regression, has a similar form to ridge regression, differing only in the regularizer: ³

$$\min : \|\mathbf{y} - X\boldsymbol{\beta}\|_2 + \lambda\|\boldsymbol{\beta}\|_1.$$

For a general uncertainty set \mathcal{U} , using the same notation as in the previous section, the robust regression formulation becomes

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^m} \max_{U \in \mathcal{U}} \|\mathbf{y} - (X + U)\boldsymbol{\beta}\|_2, \quad (1.17)$$

In the previous section, the uncertainty set was $\mathcal{U} = \{U : \|U\|_F \leq \lambda\}$. We consider a different uncertainty set here. Writing

$$U = \begin{bmatrix} | & | & \cdots & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_m \\ | & | & \cdots & | \end{bmatrix}, \quad \text{where } (\mathbf{u}_1, \dots, \mathbf{u}_m) \in \mathcal{U},$$

let the uncertainty set \mathcal{U} have the form:

$$\mathcal{U} \triangleq \left\{ (\mathbf{u}_1, \dots, \mathbf{u}_m) \mid \|\mathbf{u}_i\|_2 \leq c_i, \quad i = 1, \dots, m \right\}. \quad (1.18)$$

This is a *feature-wise uncoupled* uncertainty set: the uncertainty in different features need not satisfy any joint constraints. In contrast, the constraint $\|U\|_F \leq 1$ used in the previous section is feature-wise coupled. We revisit coupled uncertainty sets below.

Theorem 1.10. *The robust regression problem (1.17) with uncertainty set of the form (1.18) is equivalent to the following ℓ_1 regularized regression problem:*

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^m} \left\{ \|\mathbf{y} - X\boldsymbol{\beta}\|_2 + \sum_{i=1}^m c_i |\beta_i| \right\}. \quad (1.19)$$

Proof. Fix $\boldsymbol{\beta}^*$. We prove that $\max_{U \in \mathcal{U}} \|\mathbf{y} - (X + U)\boldsymbol{\beta}^*\|_2 = \|\mathbf{y} - X\boldsymbol{\beta}^*\|_2 + \sum_{i=1}^m c_i |\beta_i^*|$.

3. Again we remark that with a change of regularization parameter, this is equivalent to the more common form appearing with a square outside the norm.

The inequality

$$\max_{U \in \mathcal{U}} \|\mathbf{y} - (X + U)\boldsymbol{\beta}^*\|_2 \leq \|\mathbf{y} - X\boldsymbol{\beta}^*\|_2 + \sum_{i=1}^m |\beta_i^*| c_i,$$

follows from the triangle inequality, as in our proof in the previous section. The other inequality follows, if we take

$$\mathbf{u} \triangleq \begin{cases} \frac{\mathbf{y} - X\boldsymbol{\beta}^*}{\|\mathbf{y} - X\boldsymbol{\beta}^*\|_2} & \text{if } X\boldsymbol{\beta}^* \neq \mathbf{y}, \\ \text{any vector with unit } \ell_2 \text{ norm} & \text{otherwise;} \end{cases}$$

and let

$$\mathbf{u}_i^* \triangleq \begin{cases} -c_i \text{sgn}(\beta_i^*) \mathbf{u} & \text{if } \beta_i^* \neq 0; \\ -c_i \mathbf{u} & \text{otherwise.} \end{cases}$$

□

Taking $c_i = c$ and normalizing \mathbf{x}_i for all i , Problem (1.19) recovers the well-known Lasso (Tibshirani, 1996; Efron et al., 2004).

1.4.3.1 General Uncertainty Sets

Using this equivalence, we generalize to Lasso-like regularization algorithms in two ways: (a) to the case of arbitrary norm; and (b) to the case of coupled uncertainty sets.

Theorem 1.11. *For $\|\cdot\|_a$ an arbitrary norm in the Euclidean space, the robust regression problem*

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^m} \left\{ \max_{U \in \mathcal{U}_a} \|\mathbf{y} - (X + U)\boldsymbol{\beta}\|_a \right\};$$

where

$$\mathcal{U}_a \triangleq \left\{ (\mathbf{u}_1, \dots, \mathbf{u}_m) \mid \|\mathbf{u}_i\|_a \leq c_i, \quad i = 1, \dots, m \right\};$$

is equivalent to the following regularized regression problem

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^m} \left\{ \|\mathbf{y} - X\boldsymbol{\beta}\|_a + \sum_{i=1}^m c_i |\beta_i| \right\}.$$

We next consider feature-wise coupled uncertainty sets. This can be used to incorporate additional information about potential noise in the problem, when available, to limit the conservativeness of the worst-case formulation.

Consider the following uncertainty set:

$$\mathcal{U}' \triangleq \{(\mathbf{u}_1, \dots, \mathbf{u}_m) \mid f_j(\|\mathbf{u}_1\|_a, \dots, \|\mathbf{u}_m\|_a) \leq 0; j = 1, \dots, k\},$$

where each $f_j(\cdot)$ is a convex function. The resulting robust formulation is equivalent to a more general regularization-type problem, and moreover, it is tractable.

Theorem 1.12. *Let \mathcal{U}' be as above, and assume that the set*

$$\mathcal{Z} \triangleq \{\mathbf{z} \in \mathbb{R}^m \mid f_j(\mathbf{z}) \leq 0, j = 1, \dots, k; \mathbf{z} \geq \mathbf{0}\},$$

has non-empty relative interior. Then the robust regression problem

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^m} \left\{ \max_{U \in \mathcal{U}'} \|\mathbf{y} - (X + U)\boldsymbol{\beta}\|_a \right\},$$

is equivalent to the following regularized regression problem

$$\min_{\boldsymbol{\lambda} \in \mathbb{R}_+^k, \boldsymbol{\kappa} \in \mathbb{R}_+^m, \boldsymbol{\beta} \in \mathbb{R}^m} \left\{ \|\mathbf{y} - X\boldsymbol{\beta}\|_a + v(\boldsymbol{\lambda}, \boldsymbol{\kappa}, \boldsymbol{\beta}) \right\}; \quad (1.20)$$

$$\text{where: } v(\boldsymbol{\lambda}, \boldsymbol{\kappa}, \boldsymbol{\beta}) \triangleq \max_{\mathbf{c} \in \mathbb{R}^m} \left[(\boldsymbol{\kappa} + |\boldsymbol{\beta}|)^\top \mathbf{c} - \sum_{j=1}^k \lambda_j f_j(\mathbf{c}) \right],$$

and in particular, is efficiently solvable.

