

# Risk-Adaptive Approaches to Stochastic Optimization: A Survey

*Johannes O. Royset*

Daniel J. Epstein Department of Industrial & Systems Engineering  
University of Southern California  
royset@usc.edu

**Abstract.** Uncertainty is prevalent in engineering design, data-driven problems, and decision making broadly. Due to inherent risk-averseness and ambiguity about assumptions, it is common to address uncertainty by formulating and solving conservative optimization models expressed using measures of risk and related concepts. We survey the rapid development of risk measures over the last quarter century. From their beginning in financial engineering, we recount the spread to nearly all areas of engineering and applied mathematics. Solidly rooted in convex analysis, risk measures furnish a general framework for handling uncertainty with significant computational and theoretical advantages. We describe the key facts, list several concrete algorithms, and provide an extensive list of references for further reading. The survey recalls connections with utility theory and distributionally robust optimization, points to emerging applications areas such as fair machine learning, and defines measures of reliability.

**Keywords:** Risk-averseness, risk measure, regret measure, error measure, deviation measure, superquantile, conditional value-at-risk, regression, fairness, engineering design

**AMS Classification:** 46N10, 52B55, 65K05, 68Q32, 90C25, 91A26, 91B05, 91G70

**Date:** April 5, 2024

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Decision Making under Uncertainty</b>	<b>5</b>
2.1	Engineering Design and Control . . . . .	5
2.2	Stochastic Optimization in Statistics and Machine Learning . . . . .	7
2.3	Other Application Areas . . . . .	8
2.4	Measures of Risk . . . . .	10
2.5	Connections with Utility Theory . . . . .	13
<b>3</b>	<b>Superquantiles</b>	<b>15</b>
3.1	Equivalent Formulas . . . . .	15
3.2	Superquantiles in Optimization Models . . . . .	19
3.3	Algorithms for Superquantile Minimization . . . . .	24
3.4	Estimating Superquantiles . . . . .	29

<b>4 Risk and Regret</b>	<b>30</b>
4.1 Desirable Properties . . . . .	31
4.2 Risk Minimization . . . . .	34
4.3 Mixed Superquantiles and Law Invariance . . . . .	37
4.4 Fairness . . . . .	39
<b>5 Error and Deviation</b>	<b>40</b>
5.1 Measures of Error . . . . .	40
5.2 Measures of Deviation . . . . .	42
5.3 Superquantile Regression . . . . .	45
<b>6 Duality Theory</b>	<b>47</b>
6.1 Conjugacy . . . . .	47
6.2 Risk Envelopes and Dual Algorithms . . . . .	48
6.3 Connections with Distributionally Robust Optimization . . . . .	50
6.4 Risk Identifiers and Subgradients . . . . .	51
<b>7 Additional Examples</b>	<b>53</b>
<b>8 Computational Tools</b>	<b>56</b>
<b>9 Extensions and Challenges</b>	<b>57</b>
9.1 Measures of Reliability . . . . .	57
9.2 Dynamic and Multi-Stage Optimization . . . . .	58
9.3 Challenges and Open Problems . . . . .	58
<b>References</b>	<b>60</b>

## 1 Introduction

In many areas of applied mathematics, we aim to determine a best decision, design, estimate, or model in the presence of uncertainty about the values of key parameters and data in the problem. The uncertainty stems from our incomplete knowledge about the state of a system, the inaccuracy of approximations that might have been adopted, and the human inaptitude to predict the future. One might hope that historical observations about varying quantities can inform us and ideally eliminate the uncertainty. Despite living in the era of big data, this aspiration rarely comes to fruition. Real-world data sets tend to be noisy, biased, corrupted, or simply insufficiently large. It is therefore prudent to optimize decisions, designs, and models while accounting *conservatively* for the uncertainty. Even when our grasp of the uncertainty in an application is fairly good, human nature and our desire to avoid exceptionally “bad” outcomes motivate conservativeness in decision making and modeling broadly. Thus, financial

decisions, operational plans, military strategies, system controls, engineering designs, and statistical models are often selected conservatively.

While the various fields have approached decision making under uncertainty somewhat differently, overarching themes now emerge: the need to capture risk-averseness and ambiguity about uncertainty, promote computations and analysis, and enable explanations of algorithmic outcomes. The concept of *risk measures* provides a mathematical framework for unifying and understanding the various threads and how they connect. The framework encapsulates many common problems in robust control and optimization, distributionally robust optimization, adversarial machine learning, statistical estimation, financial risk management, utility maximization, attacker-defender games, and reliability-based design optimization. Supported by convex analysis, risk measures furnish a rich area for nonlinear analysis as well as opportunities for efficient computations. The vast literature on risk measures developed over the last 25 years is a testimony to the potency of the framework, both theoretically and practically.

In this survey, we discuss risk measures and related concepts with a focus on advances over the last quarter century. The canonical problem involves a *quantity of interest* given by a function  $f(\xi, x)$ , which depends on a *parameter*  $\xi$  and a *decision (control)*  $x$ . For example,  $f(\xi, x)$  might quantify the performance of an engineering system designed according to our decision  $x$ , given an environmental condition represented by  $\xi$ . In supervised learning,  $f(\xi, x)$  might specify the prediction error of a neural network designed according to  $x$ , given feature and label data  $\xi$ . Without fully knowing  $\xi$ , the problem is to determine a decision  $x$  such that  $f(\xi, x)$  is minimized or, alternatively,  $f(\xi, x)$  does not exceed a given threshold. The problem is ill-posed because  $\xi$  is unsettled. The uncertainty about  $\xi$  could stem from our incomplete knowledge of a “true” value, such as the demand for a product tomorrow, and/or from inherent variability in the value as exemplified by the many feature-label pairs a neural network needs to handle accurately.

There are several ways to proceed. One can estimate the value of  $\xi$  and adopt that value in the subsequent optimization and decision making. However, this fails to account for the uncertainty associated with  $\xi$ . The approach is especially problematic when  $\xi$  varies inherently and our decision  $x$  needs to perform satisfactorily under different values of  $\xi$ . Alternatively, if there is a set  $\Xi$  of possible values of  $\xi$ , then one may consider the quantity of interest in the worst case across these values, i.e.,  $\sup_{\xi \in \Xi} f(\xi, x)$ . The problem shifts to determining a decision  $x$  that minimizes this conservative quantity or makes it sufficiently low, which we refer to as *robust optimization*; see, e.g., [29, 32]. Yet another possibility is to model the uncertainty associated with  $\xi$  using a probability distribution. This brings in the vast and sophisticated tools of probability theory and statistics. In assessing a decision  $x$ , one can leverage the expected value of  $f(\cdot, x)$  computed with respect to the adopted probability distribution. Thus, the problem becomes to find a decision that is satisfactory on average. This classical approach is the central tenet of *stochastic programming*; see, e.g., [34, 307, 288].

Risk measures capture all these possibilities and many more. Through choices of probability distributions, risk measures allow us to incorporate data and other information about the possible values of  $\xi$  and their likelihoods. They reflect a wide variety of concerns and preferences a decision maker may have toward various outcomes. When restricted suitably to the class of *regular measures of risk*, they exhibit theoretically and computationally advantageous properties. In particular, convexity, linearity,

and smoothness of  $f(\xi, x)$  in  $x$  commonly carry over to the resulting optimization problem. Many risk measures also address the fact that any adopted probability distribution is an imperfect model of the uncertainty associated with  $\xi$ . Thus, they account for ambiguity about the model of uncertainty, which leads to *distributionally robust optimization problems*, as exemplified by [337, 287, 216, 307], and connections with *stochastic dominance*; see, e.g., [74].

Risk measures originated in financial engineering as an approach to quantify the reserves banks, insurance companies, and other financial institutions need to cover potential future losses; see for example the influential paper [20], the textbooks [212, 99], and the review article [100]. In this survey, we discuss the spread of risk measures *beyond* financial engineering to operations management, reliability analysis, engineering design, defense planning, statistics, and machine learning.

Among the early reviews of risk measures, [266] stands out for a concise description of the important advances taking place around the turn of the century. The tutorials [296, 169] explain key concepts and introduce connections with statistics; see also the monograph [346]. With a focus on expected utility theory and dual utility theory, [290] reviews the deep mathematical relations between risk measures and these earlier theories as well as stochastic dominance; see also the recent paper [102]. While surveying risk measures and related concepts, [276] breaks new ground by connecting concepts of risk, regret, deviation, and error and thereby relating risk management to statistics. The paper [270] emphasizes the applicability of risk measures in reliability-based engineering design. For PDE-constrained optimization and uncertainty quantification, [163] discusses the needed technical assumptions in such infinite-dimensional settings as well as algorithms. The paper [334] surveys applications for autonomous systems. Reviews of superquantile risk measures appear in [96], with a focus on applications in supply chain management, scheduling, networks, energy, and medicine, and in [175], dealing with machine learning applications. Books covering risk measures broadly include [244, 99, 307, 288].

This survey provides a succinct introduction to the area of risk measures and related concepts without devolving into the more technical aspects. We discuss the increasing interest in risk-averse approaches to statistical applications, with an updated review of the risk quadrangle proposed in [276] and refined in [271]. The survey recounts the historical development of superquantiles (a.k.a. conditional value-at-risk, average value-at-risk, tail value-at-risk, and expected shortfall) and the central role they now play in many areas of operations research, engineering, and statistics. This includes a behind-the-scene description of the derivations that led to an influential formula for superquantiles. With minimal overhead from probability theory, we describe a duality theory and connections with distributionally robust optimization. We introduce the terminology *measure of reliability* for failure probabilities, buffered failure probabilities, buffered probabilities of exceedance, and related concepts. Throughout, the survey focuses on computational aspects and implementable algorithms.

To avoid technical distractions, we mainly focus on finite-dimensional decision and uncertainty spaces, i.e.  $x$  and  $\xi$  are finite-dimensional vectors. We limit the scope to “here-and-now” decisions, but include a brief summary of multi-stage decision processes in Section 9.

The survey begins in Section 2 with examples of decision making under uncertainty and the introduction of risk measures. Section 3 discusses superquantiles and algorithms for their optimization and estimation. Section 4 reviews measures of risk broadly and their companions: measures of regret. Sec-

tion 5 shifts the focus to the statistical domain by discussing measures of error and deviation, and their application in generalized regression modeling. With their roots in convex analysis, many measures of risk are equivalently expressed by “dual formulas” that furnish additional computational possibilities as well as deep connections with distributionally robust optimization; see Section 6. Section 7 provides additional examples and Section 8 reviews computational tools for risk-adaptive optimization. The survey ends with extensions and open questions in Section 9.

## 2 Decision Making under Uncertainty

This section starts by illustrating how problems of decision making under uncertainty arise in diverse areas. We define measures of risk, give concrete examples, and make connections with utility theory.

### 2.1 Engineering Design and Control

The future responses, levels of damage, rates of deterioration, life-cycle costs, and numerous other quantities of interest for an engineering system are invariably uncertain. This is caused by unknown environmental conditions, load patterns, as well as our imperfect modeling of underlying mechanisms. Thus, design and control of engineering systems give rise to a wealth of decision making problems under uncertainty. We discuss three concrete examples, starting with a system described by explicit formulas for concreteness and ending with a high-level description of a multi-disciplinary design problem.

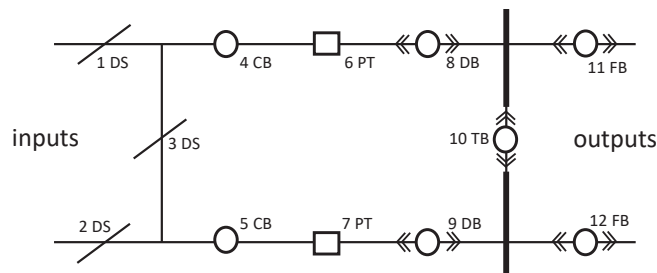


Figure 1: Outline of electrical substation (reproduced from [39]).

**2.1 Example** (electrical substation design). *Figure 1 describes a planned two-transmission-line electrical substation with 12 components of various kinds; disconnect switches (DS), circuit breakers (CB), power transformers (PT), drawout breakers (DB), tie breakers (TB), and feeder breakers (FB). In a reliability growth model, the failure rate of component  $i$  is defined as  $\lambda_i(x_i) = \alpha\beta e^{-\beta x_i}$ , where  $\alpha = 9$ ,  $\beta = 2$ , if it undergoes  $x_i$  days of testing prior to installation; see [39]. Under decision  $x_i$  and outcome  $\xi_i$  of a uniformly distributed random variable on  $[0, 1]$ , we assess the performance of component  $i$  by the number of days it falls short of the desired lifetime of  $\tau = 365$ , which becomes*

$$f_i(\xi_i, x_i) = \tau + \frac{\ln \xi_i}{\lambda_i(x_i)}$$

under an exponential failure distribution. We assess the overall system performance by identifying which component failures would trigger a disconnect between an input and an output in Figure 1. As identified in [39], there are 25 such failure modes producing the quantity of interest

$$f(\xi, x) = \max_{k=1, \dots, 25} \min_{i \in \mathbb{I}_k} f_i(\xi_i, x_i),$$

where  $\xi = (\xi_1, \dots, \xi_{12})$ ,  $x = (x_1, \dots, x_{12})$ , and the 25 sets  $\mathbb{I}_k$  are  $\{1, 2\}$ ,  $\{4, 5\}$ ,  $\{4, 7\}$ ,  $\{4, 9\}$ ,  $\{5, 6\}$ ,  $\{6, 7\}$ ,  $\{6, 9\}$ ,  $\{5, 8\}$ ,  $\{7, 8\}$ ,  $\{8, 9\}$ ,  $\{11, 12\}$ ,  $\{1, 3, 5\}$ ,  $\{1, 3, 7\}$ ,  $\{1, 3, 9\}$ ,  $\{2, 3, 4\}$ ,  $\{2, 3, 6\}$ ,  $\{2, 3, 8\}$ ,  $\{4, 10, 12\}$ ,  $\{6, 10, 12\}$ ,  $\{8, 10, 12\}$ ,  $\{5, 10, 11\}$ ,  $\{7, 10, 11\}$ ,  $\{9, 10, 11\}$ ,  $\{1, 3, 10, 12\}$ ,  $\{2, 3, 10, 11\}$ . Thus,  $f(\xi, x)$  is the number of days the overall system falls short of the desired lifetime of  $\tau = 365$  under testing decision  $x$  and outcome  $\xi$  of the inherent randomness in the components' lifetimes. An engineer would like to make a constrained decision about  $x$  such that  $f(\xi, x)$  is sufficiently low, ideally below 0, but needs to account for the uncertainty associated with  $\xi$ .

**2.2 Example** (additive manufacturing). Advances in additive manufacturing enable optimization of novel structures tailored to specific loads and boundary conditions. For example, as described in [161], we may seek to design a three-dimensional object represented by a density  $z : D \rightarrow [0, 1]$ , where  $D \subset \mathbb{R}^3$  is the physical domain. Under the assumption of linear elastic material behavior, the displacement  $u : D \rightarrow \mathbb{R}^3$  of an object with density  $z$  is determined by the weak form of the equations

$$\begin{aligned} -\operatorname{div}(E(z) : \varepsilon) &= F \text{ in } D \\ \varepsilon &= \frac{1}{2}(\nabla u + \nabla u^\top) \text{ in } D \\ \varepsilon n &= t \text{ on } \Gamma_n; \quad u = g \text{ on } \Gamma_d, \end{aligned}$$

where  $E(z)$  is the stiffness tensor (certainly dependent on  $z$ ),  $\varepsilon$  is the strain tensor,  $F$  is an external load,  $t$  is a boundary condition on the Neumann boundary  $\Gamma_n$ ,  $n$  is the outward normal of  $D$ , and  $g$  is a boundary condition on the Dirichlet boundary  $\Gamma_d$ . Typically, the material properties  $E(z)$  is uncertain, even for a fixed density  $z$ . The external load  $F$  and the boundary conditions specified by  $t$  and  $g$  are also uncertain. Thus, a solution  $u$  of the equations actually depends on the values of these uncertain parameters as well as the design decision  $z$ . A quantity of interest to be minimized could be the compliance expressed by  $\int_D F \cdot u + \int_{\Gamma_n} t \cdot u$ , which depends directly on the uncertain  $F$  and  $t$  as well as indirectly on the choice of density  $z$  and the uncertain parameters through the solution  $u$ .

**2.3 Example** (supercavitating hydrofoil). Stretching back to the pioneering work by Enrico Forlanini in 1906, naval architects have aspired to essentially lift vessels out of the water using hydrofoils and thus dramatically reduce drag and increase speed. A major challenge, however, is the inception of cavitation, i.e., water vaporization caused by the pressure on the suction side of a hydrofoil dropping below the surface tension. This causes flow unsteadiness and trigger erosion with resulting loss of material integrity. Recent efforts aim to overcome these challenges by optimizing the shape of the hydrofoil so that its hydrodynamical and structural performance is satisfactory even in the presence of uncertainty about material properties and operating conditions; see [35] and references therein. This

involves simulating the performance through numerical solution of unsteady Reynolds-averaged Navier-Stokes equations and linear elasticity equations. For instance, [35] considers five uncertain parameters  $\xi = (\xi_1, \dots, \xi_5)$ : cavitation index  $\xi_1$  (which is a surrogate for speed) and material properties represented by Young’s modulus  $\xi_2$ , Poisson’s ratio  $\xi_3$ , density  $\xi_4$ , and yield stress  $\xi_5$ . A vector  $x \in \mathbb{R}^{17}$  specifies the shape of the hydrofoil. The quantities of interest are lift force, drag-to-lift ratio, von Mises stress, and displacement at the hydrofoil tip, and thus give rise to four functions  $f_1(\xi, x), \dots, f_4(\xi, x)$ .

We refer to [110, 156] for further applications in the area of shape optimization, especially of electrical engines and elastic structures. In design of a heterogeneous lens under manufacturing uncertainty, [11] optimizes photonic nanojets. The paper [17] studies optimal control of elliptic PDEs with uncertain fractional exponents. A discussion of two-dimensional advection–diffusion equations and one-dimensional Helmholtz equations under uncertainty appears in [353]. Welded beam structures and truss bridge structures are the subjects of [40]. The paper [46] highlights the large size of many engineering problems. In the design of an acoustic cloak under uncertainty, the resulting optimization problem involves one million design variables and half a million uncertain parameters. Thus, it becomes paramount to avoid adding significant complexity while modeling uncertainty as the problem with (assumed) known parameters is already challenging.

## 2.2 Stochastic Optimization in Statistics and Machine Learning

Some problems in statistics and machine learning can also be viewed as attempting to make a decision under uncertainty. In supervised learning, we aim to find a statistical model that best predicts an unknown output based on a given input. Typically, a statistical model is specified by a vector  $c$  of coefficients, with the resulting prediction being  $g(\xi; c)$  for input  $\xi$ . If the input  $\xi = (\xi_1, \dots, \xi_n) \in \mathbb{R}^n$  and the statistical model is affine specified by the coefficients  $c = (c_0, c_1, \dots, c_n) \in \mathbb{R}^{1+n}$ , then  $g(\xi; c) = c_0 + \sum_{i=1}^n c_i \xi_i$ . But,  $g(\xi; c)$  may take many other forms, possibly utilizing neural networks. Given an input-output pair  $(\xi, \eta)$ , we would like the model’s prediction  $g(\xi; c)$  to be close in some sense to the actual output value  $\eta$ . In regression analysis, the output quantity is a scalar and “close” is commonly quantified using

$$f((\xi, \eta), c) = (\eta - g(\xi; c))^2 \quad \text{or} \quad f((\xi, \eta), c) = |\eta - g(\xi; c)|.$$

If the sign of the prediction error is important, then the quantity of interest might be  $f((\xi, \eta), c) = \eta - g(\xi; c)$ , which distinguishes between overestimating and underestimating [281, 276].