The next two corollaries are a direct application of Theorem 1.12.

Corollary 1.13. *Suppose*

$$\mathcal{U}' = \left\{ (\boldsymbol{\delta}_1, \dots, \boldsymbol{\delta}_m) \mid \|\boldsymbol{\delta}_1\|_a, \dots, \|\boldsymbol{\delta}_m\|_a \leq l \right\},$$

for arbitrary norms $\|\cdot\|_a$ and $\|\cdot\|_s$. Then the robust problem is equivalent to the regularized regression problem

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^m} \left\{ \|\mathbf{y} - X\boldsymbol{\beta}\|_a + l \|\boldsymbol{\beta}\|_s^* \right\};$$

where $\|\cdot\|_s^$ is the dual norm of $\|\cdot\|_s$.*

This corollary interprets *arbitrary* norm-based regularizers from a robust regression perspective. For example, taking both $\|\cdot\|_a$ and $\|\cdot\|_s$ to be the Euclidean norm, then \mathcal{U}' is the set of matrices with bounded Frobenius norm, and Corollary 1.13 recovers Theorem 1.9.

The next corollary considers general polytope uncertainty sets, where the column-wise norm vector of the realizable uncertainty belongs to a polytope. To illustrate the flexibility and potential use of such an uncertainty set: taking $\|\cdot\|_a$ to be the ℓ_1 norm and the polytope to be the standard simplex,

the resulting uncertainty set consists of matrices with bounded $\|\cdot\|_{2,1}$ -norm. This is the ℓ_1 -norm of the ℓ_2 -norm of the columns, and has numerous applications, including, e.g., outlier removal (Xu et al., 2010c).

Corollary 1.14. *Suppose*

$$\mathcal{U}' = \left\{ (\mathbf{u}_1, \dots, \mathbf{u}_m) \mid T\mathbf{c} \leq \mathbf{s}; \text{ where: } c_j = \|\mathbf{u}_j\|_a \right\},$$

for a given matrix T , vector \mathbf{s} , and arbitrary norm $\|\cdot\|_a$. Then the robust regression is equivalent to the following regularized regression problem with variables $\boldsymbol{\beta}$ and $\boldsymbol{\lambda}$:

$$\begin{aligned} \min_{\boldsymbol{\beta}, \boldsymbol{\lambda}} : & \quad \|\mathbf{y} - X\boldsymbol{\beta}\|_a + \mathbf{s}^\top \boldsymbol{\lambda} \\ \text{s.t.} & \quad \boldsymbol{\beta} \leq T^\top \boldsymbol{\lambda}; \\ & \quad -\boldsymbol{\beta} \leq T^\top \boldsymbol{\lambda}; \\ & \quad \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

1.4.3.2 Sparsity

In this section, we investigate the sparsity properties of robust regression, and show in particular, that Lasso is sparse *because it is robust*. This new connection between robustness and sparsity suggests that robustifying with respect to a feature-wise independent uncertainty set might be a plausible way to achieve sparsity for other problems.

We show that if there is any perturbation in the uncertainty set that makes some feature *irrelevant*, i.e., not contributing to the regression error, then the optimal robust solution puts no weight there. Thus if the features in an index set $I \subset \{1, \dots, n\}$ can be perturbed to be made irrelevant, then the solution will be supported on the complement, I^c .

To state the main theorem of this section, we introduce some notation. Given an index subset $I \subseteq \{1, \dots, n\}$, and a matrix U , let U^I denote the restriction of U to feature set I , i.e., U^I equals U on each feature indexed by $i \in I$, and is zero elsewhere. Similarly, given a feature-wise uncoupled uncertainty set \mathcal{U} , let \mathcal{U}^I be the restriction of \mathcal{U} to the feature set I , i.e., $\mathcal{U}^I \triangleq \{U^I \mid U \in \mathcal{U}\}$. Any element $U \in \mathcal{U}$ can be written as $U^I + U^{I^c}$ (here $I^c \triangleq \{1, \dots, n\} \setminus I$) with $U^I \in \mathcal{U}^I$ and $U^{I^c} \in \mathcal{U}^{I^c}$.

Theorem 1.15. *The robust regression problem*

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^m} \left\{ \max_{\Delta A \in \mathcal{U}} \|\mathbf{y} - (X + U)\boldsymbol{\beta}\|_2 \right\}, \quad (1.21)$$

has a solution supported on an index set I if there exists some perturbation

$\tilde{U} \in \mathcal{U}^{I^c}$, such that the robust regression problem

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^m} \left\{ \max_{U \in \mathcal{U}^I} \|\mathbf{y} - (X + \tilde{U} + U)\boldsymbol{\beta}\|_2 \right\}, \quad (1.22)$$

has a solution supported on the set I .

Theorem 1.15 is a special case of the following theorem with $c_j = 0$ for all $j \notin I$.

Theorem 1.15’: Let $\boldsymbol{\beta}^*$ be an optimal solution of the robust regression problem:

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^m} \left\{ \max_{U \in \mathcal{U}} \|\mathbf{y} - (X + U)\boldsymbol{\beta}\|_2 \right\}, \quad (1.23)$$

and let $I \subseteq \{1, \dots, m\}$ be such that $\beta_j^* = 0 \forall j \notin I$. Let

$$\tilde{\mathcal{U}} \triangleq \left\{ (\mathbf{u}_1, \dots, \mathbf{u}_m) \mid \|\mathbf{u}_i\|_2 \leq c_i, \quad i \in I; \quad \|\mathbf{u}_j\|_2 \leq c_j + l_j, \quad j \notin I \right\}.$$

Then, $\boldsymbol{\beta}^*$ is an optimal solution of

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^m} \left\{ \max_{U \in \tilde{\mathcal{U}}} \|\mathbf{y} - (\tilde{X} + U)\boldsymbol{\beta}\|_2 \right\}, \quad (1.24)$$

for any \tilde{X} that satisfies $\|\tilde{\mathbf{x}}_j - \mathbf{x}_j\| \leq l_j$ for $j \notin I$, and $\tilde{\mathbf{x}}_i = \mathbf{x}_i$ for $i \in I$.

In fact, we can replace the ℓ_2 norm loss by any loss function $f(\cdot)$ which satisfies the condition that if $\beta_j = 0$, X and X' only differ in the j^{th} column, then $f(\mathbf{y}, X, \boldsymbol{\beta}) = f(\mathbf{y}, X', \boldsymbol{\beta})$. This theorem thus suggests a methodology for constructing sparse algorithms by solving a robust optimization with respect to column-wise uncoupled uncertainty sets.

When we consider ℓ_2 loss, we can translate the condition of a feature being “irrelevant” to a geometric condition: orthogonality. We now use the result of Theorem 1.15 to show that robust regression has a sparse solution as long as an incoherence-type property is satisfied. This result is more in line with the traditional sparsity results, but we note that the geometric reasoning is different, now based on robustness. Specifically: we show that a feature receives zero weight, if it is “nearly” (i.e., within an allowable perturbation) orthogonal to the signal, and all relevant features.

Theorem 1.16. *Let $c_i = c$ for all i and consider ℓ_2 loss. Suppose that there exists $I \subset \{1, \dots, m\}$ such that for all $\mathbf{v} \in \text{span}(\{\mathbf{x}_i, i \in I\} \cup \{\mathbf{y}\})$, $\|\mathbf{v}\| = 1$, we have $\mathbf{v}^\top \mathbf{x}_j \leq c$, $\forall j \notin I$. Then there exists an optimal solution $\boldsymbol{\beta}^*$ that satisfies $\beta_j^* = 0$, $\forall j \notin I$.*

The proof proceeds as the previous theorem would suggest: the columns in I^c can be perturbed to be made irrelevant, and thus the optimal solution

will not be supported there; see Xu et al. (2010a) for details..

1.5 Robustness and Consistency

In this section we explore a fundamental connection between learning and robustness, by using robustness properties to re-prove the consistency of kernelized SVM, and then of Lasso. The key difference from the proofs here and those seen elsewhere (e.g., Steinwart, 2005; Steinwart and Christmann, 2008; Wainwright, 2009), is that we replace the metric entropy, VC-dimension, and stability conditions typically used, with a robustness condition. Thus we conclude that *SVM and Lasso are consistent because they are robust.*

1.5.1 Consistency of SVM

Let $\mathcal{X} \subseteq \mathbb{R}^n$ be bounded, and suppose the training samples $(\mathbf{x}_i, y_i)_{i=1}^{\infty}$ are generated according to an unknown i.i.d. distribution \mathbb{P} supported on $\mathcal{X} \times \{-1, +1\}$. The next theorem shows that our robust classifier and thus regularized SVM, asymptotically minimizes an upper-bound of the expected classification error and hinge loss, as the number of samples increases.

Theorem 1.17. *Let $K \triangleq \max_{x \in \mathcal{X}} \|x\|_2$. Then there exists a random sequence $\{\gamma_{m,c}\}$ such that:*

1. *The following bounds on the Bayes loss and the hinge loss hold uniformly for all (\mathbf{w}, b) :*

$$\begin{aligned} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}}(\mathbf{1}_{y \neq \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle + b)}) &\leq \gamma_{m,c} + c\|\mathbf{w}\|_2 + \frac{1}{m} \sum_{i=1}^m \max [1 - y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b), 0]; \\ \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}}(\max(1 - y(\langle \mathbf{w}, \mathbf{x} \rangle + b), 0)) &\leq \\ &\gamma_{m,c}(1 + K\|\mathbf{w}\|_2 + |b|) + c\|\mathbf{w}\|_2 + \frac{1}{m} \sum_{i=1}^m \max [1 - y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b), 0]. \end{aligned}$$

2. *For every $c > 0$, $\lim_{m \rightarrow \infty} \gamma_{m,c} = 0$ almost surely, and the convergence is uniform in \mathbb{P} ;*

Proof. We outline the basic idea of the proof here, and refer to Xu et al. (2009) for the technical details. We consider the testing sample set as a perturbed copy of the training sample set, and measure the magnitude of the perturbation. For testing samples that have “small” perturbations, Corollary 1.5 guarantees that $c\|\mathbf{w}\|_2 + \frac{1}{m} \sum_{i=1}^m \max [1 - y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b), 0]$ upper-bounds their total loss. Therefore, we only need to show that the

fraction of testing samples having “large” perturbations diminishes to prove the theorem. We show this using a balls and bins argument. Partitioning $\mathcal{X} \times \{-1, +1\}$, we match testing and training samples that fall in the same partition. We then use the Bretagnolle-Huber-Carol inequality for multinomial distributions to conclude that the fraction of unmatched points diminishes to zero. \square

Based on Theorem 1.17, it can be further shown that the expected classification error of the solutions of SVM converges to the Bayes risk, i.e., SVM is consistent.

1.5.2 Consistency of Lasso

In this section, we reprove the asymptotic consistency of Lasso using robustness. The basic idea of the consistency proof is as follows. We show that the robust optimization formulation can be seen to have the maximum *expected error* w.r.t. a class of probability measures. This class includes a kernel density estimator, and using this, we show that Lasso is consistent.

1.5.2.1 Robust Optimization and Kernel Density Estimation

En route to proving consistency of Lasso based on robust optimization, we discuss another result of independent interest. We link robust optimization to worst case expected utility, i.e., the worst-case expectation over a set of measures. For the proofs, and more along this direction, we refer to Xu et al. (2010b,a). Throughout this section, we use \mathcal{P} to represent the set of all probability measures (on Borel σ -algebra) of \mathbb{R}^{m+1} .

We first establish a general result on the equivalence between a robust optimization formulation and a worst-case expected utility:

Proposition 1.18. *Given a function $f : \mathbb{R}^{m+1} \rightarrow \mathbb{R}$ and Borel sets $\mathcal{Z}_1, \dots, \mathcal{Z}_n \subseteq \mathbb{R}^{m+1}$, let*

$$\mathcal{P}_n \triangleq \{\mu \in \mathcal{P} \mid \forall S \subseteq \{1, \dots, n\} : \mu(\bigcup_{i \in S} \mathcal{Z}_i) \geq |S|/n\}.$$

The following holds:

$$\frac{1}{n} \sum_{i=1}^n \sup_{(\mathbf{x}_i, y_i) \in \mathcal{Z}_i} f(\mathbf{x}_i, y_i) = \sup_{\mu \in \mathcal{P}_n} \int_{\mathbb{R}^{m+1}} f(\mathbf{x}, y) d\mu(\mathbf{x}, y).$$

This leads to the following corollary for Lasso, which states that for a given solution β , the robust regression loss over the training data is equal to the worst-case expected *generalization error*.