In a  $k$ -class classification problem, the output  $\eta \in \{1, \dots, k\}$  specifies a class and the statistical model  $g(\xi; c) = (g_1(\xi; c), \dots, g_k(\xi; c))$  produces a  $k$ -dimensional vector of probabilities representing the likelihood that input  $\xi$  corresponds to the various classes. In predicting  $\eta$  from  $\xi$ , the cross-entropy of  $g(\xi; c)$  relative to a probability mass function concentrated at  $\eta \in \{1, \dots, k\}$  becomes the quantity of interest:

$$f((\xi, \eta), c) = -\ln g_\eta(\xi; c).$$

Regardless of the specific details, when selecting a statistical model we are uncertain about the input-output  $(\xi, \eta)$  for which it should be accurate. In fact, we probably would like the statistical model to

make accurate predictions for many input-output pairs. Consequently, we are faced with the problem of selecting  $c$ , under uncertainty about  $(\xi, \eta)$ , such that a quantity of interest  $f((\xi, \eta), c)$  is “optimized.”

**2.4 Example** (support vector machine with fairness constraint). *In binary classification, we seek to predict an output  $\eta \in \{-1, 1\}$  from an input  $\xi$  using a statistical model  $g(\xi; c)$ , where  $c$  is a vector of coefficients. Support vector machines achieve this by considering the hinge-loss*

$$f_1((\xi, \eta), c) = \max \{0, 1 - \eta g(\xi; c)\}$$

as a quantity of interest. A concern in this setting is fairness; see, e.g., [130, 339]. When making mistakes, does the statistical model  $g(\cdot; c)$  exhibit a bias in certain settings such as when being applied with an input  $\xi$  corresponding to an under-represented group? To reduce bias in the statistical model, we may consider a secondary quantity of interest

$$f_2((\xi, \zeta), c) = (\zeta - \bar{\zeta})g(\xi; c),$$

where  $\zeta$  is an additional input representing sensitive attributes and  $\bar{\zeta}$  is the average across a population of attributes; see, e.g., [347]. The secondary quantity of interest can be used to quantify the covariance between attributes and predictions produced by the statistical model and thus serves as a metric of fairness. Consequently, we are faced with a problem of selecting  $c$ , under uncertainty about  $(\xi, \eta, \zeta)$ , such that two quantities of interest are satisfactory.

**2.5 Example** (surrogate building and digital twins). A computer model of a physical system, called a digital twin [145, 320, 321], as well as a coarse computer model of a high-fidelity simulation [239, 237] require calibration to achieve acceptable predictions. These models lead to surrogates that predict an output  $\eta$  from an input  $\xi$ . Given an input-output pair  $(\xi, \eta)$ , suppose that  $g(\xi; c)$  is the prediction by a surrogate with coefficients  $c$ . Thus, a quantity of interest is the prediction error  $f((\xi, \eta), c) = \eta - g(\xi; c)$ . Figure 2 shows the results of low- and high-fidelity estimates of lift force (blue asterisks) produced by a hydrofoil as described in Example 2.3. We see that the low-fidelity estimates do not match the high-fidelity estimates, but this can be addressed through calibration. With  $\xi$  being a low-fidelity estimate, we may adopt  $g(\xi; c) = c_0 + c_1\xi$  as the prediction of the corresponding high-fidelity estimate. The tuning of the coefficients  $c = (c_0, c_1)$  aims to bring the quantity of interest  $f((\xi, \eta), c) = \eta - (c_0 + c_1\xi)$  near zero, while accounting for the uncertainty (as illustrated by the spread in Figure 2) associated with  $(\xi, \eta)$ . Figure 2 shows five such surrogates fitted with varying emphasis on conservativeness; dotted lines overestimate more than the dashed lines. Surrogates are often more refined than these affine models (see, e.g., [328]), with neural networks emerging as key forms [222, 338, 147, 60]. Regardless of the sophistication, the fundamental problem remains how to select a surrogate in the presence of uncertainty.

### 2.3 Other Application Areas

Nearly all areas of human activity involve decision making under uncertainty. Financial engineering gives rise to many challenging problems in portfolio management and cashflow control; see, e.g., [212, 99].



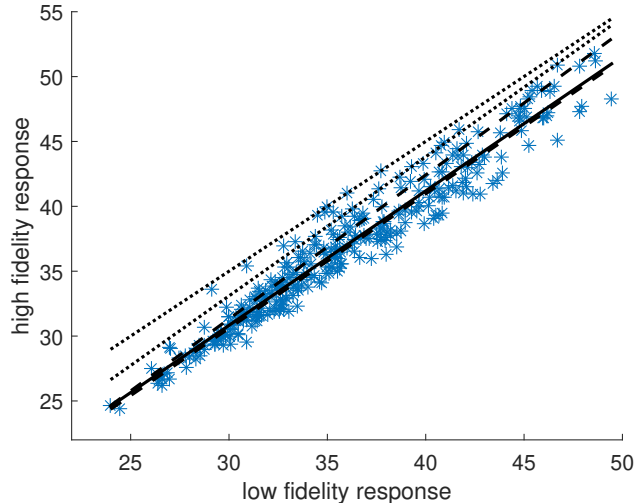


Figure 2: Low- and high-fidelity estimates of lift force (blue asterisks) produced by a hydrofoil from Example 2.3 [35] and five lines representing surrogates with varying degree of conservativeness.

The design of service networks takes place under uncertainty about future demand [141]. Demand is likewise uncertain in online and brick-and-mortar retail operations [259]. Transportation engineers account for uncertain travel times and travel patterns; see, e.g., [42, 348, 57]. Planning of humanitarian relief has to be carried out without fully knowing where and how disasters may strike [228]. The energy sector faces uncertainty in rainfall, which affects hydropower generation [299, 9], but also uncertainty at the building level [233] and in microgrids [24]. Decisions about when to charge and discharge an electrical storage unit face uncertainty about future supply and demand [105]. Military jamming missions [58] and capital investments in defense technology [319] also take place under uncertainty. We refer to [96, 200] for a summary of applications in operations management broadly. In this area, inventory control (see, e.g., [7, 129]) gives rise to many challenging problems. A classical example illustrates the setting.

**2.6 Example** (newsvendor). *In a contractual agreement, a newsvendor (a firm) places an order for the daily delivery of a fixed number of newspapers (perishable items) to meet a daily demand  $\xi$ . The newsvendor is charged  $\gamma > 0$  cents per paper ordered and sells each for  $\delta > \gamma$  cents; unsold papers cannot be returned and are worthless at the end of the day. When ordering  $x$  newspapers and  $\xi$  is the demand, the loss (expense minus income) turns out to be*

$$f(\xi, x) = \begin{cases} \gamma x - \delta x & \text{if } \xi \geq x \\ \gamma x - \delta \xi & \text{otherwise,} \end{cases}$$

*with negative values implying a profit. The loss becomes the quantity of interest that we seek to minimize while accounting for the uncertainty in  $\xi$ .*

## 2.4 Measures of Risk

As exemplified in the previous subsections, we are often faced with a quantity of interest  $f(\xi, x)$  and the need to manipulate it through the choice of  $x$ . However,  $\xi$  is unsettled and we proceed by modeling its possible values and the associated likelihoods using a probability distribution. Just as  $f(\xi, x)$  typically would be an imperfect representation of reality, the probability distribution is a *model* of the uncertainty associated with  $\xi$ . In a particular case, we construct the probability distribution using available observations of past values of  $\xi$ , expert opinions, and engineering judgement; see, e.g., [288, Section 3.D]. The probability distribution might assign each outcome in a set the same likelihood (i.e., a uniform distribution) and even concentrate at a single point, which implies certainty about the value of  $\xi$ . Regardless of the situation, we indicate the shift from a parameter  $\xi$  to a random quantity<sup>1</sup> with a probability distribution by using boldface  $\boldsymbol{\xi}$ . Outcomes of this random quantity are denoted by  $\xi$ .

For example, suppose that a quantity of interest is given by  $f(\xi, x) = (\xi - 2/3)x - 1/3$ , where  $x \in \{-1, 1\}$ , i.e., there are only two possible decisions. We model the uncertainty associated with  $\xi$  using the triangular distribution on  $[0, 2]$  with mode at 0. Thus, we are faced with the choice between the two random variables  $f(\boldsymbol{\xi}, 1)$  and  $f(\boldsymbol{\xi}, -1)$ ; the first one quantifies the randomness under the decision  $x = 1$  and the second one quantifies similarly the randomness under  $x = -1$ . Since  $\boldsymbol{\xi}$  has the assumed triangular distribution, we can compute the probability density functions of  $f(\boldsymbol{\xi}, 1)$  and  $f(\boldsymbol{\xi}, -1)$ , denoted by  $p_1$  and  $p_2$ , respectively; see Figure 3.

Throughout the article, we assume that a quantity of interest represents cost, loss, damage, displacement in excess of a threshold, or another performance metric that we wish to be *low*. In this setting, what decision is better,  $x = 1$  or  $x = -1$ ? One possible way of answering this question is to consider the expectations  $\mathbb{E}[f(\boldsymbol{\xi}, 1)]$  and  $\mathbb{E}[f(\boldsymbol{\xi}, -1)]$ , but they both come out as  $-1/3$ . Figure 3 highlights a significant difference between the decisions, however:  $f(\boldsymbol{\xi}, 1)$  has values as high as 1, while  $f(\boldsymbol{\xi}, -1)$  does not exceed  $1/3$ . Thus, from a worst-case point of view,  $x = -1$  appears better.

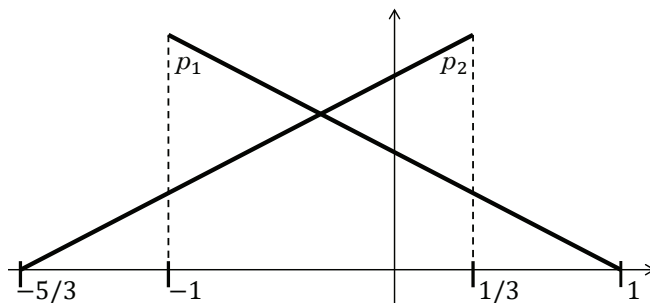


Figure 3: Probability density functions  $p_1$  and  $p_2$  for random variables  $f(\boldsymbol{\xi}, 1)$  and  $f(\boldsymbol{\xi}, -1)$ , respectively.

Expectations as well as worst-case outcomes can be viewed as defining measures of risk. A risk measure converts a random variable, say  $f(\boldsymbol{\xi}, 1)$ , into a scalar, which furnishes a basis for comparison

---

<sup>1</sup>As usual, a random quantity  $\boldsymbol{\xi}$  is a measurable mapping from some underlying probability space  $(\Omega, \mathcal{A}, \mathbb{P})$  to a space of outcomes such as  $\mathbb{R}$ , which makes  $\boldsymbol{\xi}$  a random variable, or to  $\mathbb{R}^m$ , which makes  $\boldsymbol{\xi}$  a random vector.

with the scalar obtained from, say  $f(\boldsymbol{\xi}, -1)$ . Thus, risk measures facilitate decision making when the quantity of interest is uncertain.

A risk measure assesses a random variable, which might have a complicated definition in terms of other random variables and decision variables. This is certainly the case when the random variable represents a real-world quantity of interest as we have seen above. Even in the simple case of Figure 3, risk measures assessed the random variables  $f(\boldsymbol{\xi}, 1)$  and  $f(\boldsymbol{\xi}, -1)$ , which in turn are defined by the triangularly distributed random variable  $\boldsymbol{\xi}$ . Regardless of the circumstances, it is sometimes beneficial to “hide” this complexity and we often write  $\boldsymbol{\xi}$  for a generic random variable that we seek to assess using a risk measure.

**2.7 Definition** (risk measure). *A measure of risk (also called risk measure)  $\mathcal{R}$  assigns to a random variable  $\boldsymbol{\xi}$  a number  $\mathcal{R}(\boldsymbol{\xi}) \in [-\infty, \infty]$  as a quantification of its risk, with this number being  $\xi$  if  $\text{prob}\{\boldsymbol{\xi} = \xi\} = 1$ .*

Since we consistently prefer lower values to higher values, the choice between two random variables  $\boldsymbol{\xi}_1$  and  $\boldsymbol{\xi}_2$  now reduces to comparing the two numbers  $\mathcal{R}(\boldsymbol{\xi}_1)$  and  $\mathcal{R}(\boldsymbol{\xi}_2)$  and selecting the random variable with the lower number.

The requirement  $\mathcal{R}(\boldsymbol{\xi}) = \xi$  for a random variable  $\boldsymbol{\xi}$  concentrated at  $\xi$  stems from the translation equivariance condition of [20] and it is brought out explicitly in [266]. The requirement is certainly reasonable: a quantity of interest that has a specific value with certainty should indeed be assessed as that value. Subsection 4.1 associates additional properties with risk measures. However, we prefer Definition 2.7 as a means to introduce the concept without imposing potentially limiting requirements.

Since a risk measure simply converts a random variable into a scalar, with a minor requirement for constant random variables, a large number of possibilities exist.

**Expectation.** The choice  $\mathcal{R}(\boldsymbol{\xi}) = \mathbb{E}[\boldsymbol{\xi}]$ , the expected value, defines a risk measure, but one that is insensitive to the possibility of high values as long as they are offset by low values. Thus, it is referred to as *risk-neutral*. This risk measure is meaningful in settings where a decision is implemented repeatedly and the focus is on good average performance, with little concern about fluctuations in performance. Moreover, one should be quite confident that the assumed probability distributions are accurate, at least in the sense that they produce the correct expected values. Stochastic programming originated with a focus on expectation minimization; see, e.g., [336] for early developments. In statistics, expectation minimization arises in classification, regression, and other contexts.

Many real-life situations involve a decision that will only be executed once or at most a few times. Then, the average performance across many (hypothetical) instances is less important than guarantees about reasonable performance in the actual instances that will take place. This is often the case for major events such as an earthquake, a mars landing, or certain military missions. Even if a decision will be put to test often, we may not accept poor performance even though it is “countered” by excellent performance other times. This is common for engineering systems such as driverless cars where the performance needs to be sufficiently good nearly all the time.

We may also suspect that our models, producing the quantity of interest and probability distributions representing the uncertainty, are incomplete. The performance that we predict, even under assumed known conditions, could be erroneous. Thus, we may consider “worst cases” as a means to compensate for modeling deficiencies. These situations trigger the need for *risk-averse* measures of risk, which is the main focus of this article. They allow us to address decision makers that are willing to forego good average performance to achieve reduced possibility and severity of poor outcomes.

Generally, for a random variable  $\xi$  with support<sup>2</sup>  $\Xi$  let  $\sup \xi = \sup \Xi$ , which can be interpreted as the largest possible value of  $\xi$ .

**Worst-case risk.** The choice  $\mathcal{R}(\xi) = \sup \xi$  represents a conservative outlook, possibly overly so because  $\mathcal{R}(\xi) = \infty$  when  $\xi$  has, for example, a normal or exponential distribution. This risk measure ignores the probabilities of the various outcomes of  $\xi$  and utilizes only the support  $\Xi$ . Thus, it is especially suitable when there is little information for building a probability distribution. A worst-case risk measure underpins robust optimization [29, 32], semiinfinite programming [255, Chapter 3], robust control [255, Chapter 4], adversarial training [204], and diametrical risk minimization [225].

**Mean-plus-standard-deviation risk.** A natural choice motivated by statistical confidence intervals is to adopt a mean-plus-standard-deviation risk measure defined by  $\mathcal{R}(\xi) = \mathbb{E}[\xi] + \lambda \text{std}(\xi)$ , where  $\text{std}(\xi)$  is the standard deviation of  $\xi$  and  $\lambda$  is a positive constant. While mean-plus-standard-deviation risk measures can be traced back to early efforts in portfolio management [210], where in fact the variance was used instead of the standard deviation, their close ties to normal distributions are problematic. In particular, the two random variables from Figure 3 have the same mean and standard deviation and thus the same mean-plus-standard-deviation risk. This is caused by an equal treatment of outcomes above the mean and below the mean in the calculation of standard deviations, which is counter to many decision makers’ greater concern about variability above the mean than below the mean.

Beyond mean and standard deviations, quantiles and superquantiles are central quantities for summarizing a random variable. Most importantly, they treat high values differently than low values. For  $\alpha \in (0, 1)$ , the  $\alpha$ -*quantile* of a random variable  $\xi$  with cumulative distribution function<sup>3</sup>  $P$  is given by

$$Q(\alpha) = \min \{ \xi \in \mathbb{R} \mid P(\xi) \geq \alpha \} \tag{2.1}$$

and the  $\alpha$ -*superquantile* of an integrable random variable  $\xi$  is given by

$$\bar{Q}(\alpha) = Q(\alpha) + \frac{1}{1-\alpha} \mathbb{E}[\max\{0, \xi - Q(\alpha)\}]. \tag{2.2}$$

We also define  $\bar{Q}(0) = \mathbb{E}[\xi]$  and  $\bar{Q}(1) = \sup \xi$ , which indeed are the limits of  $\bar{Q}(\alpha)$  as  $\alpha \searrow 0$  and  $\alpha \nearrow 1$ , respectively; see, e.g., [269, Theorem 2].

---

<sup>2</sup>The support of  $\xi$  is the smallest closed subset of  $\mathbb{R}$  containing outcomes of  $\xi$  occurring with probability one.

<sup>3</sup>For a random variable  $\xi$  defined on  $(\Omega, \mathcal{A}, \mathbb{P})$ , its cumulative distribution function is  $P(\xi) = \mathbb{P}(\{\omega \in \Omega \mid \xi(\omega) \leq \xi\})$ .

**Quantile risk.** The choice of risk measure  $\mathcal{R}(\boldsymbol{\xi}) = Q(\alpha)$ , the  $\alpha$ -quantile of  $\boldsymbol{\xi}$ , is widely used in financial engineering (see, e.g., [20]) under the name “value-at-risk” with typical  $\alpha$  near 1. In Figure 3,  $Q(0.8) = 0.106$  for density function  $p_1$  and  $Q(0.8) = 0.122$  for  $p_2$ . Thus,  $p_1$  (corresponding to decision  $x = 1$ ) is better according to this measure of risk. The conclusion is problematic because  $p_1$  extends further to the right in Figure 3 than  $p_2$ , which means that worst outcomes are possible under decision  $x = 1$  than under  $x = -1$ . In fact, this measure of risk ignores the *magnitude* of outcomes above  $Q(\alpha)$  and thus fails to reflect the severity of poor outcomes.