Corollary 1.19. *Given $\mathbf{y} \in \mathbb{R}^n$, $X \in \mathbb{R}^{n \times m}$, the following equation holds for any $\boldsymbol{\beta} \in \mathbb{R}^m$,*

$$\|\mathbf{y} - X\boldsymbol{\beta}\|_2 + \sqrt{nc_n}(\|\boldsymbol{\beta}\|_1 + 1) = \sup_{\mu \in \hat{\mathcal{P}}(n)} \sqrt{n \int_{\mathbb{R}^{m+1}} (y' - \mathbf{x}'^\top \boldsymbol{\beta})^2 d\mu(\mathbf{x}', y')}. \quad (1.25)$$

Where we let x_{ij} and u_{ij} be the (i, j) -entry of X and U , respectively, and

$$\hat{\mathcal{P}}(n) \triangleq \bigcup_{\|\boldsymbol{\sigma}\|_2 \leq \sqrt{nc_n}; \forall i: \|\mathbf{u}_i\|_2 \leq \sqrt{nc_n}} \mathcal{P}_n(X, U, \mathbf{y}, \boldsymbol{\sigma});$$

$$\begin{aligned} \mathcal{P}_n(X, U, \mathbf{y}, \boldsymbol{\sigma}) \triangleq & \{ \mu \in \mathcal{P} | \mathcal{Z}_i = [y_i - \sigma_i, y_i + \sigma_i] \times \prod_{j=1}^m [x_{ij} - u_{ij}, x_{ij} + u_{ij}]; \\ & \forall S \subseteq \{1, \dots, n\} : \mu(\bigcup_{i \in S} \mathcal{Z}_i) \geq |S|/n \}. \end{aligned}$$

The proof of consistency relies on showing that this set $\hat{\mathcal{P}}(n)$ of distributions contains a kernel density estimator. Recall the basic definition: The *kernel density estimator* for a density h in \mathbb{R}^d , originally proposed in Rosenblatt (1956) and Parzen (1962), is defined by

$$h_n(\mathbf{x}) = (nc_n^d)^{-1} \sum_{i=1}^n K\left(\frac{\mathbf{x} - \hat{\mathbf{x}}_i}{c_n}\right),$$

where $\{c_n\}$ is a sequence of positive numbers, $\hat{\mathbf{x}}_i$ are i.i.d. samples generated according to h , and K is a Borel measurable function (kernel) satisfying $K \geq 0$, $\int K = 1$. See Devroye and Györfi (1985); Scott (1992) and references therein for detailed discussions. A celebrated property of a kernel density estimator is that it converges in \mathcal{L}^1 to h when $c_n \downarrow 0$ and $nc_n^d \uparrow \infty$ (Devroye and Györfi, 1985).

1.5.2.2 Density Estimation and Consistency of Lasso

We now use robustness of Lasso to prove its consistency. Throughout, we use c_n to represent the robustness level c where there are n samples. We take c_n to zero as n grows.

Recall the standard generative model in statistical learning: let \mathbb{P} be a probability measure with bounded support that generates i.i.d. samples (y_i, \mathbf{x}_i) , and has a density $f^*(\cdot)$. Denote the set of the first n samples by

\mathcal{S}_n . Define

$$\begin{aligned}\beta(c_n, \mathcal{S}_n) &\triangleq \arg \min_{\beta} \left\{ \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \beta)^2} + c_n \|\beta\|_1 \right\} \\ &= \arg \min_{\beta} \left\{ \frac{\sqrt{n}}{n} \sqrt{\sum_{i=1}^n (y_i - \mathbf{x}_i^\top \beta)^2} + c_n \|\beta\|_1 \right\}; \\ \beta(\mathbb{P}) &\triangleq \arg \min_{\beta} \left\{ \sqrt{\int_{y, \mathbf{x}} (y - \mathbf{x}^\top \beta)^2 d\mathbb{P}(y, \mathbf{x})} \right\}.\end{aligned}$$

In words, $\beta(c_n, \mathcal{S}_n)$ is the solution to Lasso with the tradeoff parameter set to $c_n \sqrt{n}$, and $\beta(\mathbb{P})$ is the “true” optimal solution. We establish that $\beta(c_n, \mathcal{S}_n) \rightarrow \beta(\mathbb{P})$ using robustness.

Theorem 1.20. *Let $\{c_n\}$ be such that $c_n \downarrow 0$ and $\lim_{n \rightarrow \infty} n(c_n)^{m+1} = \infty$. Suppose there exists a constant H such that $\|\beta(c_n, \mathcal{S}_n)\|_2 \leq H$ for all n . Then,*

$$\lim_{n \rightarrow \infty} \sqrt{\int_{y, \mathbf{x}} (y - \mathbf{x}^\top \beta(c_n, \mathcal{S}_n))^2 d\mathbb{P}(y, \mathbf{x})} = \sqrt{\int_{y, \mathbf{x}} (y - \mathbf{x}^\top \beta(\mathbb{P}))^2 d\mathbb{P}(y, \mathbf{x})},$$

almost surely.

We give an outline of the proof, and refer to Xu et al. (2010a) for the details. In Section 1.4.3 we showed that Lasso is a special case of robust optimization. Then in Section 1.5.2.1, we proved that robust optimization is equivalent to a worst-case expectation. The proof follows by showing that the sets \mathcal{P}_n in the worst-case expectation equivalent to Lasso, contain a kernel density estimator. Since these sets shrink, consistency follows.

The assumption that $\|\beta(c_n, \mathcal{S}_n)\|_2 \leq H$ can be removed. As in Theorem 1.20, the proof technique rather than the result itself is of interest. We refer the interested reader to Xu et al. (2010a).

1.6 Robustness and Generalization

We have already seen that regularized regression and regularized SVMs are a special case of robust optimization, and hence exhibit robustness to perturbed data. This robustness was used above to show that ridge regression has a Lipschitz solution, that Lasso is sparse, and SVM and Lasso are consistent. In this section, we show that robustness can be used to control the estimation of the risk (i.e., generalization error) of learning algorithms.

The results we describe are based on Xu and Mannor (2010b).

Several approaches have been proposed to bound the deviation of the risk from its empirical measurement, among which methods based on uniform convergence and stability are most widely used (e.g., Vapnik and Chervonenkis, 1991; Evgeniou et al., 2000; Alon et al., 1997; Bartlett, 1998; Bartlett and Mendelson, 2002; Bartlett et al., 2005; Bousquet and Elisseeff, 2002; Poggio et al., 2004; Mukherjee et al., 2006, and many others). We provide a new, robustness-driven approach to proving generalization bounds.

Whereas in the past sections, “robustness” was defined directly in terms of robust optimization, we abstract this definition here. Because we consider abstract algorithms in this section, we introduce some necessary notations, different from previous sections. We use \mathcal{Z} and \mathcal{H} to denote the set from which each sample is drawn, and the hypothesis set, respectively. Throughout this section we use $\mathbf{s} \in \mathcal{Z}^m$ to denote the training sample set consisting of m training samples (s_1, \dots, s_m) . A learning algorithm \mathcal{A} is thus a mapping from \mathcal{Z}^m to \mathcal{H} . We use $\mathcal{A}_{\mathbf{s}}$ to represent the hypothesis learned, given training set \mathbf{s} . For each hypothesis $h \in \mathcal{H}$ and point $z \in \mathcal{Z}$, there is an associated loss $l(h, z)$, which is non-negative and upper-bounded uniformly by a scalar M . In the special case of supervised learning, the sample space can be decomposed as $\mathcal{Z} = \mathcal{Y} \times \mathcal{X}$, and the goal is to learn a mapping from \mathcal{X} to \mathcal{Y} , i.e., to predict the y -component given x -component. We hence use $\mathcal{A}_{\mathbf{s}}(x)$ to represent the predicted y -component (label) of $x \in \mathcal{X}$ when \mathcal{A} is trained on \mathbf{s} . We call \mathcal{X} the input space and \mathcal{Y} the output space. We use $|\cdot|_x$ and $|\cdot|_y$ to denote the x -component and y -component of a point. For example, $s_{i|x}$ is the x -component of s_i . Finally, we use $\mathcal{N}(\epsilon, T, \rho)$ to denote the ϵ -covering number of a space T equipped with a metric ρ (see van der Vaart and Wellner, 2000, for a precise definition.).

The following definition says that an algorithm is called robust, if we can partition the sample set into finite subsets, such that if a new sample falls into the same subset as a training sample, then the loss of the former is close to the loss of the latter.

Definition 1.21. *Algorithm \mathcal{A} is $(K, \epsilon(\mathbf{s}))$ robust if \mathcal{Z} can be partitioned into K disjoint sets, denoted by $\{C_i\}_{i=1}^K$, such that $\forall \mathbf{s} \in \mathbf{s}$,*

$$s, z \in C_i, \implies |l(\mathcal{A}_{\mathbf{s}}, s) - l(\mathcal{A}_{\mathbf{s}}, z)| \leq \epsilon(\mathbf{s}). \quad (1.26)$$

1.6.1 Generalization Properties of Robust Algorithms

In this section we use the above definition to derive PAC bounds for robust algorithms. Let the sample set \mathbf{s} consist of m i.i.d. samples generated by an unknown distribution μ . Let $\hat{l}(\cdot)$ and $l_{\text{emp}}(\cdot)$ denote the expected error and

the training error, respectively, i.e.,

$$\hat{l}(\mathcal{A}_{\mathbf{s}}) \triangleq \mathbb{E}_{z \sim \mu} l(\mathcal{A}_{\mathbf{s}}, z); \quad l_{\text{emp}}(\mathcal{A}_{\mathbf{s}}) \triangleq \frac{1}{m} \sum_{s_i \in \mathbf{s}} l(\mathcal{A}_{\mathbf{s}}, s_i).$$

Theorem 1.22. *If \mathbf{s} consists of m i.i.d. samples, the loss function $l(\cdot, \cdot)$ is upper bounded by M , and \mathcal{A} is $(K, \epsilon(\mathbf{s}))$ -robust, then for any $\delta > 0$, with probability at least $1 - \delta$,*

$$\left| \hat{l}(\mathcal{A}_{\mathbf{s}}) - l_{\text{emp}}(\mathcal{A}_{\mathbf{s}}) \right| \leq \epsilon(\mathbf{s}) + M \sqrt{\frac{2K \ln 2 + 2 \ln(1/\delta)}{m}}.$$

Proof. The proof follows by partitioning the set and using inequalities for multinomial random variables, á la the Bretagnolle-Huber-Carol inequality. \square

Theorem 1.22 requires that we fix a K *a priori*. However, it is often worthwhile to consider adaptive K . For example, in the large-margin classification case, typically the margin is known only after \mathbf{s} is realized. That is, the value of K depends on \mathbf{s} . Because of this dependency, we need a generalization bound that holds uniformly for all K .

Corollary 1.23. *If \mathbf{s} consists of m i.i.d. samples, and \mathcal{A} is $(K, \epsilon_K(\mathbf{s}))$ robust for all $K \geq 1$, then for any $\delta > 0$, with probability at least $1 - \delta$,*

$$\left| \hat{l}(\mathcal{A}_{\mathbf{s}}) - l_{\text{emp}}(\mathcal{A}_{\mathbf{s}}) \right| \leq \inf_{K \geq 1} \left[\epsilon_K(\mathbf{s}) + M \sqrt{\frac{2K \ln 2 + 2 \ln \frac{K(K+1)}{\delta}}{m}} \right].$$

If $\epsilon(\mathbf{s})$ does not depend on \mathbf{s} , we can sharpen the bound given in Corollary 1.23.