**Superquantile risk.** The choice of risk measure  $\mathcal{R}(\boldsymbol{\xi}) = \bar{Q}(\alpha)$ , the  $\alpha$ -superquantile of  $\boldsymbol{\xi}$ , covers the whole range of possibilities from the risk-neutral  $\mathcal{R}(\boldsymbol{\xi}) = \mathbb{E}[\boldsymbol{\xi}]$  to the worst-case risk measure  $\mathcal{R}(\boldsymbol{\xi}) = \sup \boldsymbol{\xi}$  by adjusting  $\alpha$  from 0 to 1. In Figure 3,  $p_1$  and  $p_2$  produce superquantiles  $\bar{Q}(0.8) = 0.404$  and  $\bar{Q}(0.8) = 0.230$ , respectively. An assessment of the two decisions underpinning Figure 3 based on 0.8-superquantiles thus gives a decisive advantage to  $p_2$  and the decision  $x = -1$ . This reversal compared to the conclusion reached in the previous paragraph is caused by the fact that superquantiles account for the magnitude of poor outcomes and therefore flag decision  $x = 1$  due to the possibility of exceptionally poor outcomes under  $p_1$ .

These risk measures are just examples, with many others given below. We already see that the choice of risk measure profoundly affects a decision and thus should, ideally, reflect a decision maker’s preferences and any model ambiguity that might be present. A risk measure should also be computationally tractable because any complication comes on top of the challenges already associated with a quantity of interest. Section 4 reviews desirable properties for risk measures that both promote computations and reflect typical preferences among decision makers.

In reliability analysis, one often considers the probability of a random variable  $\boldsymbol{\xi}$  exceeding a threshold  $\tau$ , i.e.,  $\text{prob}\{\boldsymbol{\xi} > \tau\}$ . Such probabilities do *not* define measures of risk because they fail to produce  $\xi$  when the random variable takes the value  $\xi$  with probability one. For this reason, we prefer to treat such expressions separately in Subsection 9.1 under the name *measures of reliability*.

## 2.5 Connections with Utility Theory

How to make a choice between two random variables has a long history, especially in economics. The classical prescription is due to von Neumann and Morgenstern [330] and their *expected utility theory*. Faced with a choice between random variables  $\boldsymbol{\xi}_1$  and  $\boldsymbol{\xi}_2$ , adopt a utility function  $u : \mathbb{R} \rightarrow \mathbb{R}$ , and make the choice using  $\mathbb{E}[u(\boldsymbol{\xi}_1)]$  and  $\mathbb{E}[u(\boldsymbol{\xi}_2)]$ , with higher numbers being preferred. For brief summaries of expected utility theory, we refer to [270, 1, 102]. With our focus on achieving lower values (of cost, damage, shortfall relative to a target, etc.), there is a mismatch with the orientation of expected utility theory. This is easily rectified by reversing the sign of utility functions and thus producing *disutility functions*; see, e.g., [267]. For random variables  $\boldsymbol{\xi}_1$  and  $\boldsymbol{\xi}_2$  representing quantities of interest oriented toward lower values, we can define a disutility function  $v : \mathbb{R} \rightarrow \mathbb{R}$  from a utility function  $u : \mathbb{R} \rightarrow \mathbb{R}$  by setting  $v(\xi) = -u(-\xi)$ . Then, we may deem  $\boldsymbol{\xi}_1$  preferable over  $\boldsymbol{\xi}_2$  if  $\mathbb{E}[v(\boldsymbol{\xi}_1)] \leq \mathbb{E}[v(\boldsymbol{\xi}_2)]$ .

In stochastic programming models (see, e.g., [34, 307, 288]), one often seeks to match a quantity of interest  $f(\xi, x)$  with some target level, say 0. This produces the constraint  $f(\xi, x) = 0$ , which rarely can be satisfied due to the uncertainty associated with  $\xi$ . Simple recourse models [336] circumvent this difficulty by removing the constraint, assuming that the values of the uncertain parameters are governed by the random vector  $\boldsymbol{\xi}$ , and adding the term  $\mathbb{E}[v(f(\boldsymbol{\xi}, x))]$  to the objective function. Here,  $v : \mathbb{R} \rightarrow \mathbb{R}$  can be viewed as a disutility function, often of the form  $v(\eta) = \max\{\delta\eta, \gamma\eta\}$  with  $\delta \leq 0 < \gamma$  being fixed parameters; see, e.g., [288, Section 3.H].

While expected utility theory and related approaches such as prospect theory [142] are widely adopted tools for risk-averse decision making, they do not automatically lead to measures of risk because  $\mathbb{E}[v(\boldsymbol{\xi})] \neq v(\xi)$  for a random variable  $\boldsymbol{\xi}$  concentrated at  $\xi$ ; see [77, 267] for efforts to bridge this gap. These theories have also come under criticism [101, 102]. The main concern is that the “right” utility function (or disutility function) is challenging to determine in practice, with more intricate functions resulting in difficulties with explaining how one reached a particular decision. There may also be no utility function that fully captures a decision maker’s preference; she may be thinking about possible recourse actions that cannot be captured by utility functions alone. These concerns worsen when there are multiple stakeholders, which is common in many engineering applications. Moreover, there might be ambiguity about the underlying models and probability distributions making expected disutility less meaningful; see [102] for a detailed discussion.

**2.8 Example** (disutility functions and beyond). *Suppose that a firm faces a choice between two uncertain costs in the future described by the random variables  $\boldsymbol{\xi}_1$  and  $\boldsymbol{\xi}_2$  distributed uniformly between  $[-3/2, 1]$  and  $[-8, 2]$ , respectively. (A negative cost means that the firm receives money.) The comparison between  $\mathbb{E}[\boldsymbol{\xi}_1] = -1/4$  and  $\mathbb{E}[\boldsymbol{\xi}_2] = -3$  leads to the conclusion that the better choice is to select  $\boldsymbol{\xi}_2$ ; on average it results in a lower cost. However, this may not be the right decision if the firm is concern about the possibility of a high cost, i.e., it is risk-averse.*

**Detail.** Following expected (dis)utility theory, we may adopt a disutility function  $v(\xi) = \max\{0, \xi\}/(1-\alpha)$ , where  $\alpha \in (0, 1)$  is a fixed parameter. The disutility function reflects a preference for outcomes  $\xi \in (-\infty, 0]$ , with increasing displeasure associated with positive outcomes. This risk-averse perspective leads to a comparison between  $\mathbb{E}[\max\{0, \boldsymbol{\xi}_1\}/(1-\alpha)]$  and  $\mathbb{E}[\max\{0, \boldsymbol{\xi}_2\}/(1-\alpha)]$ . These quantities can be thought of as the expected levels of displeasure felt by the firm when facing the random costs  $\boldsymbol{\xi}_1$  and  $\boldsymbol{\xi}_2$ , respectively. Regardless of  $\alpha$ , the comparison produces a tie between the two choices because  $\mathbb{E}[\max\{0, \boldsymbol{\xi}_1\}] = \mathbb{E}[\max\{0, \boldsymbol{\xi}_2\}] = 1/5$ .

Neither of these two approaches for assessing the merit of adopting one cost over the other considers the possibility of mitigating actions by the firm. It turns out that such deeper considerations may change the decision. Let us suppose that the costs  $\boldsymbol{\xi}_1$  and  $\boldsymbol{\xi}_2$  are given in present money and that displeasure is quantified as above. If the firm is more active and invests  $\gamma$  amount of money in a risk-free asset (bonds or bank deposit) now, then the future displeasure, as perceived now, is reduced from  $\mathbb{E}[\max\{0, \boldsymbol{\xi}_i\}/(1-\alpha)]$  to  $\mathbb{E}[\max\{0, \boldsymbol{\xi}_i - \gamma\}/(1-\alpha)]$  as  $\gamma$  will be available at the future point in time to offset the cost  $\boldsymbol{\xi}_i$ . The upfront expense  $\gamma$  also needs to be considered and the goal becomes to select

the investment  $\gamma$  such that

$$\gamma + \frac{1}{1-\alpha} \mathbb{E}[\max\{0, \xi_i - \gamma\}] \text{ is minimized.} \quad (2.3)$$

As we see in Theorem 3.1 below, the resulting minimum value is the  $\alpha$ -superquantile of  $\xi_i$ . A comparison between the 0.8-superquantile of  $\xi_1$ , which is  $3/4$ , and the 0.8-superquantile of  $\xi_2$ , which is 1, reveals that  $\xi_1$  is preferred. When accounting for the possibility of mitigating future displeasure by investing in a risk-free asset, the advantage tilts decisively toward the first cost. We see that a superquantile measure of risk inherently incorporates in its assessment of a random variable the possibility of such mitigation.  $\square$

The example illustrates the difference between making decisions based on expected (dis)utility theory and based on superquantile risk. The latter turns out to be deeply rooted in *dual utility theory* [344], which relies on axioms parallel to those of expected utility theory; see [75, 267] and our discussion in Subsection 4.3. Moreover, superquantile risk has the following axiomatic justification [332]: For a real-valued risk measure  $\mathcal{R}$ , defined on the integrable random variables, with  $\mathcal{R}(\xi) = 1$  for constant random variables  $\xi$  concentrated at 1, we have that

$$\begin{aligned} \mathcal{R} \text{ satisfies the monotonicity, law invariance, prudence, and no-reward-for-concentration axioms} \\ \iff \mathcal{R} \text{ is a superquantile risk measure given by } \bar{Q}(\alpha) \text{ for some } \alpha \in (0, 1). \end{aligned}$$

Thus, superquantile risk measures are the *only* risk measures with the listed axiomatic properties. (We define monotonicity and law invariance in Subsection 4.1, and prudence relates to lower semicontinuity as also defined in that section. For details about prudence and no-reward-for-concentration, we refer to [332, Sections 2.1-2.2].) This theoretical foundation and their practical usage support the claim that superquantile risk measures are “currently the most important risk measure in banking practice” [332].

### 3 Superquantiles

In the vast landscape of risk measures, superquantiles emerge as central. They capture risk-averseness and one-sided concerns about high values. They span the range of preferences from the risk-neutral perspective (by setting  $\alpha = 0$ ) to a focus on worst-case outcomes (by setting  $\alpha = 1$ ). With the dependence on a single parameter ( $\alpha$ ), superquantiles are easy to explain to decision makers. In the following, we identify three other features: (i) superquantiles have mathematical properties that facilitate computational optimization, (ii) they capture ambiguity about the adopted probability distribution and thus connect with distributionally robust optimization, and (iii) they furnish the fundamental building blocks for many “reasonable” measures of risk. Thus, we proceed with a comprehensive review of superquantiles; Sections 4 and 7 cover more general risk measures.

#### 3.1 Equivalent Formulas

Superquantiles can at least be traced back to the concept of *optimized certainty equivalents* as developed in [30] for random variables oriented toward high value as preferable. Converting the concept into the

present orientation, this pioneering paper effectively defines the quantity

$$\min_{\gamma \in \mathbb{R}} \gamma + \mathbb{E}[v(\boldsymbol{\xi} - \gamma)],$$

where  $v : \mathbb{R} \rightarrow \mathbb{R}$  is an increasing, strictly convex, and twice continuously differentiable disutility function, normalized with  $v(0) = 0$  and derivative  $v'(0) = 1$ . While the increasing and smoothness properties as well as strict convexity are violated by  $v(\xi) = \max\{0, \xi\}/(1 - \alpha)$ , this particular choice brings us to the minimization formula for an  $\alpha$ -superquantile used in conjunction with (2.3) and formally expressed as (3.2) below; see [31] for details about these connections.

Averages beyond a quantile, which are closely related to (2.2), are mentioned in [19, 20] under the name tail value-at-risk; see also [85]. Using the term conditional value-at-risk (CVaR), [274] also starts from an average beyond a quantile and derives the equivalence below between (3.1) and (3.2) under the assumption that the random variable  $\boldsymbol{\xi}$  is continuously distributed. The equivalence between (3.1) and (3.2) for general, integrable random variables is confirmed in [240]; see also [275], which further establishes a connection with the modified tail expectation (3.4). This modified tail expectation accounts for random variables with a positive probability of taking a value exactly at a quantile and is also the starting point for [3]. That paper proceeds by establishing equivalence with the formula (3.3) from [2] involving an integral of quantiles across different probability levels as well as with the minimization formula (3.2) from [274]. Integrals of quantiles have a long history in statistics and underpin Lorenz curves [198]; see, for example, the discussion in Section 3 of [269] for additional connections. Further insight stems from [67], which expresses a superquantile of a continuous random variable as the worst-case expectation over a family of probability distributions. The situation for general distributions is hinted to in [67], but brought out more clearly in [278]. In summary, the flurry of activity around the turn of the century produced the following equivalent formulas for an  $\alpha$ -superquantile.

**3.1 Theorem** (equivalent formulas for superquantiles). *For  $\alpha \in (0, 1)$  and an integrable random variable  $\boldsymbol{\xi}$  with cumulative distribution function  $P$  and quantile function  $Q$ , the following hold:*

$$Q(\alpha) + \frac{1}{1 - \alpha} \mathbb{E}[\max\{0, \boldsymbol{\xi} - Q(\alpha)\}] \tag{3.1}$$

$$= \min_{\gamma \in \mathbb{R}} \gamma + \frac{1}{1 - \alpha} \mathbb{E}[\max\{0, \boldsymbol{\xi} - \gamma\}] \tag{3.2}$$

$$= \frac{1}{1 - \alpha} \int_{\alpha}^1 Q(\beta) d\beta \tag{3.3}$$

$$= \text{expectation of the } \alpha\text{-tail distribution of } \boldsymbol{\xi}, \tag{3.4}$$

where the  $\alpha$ -tail distribution is defined as having  $P^{[\alpha]}(\xi) = \max\{0, P(\xi) - \alpha\}/(1 - \alpha)$  as its cumulative distribution function. Thus, any of these formulas can be taken as the definition of the  $\alpha$ -superquantile of  $\boldsymbol{\xi}$ , which we denote by  $\bar{Q}(\alpha)$  or  $\text{s-rsk}_{\alpha}(\boldsymbol{\xi})$  to highlight the dependence on  $\boldsymbol{\xi}$ .

If  $\boldsymbol{\xi}$  is square integrable, then

$$\bar{Q}(\alpha) = \operatorname{argmin}_{\gamma \in \mathbb{R}} \left\{ \gamma + \frac{1}{1 - \alpha} \int_0^1 \max\{0, \bar{Q}(\beta) - \gamma\} d\beta \right\}. \tag{3.5}$$



If  $\boldsymbol{\xi}$  has finite support  $\Xi$  of cardinality  $s$  and corresponding probabilities  $p_\xi, \xi \in \Xi$ , then

$$\bar{Q}(\alpha) = \max_{\bar{p} \in \Delta_\alpha} \sum_{\xi \in \Xi} \xi \bar{p}_\xi, \quad \text{where } \Delta_\alpha = \left\{ \bar{p} \in \mathbb{R}^s \mid 0 \leq \bar{p}_\xi \leq \frac{p_\xi}{1-\alpha}, \xi \in \Xi, \sum_{\xi \in \Xi} \bar{p}_\xi = 1 \right\}. \quad (3.6)$$

**Proof.** While originally developed by several researchers as discussed before the theorem, the four equivalences are concisely summarized in [269]; see the discussion around equation (3.4) in that reference for the equivalence between (3.3) and (3.4) above. Section 4 and Theorem 2 of [269] confirm the other equivalences.

Theorem 7 in [269] establishes the argmin formula (3.5). A proof of the last expression (3.6) appears in [278]; see also our discussion in Section 6.  $\square$

In general, the  $\alpha$ -superquantile  $\bar{Q}(\alpha)$  of a random variable  $\boldsymbol{\xi}$  is equal to neither  $\mathbb{E}[\boldsymbol{\xi} \mid \boldsymbol{\xi} \geq Q(\alpha)]$  nor  $\mathbb{E}[\boldsymbol{\xi} \mid \boldsymbol{\xi} > Q(\alpha)]$  and this sometimes causes confusion. The discrepancy is reflected in (3.4), where any probability atom at the  $\alpha$ -quantile  $Q(\alpha)$  of  $\boldsymbol{\xi}$  is carefully “split.” Nevertheless, in the absence of such an atom,  $\bar{Q}(\alpha) = \mathbb{E}[\boldsymbol{\xi} \mid \boldsymbol{\xi} \geq Q(\alpha)] = \mathbb{E}[\boldsymbol{\xi} \mid \boldsymbol{\xi} > Q(\alpha)]$ . In particular, if  $\boldsymbol{\xi}$  has a density function  $p$ , then

$$\begin{aligned} \bar{Q}(\alpha) &= Q(\alpha) + \frac{1}{1-\alpha} \int_{-\infty}^{\infty} \max\{0, \xi - Q(\alpha)\} p(\xi) d\xi \\ &= Q(\alpha) + \frac{1}{1-\alpha} \int_{Q(\alpha)}^{\infty} \xi p(\xi) d\xi - \frac{Q(\alpha)}{1-\alpha} \int_{Q(\alpha)}^{\infty} p(\xi) d\xi \\ &= \frac{1}{1-\alpha} \int_{Q(\alpha)}^{\infty} \xi p(\xi) d\xi, \end{aligned} \quad (3.7)$$

which coincides with  $\mathbb{E}[\boldsymbol{\xi} \mid \boldsymbol{\xi} \geq Q(\alpha)]$  and  $\mathbb{E}[\boldsymbol{\xi} \mid \boldsymbol{\xi} > Q(\alpha)]$ .

These finer points may sometimes be glossed over. Based on our experience with practitioners, we recommend adopting (3.1) as the definition of a superquantile, with the supporting remark:

$$\alpha\text{-superquantile of } \boldsymbol{\xi} = \text{average of the worst } (1-\alpha)100\% \text{ outcomes of } \boldsymbol{\xi}.$$

The word “worst” is intuitively understood by the practitioner and is sufficiently ambiguous to provide cover for the mathematician.

The importance of (3.2) in computations emerges in Subsection 3.2. It is apparent from (3.1) and (3.2) that the  $\alpha$ -quantile  $Q(\alpha)$  is a minimizer of the optimization problem over  $\gamma \in \mathbb{R}$  in the latter formula. However, it may not be the only minimizer. There are multiple minimizers if the cumulative distribution function  $P$  for  $\boldsymbol{\xi}$  has a “flat stretch” to the right of  $Q(\alpha)$  and  $P(Q(\alpha)) = \alpha$ . For any  $\alpha \in (0, 1)$ , we obtain that

$$[Q(\alpha), Q^+(\alpha)] = \operatorname{argmin}_{\gamma \in \mathbb{R}} \left\{ \gamma + \frac{1}{1-\alpha} \mathbb{E}[\max\{0, \boldsymbol{\xi} - \gamma\}] \right\}, \quad (3.8)$$

where  $Q^+(\alpha) = \sup\{\gamma \in \mathbb{R} \mid P(\gamma) \leq \alpha\}$ ; see, e.g., [269, Equation (4.1)].

The insight leading to (3.2) in [274], as recounted in [269], was that an integrable random variable  $\boldsymbol{\xi}$  can be associated with the convex function given by  $e(\gamma) = \mathbb{E}[\max\{\gamma, \boldsymbol{\xi}\}]$  whose subgradients recover

the cumulative distribution function of the random variable. The conjugate function  $e^*$  of the convex function  $e$ , as given by  $e^*(\alpha) = \sup_{\gamma \in \mathbb{R}} \alpha\gamma - e(\gamma)$ , turns out to furnish (3.2) after a scaling with  $\alpha - 1$ . Specifically, for  $\alpha \in (0, 1)$ , one has

$$\frac{e^*(\alpha)}{\alpha - 1} = \min_{\gamma \in \mathbb{R}} \gamma + \frac{1}{1 - \alpha} \mathbb{E}[\max\{0, \xi - \gamma\}]$$

as seen in [269, Theorem 2]. Similar relations between conjugate pairs and their connections with cumulative distribution functions and quantile functions were examined independently in [232]; see also [72].

The argmin-formula (3.5) in Theorem 3.1 may at first appear less useful as it requires the knowledge of all superquantiles to compute one of them. However, it is the linchpin for superquantile regression (see [273] and our discussion in Subsection 5.3) and allows us to estimate one superquantile from a finite number of other ones via numerical integration. Additional argmin-formulas for superquantiles appear in [160].

As we see in Section 6, the formula (3.6), referred to as the *dual formula*, holds much beyond finite distributions. Still, recording this special case is useful as it avoids all technical overhead while addressing important applications in data-driven optimization and learning. The key insight from (3.6) is that two seemingly different decision makers will make the same assessment of a random variable. Specifically, for  $\alpha \in (0, 1)$  and random variables with finite support, consider the two individuals:

Mr. Averse has full confidence in the assumed probability distributions of random variables, but is inherently risk-averse and makes decision by comparing  $\alpha$ -superquantiles of the random variables.

Ms. Ambiguous is risk-neutral and makes decisions based on expectations, but is suspicious about the assumed probability distributions. She computes an expectation using the worst-case probability distribution obtained by scaling the assumed probabilities with factors between 0 and  $1/(1 - \alpha)$ .

In effect, Mr. Averse computes the left-hand side of (3.6) and Ms. Ambiguous computes the right-hand side. Thus, they reach the same assessment. We conclude that superquantiles can be used to address a decision maker's inherent risk-averseness as well as ambiguity about probability distributions, for example due to lack of data, fear of contamination, or adversarial interference.

It follows almost immediately from the definitions that the superquantiles  $\bar{Q}(\alpha)$  of an integrable random variable  $\xi$  tend to the worst-case  $\sup \xi = \bar{Q}(1)$  as  $\alpha \nearrow 1$ . A useful quantification of the difference  $\bar{Q}(1) - \bar{Q}(\alpha)$  appears in [14]. For further properties of superquantiles, we refer to [269].

With its origin in financial engineering, superquantiles are also known as *tail value-at-risk*, *expected shortfall*, *average value-at-risk*, and *conditional value-at-risk* (with CVaR as an abbreviation). Despite slight difference in definitions originally, we now accept these as synonyms for the quantity here called superquantiles. The name superquantile stems from [268] and was motivated by a need for making this fundamental concept free from dependence on financial terminology.

We give explicit formulas for superquantiles in three cases next; see [223] for many other cases involving common probability distributions.

**3.2 Example** (triangular distribution). Given  $\alpha \in (0, 1)$ , the  $\alpha$ -quantiles and  $\alpha$ -superquantiles of a random variable  $\xi$  with the triangular density function  $p_1$  in Figure 3 are

$$Q(\alpha) = 1 - 2\sqrt{1 - \alpha} \quad \text{and} \quad \bar{Q}(\alpha) = 1 - \frac{4}{3}\sqrt{1 - \alpha}.$$

Moreover,  $\bar{Q}(0) = -1/3$  and  $\bar{Q}(1) = 1$ .

**Detail.** Let  $P$  be the cumulative distribution function of  $\xi$ . Since  $p_1(\xi) = -\xi/2 + 1/2$  for  $\xi \in [-1, 1]$ , the solution of the equation  $P(\xi) = \int_{-1}^{\xi} p_1(\eta) d\eta = \alpha$  is  $Q(\alpha) = 1 - 2\sqrt{1 - \alpha}$ . This formula and (3.7) give  $\bar{Q}(\alpha)$ .  $\square$

**3.3 Example** (normal distribution). Given  $\alpha \in (0, 1)$ , the  $\alpha$ -quantiles and  $\alpha$ -superquantiles of a random variable that is normally distributed with mean  $\mu$  and variance  $\sigma^2$  are

$$Q(\alpha) = \mu + \sigma\Phi^{-1}(\alpha) \quad \text{and} \quad \bar{Q}(\alpha) = \mu + \frac{\sigma\varphi(\Phi^{-1}(\alpha))}{1 - \alpha},$$

where  $\varphi$  is the standard normal density function and  $\Phi^{-1}(\alpha)$  is the corresponding  $\alpha$ -quantile given by the standard normal cumulative distribution function  $\Phi$ . Moreover,  $\bar{Q}(0) = \mu$  and  $\bar{Q}(1) = \infty$ .

**3.4 Example** (finite distribution). For  $\alpha \in [0, 1]$ , the  $\alpha$ -superquantile of a finitely distributed random variable with values  $\xi_1 < \xi_2 < \dots < \xi_r$ , which occur with probabilities  $p_1, p_2, \dots, p_r$ , respectively, is given by

$$\bar{Q}(\alpha) = \begin{cases} \sum_{j=1}^r p_j \xi_j & \text{if } \alpha = 0 \\ \frac{1}{1-\alpha} \left( \left( \sum_{j=1}^i p_j \right) - \alpha \right) \xi_i + \sum_{j=i+1}^r p_j \xi_j & \text{if } \sum_{j=1}^{i-1} p_j < \alpha \leq \sum_{j=1}^i p_j < 1 \\ \xi_r & \text{if } \alpha > 1 - p_r. \end{cases}$$

## 3.2 Superquantiles in Optimization Models

Some of the formulas in Theorem 3.1 might leave the impression that superquantiles are more complicated to compute and optimize than quantiles. However, this is not the case. The formula (3.2) decouples superquantiles from quantiles and this turns out to be especially important in optimization models.

Consider a quantity of interest  $f : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$  and an  $m$ -dimensional random vector  $\xi$ . Suppose that for each  $x \in \mathbb{R}^n$ ,  $f(\xi, x)$  is an integrable random variable. With

$$\text{s-rsk}_\alpha(f(\xi, x)) = \alpha\text{-superquantile of } f(\xi, x),$$

we may seek a decision  $x$  that minimizes  $\text{s-rsk}_\alpha(f(\xi, x))$ . If the minimum value turns out to be  $\tau$ , then the obtained decision  $x^*$  has the guarantee that  $f(\xi, x^*) \leq \tau$  on average across the worst  $(1 - \alpha)100\%$

outcomes. Moreover, the alternative formula (3.6) leads to the insight that  $x^*$  is a decision that minimizes the worst-case expected value of the quantity of interest across a set of probability distributions “near” the nominal one. The resulting decision would typically be rather different than those obtained by minimizing  $\mathbb{E}[f(\boldsymbol{\xi}, x)]$  under the nominal probability distribution, which pay no particular attention to the possibility of high values of  $f(\boldsymbol{\xi}, x)$  or deviations from the nominal distribution.

Superquantiles preserve the important convexity property as recognized in [274]; here we recall [288, Proposition 3.10].

**3.5 Proposition** (convexity of superquantile functions). *For a random vector  $\boldsymbol{\xi}$  with support  $\Xi \subset \mathbb{R}^m$  and  $f : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$ , suppose that*

- (a)  $f(\boldsymbol{\xi}, x)$  is integrable for all  $x \in \mathbb{R}^n$
- (b)  $f(\boldsymbol{\xi}, \cdot)$  is convex for all  $\boldsymbol{\xi} \in \Xi$ .

*Then, for any  $\alpha \in [0, 1)$ , the function  $x \mapsto \text{s-rsk}_\alpha(f(\boldsymbol{\xi}, x))$  is convex and real-valued.*

If we ignore the uncertainty associated with  $\boldsymbol{\xi}$  and simply consider minimizing  $f(\hat{\boldsymbol{\xi}}, x)$  for some nominal value  $\hat{\boldsymbol{\xi}}$ , then we would face a convex problem as long as  $f(\hat{\boldsymbol{\xi}}, \cdot)$  is convex. The proposition asserts that  $x \mapsto \text{s-rsk}_\alpha(f(\boldsymbol{\xi}, x))$ , which treats uncertainty much more comprehensively, is convex too when  $f(\boldsymbol{\xi}, \cdot)$  is convex for all  $\boldsymbol{\xi} \in \Xi$ . Thus, optimization under uncertainty carried out in this manner does not muddle up convexity that might be present in  $f$ .

We next consider computational approaches. Given a feasible set  $X \subset \mathbb{R}^n$  and  $\alpha \in (0, 1)$ , suppose that we seek to solve the optimization problem

$$\underset{x \in X}{\text{minimize}} \quad \text{s-rsk}_\alpha(f(\boldsymbol{\xi}, x)). \quad (3.9)$$

The formula (3.2) enables us to reformulated this problem equivalently as

$$\underset{x \in X, \gamma \in \mathbb{R}}{\text{minimize}} \quad \mathbb{E} \left[ \gamma + \frac{1}{1 - \alpha} \max \{0, f(\boldsymbol{\xi}, x) - \gamma\} \right], \quad (3.10)$$

which brings us back to minimizing an expectation function for which there are many algorithms including (stochastic) subgradient type methods (cf. Subsection 3.3). Moreover, if  $\boldsymbol{\xi}$  is finitely distributed with support  $\Xi$  and corresponding probabilities  $\{p_\xi > 0, \boldsymbol{\xi} \in \Xi\}$ , then the problem simplifies to

$$\underset{x \in X, \gamma \in \mathbb{R}}{\text{minimize}} \quad \gamma + \frac{1}{1 - \alpha} \sum_{\boldsymbol{\xi} \in \Xi} p_\xi \max \{0, f(\boldsymbol{\xi}, x) - \gamma\}, \quad (3.11)$$

which in turn can be reformulated as

$$\underset{x \in X, \gamma \in \mathbb{R}, z \in \mathbb{R}^s}{\text{minimize}} \quad \gamma + \frac{1}{1 - \alpha} \sum_{\boldsymbol{\xi} \in \Xi} p_\xi z_\xi \quad \text{subject to} \quad f(\boldsymbol{\xi}, x) - \gamma \leq z_\xi, \quad 0 \leq z_\xi \quad \forall \boldsymbol{\xi} \in \Xi, \quad (3.12)$$

where  $z = (z_\xi, \boldsymbol{\xi} \in \Xi) \in \mathbb{R}^s$  are additional variables and  $s$  is the cardinality of  $\Xi$ . Both (3.11) and (3.12) are convex as long as  $X$  is convex and  $f(\boldsymbol{\xi}, \cdot)$  is convex for all  $\boldsymbol{\xi} \in \Xi$ . In fact, if  $f(\boldsymbol{\xi}, \cdot)$  is affine for all  $\boldsymbol{\xi} \in \Xi$  and  $X$  is polyhedral, then (3.12) is a linear optimization problem.

Smoothness is also preserved. If  $f(\xi, \cdot)$  is continuously differentiable for all  $\xi \in \Xi$ , then (3.12) involves inequality constraints of the kind commonly addressed by nonlinear programming algorithms.

In light of (3.8), we see that the  $\gamma$ -portion of any minimizer  $(x^*, \gamma^*, z^*)$  from (3.12) is the  $\alpha$ -quantile of  $f(\xi, x^*)$  or possibly a (slightly) larger quantity as specified by (3.8). The minimum value from (3.12) is of course the  $\alpha$ -superquantile of  $f(\xi, x^*)$ . The advantages of (3.12) are therefore clear, but the reformulation only applies to finite distributions.

A superquantile appearing as a constraint is treated analogously. For  $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$  and  $\tau \in \mathbb{R}$ , the problem

$$\underset{x \in X}{\text{minimize}} \quad f_0(x) \quad \text{subject to} \quad \text{s-rsk}_\alpha(f(\xi, x)) \leq \tau \quad (3.13)$$

is equivalently stated as

$$\underset{x \in X, \gamma \in \mathbb{R}, z \in \mathbb{R}^s}{\text{minimize}} \quad f_0(x) \quad \text{subject to} \quad \gamma + \frac{1}{1-\alpha} \sum_{\xi \in \Xi} p_\xi z_\xi \leq \tau \quad (3.14)$$

$$f(\xi, x) - \gamma \leq z_\xi, \quad 0 \leq z_\xi \quad \forall \xi \in \Xi,$$

again provided that  $\xi$  has a finite distribution with  $s$  outcomes. Following a similar pattern, we achieve formulations for multiple quantities of interest as well.

The convexity property associated with superquantiles stands in sharp contrast to the situation for quantiles. The following illustration is taken from [288, Example 3.12].

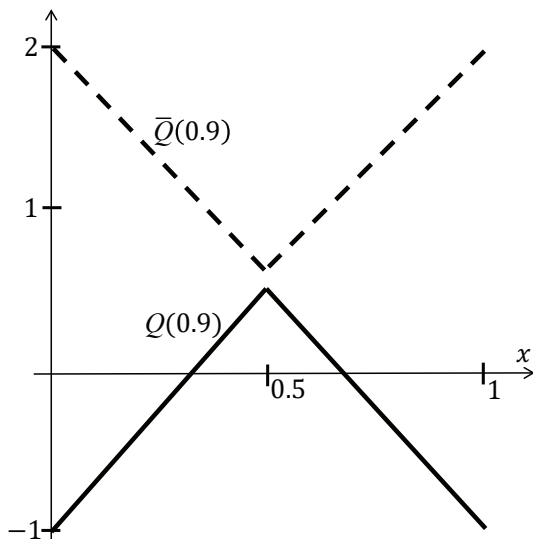


Figure 4: Quantile  $Q(0.9)$  (solid line) and superquantile  $\bar{Q}(0.9)$  (dashed) of  $f(\xi, x)$  in Example 3.6 as functions of  $x$ .

**3.6 Example** (quantile vs superquantile minimization). *Suppose that we need to allocate funds between two financial instruments, both associated with uncertainty. The first instrument requires, with*

probability 0.1, a payment of 2 dollars per dollar committed and, with probability 0.9, yields one dollar per dollar committed. Let the random variable  $\xi_1$  model these losses; it takes the value 2 with probability 0.1 and the value  $-1$  with probability 0.9. The second instrument is modeled by the random variable  $\xi_2$ , which has the same probability distribution as that of  $\xi_1$  but the two random variables are statistically independent. Suppose that we need to allocate 1 million dollars between these two instruments. Since the instruments appear equally unappealing, we might be led to believe that any allocation is fine. This would be a big mistake, but one that remains hidden for an analyst examining quantiles. Superquantiles reveal that the best strategy would be to allocate half a million to each instrument.

**Detail.** Let  $x \in [0, 1]$  be the fraction of our million dollars allocated to the first instrument, the remainder is allocated to the second instrument. Then, the quantity of interest describing our loss is the random variable  $f(\xi, x) = x\xi_1 + (1 - x)\xi_2$ . Since there are only four possible outcomes of  $\xi = (\xi_1, \xi_2)$ , we find that

$$f(\xi, x) = \begin{cases} -1 & \text{with probability 0.81} \\ 2 - 3x & \text{with probability 0.09} \\ 3x - 1 & \text{with probability 0.09} \\ 2 & \text{with probability 0.01.} \end{cases}$$

With  $\alpha = 0.9$ , the  $\alpha$ -quantile of  $f(\xi, x)$  becomes

$$Q(0.9) = \begin{cases} 3x - 1 & \text{if } x \in [0, 1/2] \\ 2 - 3x & \text{otherwise,} \end{cases}$$

which is depicted with a solid line in Figure 4. By (3.1), the 0.9-superquantile of  $f(\xi, x)$  becomes

$$\bar{Q}(0.9) = \begin{cases} -2.7x + 2 & \text{if } x \in [0, 1/2] \\ 2.7x - 0.7 & \text{otherwise,} \end{cases}$$

which is also shown in Figure 4 with a dashed line. As a function of  $x$ , the superquantiles define a convex function but the quantiles do not. The minimization of superquantiles has  $x^* = 0.5$  as minimizer, but the minimization of quantiles results in  $x = 0$  and  $x = 1$ . The decision  $x^*$  involves *hedging*, a well-known strategy in finance to reduce risk. The decision lowers the probability of a loss of 2 million dollars from 0.1 to 0.01 compared to the choice  $x = 0$  or  $x = 1$ . This comes at the expense of reducing the probability of a loss of  $-1$  million dollars from 0.9 to 0.81. The decision  $x^*$  also results in the possibility of a loss of 0.5 million dollars (with probability 0.18), but this might be much more palatable than a 2-million-dollar loss. The minimum value of the 0.9-superquantiles is 0.65 million dollars, which is a more reasonable and conservative assessment of the uncertain future loss than the wildly optimistic  $-1$  million provided by the minimum value of the 0.9-quantiles.  $\square$

The advantages of superquantiles over quantiles are not limited to convexity properties and hedging strategies as indicated in the previous example. Using [288, Example 3.13], we next illustrate that superquantiles reveal the magnitude of “typical” poor outcomes and not only their likelihood.

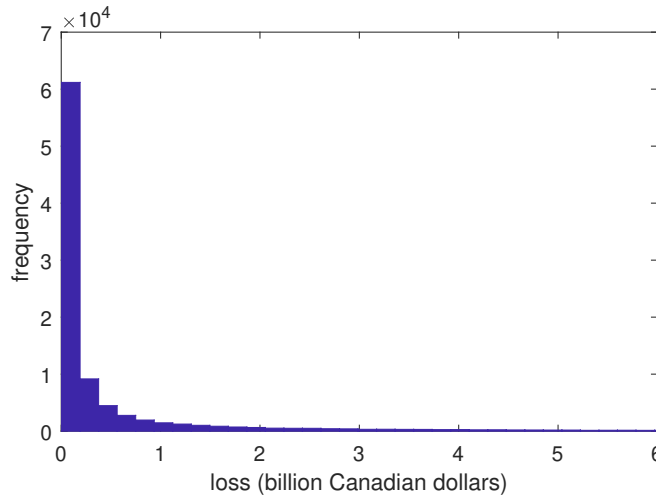


Figure 5: Histogram of 100000 outcomes of the cumulative loss due to earthquake damage in Vancouver.

**3.7 Example** (revealing vulnerabilities). *Superquantiles have the ability to reveal vulnerabilities that might remain hidden if the focus is on quantiles. From a study of the cumulative loss due to earthquake damage during the next 50 years in the greater Vancouver region in Canada [208], we obtain 100000 loss values that one can view as outcomes for a random variable modeling the cumulative loss; see the histogram of Figure 5. The cumulative loss is most likely less than one billion Canadian dollars. In fact, the 0.9-quantile of the cumulative loss is 4.39 billions. Although sizable, it seems that the region is quite resilient to earthquakes. It turns out that this conclusion is flawed and a study of superquantiles paints a gloomier picture.*

**Detail.** The largest outcome is actually 373 billions; Figure 5 should have been extended much to the right. Superquantiles quantify this long right-tail. Specifically, the 0.9-superquantile of the cumulative loss is 28.92 billions. This means that when the cumulative loss exceeds 4.39 billions, it does so substantially. The region might have a significant vulnerability after all.  $\square$

The advantageous properties of a superquantile, both computationally and from a modeling perspective, have brought this measure of risk to the forefront. From its start in financial applications [274, 168] (see [12, 218] for examples of recent work in that area), superquantiles have been used in retail operations [199], in planning military jamming missions [58], in supporting capital investment decisions [319], and in inventory control [7, 129]; see [96, 200] for a summary from the area of operations management broadly. Applications of superquantiles in the energy sector are exemplified by [111, 343]. Superquantiles appear also in the control of PDEs [164, 108] and design of physical systems [285, 35, 44, 45].