Corollary 1.24. *If \mathbf{s} consists of m i.i.d. samples, and \mathcal{A} is (K, ϵ_K) robust for all $K \geq 1$, then for any $\delta > 0$, with probability at least $1 - \delta$,*

$$\left| \hat{l}(\mathcal{A}_{\mathbf{s}}) - l_{\text{emp}}(\mathcal{A}_{\mathbf{s}}) \right| \leq \inf_{K \geq 1} \left[\epsilon_K + M \sqrt{\frac{2K \ln 2 + 2 \ln \frac{1}{\delta}}{m}} \right].$$

1.6.2 Examples of Robust Algorithms

In this section we provide some examples of robust algorithms. For the proofs of the examples, we refer to Xu and Mannor (2010b) and Xu and Mannor (2010a). Our first example is Majority Voting (MV) classification (e.g., Section 6.3 of Devroye et al., 1996) that partitions the input space \mathcal{X} and labels each partition set according to a majority vote of the training samples

belonging to it.

Example 1.25 (Majority Voting). *Let $\mathcal{Y} = \{-1, +1\}$. Partition \mathcal{X} to $\mathcal{C}_1, \dots, \mathcal{C}_K$, and use $\mathcal{C}(x)$ to denote the set to which x belongs. A new sample $x_a \in \mathcal{X}$ is labeled by*

$$\mathcal{A}_{\mathbf{s}}(x_a) \triangleq \begin{cases} 1, & \text{if } \sum_{s_i \in \mathcal{C}(x_a)} \mathbf{1}(s_i|y = 1) \geq \sum_{s_i \in \mathcal{C}(x_a)} \mathbf{1}(s_i|y = -1); \\ -1, & \text{otherwise.} \end{cases}$$

If the loss function is the prediction error $l(\mathcal{A}_{\mathbf{s}}, z) = \mathbf{1}_{z|y \neq \mathcal{A}_{\mathbf{s}}(z|_x)}$, then MV is $(2K, 0)$ robust.

The MV algorithm has a natural partition of the sample space that makes it robust. Another class of robust algorithms is those that have approximately the same testing loss for testing samples that are close (in the sense of geometric distance) to each other, since we can partition the sample space with norm balls, as in the standard definition of covering numbers (van der Vaart and Wellner, 2000). The next theorem states that an algorithm is robust if two samples being close implies that they have similar testing error. Thus, in particular, this means that robustness is weaker than uniform stability (Bousquet and Elisseeff, 2002).

Theorem 1.26. *Fix $\gamma > 0$ and metric ρ of \mathcal{Z} . Suppose \mathcal{A} satisfies*

$$|l(\mathcal{A}_{\mathbf{s}}, z_1) - l(\mathcal{A}_{\mathbf{s}}, z_2)| \leq \epsilon(\mathbf{s}), \quad \forall z_1, z_2 : z_1 \in \mathbf{s}, \rho(z_1, z_2) \leq \gamma,$$

and $\mathcal{N}(\gamma/2, \mathcal{Z}, \rho) < \infty$. Then \mathcal{A} is $(\mathcal{N}(\gamma/2, \mathcal{Z}, \rho), \epsilon(\mathbf{s}))$ -robust.

Theorem 1.26 leads to the next example: if the testing error given the output of an algorithm is Lipschitz continuous, then the algorithm is robust.

Example 1.27 (Lipschitz continuous functions). *If \mathcal{Z} is compact w.r.t. metric ρ , and $l(\mathcal{A}_{\mathbf{s}}, \cdot)$ is Lipschitz continuous with Lipschitz constant $c(\mathbf{s})$, i.e.,*

$$|l(\mathcal{A}_{\mathbf{s}}, z_1) - l(\mathcal{A}_{\mathbf{s}}, z_2)| \leq c(\mathbf{s})\rho(z_1, z_2), \quad \forall z_1, z_2 \in \mathcal{Z},$$

then \mathcal{A} is $(\mathcal{N}(\gamma/2, \mathcal{Z}, \rho), c(\mathbf{s})\gamma)$ -robust for all $\gamma > 0$.

Theorem 1.26 also implies that SVM, Lasso, feed-forward neural network and PCA are robust, as stated in Example 1.28 to Example 1.31.

Example 1.28 (Support Vector Machines). *Let \mathcal{X} be compact. Consider the standard SVM formulation (Cortes and Vapnik, 1995; Schölkopf and Smola,*

2002), as discussed in Sections 1.3 and 1.4.

$$\begin{aligned} \min_{\mathbf{w}, d} \quad & c\|\mathbf{w}\|_{\mathcal{H}}^2 + \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & 1 - s_{i|y}[\langle \mathbf{w}, \phi(s_{i|x}) \rangle + d] \leq \xi_i, \quad i = 1, \dots, m; \\ & \xi_i \geq 0, \quad i = 1, \dots, m. \end{aligned}$$

Here $\phi(\cdot)$ is a feature mapping, $\|\cdot\|_{\mathcal{H}}$ is its RKHS kernel, and $k(\cdot, \cdot)$ is the kernel function. Let $l(\cdot, \cdot)$ be the hinge-loss, i.e., $l((w, d), z) = [1 - z_{|y}(\langle w, \phi(z_{|x}) \rangle + d)]^+$, and define $f_{\mathcal{H}}(\gamma) \triangleq \max_{\mathbf{a}, \mathbf{b} \in \mathcal{X}, \|\mathbf{a} - \mathbf{b}\|_2 \leq \gamma} (k(\mathbf{a}, \mathbf{a}) + k(\mathbf{b}, \mathbf{b}) - 2k(\mathbf{a}, \mathbf{b}))$. If $k(\cdot, \cdot)$ is continuous, then for any $\gamma > 0$, $f_{\mathcal{H}}(\gamma)$ is finite, and SVM is $(2\mathcal{N}(\gamma/2, \mathcal{X}, \|\cdot\|_2), \sqrt{f_{\mathcal{H}}(\gamma)/c})$ robust.

Example 1.29 (Lasso). Let \mathcal{Z} be compact and the loss function be $l(A_{\mathbf{s}}, z) = |z_{|y} - A_{\mathbf{s}}(z_{|x})|$. Lasso (Tibshirani, 1996), which is the following regression formulation:

$$\min_{\mathbf{w}} : \quad \frac{1}{m} \sum_{i=1}^m (s_{i|y} - \mathbf{w}^\top s_{i|x})^2 + c\|\mathbf{w}\|_1,$$

is $(\mathcal{N}(\gamma/2, \mathcal{Z}, \|\cdot\|_\infty), (Y(\mathbf{s})/c + 1)\gamma)$ -robust for all $\gamma > 0$, where $Y(\mathbf{s}) \triangleq \frac{1}{n} \sum_{i=1}^n s_{i|y}^2$.

Example 1.30 (Feed-forward Neural Networks). Let \mathcal{Z} be compact and the loss function be $l(A_{\mathbf{s}}, z) = |z_{|y} - A_{\mathbf{s}}(z_{|x})|$. Consider the d -layer neural network (trained on \mathbf{s}), which is the following predicting rule given an input $x \in \mathcal{X}$

$$\begin{aligned} x^0 & := z_{|x} \\ \forall v = 1, \dots, d-1 : \quad x_i^v & := \sigma\left(\sum_{j=1}^{N_{v-1}} w_{ij}^{v-1} x_j^{v-1}\right); \quad i = 1, \dots, N_v; \\ A_{\mathbf{s}}(x) & := \sigma\left(\sum_{j=1}^{N_{d-1}} w_j^{d-1} x_j^{d-1}\right); \end{aligned}$$

If there exists α and β such that the d -layer neural network satisfying that $|\sigma(a) - \sigma(b)| \leq \beta|a - b|$, and $\sum_{j=1}^{N_v} |w_{ij}^v| \leq \alpha$ for all v, i , then it is $(\mathcal{N}(\gamma/2, \mathcal{Z}, \|\cdot\|_\infty), \alpha^d \beta^d \gamma)$ -robust, for all $\gamma > 0$.

We remark that in Example 1.30, the number of hidden units in each layer has no effect on the robustness of the algorithm and consequently the bound on the testing error. This indeed agrees with Bartlett (1998), where the author showed (using a different approach based on fat-shattering dimension) that for neural networks, the weight plays a more important role

than the number of hidden units.

The next example considers an unsupervised learning algorithm, namely the principal component analysis algorithm. We show that it is robust if the sample space is *bounded*. This does not contradict the well known fact that the principal component analysis is sensitive to outliers which are far away from the origin.

Example 1.31 (Principal Component Analysis (PCA)). *Let $\mathcal{Z} \subset \mathbb{R}^m$ be such that $\max_{z \in \mathcal{Z}} \|z\|_2 \leq B$. If the loss function is $l((w_1, \dots, w_d), z) = \sum_{k=1}^d (w_k^\top z)^2$, then finding the first d principal components, which solves the following optimization problem over d vectors $w_1, \dots, w_d \in \mathbb{R}^m$,*

$$\begin{aligned} \max_{\mathbf{w}_1, \dots, \mathbf{w}_d} \quad & \sum_{i=1}^m \sum_{k=1}^d (\mathbf{w}_k^\top s_i)^2 \\ \text{s.t.} \quad & \|\mathbf{w}_k\|_2 = 1, \quad k = 1, \dots, d; \\ & \mathbf{w}_i^\top \mathbf{w}_j = 0, \quad i \neq j. \end{aligned}$$

is $(\mathcal{N}(\gamma/2, \mathcal{Z}, \|\cdot\|_2), 2d\gamma B)$ -robust.

1.7 Conclusion

The purpose of this chapter has been to hint at the wealth of applications and uses of robust optimization in machine learning. Broadly speaking, there are two main methodological frameworks developed here: robust optimization used as a way to make an optimization-based machine learning algorithm robust to noise; and robust optimization as itself a fundamental tool for analyzing properties of machine learning algorithms, and for constructing algorithms with special properties. The properties we have discussed here include sparsity, consistency and generalization. There are many directions of interest future work can pursue. We highlight two that we consider of particular interest and promise. The first is learning in the high dimensional setting, where the dimensionality of the models (or parameter space) is of the same order of magnitude as the number of training samples available. Hidden structure, like sparsity or low-rank, have offered ways around the challenges of this regime. Robustness and robust optimization may offer clues as to how to develop new tools and new algorithms for this setting. A second direction of interest is the design of uncertainty sets for robust optimization, from data. Constructing uncertainty sets from data is a central problem in robust optimization, that has not been adequately addressed, and machine learning methodology may be able to provide a way forward.

References

- N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haussler. Scale-sensitive dimension, uniform convergence, and learnability. *Journal of the ACM*, 44(4):615–631, 1997.
- M. Anthony and P. L. Bartlett. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, 1999.
- P. L. Bartlett. The sample complexity of pattern classification with neural networks: The size of the weight is more important than the size of the network. *IEEE Transactions on Information Theory*, 44(2):525–536, 1998.
- P. L. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, November 2002.
- P. L. Bartlett, O. Bousquet, and S. Mendelson. Local Rademacher complexities. *The Annals of Statistics*, 33(4):1497–1537, 2005.
- A. Ben-Tal and A. Nemirovski. Robust solutions of linear programming problems contaminated with uncertain data. *Mathematical Programming, Serial A*, 88:411–424, 2000.
- A. Ben-Tal, L. El Ghaoui, and A. Nemirovski. *Robust Optimization*. Princeton University Press, 2009.
- D. Bertsimas and M. Sim. The price of robustness. *Operations Research*, 52(1):35–53, January 2004.
- D. Bertsimas, D. B. Brown, and C. Caramanis. Theory and applications of robust optimization. Submitted, available from <http://users.ece.utexas.edu/~cmcaram>, 2010.
- C. Bhattacharyya, L. R. Grate, M. I. Jordan, L. El Ghaoui, and I. S. Mian. Robust sparse hyperplane classifiers: Application to uncertain molecular profiling data. *Journal of Computational Biology*, 11(6):1073–1089, 2004.
- P. Bickel, Y. Ritov, and A. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics*, 37:1705–1732, 2009.
- B. E. Boser, I. M. Guyon, and V. N. Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the Fifth Annual ACM Workshop on Computational Learning Theory*, pages 144–152, New York, NY, 1992.
- O. Bousquet and A. Elisseeff. Stability and generalization. *Journal of Machine Learning Research*, 2:499–526, 2002.
- E. J. Candès and T. Tao. Near-optimal signal recovery from random projections: universal encoding strategies. *IEEE Transactions on Information*