Superquantiles are increasingly being used in machine learning. For computer vision problems, [185, 345] report promising empirical results. The paper [254] shows the potential in federated learning

with simulations in the areas of character recognition and sentiment analysis; see [174] for distributed learning on mobile devices involving heterogeneous distributions. Modeling using superquantiles also emerges as a tool to address fairness in machine learning [339, 102] and, after slight adjustments, also outliers [196]. Superquantiles in reinforcement learning are discussed in [217, 309, 325]. Additional recent efforts include [311, 61]; see also [93, 149] for closely related top  $k$ -approaches, with promising numerical results, and [188] for a related tilted loss approach. We refer to the recent review article [175] for a summary of superquantile applications in machine learning. A general statistical estimation point of view, especially in the broad context of distributional shifts, appears in [80].

Connections between superquantiles and distributionally robust optimization with a Wasserstein ambiguity set emerge in [135]; see also our discussion in Subsection 6.3. Superquantiles relate to the classical newsvendor problem from Example 2.6 [119] as well as support vector machines for binary classification [318, 120, 122].

For applications in the area of two-stage stochastic mixed-integer programming, we refer to [298, 326] and the paper [94], which focuses on fixed-charge transportation problems; see also [133] for applications to energy storage and transportation, [323] for routing hazardous material, and [227] for modeling of disaster relief. Superquantiles even define norms as discussed in [235, 206, 121] and quantify the distance between cumulative distribution functions [236].

### 3.3 Algorithms for Superquantile Minimization

The basic approach for minimizing a superquantile risk over a feasible set  $X$  takes us from (3.9) through a reformulation to (3.12), with a similar treatment when a superquantile appears as a constraint; see (3.13) and (3.14). The approach is appealing because it can leverage any state-of-the-art algorithm for solving the resulting problem. However, it breaks down if the random vector  $\boldsymbol{\xi}$  does not have a finite distribution or if its finite distribution involves a massively large number of outcomes. In this subsection, we survey more advanced algorithms for superquantile minimization.

**Stochastic Subgradient Methods.** The superquantile minimization problem (3.9) can always be written as (3.10) using the formula (3.2). Thus, we are back in the familiar domain of expectation minimization. Let

$$\psi(\boldsymbol{\xi}, (x, \gamma)) = \gamma + \frac{1}{1 - \alpha} \max \{0, f(\boldsymbol{\xi}, x) - \gamma\},$$

which thus defines the integrand in (3.10), and  $\varphi(z) = \mathbb{E}[\psi(\boldsymbol{\xi}, z)]$ . We can then apply a standard stochastic subgradient method to  $\varphi$ , among which the SGD<sup>4</sup> method is well known.

#### SGD Method.

**Data.**  $z^0 \in X \times \mathbb{R}$  and step sizes  $\lambda^\nu \in (0, \infty)$ ,  $\nu = 0, 1, 2, \dots$

**Step 0.** Set iteration counter  $\nu = 0$ .

---

<sup>4</sup>SGD stands for “stochastic gradient descent” but this is doubly misleading as it is neither a descent method nor involves gradients only;  $\psi(\boldsymbol{\xi}, \cdot)$  is nonsmooth and thus we need to consider subgradients.



**Step 1.** Generate an observation  $\xi^\nu$  according to the distribution of  $\xi$ .

**Step 2.** Compute a subgradient  $v$  of  $\psi(\xi^\nu, \cdot)$  at the point  $z^\nu$  and set

$$z^{\nu+1} \in \text{prj}_{X \times \mathbb{R}}(z^\nu - \lambda^\nu v).$$

**Step 3.** Replace  $\nu$  by  $\nu + 1$  and go to Step 1.

Step 2 requires the projection of  $z^\nu - \lambda^\nu v$  onto the set  $X \times \mathbb{R}$ , which thus needs to be relatively easy to accomplish. There are several possible convergence results for this algorithm. Regardless, the algorithm is inherently random and we need to view the observations  $\{\xi^0, \xi^1, \dots\}$  obtained in Step 1 as outcomes of the independent random vectors  $\{\xi^0, \xi^1, \dots\}$  distributed as  $\xi$ . Then, the iterates  $\{z^\nu\}_{\nu=1}^\infty$  produced by the algorithm and their averages become random vectors as indicated by switching to bold face:  $\{z^\nu\}_{\nu=1}^\infty$ . In the convex case, the average of the first  $\bar{\nu}$  iterates is a near-minimizer with a specific tolerance that is guaranteed in expectation.

**3.8 Theorem** (SGD method). *For a random vector  $\xi$  with support  $\Xi \subset \mathbb{R}^m$ , a function  $f : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$ , and a nonempty, closed, and convex set  $X \subset \mathbb{R}^n$ , suppose that*

- (a)  $f(\xi, x)$  is integrable for all  $x \in X$
- (b)  $f(\xi, \cdot)$  is convex for all  $\xi \in \Xi$
- (c) for some  $\beta \in \mathbb{R}$ ,

$$\sup \left\{ \mathbb{E}[\|v(\xi)\|_2^2] \mid v(\xi) \text{ integrable, } v(\xi) \in \partial_z \psi(\xi, z) \ \forall \xi \in \Xi, z \in X \times \mathbb{R} \right\} \leq \beta^2.$$

If the SGD method has generated  $\{z^\nu, \nu = 1, 2, \dots, \bar{\nu}\}$  using step size  $\lambda^\nu = \lambda = \|z^0 - z^*\|_2 / (\beta\sqrt{\bar{\nu} + 1})$  for all  $\nu$ , then

$$\mathbb{E}[\varphi(\bar{z})] - \inf_{z \in X \times \mathbb{R}} \varphi(z) \leq \frac{\beta \|z^0 - z^*\|_2}{\sqrt{\bar{\nu} + 1}}$$

where  $z^* \in \text{argmin}_{z \in X \times \mathbb{R}} \varphi(z)$  and  $\bar{z} = \frac{1}{\bar{\nu} + 1} (\sum_{\nu=1}^{\bar{\nu}} z^\nu + z^0)$ .

A slight adjustment of the arguments in Sections 2.1 and 2.2 of [219] leads to the theorem; see [288, Section 3.G] for details. These arguments stem originally from [220, 221]. With the importance of expectation minimization in machine learning, there is a rapid development of closely related algorithms under milder assumptions, which we do not attempt to review systematically; see [63, 81] for the setting of weakly convex functions, [292] for functions that are differentiable in the generalized sense of V. Norkin, and [349, 64] for modifications of Goldstein's subgradient method. There is also an extensive literature that only relies on function evaluations; see, e.g., [192]. The monograph [177] provides an in-depth treatment, while the overview article [109] offers accessible proofs. For estimates of the actual minimum value in the context of stochastic subgradient type algorithms, we refer to [178]. Theoretical justification for adapting the SGD method to (3.10), with some empirical evidence, appears in [311].

While (3.10) indeed is just an expectation minimization problem over  $z = (x, \gamma)$ , it has a rather specific structure which may cause difficulties [61]. For example, if  $\alpha$  is near one, then the optimal  $\gamma$  for any fixed  $x$  tends to be high; in fact it can be taken as the  $\alpha$ -quantile of  $f(\boldsymbol{\xi}, x)$ ; see Theorem 3.1. Thus, there will be “few” outcomes of  $f(\boldsymbol{\xi}, x)$  that exceed  $\gamma$ . A (small) minibatch  $\{\xi^1, \dots, \xi^\nu\}$  may produce  $\max\{0, f(\xi^i, x) - \gamma\} = 0$  for all  $i = 1, \dots, \nu$  and then a zero (sub)gradient with respect to  $x$ . This can be addressed, in part, by keeping  $\gamma$  relatively low and/or increasing the minibatch size [345]. Regardless, the main advantage of minimizing superquantiles via (3.10) is the possibility of leveraging the vast computational infrastructure for expectation minimization.

Another possibility, viable when the support of  $\boldsymbol{\xi}$  has low cardinality (say, less than 5,000), is to minimize a superquantile directly without passing through the formula (3.2). This would require a subgradient of  $x \mapsto \text{s-rsk}_\alpha(f(\boldsymbol{\xi}, x))$ . Proposition 2 of [175] furnishes the following convenient expression.

**3.9 Proposition** (subgradients of superquantile functions). *For  $\alpha \in (0, 1)$ ,  $f : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$ , and a random vector  $\boldsymbol{\xi}$  with finite support  $\Xi \subset \mathbb{R}^m$  of cardinality  $s$  and associated probabilities  $p_\xi = 1/s$  for all  $\xi \in \Xi$ , suppose that  $f(\xi, \cdot)$  is continuously differentiable for all  $\xi \in \Xi$ . Consider the function given by*

$$\varphi(x) = \text{s-rsk}_\alpha(f(\boldsymbol{\xi}, x)).$$

*Then, then the set of subgradient<sup>5</sup> of  $\varphi$  at  $x \in \mathbb{R}^n$  is*

$$\partial\varphi(x) = \frac{1}{s(1-\alpha)} \sum_{\xi \in \Xi_+(x)} \nabla_x f(\xi, x) + \frac{\sigma - \alpha}{1 - \alpha} \text{con} \{ \nabla_x f(\xi, x), \xi \in \Xi_0(x) \},$$

*where  $\sigma = 1 - |\Xi_+(x)|/s$ ,  $\text{con } A$  is the convex hull of the set  $A$ , and*

$$\Xi_+(x) = \{ \xi \in \Xi \mid f(\xi, x) > Q(\alpha) \} \quad \text{and} \quad \Xi_0(x) = \{ \xi \in \Xi \mid f(\xi, x) = Q(\alpha) \},$$

*with  $Q(\alpha)$  being the  $\alpha$ -quantile of  $f(\boldsymbol{\xi}, x)$  and  $|\Xi_+(x)|$  being the cardinality of  $\Xi_+(x)$ .*

At least in the convex case, the proposition furnishes the necessary ingredient for a subgradient method (see, e.g., [288, Subsection 2.I]) to minimize a superquantile. If  $f(\xi, \cdot)$  is nonconvex but smooth with Lipschitz continuous gradients  $\nabla_x f(\xi, \cdot)$  for each  $\xi$  in the finite set  $\Xi$ , then  $\varphi$  is weakly convex and one can leverage developments in [65]; see also references therein. The calculation of a subgradient requires a quantile of  $f(\boldsymbol{\xi}, x)$ , which essentially amounts to sorting of the  $s$  values  $f(\xi, x), \xi \in \Xi$ . Thus, the worst-case complexity of computing a subgradient is of order  $O(s \ln s)$ . This means that a subgradient method for minimizing a superquantile has a relatively high, per-iteration computational cost if  $s$  is large. For a more detailed complexity analysis and comparisons, we refer to [185] and references therein.

**Primal Smoothing Methods.** Even if the quantity of interest  $f(\xi, x)$  is continuously differentiable in  $x$  for every  $\xi \in \Xi$ , the function  $x \mapsto \text{s-rsk}_\alpha(f(\boldsymbol{\xi}, x))$  is only exceptionally smooth. (The cases  $\alpha = 0$  and  $\Xi$  being a singleton are examples of such exceptions.) The situation is similar as we pass to (3.10)

---

<sup>5</sup>Subgradients are defined in a general (Mordukhovich) sense; see, e.g., [288, Section 4.I].

after invoking the formula (3.2); the max-term is typically nonsmooth. The simple form of this term, however, makes it easy to approximate using a continuously differentiable function; see [10, 311] for a quadratic approximation and [25, 164, 345] for exponential smoothing. The latter type of smoothing has the advantage of preserving any order of smoothness in  $f(\xi, \cdot)$ . Concretely, exponential smoothing (see, e.g., [288, Example 4.16]) amounts to replacing  $\gamma \mapsto \max\{0, \gamma\}$  by the function  $h_\theta : \mathbb{R} \rightarrow \mathbb{R}$  defined as

$$h_\theta(\gamma) = \frac{1}{\theta} \ln(1 + \exp(\theta\gamma)),$$

where  $\theta > 0$  is a parameter. The error caused by smoothing is controlled by

$$0 \leq h_\theta(\gamma) - \max\{0, \gamma\} \leq \frac{\ln 2}{\theta} \quad \forall \gamma \in \mathbb{R}.$$

The function  $h_\theta$  is differentiable any number of times, with easily accessible derivative formulas; see [288, Example 4.16]. Thus, there is strong theoretical backing for approximating (3.10) by

$$\underset{x \in X, \gamma \in \mathbb{R}}{\text{minimize}} \mathbb{E} \left[ \gamma + \frac{1}{1 - \alpha} h_\theta(f(\xi, x) - \gamma) \right], \quad (3.15)$$

which involves a continuously differentiable objective function provided that  $f(\xi, \cdot)$  is continuously differentiable and mild additional assumptions hold; see, e.g., [307, Subsection 9.2.5]. We refer to [284] for general tools to justify such approximations and to [47, 164] for many related smoothing schemes. Regardless, we are faced with the expectation minimization problem (3.15) and can leverage a vast array of existing algorithm, which may perform better practically and theoretically when applied to a smooth problem than to the (potentially) nonsmooth problem (3.10).

**Proximal Composite Method.** Minimizing a superquantile of  $f(\xi, x)$  over  $x \in X$  when  $\xi$  is finitely distributed with support  $\Xi$  is equivalent to (3.11) as achieved by the formula (3.2). This reformulation is well structured because it can be written in terms of the real-valued convex function  $h : \mathbb{R}^{1+s} \rightarrow \mathbb{R}$  given by

$$h(\gamma, u) = \gamma + \frac{1}{1 - \alpha} \sum_{\xi \in \Xi} \max\{0, u_\xi - \gamma\}, \quad u = (u_\xi, \xi \in \Xi),$$

and the mapping  $F : \mathbb{R}^{1+n} \rightarrow \mathbb{R}^{1+s}$ , where  $F(\gamma, x) = (\gamma, (f(\xi, x), \xi \in \Xi))$ ;  $s$  is the cardinality of  $\Xi$ . Thus, (3.11) is equivalently expressed as minimizing  $h(F(\gamma, x))$  over  $\gamma \in \mathbb{R}$  and  $x \in X$ . The main advantage of this perspective is that the potentially difficult functions  $f(\xi, \cdot)$ ,  $\xi \in \Xi$ , are separated from the remaining parts of the problem represented by  $h$  and  $X$ . If  $f(\xi, \cdot)$  is smooth, then we may linearize it and potentially avoid evaluating it a large number of times. The following algorithm and its convergence result is a modification of [186]; see [288, Section 6.F] for details.

**Proximal Composite Method.**

**Data.**  $z^0 \in \mathbb{R} \times X$ ,  $\tau \in (1, \infty)$ ,  $\sigma \in (0, 1)$ ,  $\bar{\lambda} \in (0, \infty)$ ,  $\lambda^0 \in (0, \bar{\lambda}]$ .

**Step 0.** Set  $\nu = 0$ .

**Step 1.** Compute

$$\bar{z}^\nu \in \operatorname{argmin}_{z \in \mathbb{R} \times X} \left\{ h(F(z^\nu) + \nabla F(z^\nu)(z - z^\nu)) + \frac{1}{2\lambda^\nu} \|z - z^\nu\|_2^2 \right\}.$$

If  $\bar{z}^\nu = z^\nu$ , then Stop.

**Step 2.** If

$$h(F(z^\nu)) - h(F(\bar{z}^\nu)) \geq \sigma \left( h(F(z^\nu)) - h(F(z^\nu) + \nabla F(z^\nu)(\bar{z}^\nu - z^\nu)) \right),$$

then set  $\lambda^{\nu+1} = \min\{\tau\lambda^\nu, \bar{\lambda}\}$  and go to Step 3.

Else, replace  $\lambda^\nu$  by  $\lambda^\nu/\tau$  and go to Step 1.

**Step 3.** Set  $z^{\nu+1} = \bar{z}^\nu$ , replace  $\nu$  by  $\nu + 1$ , and go to Step 1.

**3.10 Theorem** (proximal composite method). *For closed convex  $X \subset \mathbb{R}^n$  and twice continuously differentiable  $f(\xi, \cdot)$ ,  $\xi \in \Xi$ , suppose that the proximal composite method has generated  $\{z^\nu\}_{\nu=1}^\infty$  with a cluster point  $(\gamma^*, x^*)$ . Then,  $(\gamma^*, x^*)$  satisfies a necessary optimality condition<sup>6</sup> for (3.11):*

$$\exists y = (y_\xi, \xi \in \Xi) \in \mathbb{R}^s \text{ such that } (0, y) \in \partial h(F(\gamma^*, x^*)) \text{ and } - \sum_{\xi \in \Xi} y_\xi \nabla_x f(\xi, x^*) \in N_X(x^*).$$

The advantage of this approach is that any convex optimization algorithm can be brought in to solve the subproblem in Step 1. If  $X$  is polyhedral, then one may even solve the subproblem as a convex quadratic problem after introducing auxiliary variables. Step 2 tests whether the present linear approximation is sufficiently accurate and decreases  $\lambda^\nu$  if it is not.

**Dual Algorithms.** The dual formula (3.6) gives rise to the following approach. As recognized in [173, 175, 253], we find that

$$\text{s-rsk}_\alpha(f(\xi, x)) = \max_{\bar{p} \in \Delta_\alpha} \sum_{\xi \in \Xi} \bar{p}_\xi f(\xi, x) \approx \max_{\bar{p} \in \Delta_\alpha} \sum_{\xi \in \Xi} \bar{p}_\xi f(\xi, x) - \varepsilon \varphi(\bar{p}),$$

where  $\varepsilon > 0$  and  $\varphi$  is a nonnegative strongly convex real-valued function. Proposition 3 of [175] asserts that the approximation error vanishes as  $\varepsilon \rightarrow 0$  and the approximation is continuously differentiable if  $f(\xi, \cdot)$  is continuously differentiable for all  $\xi \in \Xi$ . Thus, one can minimize a superquantile approximately by minimizing a smooth approximation. A challenge is that each function and gradient computation of the approximating objective function requires the full set  $\Xi$ , which might be of high cardinality, especially in learning applications. We refer to [253] for further discussion of such issues. This dual smoothing approach can also be linked to primal smoothing methods; see [175, Corollary 6]. A related possibility is mentioned in [288, Example 8.32].