- Theory*, 52:5406–5425, 2006.
- E. J. Candès and T. Tao. The Dantzig selector: Statistical estimation when p is much larger than n . *The Annals of Statistics*, 35(6):2313–2351, 2007.
- E. J. Candès and T. Tao. Reflections on compressed sensing. *IEEE Information Theory Society Newsletter*, 58(4):20–23, 2008.
- E. J. Candès, J. Romberg, and T. Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on Information Theory*, 52(2):489–509, 2006.
- C. Caramanis and S. Mannor. Learning in the limit with adversarial disturbances. In *Proceedings of The 21st Annual Conference on Learning Theory*, 2008.
- S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20(1):33–61, 1999.
- C. Cortes and V. N. Vapnik. Support vector networks. *Machine Learning*, 20:1–25, 1995.
- A. d’Aspremont, L. El Ghaoui, M. I. Jordan, and G. R. Lanckriet. A direct formulation for sparse PCA using semidefinite programming. *SIAM Review*, 49(3):434–448, 2007.
- L. Devroye and L. Györfi. *Nonparametric Density Estimation: the l_1 View*. John Wiley & Sons, 1985.
- L. Devroye, L. Györfi, and G. Lugosi. *A Probabilistic Theory of Pattern Recognition*. Springer, New York, 1996.
- D. L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52(4):1289–1306, 2006.
- G. E. Dullerud and F. Paganini. *A Course in Robust Control Theory: A Convex Approach*, volume 36 of *Texts in Applied Mathematics*. Springer-Verlag, New York, 2000.
- B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *The Annals of Statistics*, 32(2):407–499, 2004.
- L. El Ghaoui and H. Le Bret. Robust solutions to least-squares problems with uncertain data. *SIAM Journal on Matrix Analysis and Applications*, 18:1035–1064, 1997.
- T. Evgeniou, M. Pontil, and T. Poggio. Regularization networks and support vector machines. In A. J. Smola, P. L. Bartlett, B. Schölkopf, and D. Schuurmans, editors, *Advances in Large Margin Classifiers*, pages 171–203, Cambridge, MA, 2000. MIT Press.
- A. Globerson and S. Roweis. Nightmare at test time: Robust learning by feature deletion. In *Proceedings of the 23rd International Conference on*

- Machine Learning*, pages 353–360, New York, NY, USA, 2006. ACM Press.
- A. Hoerl. Application of ridge analysis to regression problems. *Chemical Engineering Progress*, 58:54–59, 1962.
- P. Kall and S. W. Wallace. *Stochastic Programming*. John Wiley & Sons, 1994.
- S.-J. Kim, A. Magnani, and S. Boyd. Robust fisher discriminant analysis. In *Advances in Neural Information Processing Systems*, number 1659-1666, 2005.
- 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, 2003.
- S. Mukherjee, P. Niyogi, T. Poggio, and R. Rifkin. Learning theory: Stability is sufficient for generalization and necessary and sufficient for consistency of empirical risk minimization. *Advances in Computational Mathematics*, 25(1-3):161–193, 2006.
- A. Nilim and L. El Ghaoui. Robust control of Markov decision processes with uncertain transition matrices. *Operations Research*, 53(5):780–798, September 2005.
- E. Parzen. On the estimation of a probability density function and the mode. *The Annals of Mathematical Statistics*, 33:1065–1076, 1962.
- T. Poggio, R. Rifkin, S. Mukherjee, and P. Niyogi. General conditions for predictivity in learning theory. *Nature*, 428(6981):419–422, 2004.
- A. Prékopa. *Stochastic Programming*. Kluwer, 1995.
- M. Rosenblatt. Remarks on some nonparametric estimates of a density function. *The Annals of Mathematical Statistics*, 27:832–837, 1956.
- B. Schölkopf and A. J. Smola. *Learning with Kernels*. MIT Press, 2002.
- D. W. Scott. *Multivariate Density Estimation: Theory, Practice and Visualization*. John Wiley & Sons, New York, 1992.
- 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, July 2006.
- A. L. Soyster. Convex programming with set-inclusive constraints and applications to inexact linear programming. *Operations Research*, 21:1154–1157, 1973.
- I. Steinwart. Consistency of support vector machines and other regularized kernel classifiers. *IEEE Transactions on Information Theory*, 51(1):128–142, 2005.

- I. Steinwart and A. Christmann. *Support Vector Machines*. Springer, New York, 2008.
- 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 and V. Arsenin. *Solutions of Ill-Posed Problems*. Wiley, New York, 1977.
- T. Trafalis and R. Gilbert. Robust support vector machines for classification and computational issues. *Optimization Methods and Software*, 22(1):187–198, February 2007.
- A. W. van der Vaart and J. A. Wellner. *Weak Convergence and Empirical Processes*. Springer-Verlag, New York, 2000.
- V. N. Vapnik and A. Chervonenkis. The necessary and sufficient conditions for consistency in the empirical risk minimization method. *Pattern Recognition and Image Analysis*, 1(3):260–284, 1991.
- V. N. Vapnik and A. Lerner. Pattern recognition using generalized portrait method. *Automation and Remote Control*, 24:744–780, 1963.
- M. Wainwright. Sharp thresholds for noisy and high-dimensional recovery of sparsity using ℓ_1 -constrained quadratic programming (Lasso). *IEEE Transactions on Information Theory*, 55:2183–2202, 2009.
- H. Xu and S. Mannor. Robustness and generalization. ArXiv: 1005.2243, 2010a.
- H. Xu and S. Mannor. Robustness and generalizability. In *Proceeding of the Twenty-third Annual Conference on Learning Theory*, pages 503–515, 2010b.
- H. Xu, C. Caramanis, and S. Mannor. Robustness and regularization of support vector machines. *Journal of Machine Learning Research*, 10(Jul): 1485–1510, 2009.
- H. Xu, C. Caramanis, and S. Mannor. Robust regression and Lasso. *IEEE Transactions on Information Theory*, 56(7):3561–3574, 2010a.
- H. Xu, C. Caramanis, and S. Mannor. A distributional interpretation to robust optimization. submitted, 2010b.
- H. Xu, C. Caramanis, and S. Sanghavi. Robust PCA via outlier pursuit. To appear *Advances in Neural Information Processing Systems*, 2010c.
- K. Zhou, J. Doyle, and K. Glover. *Robust and Optimal Control*. Prentice-Hall, 1996.