**Other Algorithmic Approaches.** By specializing the classical L-shaped method (see, e.g., [288, Section 5.H]), the paper [170] achieves a decomposition algorithm for minimizing superquantiles involving

<sup>6</sup>The normal cone to a set  $X$  at a point  $x$  is denoted by  $N_X(x)$ ; see, e.g., [288, Section 4.G].

quantities of interest that are affine in  $x$ ; see [89, 90] for further improvements. In the limited setting of simplex constraints and affine quantities of interest, [190] develops a three-phase algorithm that starts with a gradient descent heuristic applied to (3.11), continues with steps akin to the subgradient method, and ends with solving (3.12) from (hopefully) a good starting point using the simplex method. Simple active-set strategies to reduce the number of constraints in (3.12) generally appear highly beneficial [25, 40].

Essentially all the algorithms discussed in this subsection require that we can evaluate the quantity of interest  $f(\xi, x)$  as well as its gradients or subgradients with respect to  $x$ . If only function values are available, one can resort to black-box optimization [22] and other zeroth-order methods, e.g., [192]. If function values are computationally costly to compute (possibly even with noise), then one might resort to surrogates; cf. Example 2.5.

### 3.4 Estimating Superquantiles

Except when a random variable  $\xi$  is finitely distributed or follows some other standard distribution (cf. Examples 3.2, 3.3, and 3.4 as well as [223]), there is no explicit formula for its  $\alpha$ -superquantile  $\bar{Q}(\alpha)$ . If  $\xi$  is square integrable, then we have the bound from [273, Proposition 1] for any  $\alpha \in [0, 1)$ :

$$\mathbb{E}[\xi] \leq \bar{Q}(\alpha) \leq \min \left\{ \mathbb{E}[\xi] + \frac{\text{std}(\xi)}{\sqrt{1-\alpha}}, \sup \xi \right\}. \quad (3.16)$$

Much more accurate estimates of superquantiles are available via sampling. Suppose that  $\xi_1, \xi_2, \dots$  are independent random variables with the same distribution as  $\xi$ , which is assumed to be integrable. Then, for  $\alpha \in (0, 1)$ , the estimator

$$\bar{\mathbf{Q}}^\nu(\alpha) = \min_{\gamma \in \mathbb{R}} \gamma + \frac{1}{\nu(1-\alpha)} \sum_{i=1}^{\nu} \max\{0, \xi_i - \gamma\} \quad (3.17)$$

is strongly consistent, i.e.,  $\bar{\mathbf{Q}}^\nu(\alpha) \rightarrow \bar{Q}(\alpha)$  as  $\nu \rightarrow \infty$  almost surely ([306, Section 6.5.1] and [340]); see also [288, Example 8.57]. As that example shows, the consistency carries over in the sense of epi-convergence to the case when  $\xi$  is replaced by a quantity of interest  $f(\xi, x)$  that is continuous in  $x$ . Thus, passing from (3.10) to (3.11) by generating a finite sample independently is fundamentally sound, with cluster points of minimizers of the sample average approximation (3.11) being minimizers of the actual problem (3.10) almost surely. For superquantiles and many other risk measures, [302] furnishes a comprehensive treatment of strong consistency and [126] develops asymptotics for minimum values and minimizers obtained through solving sample average approximations as well as associated hypothesis testing methodology. Asymptotics for  $\bar{\mathbf{Q}}^\nu(\alpha)$  appeared already in [306, Section 6.5.1].

There is a wealth of concentration inequalities for bounding the probability that a sample average deviates from the mean with some  $\varepsilon$  for a fixed sample size. Hoeffding's inequality addresses bounded random variables, but extensions to subgaussian random variables and beyond are available; see, for example, [36]. These concentration inequalities can be made to hold uniformly in some sense across values of the auxiliary variable  $\gamma$  in the optimization problems giving  $\bar{\mathbf{Q}}^\nu(\alpha)$  and  $\bar{Q}(\alpha)$ ; see [307, Section

9.2.11] for a general discussion. An early effort to capitalize on these possibilities is [37], which assumes that the support  $\Xi$  of  $\xi$  is a bounded interval  $[0, \beta]$ . Then, one has for every  $\varepsilon \in (0, \infty)$  and  $\alpha \in (0, 1)$  that

$$\text{prob}\left\{\bar{Q}^\nu(\alpha) \geq \bar{Q}(\alpha) + \varepsilon\right\} \leq \exp\left(-2(\alpha\varepsilon/\beta)^2\nu\right) \quad (3.18)$$

and, provided that  $\xi$  is continuously distributed, one also has

$$\text{prob}\left\{\bar{Q}^\nu(\alpha) \leq \bar{Q}(\alpha) - \varepsilon\right\} \leq 3 \exp\left(-\alpha(\varepsilon/\beta)^2\nu/5\right).$$

The square dependence of  $\alpha$  in (3.18) is improved to linear dependence in [335]; see also [106] for a discussion of Berry–Essen bounds, the law of iterated logarithm, and large deviation results. Recent efforts for bounded random variables also include [8]. For unbounded subgaussian continuously distributed random variables, [155] achieves concentration bounds as long as the cumulative distribution functions are strictly increasing; see also [26] for a slight tightening and [258] for a treatment of random variables with only finite  $p$ th moment for  $p \in (1, 2]$ . Deviating from (3.17), [322] considers a more complicated estimator and achieves a strict improvement compared to [37], with strong empirical performance, but still in the setting of bounded random variables. The high-water mark for (3.17) appears to be the recent PAC-Bayesian bound in [213]. The paper [137] reviews developments up to 2014. Uniform bounds holding across a set of decisions appear in [181, 151]. The paper [114] examines central limit theorems for superquantiles and related conditional expectations.

Going beyond independent samples, [203] establishes strong laws of large numbers and associated convergence rates for  $\alpha$ -mixing<sup>7</sup> sequences and [202] considers a kernel estimator also for  $\alpha$ -mixing sequences. In the context of computationally expensive models of physical systems, [132, 131] discuss reduced-order models to guide a choice of importance sampling and reduce the sample size needed in (3.17). Similarly, [45] explores variance reduction via the cross-entropy method, [113] examines efficient multi-level nested Monte Carlo simulations, [179, 180] leverage multi-fidelity simulations and polynomial chaos expansion, and [48] uses metamodels based on stochastic kriging with further extensions in [83, 150] that also bring in extreme value theory. The paper [91] discusses how to construct a finite distribution that results in an accurate approximation of an actual problem involving superquantile minimization; see also [18] for other scenario reduction techniques.

## 4 Risk and Regret

As seen in Subsection 2.4, there are many possible measures of risk with superquantiles furnishing main examples. In this section, we discuss what constitutes a “good” measure of risk, both from modeling and computing points of view. We introduce measures of regret as key building blocks for constructing measures of risk and discuss how superquantiles generate many, if not most, meaningful measures of risk. The section ends with a review of the role risk measures play in achieving *fairness*.

---

<sup>7</sup>The  $\alpha$  here is of course unrelated to the  $\alpha$  specifying which superquantile is being considered.

## 4.1 Desirable Properties

The purpose of a risk measure is to model risk-averseness in a meaningful way and especially avoid embarrassing paradoxes that might discredit the modeling effort in the eyes of a decision maker. We would like to achieve this without introducing excessive computational complexity. So what properties would we like a measure of risk to possess? To formally answer this question, we view risk measures as functionals defined on a space of random variables.

For a given probability space  $(\Omega, \mathcal{A}, \mathbb{P})$ , we consider all square-integrable random variables:

$$\mathcal{L}^2 = \{\xi : \Omega \rightarrow \mathbb{R} \mid \mathbb{E}[\xi^2] < \infty\}.$$

As seen in [294, 52, 99, 70], one may consider alternative spaces, with the choice possibly depending on the risk measures of interest. For a detailed discussion of suitable spaces on which to define risk measures, see the papers [248, 102]. In this survey, we follow [276] and limit the scope to  $\mathcal{L}^2$ , which simplifies the exposition significantly.

We equip  $\mathcal{L}^2$  with the standard *norm*  $\|\xi\|_{\mathcal{L}^2} = (\mathbb{E}[\xi^2])^{1/2}$ . For  $\{\xi, \xi^\nu \in \mathcal{L}^2, \nu = 1, 2, \dots\}$ , the random variables  $\xi^\nu$  *converge* to the random variable  $\xi$ , written  $\xi^\nu \rightarrow \xi$ , when  $\|\xi^\nu - \xi\|_{\mathcal{L}^2} \rightarrow 0$  or, equivalently, when  $\mathbb{E}[(\xi^\nu - \xi)^2] \rightarrow 0$ . A subset  $\mathcal{C} \subset \mathcal{L}^2$  is *closed* if  $\{\xi^\nu \in \mathcal{C}, \nu = 1, 2, \dots\}$  and  $\xi^\nu \rightarrow \xi$  imply  $\xi \in \mathcal{C}$ . The set  $\mathcal{C}$  is *convex* if  $(1 - \lambda)\xi_0 + \lambda\xi_1 \in \mathcal{C}$  for all  $\lambda \in [0, 1]$  and  $\xi_0, \xi_1 \in \mathcal{C}$ . We recall that a random variable  $\xi \in \mathcal{L}^2$  is *constant* if  $\text{prob}\{\xi = \alpha\} = \mathbb{P}(\{\omega \in \Omega \mid \xi(\omega) = \alpha\}) = 1$  for some  $\alpha \in \mathbb{R}$ . The constant random variable with value 0 is denoted by  $\mathbf{0}$ . The underlying probability space  $(\Omega, \mathcal{A}, \mathbb{P})$  is *finite* if the cardinality of  $\Omega$  is finite, which means that any random variable in  $\mathcal{L}^2$  has a finite number of possible values, i.e., is finitely distributed. Throughout, we adopt the usual extended real-valued arithmetic rules such as  $0 \cdot \infty = 0$ ,  $0 \cdot (-\infty) = 0$ ,  $-\infty + \infty = \infty$ , and  $\infty - \infty = \infty$ ; see, e.g., [282, Section 1.E] or [288, Section 1.D].

A functional  $\mathcal{F} : \mathcal{L}^2 \rightarrow [-\infty, \infty]$  may satisfy any of the following properties:

<i>Constancy:</i>	$\mathcal{F}(\xi) = \alpha$ whenever $\xi$ is constant with value $\alpha \in \mathbb{R}$ .
<i>Averseness:</i>	$\mathcal{F}(\xi) > \mathbb{E}[\xi]$ for all nonconstant $\xi$ .
<i>Convexity:</i>	$\mathcal{F}((1 - \lambda)\xi_0 + \lambda\xi_1) \leq (1 - \lambda)\mathcal{F}(\xi_0) + \lambda\mathcal{F}(\xi_1) \quad \forall \xi_0, \xi_1, \lambda \in [0, 1]$ .
<i>Lower semicontinuity:</i>	$\{\mathcal{F} \leq \alpha\} = \{\xi \in \mathcal{L}^2 \mid \mathcal{F}(\xi) \leq \alpha\}$ is closed $\forall \alpha \in \mathbb{R}$ .
<i>Positive homogeneity:</i>	$\mathcal{F}(\lambda\xi) = \lambda\mathcal{F}(\xi) \quad \forall \lambda \in [0, \infty)$ .
<i>Monotonicity:</i>	$\mathcal{F}(\xi_0) \leq \mathcal{F}(\xi_1)$ when $\xi_0(\omega) \leq \xi_1(\omega)$ for $\mathbb{P}$ -almost every $\omega \in \Omega$ .
<i>Law invariance:</i>	$\mathcal{F}(\xi_0) = \mathcal{F}(\xi_1)$ when $\xi_0, \xi_1$ have the same distribution.

The constancy requirement is certainly reasonable and in fact a prerequisite for a measure of risk; see Definition 2.7. The averseness property excludes the possibility  $\mathcal{F}(\xi) = \mathbb{E}[\xi]$ , which is better treated separately. The purpose of a risk measure is after all to model risk-averseness.

The convexity property is key to make sure that a risk measure is computationally attractive. As we see from the comparison between quantiles and superquantiles in Example 3.6, it is also important in promoting hedging. Generally, the measure of risk  $\mathcal{R}(\xi) = Q(\alpha)$ , where  $Q(\alpha)$  is the  $\alpha$ -quantile of  $\xi$ ,

fails the convexity requirement. For example, consider two statistically independent random variables  $\xi_0, \xi_1$  with common density function value 0.9 on  $[-1, 0]$ , value 0.05 on  $(0, 2]$ , and value zero otherwise. Then, the 0.9-quantile is 0 for both  $\xi_0$  and  $\xi_1$ , but the 0.9-quantile for  $\xi_0/2 + \xi_1/2$  is approximately 0.25.

Lower semicontinuity (lsc) is a technical condition beneficial in convex analysis (as seen in Subsection 6). We know even from finite-dimensional optimization that the minimum of a function over a compact set may not be attained if the function is not lsc. The lsc property holds for example when  $\mathcal{F}$  is real-valued and convex and either  $\mathcal{F}$  is monotone or the probability space is finite. In fact, then  $\mathcal{F}$  is continuous; see [276, Equation (3.5)] for a summary of this claim which in turn relies on the fundamental Proposition 3.1 about continuity of convex and monotone functionals in [294].

Positive homogeneity implies a certain invariance to scaling. If we convert the quantity of interest from dollar to euro, then the associated risk should not fundamentally change. However, this property fails for the variance  $\mathcal{F}(\xi) = \text{var}(\xi) = (\text{std}(\xi))^2$ .

Monotonicity is a natural requirement because risk would typically be deemed less for  $\xi_0$  than for  $\xi_1$  when, for every pair of outcomes  $\{\xi_0(\omega), \xi_1(\omega)\}$ , the random variable  $\xi_0$  never comes out worse than  $\xi_1$ . The mean-plus-standard-deviation risk measure  $\mathcal{R}(\xi) = \mathbb{E}[\xi] + \lambda \text{std}(\xi)$  fails this requirement because the standard deviation does not distinguish between variability above and below the mean. Thus, a random variable with a high variability below the mean is deemed “high risk,” even though low values represent no “real” risk under our orientation. For example, consider the random variables  $\xi_0$  and  $\xi_1$  and their joint distribution that assigns probability 1/2 to the outcome  $(0, 0)$  and probability 1/2 to the outcome  $(-2, -1)$ . Thus, for every outcome the random variables have either the same value or  $\xi_0$  has a value below that of  $\xi_1$ . Still, with  $\lambda = 2$ , we obtain  $\mathbb{E}[\xi_0] + \lambda \text{std}(\xi_0) = 1$  and  $\mathbb{E}[\xi_1] + \lambda \text{std}(\xi_1) = 1/2$ .

In elementary probability theory, random variables are thought of as being “fully” described by their distributions. However, random variables are (measurable) functions from a sample space (here denoted by  $\Omega$ ) to the reals, with many random variables potentially having the same distribution. Thus, the law invariance property is not automatically satisfied but certainly appears reasonably in many situations.

For  $\alpha \in (0, 1]$ , the superquantile risk measure  $\text{s-rsk}_\alpha$  satisfies all the seven properties above [276];  $\text{s-rsk}_0 = \mathbb{E}[\cdot]$  misses the averseness requirement. It is also real-valued when  $\alpha \in [0, 1)$  by (3.16).

Constancy, convexity, and lsc properties together imply that

$$\mathcal{F}(\xi + \alpha) = \mathcal{F}(\xi) + \alpha \quad \forall \xi \in \mathcal{L}^2, \alpha \in \mathbb{R}.$$

Thus,  $\mathcal{F}$  quantifies, for example, risk in a manner that is translation invariant; adding a constant amount to an uncertain future cost changes the risk by that amount.

Pioneering works in finance [153, 19, 20, 68] identify many of these properties as desirable. Although originally expressed slightly differently, [20] defines a real-valued risk measure to be *coherent* if it satisfies the constancy, convexity, positive homogeneity, and monotonicity properties. The initial development in [20] was for finite probability spaces, with an extension to spaces of bounded random variables in [68]. The coherency of the quantity in (3.2) (which we now of course call a superquantile) was first established in [240] and, as the complete picture about the equivalences in Theorem 3.1 emerged, also in [3, 275].



Further developments under the names *convex measures of risk* and *convex risk functions* relaxed the positive homogeneity condition [97, 294]. We adopt the concept of *regularity* from [276] (with refinements in [271]), which is also used in [166].

**4.1 Definition** (regular measure of risk). *A regular measure of risk  $\mathcal{R}$  is a functional from  $\mathcal{L}^2$  to  $(-\infty, \infty]$  that is lsc, convex and also satisfies the constancy and averseness properties.*

As we see from the following development, these requirements align well with convex analysis while permitting a wide set of risk measures including those that may assign  $\infty$  to some random variables and that may lack positive homogeneity and monotonicity.

While measures of risk take a central role, there are other classes of functionals as well. In Subsection 2.5, we discuss utility and disutility functions. It turns out that particular disutility functions, which can be traced back to studies of optimized certainty equivalents [30] (see also [31]), emerge as important building blocks for regular measures of risk. Here, we adopt the definition in [271].

**4.2 Definition** (regular measure of regret). *A regular measure of regret  $\mathcal{V}$  is a functional from  $\mathcal{L}^2$  to  $(-\infty, \infty]$  that is lsc and convex, with*

$$\mathcal{V}(\mathbf{0}) = 0, \quad \text{but } \mathcal{V}(\boldsymbol{\xi}) > \mathbb{E}[\boldsymbol{\xi}] \quad \forall \boldsymbol{\xi} \neq \mathbf{0}.$$

*The quantity  $\mathcal{V}(\boldsymbol{\xi})$  is the regret of  $\boldsymbol{\xi}$ .*

We note that the definition of a regular measure of regret in [276] includes a limiting condition that is shown to be superfluous in [271].

**4.3 Example** (measures of regret and risk). *We have the following examples of regular measures of regret  $\mathcal{V}$  and regular measures of risk  $\mathcal{R}$ ; see, e.g., [288, Examples 8.8 and 8.10]:*

(a) Penalty regret and superquantile risk with  $\alpha \in (0, 1)$ :

$$\mathcal{V}(\boldsymbol{\xi}) = \frac{1}{1-\alpha} \mathbb{E}[\max\{0, \boldsymbol{\xi}\}] \quad \mathcal{R}(\boldsymbol{\xi}) = \text{s-rsk}_\alpha(\boldsymbol{\xi}).$$

(b) Worst-case regret and risk:

$$\mathcal{V}(\boldsymbol{\xi}) = \begin{cases} 0 & \text{if } \sup \boldsymbol{\xi} \leq 0 \\ \infty & \text{otherwise} \end{cases} \quad \mathcal{R}(\boldsymbol{\xi}) = \sup \boldsymbol{\xi}.$$

(c) Moment-based regret and mean-plus-standard-deviation risk with  $\lambda > 0$ :

$$\mathcal{V}(\boldsymbol{\xi}) = \mathbb{E}[\boldsymbol{\xi}] + \lambda \sqrt{\mathbb{E}[\boldsymbol{\xi}^2]} \quad \mathcal{R}(\boldsymbol{\xi}) = \mathbb{E}[\boldsymbol{\xi}] + \lambda \text{std}(\boldsymbol{\xi}).$$

(d) Mixed regret and risk with  $\lambda_1, \dots, \lambda_q > 0$  and  $\sum_{i=1}^q \lambda_i = 1$ :

$$\mathcal{V}(\boldsymbol{\xi}) = \inf_{\gamma_1, \dots, \gamma_q} \left\{ \sum_{i=1}^q \lambda_i \mathcal{V}_i(\boldsymbol{\xi} - \gamma_i) \mid \sum_{i=1}^q \lambda_i \gamma_i = 0 \right\} \quad \mathcal{R}(\boldsymbol{\xi}) = \sum_{i=1}^q \lambda_i \mathcal{R}_i(\boldsymbol{\xi})$$

when  $\mathcal{V}_1, \dots, \mathcal{V}_q$  ( $\mathcal{R}_1, \dots, \mathcal{R}_q$ ) are regular measures of regret (risk).

Several additional examples appear in Section 7. The distinction between regular measures of regret and risk emerges from the assessment of constant random variables. While  $\mathcal{V}(\mathbf{0}) = \mathcal{R}(\mathbf{0}) = 0$ , a constant random variable  $\boldsymbol{\xi}$  with value  $\alpha \neq 0$  is treated differently:

$$\mathcal{V}(\boldsymbol{\xi}) > \alpha, \quad \text{but} \quad \mathcal{R}(\boldsymbol{\xi}) = \alpha.$$

This seemingly minor discrepancy results in profoundly different roles, with regular measures of regret providing a means to construct regular measures of risk and thereby extending the useful relationship between  $\frac{1}{1-\alpha}\mathbb{E}[\max\{0, \boldsymbol{\xi}\}]$  and  $\text{s-rsk}_\alpha(\boldsymbol{\xi})$  in the formula (3.2) for superquantiles.

**4.4 Theorem** (regret-risk). *For a regular measure of regret  $\mathcal{V}$ , a regular measure of risk  $\mathcal{R}$  emerges from*

$$\mathcal{R}(\boldsymbol{\xi}) = \min_{\gamma \in \mathbb{R}} \gamma + \mathcal{V}(\boldsymbol{\xi} - \gamma).$$

*For every regular measure of risk  $\mathcal{R}$  there exists a regular measure of regret  $\mathcal{V}$ , not necessarily unique, that constructs  $\mathcal{R}$  through this minimization formula.*

**Proof.** Under slightly more restrictive assumptions, [276] establishes the first conclusion. The present form is taken from [271, Theorem 2.2]; see [288, Theorem 8.9] for the existence claim.  $\square$

The theorem generalizes the situation for superquantiles as established by the formula (3.2): a regular measure of risk finds a “best” way of reducing displeasure with a mix of uncertain outcomes by optimizing a scalar  $\gamma$ . Thus, a regular measure of risk provides a deeper assessment of a random variable compared to a regular measure of regret. We note that by writing “min” instead of “inf” we assert that the minimum value in the theorem is attained for some  $\gamma$ .

As seen from [288, Example 8.10], which in part relies on [276], the regular measures of regret and risk in Example 4.3(a,b,c) pair in the sense of Theorem 4.4. Moreover, if  $\mathcal{V}_i$  and  $\mathcal{R}_i$  pair by satisfying the formula in Theorem 4.4, then  $\mathcal{V}$  and  $\mathcal{R}$  in Example 4.3(d) also satisfy that formula.

Generally, Theorem 4.4 provides a path to constructing new regular measures of risk from regular measures of regret, which in turn might be motivated by disutility functions.

## 4.2 Risk Minimization

A regular measure of risk offers key computational advantages when implemented in minimization problems. Suppose that  $f : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$  represents a quantity of interest,  $\boldsymbol{\xi}$  is an  $m$ -dimensional random vector,  $X \subset \mathbb{R}^n$ ,  $\mathcal{R}$  is a regular measure of risk, and  $f(\boldsymbol{\xi}, x) \in \mathcal{L}^2$  for all  $x \in \mathbb{R}^n$ . Then, the problem

$$\underset{x \in X}{\text{minimize}} \quad \varphi(x) = \mathcal{R}(f(\boldsymbol{\xi}, x)) \tag{4.1}$$

is convex provided that  $X$  is convex and  $f(\boldsymbol{\xi}, \cdot)$  is affine for all  $\boldsymbol{\xi} \in \Xi$ , the support of  $\boldsymbol{\xi}$ ; see, e.g., [288, Example 8.11] for a proof. Thus, a regular measure of risk does not disrupt the convexity of the affine functions  $f(\boldsymbol{\xi}, \cdot)$ .

For any regular measure of regret  $\mathcal{V}$  that pairs with  $\mathcal{R}$  via Theorem 4.4, the problem (4.1) is equivalent to

$$\underset{x \in X, \gamma \in \mathbb{R}}{\text{minimize}} \quad \gamma + \mathcal{V}(f(\boldsymbol{\xi}, x) - \gamma), \quad (4.2)$$

which might be computationally more attractive; cf. the situation for superquantiles in (3.9) and (3.10).

If the regular measure of risk in (4.1) is also monotone, then (4.1) is a convex problem when  $X$  is convex and  $f(\xi, \cdot)$  is convex for all  $\xi \in \Xi$ ; see, e.g., [288, Example 8.13]. Such monotonicity holds if  $\mathcal{V}$  is monotone as can be seen directly. Then, (4.2) is also a convex problem when  $X$  is convex and  $f(\xi, \cdot)$  is convex for all  $\xi \in \Xi$ .

The seemingly complicated situation for mixed risk  $\mathcal{R}$  from Example 4.3(d) simplifies even further. In that case, (4.1) is equivalent to

$$\underset{x \in X, \gamma_0, \gamma_1, \dots, \gamma_q}{\text{minimize}} \quad \gamma_0 + \sum_{i=1}^q \lambda_i \mathcal{V}_i(f(\boldsymbol{\xi}, x) - \gamma_0 - \gamma_i) \quad \text{subject to} \quad \sum_{i=1}^q \lambda_i \gamma_i = 0,$$

where  $\mathcal{V}_1, \dots, \mathcal{V}_q$  are the regular measures of regret that pair with the regular measures of risk  $\mathcal{R}_1, \dots, \mathcal{R}_q$ , which in turn produce the mixed risk  $\mathcal{R}$ . Reformulations of this kind extend to situations with multiple quantities of interest, possibly some emerging as constraints.

Optimality conditions are important in algorithmic development and to interpret solutions of optimization problems. For a discussion of this subject, we refer to [294, 280, 165]; see also the subgradient formulas in Subsection 6.4. Algorithms for computing stationary points of nonconvex, nonsmooth risk-minimization problems emerge from [195].

**Risk Measures from Expectation Regret.** A main approach to constructing a regular measure of risk is to start with a lsc convex function  $v : \mathbb{R} \rightarrow (-\infty, \infty]$  that also satisfies the properties

$$v(0) = 0, \quad \text{but} \quad v(\xi) > \xi \quad \forall \xi \in \mathbb{R} \setminus \{0\}.$$

The value  $v(\xi)$  could reflect a decision maker's displeasure with an outcome  $\xi$ . The expectation theorem in [276] shows that the functional on  $\mathcal{L}^2$  defined by  $\mathcal{V}(\boldsymbol{\xi}) = \mathbb{E}[v(\boldsymbol{\xi})]$  is a regular measure of regret and thus defines a regular measure of risk via Theorem 4.4:  $\mathcal{R}(\boldsymbol{\xi}) = \min_{\gamma \in \mathbb{R}} \gamma + \mathbb{E}[v(\boldsymbol{\xi} - \gamma)]$ . Superquantiles fall within this class with  $v(\xi) = \max\{0, \xi\}/(1 - \alpha)$  for  $\alpha \in (0, 1)$ . If  $v$  is not bounded from above by a quadratic function, then these regular measures of risk and regret may not be real-valued (on  $\mathcal{L}^2$ ). Generally, one can view  $v$  as a normalized disutility function; see [31, 276, 267]. Minimization of regular measures of risk constructed in this manner results in expectation minimization problems after passing from (4.1) to (4.2). Thus, one can leverage the vast array of algorithms for such problems including the SGD method (cf. Subsection 3.3), the broad theory of  $M$ -estimators [134], and also tackle extensions in the context of PDE-constrained optimization using local approximations of function values and gradients [354] or low-rank tensor approximations of random fields [16].

**Approximations using Epi-Regularization.** As exemplified by superquantile risk and penalty regret (see Example 4.3(a)), a regular measure of regret  $\mathcal{V}$  when applied to a quantity of interest  $f(\xi, x)$

might result in a nonsmooth function  $x \mapsto \mathcal{V}(f(\boldsymbol{\xi}, x))$ , which is less ideal when solving (4.2). Likewise,  $x \mapsto \mathcal{R}(f(\boldsymbol{\xi}, x))$  is often nonsmooth. However, it is possible to construct approximations using epi-regularization [166]; see also [38]. Specifically, one can replace a regular measure of risk  $\mathcal{R}$  by the approximation

$$\mathcal{R}_\varepsilon^\Phi(\boldsymbol{\xi}) = \inf \{ \mathcal{R}(\boldsymbol{\xi} - \boldsymbol{\eta}) + \varepsilon \Phi(\boldsymbol{\eta}/\varepsilon) \mid \boldsymbol{\eta} \in \mathcal{L}^2 \},$$

where  $\varepsilon \in (0, \infty)$  and  $\Phi : \mathcal{L}^2 \rightarrow (-\infty, \infty]$  is lsc and convex, with  $\Phi(\boldsymbol{\xi}) < \infty$  for at least one  $\boldsymbol{\xi} \in \mathcal{L}^2$ . Since there is much flexibility, we can choose  $\Phi$  such that  $\mathcal{R}_\varepsilon^\Phi$  has desirable properties. For example,  $\Phi(\boldsymbol{\xi}) = \frac{1}{2}\mathbb{E}[\boldsymbol{\xi}^2]$  produces in the case of  $\alpha$ -superquantiles (i.e.,  $\mathcal{R} = \text{s-rsk}_\alpha$ ) with  $\alpha \in (0, 1)$  that

$$\mathcal{R}_\varepsilon^\Phi(\boldsymbol{\xi}) = \min_{\gamma \in \mathbb{R}} \gamma + \frac{1}{1-\alpha} \mathbb{E}[v_\varepsilon(\boldsymbol{\xi} - \gamma)], \quad \text{where } v_\varepsilon(\xi) = \begin{cases} 0 & \text{if } \xi \leq 0 \\ \frac{1}{2\varepsilon}\xi^2 & \text{if } \xi \in (0, \frac{\varepsilon}{1-\alpha}) \\ \frac{1}{1-\alpha}(\xi - \frac{\varepsilon}{2(1-\alpha)}) & \text{otherwise.} \end{cases}$$

Since  $v_\varepsilon$  is continuously differentiable, we may approximate  $\mathcal{V}$  in (4.2) by the functional  $\mathcal{V}_\varepsilon(\boldsymbol{\xi}) = \mathbb{E}[v_\varepsilon(\boldsymbol{\xi})]$  and reap computational benefits. The approximation error will be small when  $\varepsilon$  is small because then  $v_\varepsilon$  accurately approximates  $\xi \mapsto \max\{0, \xi\}/(1-\alpha)$ ; see [166] for details. However, the functional  $\mathcal{V}_\varepsilon$  is not a regular measure of regret;  $\mathcal{V}_\varepsilon(\boldsymbol{\xi})$  might not be strictly larger than  $\mathbb{E}[\boldsymbol{\xi}]$ .

**Approximations using Sampling.** Almost always in practice, the “true” probability distribution of a random variable  $\boldsymbol{\xi}$  is unknown but one might have the ability to generate a sample  $\xi_1, \dots, \xi_\nu$ , which produces an empirical distribution. Let  $\boldsymbol{\xi}^\nu$  be a random variable with this empirical distribution, i.e., each of the values  $\xi_i$ ,  $i = 1, \dots, \nu$ , occurs with probability  $1/\nu$ . This means that while the actual risk  $\mathcal{R}(\boldsymbol{\xi})$  might not be computable because the distribution of  $\boldsymbol{\xi}$  is unknown, that of  $\boldsymbol{\xi}^\nu$  tends to be available but of course is also a bit different. From [302] we know that if  $\mathcal{R}$  is a law-invariant, real-valued, monotone, and regular risk measure, then  $\mathcal{R}(\boldsymbol{\xi}^\nu) \rightarrow \mathcal{R}(\boldsymbol{\xi})$  as the sample size  $\nu \rightarrow \infty$  almost surely<sup>8</sup>. In fact, this also holds beyond random variables defined on  $\mathcal{L}^2$ ; see [302] for details and [69] for recent refinements. (We note that the results from [302] and those in the following theorem rely on the minor technical assumption that the probability space  $(\Omega, \mathcal{A}, \mathbb{P})$  is atomless and complete.)

**4.5 Theorem** (epi-convergence under sampling). *For a law-invariant, monotone, and regular measure of risk  $\mathcal{R} : \mathcal{L}^2 \rightarrow \mathbb{R}$ , consider  $f : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$ , an  $m$ -dimensional random vector  $\boldsymbol{\xi}$  with support  $\Xi$ , and its sample-based approximation  $\boldsymbol{\xi}^\nu$  (see the previous paragraph). Suppose that the following hold:*

- (a)  $f(\boldsymbol{\xi}, x) \in \mathcal{L}^2$  for all  $x \in \mathbb{R}^n$ .
- (b)  $f$  is a random lsc function<sup>9</sup>.
- (c) For every  $x \in \mathbb{R}^n$ , there is neighborhood  $N$  and  $g : \mathbb{R}^m \rightarrow [0, \infty)$  such that  $g(\boldsymbol{\xi}) \in \mathcal{L}^2$  and  $f(\boldsymbol{\xi}, x) \geq g(\boldsymbol{\xi})$  for all  $x \in N$  and  $\boldsymbol{\xi} \in \Xi$ .

<sup>8</sup>Since the sample that generates  $\boldsymbol{\xi}^\nu$  is random, the “almost surely” refers to the fact that the convergence holds for all samples  $\{\xi_1, \xi_2, \dots\}$  in a set of probability one.

<sup>9</sup>A sufficient (but not necessary) condition for  $f$  to be random lsc is that it is measurable in its first argument and continuous in its second argument; see, e.g., [288, Section 8.G].

Consider the functions  $\varphi, \varphi^\nu : \mathbb{R}^n \rightarrow \mathbb{R}$  given by  $\varphi(x) = \mathcal{R}(f(\boldsymbol{\xi}, x))$  and  $\varphi^\nu(x) = \mathcal{R}(f(\boldsymbol{\xi}^\nu, x))$ . Then,  $\varphi$  is lsc and  $\varphi^\nu$  epi-converges to  $\varphi$  as  $\nu \rightarrow \infty$  almost surely.

Moreover, if  $f(\xi, \cdot)$  is convex for all  $\xi \in \Xi$ , then assumptions (b,c) can be dropped and the conclusion strengthened to  $\varphi^\nu$  converging to  $\varphi$  uniformly on compact sets as  $\nu \rightarrow \infty$  almost surely.

**Proof.** The result appears in [302], which also handles integrable random variables.  $\square$

The main consequence of  $\varphi^\nu$  epi-converging to  $\varphi$  is that all cluster points of a sequence of minimizers of  $\varphi^\nu$  are minimizers of  $\varphi$ . For an introduction to epi-convergence, we refer to [288, Section 4.C]. Thus, under the assumptions of the theorem, we are on solid footing if we approximate a distribution by an empirical distribution in a risk-minimization problem. For the effect of more general changes to a probability distribution, we refer to [247, 152, 56, 86]. In [245], we find expressions for the asymptotic distributions of sample-based approximations of law-invariant risk measures; far reaching extensions emerge in [73] for nested expectations and, recently, in [71] for kernels, wavelets, and other estimators. Asymptotics for minimum values and minimizers appear in [126]. We also have uniform error bounds for broad classes of risk measures [181, 151, 183]. In the specific setting of a derivative-free stochastic mirror descent algorithm and spectral risk measures, [136] develops generalization errors of the expectation and high-probability kinds. The paper [143] establishes convergence rates of stochastic subgradient methods for mean-plus-semideviation risk measures, while [113] obtains rates for multi-level nested simulations of risk measures.

### 4.3 Mixed Superquantiles and Law Invariance

We can extend the mixing of risk measures in Example 4.3(d) to integrals over an uncountable collection of risk measures, especially superquantiles, to produce new regular measures of risk as pioneered by Kusuoka [171]. This subsection confirms that a large class of reasonable risk measures stem from superquantiles in this manner, which further highlights the centrality of superquantiles in decision making under uncertainty.

**4.6 Definition** (mixed superquantile measure of risk). A weighting measure<sup>10</sup>  $\lambda$  produces a mixed superquantile measure of risk  $\mathcal{R} : \mathcal{L}^2 \rightarrow (-\infty, \infty]$  given by

$$\mathcal{R}(\boldsymbol{\xi}) = \int_0^1 \bar{Q}(\beta) d\lambda(\beta), \quad (4.3)$$

where  $\bar{Q}(\beta)$  is the  $\beta$ -superquantile of  $\boldsymbol{\xi}$ .

If the weighting measure is supported on  $\{\alpha_i \in [0, 1), i = 1, \dots, q\}$ , then the resulting mixed superquantile measure of risk reduces to Example 4.3(d) with  $\mathcal{R}_i = \text{s-rsk}_{\alpha_i}$ ,  $i = 1, \dots, q$ . A common choice is  $\alpha_1 = 0$  and  $\alpha_2 = \alpha$ , which produces  $\mathcal{R}(\boldsymbol{\xi}) = \lambda_1 \mathbb{E}[\boldsymbol{\xi}] + \lambda_2 \text{s-rsk}_\alpha(\boldsymbol{\xi})$ , where  $\lambda_1, \lambda_2$  are nonnegative weights summing to 1; see for example [102] for its use in machine learning and Section 7.

<sup>10</sup>A weighting measure is a probability measure on the Borel subsets of  $[0, 1)$ .

The definition of mixed superquantile measures of risk allows for “averaging” more than a finite number of superquantiles. Given  $\alpha \in [0, 1)$ , one possibility places zero weight on superquantiles  $\bar{Q}(\beta)$  for  $\beta \in [0, \alpha)$  and weight  $1/(1 - \alpha)$  on  $\bar{Q}(\beta)$  for  $\beta \in [\alpha, 1)$ . This produces the  $\alpha$ -second-order superquantile of a random variable  $\xi$  [272]:

$$\bar{\bar{Q}}(\alpha) = \frac{1}{1 - \alpha} \int_{\alpha}^1 \bar{Q}(\beta) d\beta. \quad (4.4)$$

Second-order superquantiles define measures of risk, which are more conservative than those based on quantiles and superquantiles. We see that  $Q(\alpha) \leq \bar{Q}(\alpha) \leq \bar{\bar{Q}}(\alpha)$  regardless of  $\alpha \in (0, 1)$ . In addition to producing new measures of risk, second-order superquantiles also underpin the argmin-formula for superquantiles (3.5); see [272, 160].

Properties of mixed superquantile risk measures are traced back to [2, 278, 279]; here we summarize key facts as given in [272].

**4.7 Proposition** (mixed superquantile properties). *A mixed superquantile risk measure  $\mathcal{R}$  as defined in (4.3) is well-defined, monotone, and positively homogeneous. It is regular if  $\lambda(\{0\}) < 1$ , but lacking averseness if  $\lambda(\{0\}) = 1$ . It is real-valued on  $\mathcal{L}^2$  whenever  $\lambda$  satisfies  $\int_0^1 (1 - \beta)^{-1/2} d\lambda(\beta) < \infty$  and, regardless of the weighting measure, has  $\mathcal{R}(\xi) < \infty$  whenever  $\sup \xi < \infty$ .*

*In terms of the quantile function  $Q$  of  $\xi$ , one has the alternative expression*

$$\mathcal{R}(\xi) = \int_0^1 Q(\beta) \varphi(\beta) d\beta, \text{ where } \varphi(\beta) = \int_{0 \leq \alpha < \beta} \frac{1}{1 - \alpha} d\lambda(\alpha), \beta \in [0, 1]. \quad (4.5)$$

*The risk profile function  $\varphi$  is right-continuous and nondecreasing on  $[0, 1]$  with  $\varphi(0) = 0$  and satisfies  $\int_0^1 (1 - \alpha) d\varphi(\alpha) = 1$ . Conversely, any  $\varphi$  with these properties arises from a unique weighting measure  $\lambda$  given by  $d\lambda(\alpha) = (1 - \alpha) d\varphi(\alpha)$ .*

The alternative expression in (4.5) makes a deep connection with the integral formula (3.3) for a (single) superquantile. The latter averages all quantiles above  $\alpha$ . In contrast, (4.5) potentially weighs the quantiles differently using a risk profile function  $\varphi$ . The resulting risk measures are also known as spectral risk measures [2]. They connect fundamentally with distortion functionals [241] common in insurance applications. There are further connections with dual utility theory [344]; see also [75].

We refer to [160] for ways to compute mixed superquantile risk using numerical integration. The paper [329] exemplifies recent efforts to use mixed superquantile risk measures in reinforcement learning. Since the risk profile function models a decision maker’s preferences, it is often unknown and [333] examines worst-case models over a class of such functions; see also [128]. The paper [187] examines worst-case values of law-invariant measures of risk under only partial information about random variables.

While the class of mixed superquantile risk measures is clearly large, the pioneering contribution in [171] was to show that *every* real-valued, law-invariant, positively homogeneous, monotone, regular measure of risk  $\mathcal{R}$  can be written in the form

$$\mathcal{R}(\xi) = \sup_{\lambda \in \Lambda} \int_0^1 \bar{Q}(\beta) d\lambda(\beta)$$

for some set  $\Lambda$  of weighting measures. This is the *Kusuoka representation* of  $\mathcal{R}$ . (The paper [171] showed the existence of such  $\Lambda$  for bounded random variables, with extensions to  $p$ -integrable random variables appearing in [244], refinements about uniqueness being established by [303], and insight in the case of atomic probability spaces emerging from [229].) Thus, superquantiles are the fundamental building blocks of a large class of meaningful measures of risk.

An immediate consequence of the Kusuoka representation is the following. For two random variables  $\xi_0$  and  $\xi_1$  and *any* real-valued, law-invariant, positively homogeneous, monotone, regular measure of risk  $\mathcal{R}$ , one has that

$$\text{s-rsk}_\alpha(\xi_0) \leq \text{s-rsk}_\alpha(\xi_1) \quad \forall \alpha \in [0, 1) \quad \implies \quad \mathcal{R}(\xi_0) \leq \mathcal{R}(\xi_1).$$

Thus, if  $\xi_0$  dominates  $\xi_1$  in the sense of the left-hand side of the implication, then it becomes less critical to determine the “right” measure of risk; most meaningful risk measures will prefer  $\xi_0$  over  $\xi_1$ . (This notion of dominance is equivalent to second-order stochastic dominance; see [72, 269].) Even without such complete dominance, we can plot  $\text{s-rsk}_\alpha(\xi_0)$  and  $\text{s-rsk}_\alpha(\xi_1)$  as functions of  $\alpha$  to highlight the pros and cons with each decision (random variable). In the context of machine learning, [102] generates such plots for the purpose of comparing different statistical models and their resulting errors.

#### 4.4 Fairness

Fairness is an emerging area for application of risk measures. As pioneered in [339] (see also [184, 102]), one can even out the performance of a classifier across various subgroups by replacing the usual expectation with another risk measure during the training. Specifically, consider a quantity of interest  $f((\xi, \eta), c)$  representing loss or estimation error, where  $c \in \mathbb{R}^n$  is a vector of coefficients specifying a statistical model (e.g., a neural network) and  $(\xi, \eta) \in \mathbb{R}^m$  is a vector consisting of input (features)  $\xi$  and output (labels)  $\eta$ . An algorithm for supervised learning may seek to find  $c$  that minimizes  $\mathbb{E}[f((\xi, \eta), c)]$ , where  $(\xi, \eta)$  is a random vector with some assumed probability distribution typically taken as the empirical distribution of a training data set. The choice of expectation as risk measure may produce highly uneven performance of the resulting statistical model across subgroups. For certain types of outcomes of  $(\xi, \eta)$ , for instance those corresponding to male customers, the statistical model may over-estimate the output  $\eta$ , while for other outcomes, corresponding to female customers, the model may under-estimate the output. Biases of these kinds are problematic and sometimes illegal.

Suppose that we augment the input-output vector by also including an input  $\zeta$ , which represents sensitive characteristics such as gender. Thus, the random vector  $(\xi, \eta, \zeta)$  has now three parts. As proposed in [339], we may change the training problem from that of minimizing  $\mathbb{E}[f((\xi, \eta, \zeta), c)]$  to solving

$$\underset{c \in \mathbb{R}^n}{\text{minimize}} \mathcal{R}(g(\zeta, c)), \quad \text{where } g(\zeta, c) = \mathbb{E}[f((\xi, \eta, \zeta), c) \mid \zeta = \zeta] \quad (4.6)$$

and  $\mathcal{R}$  is a measure of risk that acts on a new quantity of interest  $g(\zeta, c)$ , which is uncertain because the value of  $\zeta$  is governed by the distribution of  $\zeta$ . For a fixed  $\zeta$  (say “male”) and coefficient vector  $c$  describing a statistical model,  $g(\zeta, c)$  is the conditional expectation of the original quantity of interest, given  $\zeta = \zeta$ . Thus, it quantifies the average performance of the statistical model when applied to

subjects with sensitive feature  $\zeta$ . If  $g(\zeta, c)$  is much higher for some values of  $\zeta$  than others, then the statistical model specified by  $c$  might be perceived as unfair. The role of the risk measure  $\mathcal{R}$  is to assess the random variable  $g(\zeta, c)$  and to safeguard against high values. Thus, a solution  $c^*$  of (4.6) tends to specify a statistical model that performs reasonably well regardless of sensitive characteristics. We avoid an upper tail in the distribution of  $g(\zeta, c^*)$  that extends far to the right. Since quantities of interest in machine learning are often nonnegative, the left tail may also be “light.” This means that  $g(\zeta, c^*)$  has little variability, which might be our goal. (The paper [339] defines *perfect fairness* of a statistical model given by  $c$  as having  $g(\zeta, c) = g(\zeta', c)$  for all possible values  $\zeta, \zeta'$ .)

In certain settings, one is not allowed (by law or company policy) to use sensitive data during the training and testing of an artificial intelligence system [130]. This precludes a direct application of the above approach. Still, by switching from expected loss to a risk measure, we can control the upper tail of the resulting loss distribution across all individuals. This safeguards against statistical models with poor performance for some individuals and excellent performance for others. Consequently, even if the sensitive characteristic for each individual is unknown, we can achieve more consistent performance by using risk measures; see also [130] for a related approach based on distributionally robust optimization.

## 5 Error and Deviation

The broad approaches to decision making based on measures of risk and regret as developed in the previous section extend to statistical concepts such as mean-squared error, standard deviation, and generalized regression. Already in the pioneering work [211] on portfolio optimization, we find a discussion of semivariation, semideviation, and expected absolute deviation as means of quantifying “variability” in a random variable and the “skewness” of distributions; see, e.g., [312, 157, 230, 232] for related quantities and computational methods. The first systematic studies of general classes of *measures of deviation* and *measures of error* appear to be [278, 279, 281], which approach the concepts axiomatically. Further connections with measures of risk and regret appear in [276]; see [271] for technical refinements. This produces quadrangles of risk, regret, error, and deviation measures with new insights and computational possibilities.

### 5.1 Measures of Error

We start by examining how to quantify the “nonzeroness” of a random variable and follow [281, 276, 271].

**5.1 Definition** (regular measure of error and its statistic). *A regular measure of error  $\mathcal{E}$  is a functional from  $\mathcal{L}^2$  to  $[0, \infty]$  that is lsc and convex, with*

$$\mathcal{E}(\mathbf{0}) = 0, \quad \text{but } \mathcal{E}(\boldsymbol{\xi}) > 0 \quad \forall \boldsymbol{\xi} \neq \mathbf{0}.$$

*The quantity  $\mathcal{E}(\boldsymbol{\xi})$  is the error of  $\boldsymbol{\xi}$ . The corresponding statistic is  $\mathcal{S}(\boldsymbol{\xi}) = \operatorname{argmin}_{\gamma \in \mathbb{R}} \mathcal{E}(\boldsymbol{\xi} - \gamma)$ .*

Measures of error are important in regression analysis where the goal is to minimize the “nonzeroness” of certain random variables representing the “residual” difference between observations and



model predications. The statistic of  $\boldsymbol{\xi}$  is the set of scalars, possibly a single number, that best approximate a random variable  $\boldsymbol{\xi}$  in the sense of the corresponding regular measure of error. (Inherited from [281], we use “statistic” in a somewhat different meaning than in the Statistics literature.) We note that the definition in [276] includes a limiting condition that is shown to be superfluous in [271].

**5.2 Example** (regular measures of error). *We have the following examples of regular measures of error  $\mathcal{E}$  and corresponding statistics  $\mathcal{S}$ ; see, e.g., [288, Example 8.15] for justifications of regularity.*

(a) Koenker-Bassett error and quantile statistic<sup>11</sup> with  $\alpha \in (0, 1)$ :

$$\mathcal{E}(\boldsymbol{\xi}) = \frac{1}{1-\alpha} \mathbb{E}[\max\{0, \boldsymbol{\xi}\}] - \mathbb{E}[\boldsymbol{\xi}] \quad \mathcal{S}(\boldsymbol{\xi}) = [Q(\alpha), Q^+(\alpha)].$$

(b) Worst-case error and statistic:

$$\mathcal{E}(\boldsymbol{\xi}) = \begin{cases} -\mathbb{E}[\boldsymbol{\xi}] & \text{if } \sup \boldsymbol{\xi} \leq 0 \\ \infty & \text{otherwise} \end{cases} \quad \mathcal{S}(\boldsymbol{\xi}) = \begin{cases} \{\sup \boldsymbol{\xi}\} & \text{if } \sup \boldsymbol{\xi} < \infty \\ \emptyset & \text{otherwise.} \end{cases}$$

(c)  $\mathcal{L}^2$ -error and statistic with  $\lambda > 0$ :

$$\mathcal{E}(\boldsymbol{\xi}) = \lambda \sqrt{\mathbb{E}[\boldsymbol{\xi}^2]} = \lambda \|\boldsymbol{\xi}\|_{\mathcal{L}^2} \quad \mathcal{S}(\boldsymbol{\xi}) = \{\mathbb{E}[\boldsymbol{\xi}]\}.$$

(d) Mixed error and statistic with  $\lambda_1, \dots, \lambda_q > 0$  and  $\sum_{i=1}^q \lambda_i = 1$ :

$$\mathcal{E}(\boldsymbol{\xi}) = \inf_{\gamma_1, \dots, \gamma_q} \left\{ \sum_{i=1}^q \lambda_i \mathcal{E}_i(\boldsymbol{\xi} - \gamma_i) \mid \sum_{i=1}^q \lambda_i \gamma_i = 0 \right\} \quad \mathcal{S}(\boldsymbol{\xi}) = \sum_{i=1}^q \lambda_i \mathcal{S}_i(\boldsymbol{\xi}),$$

where  $\mathcal{S}_i(\boldsymbol{\xi}) = \operatorname{argmin}_{\gamma \in \mathbb{R}} \mathcal{E}_i(\boldsymbol{\xi} - \gamma)$  and  $\mathcal{E}_i$  is a regular measure of error.

In regression analysis, we seek to predict (or forecast) the value of a random variable  $\boldsymbol{\eta}$  from the values of some other random variables  $\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n$ . For instance, these other random variables could be input to a system, which we observe or perhaps even control, and  $\boldsymbol{\eta}$  is the output of the system, which we hope to predict. Let  $\boldsymbol{\xi} = (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n)$  and suppose that  $\boldsymbol{\eta}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n \in \mathcal{L}^2$ . We may seek a “best” statistical model of the form  $\gamma + \langle c, \boldsymbol{\xi} \rangle$  using a regular measure of error  $\mathcal{E}$ . This leads to the *generalized regression problem*

$$\underset{(\gamma, c) \in \mathbb{R}^{1+n}}{\text{minimize}} \quad \mathcal{E}(\boldsymbol{\eta} - \gamma - \langle c, \boldsymbol{\xi} \rangle). \quad (5.1)$$

A minimizer  $(\gamma^*, c^*)$  of the problem defines the random variable  $\gamma^* + \langle c^*, \boldsymbol{\xi} \rangle$ , which then makes the “nonzeroness” of  $\boldsymbol{\eta} - (\gamma + \langle c, \boldsymbol{\xi} \rangle)$  as low as possible in the sense of the selected measure of error. Thus,  $\gamma^* + \langle c^*, \boldsymbol{\xi} \rangle$  is a best possible approximation of  $\boldsymbol{\eta}$  in this sense. Since  $\mathcal{E}$  is convex, it follows that the function  $(\gamma, c) \mapsto \mathcal{E}(\boldsymbol{\eta} - \gamma - \langle c, \boldsymbol{\xi} \rangle)$  is convex. Thus, (5.1) is computationally appealing.

---

<sup>11</sup>For the definition of  $Q^+(\alpha)$ , see (3.8).

The choice in Example 5.2(c) leads to *least-squares regression* regardless of  $\lambda > 0$ . Since the measure of error can be squared without affecting the set of minimizers, in this case one has

$$\operatorname{argmin}_{(\gamma,c) \in \mathbb{R}^{1+n}} \mathcal{E}(\boldsymbol{\eta} - \gamma - \langle c, \boldsymbol{\xi} \rangle) = \operatorname{argmin}_{(\gamma,c) \in \mathbb{R}^{1+n}} \mathbb{E} \left[ (\boldsymbol{\eta} - \gamma - \langle c, \boldsymbol{\xi} \rangle)^2 \right].$$

In fact,  $\boldsymbol{\eta} \mapsto \mathbb{E}[\boldsymbol{\eta}^2]$  is also a regular measure of error and we could just as well have adopted it from the start.

There is a long tradition for considering measures of error beyond least-squares, especially in robust statistics; see [138]. The Koenker-Bassett error in Example 5.2(a) leads to *quantile regression*. Figure 2 illustrates statistical models for predicting lift force using least-squares regression (solid line) and quantile regression with  $\alpha = 0.5, 0.75$  (dashed lines) and  $\alpha = 0.95, 0.995$  (dotted lines). While the lines are quite similar, quantile regression with higher values of  $\alpha$  results in more conservative estimates in the sense that the values of  $\gamma + c\boldsymbol{\xi}$  tend to exceed those of  $\boldsymbol{\eta}$ .

As eluded to in Subsection 2.2, one can go much beyond linear statistical models in (5.1). We refer to [281] for several refinements and to [125] for a discussion of the connection between regular measures of error and inter-regenerative relationships via log-likelihood and entropy maximization. Innovative thinking about alternative regression approaches based on constrained residuals, leveraging superquantiles, appears in [324].

## 5.2 Measures of Deviation

We next quantify “nonconstancy” of a random variable. This can be accomplished using the standard deviation, but there are many other possibilities as well; see the axiomatic development in [278, 279]. Here, we adopt the definition in [276].

**5.3 Definition** (regular measure of deviation). *A regular measure of deviation  $\mathcal{D}$  is a functional from  $\mathcal{L}^2$  to  $[0, \infty]$  that is lsc and convex, with*

$$\mathcal{D}(\boldsymbol{\xi}) = 0 \text{ if } \boldsymbol{\xi} \text{ is constant; } \quad \mathcal{D}(\boldsymbol{\xi}) > 0 \text{ otherwise.}$$

*The quantity  $\mathcal{D}(\boldsymbol{\xi})$  is the deviation of  $\boldsymbol{\xi}$ .*

**5.4 Example** (regular measures of deviation). *We have the following examples of regular measures of deviation  $\mathcal{D}$ ; see, e.g., [288, Example 8.18] for supporting arguments.*

(a) Superquantile deviation with  $\alpha \in (0, 1)$ :

$$\mathcal{D}(\boldsymbol{\xi}) = \text{s-rsk}_\alpha(\boldsymbol{\xi}) - \mathbb{E}[\boldsymbol{\xi}].$$

(b) Worst-case deviation:

$$\mathcal{D}(\boldsymbol{\xi}) = \sup \boldsymbol{\xi} - \mathbb{E}[\boldsymbol{\xi}].$$

(c) Scaled standard deviation with  $\lambda > 0$ :

$$\mathcal{D}(\boldsymbol{\xi}) = \lambda \text{std}(\boldsymbol{\xi}).$$

(d) Mixed deviation with  $\lambda_1, \dots, \lambda_q > 0$  and  $\sum_{i=1}^q \lambda_i = 1$ :

$$\mathcal{D}(\boldsymbol{\xi}) = \sum_{i=1}^q \lambda_i \mathcal{D}_i(\boldsymbol{\xi}), \quad \text{for regular measures of deviation } \mathcal{D}_1, \dots, \mathcal{D}_q.$$

The regular measures of deviation in Example 5.4(a,b,c) share the common property that they are obtained from the regular measures of risk in Example 4.3(a,b,c), respectively, by subtracting  $\mathbb{E}[\boldsymbol{\xi}]$ . Similarly, the regular measures of error in Example 5.2(a,b,c) can be constructed from the regular measures of regret in Example 4.3(a,b,c), respectively, by subtracting  $\mathbb{E}[\boldsymbol{\xi}]$ . These connections hold in general, as the next theorem (adapted from [271]) asserts.

**5.5 Theorem** (expectation translations). *Every regular measure of deviation  $\mathcal{D}$  defines a regular measure of risk  $\mathcal{R}$  and vice versa through the relations:*

$$\mathcal{R}(\boldsymbol{\xi}) = \mathcal{D}(\boldsymbol{\xi}) + \mathbb{E}[\boldsymbol{\xi}] \quad \text{and} \quad \mathcal{D}(\boldsymbol{\xi}) = \mathcal{R}(\boldsymbol{\xi}) - \mathbb{E}[\boldsymbol{\xi}].$$

*Similarly, every regular measure of error  $\mathcal{E}$  defines a regular measure of regret  $\mathcal{V}$  and vice versa through the relations:*

$$\mathcal{V}(\boldsymbol{\xi}) = \mathcal{E}(\boldsymbol{\xi}) + \mathbb{E}[\boldsymbol{\xi}] \quad \text{and} \quad \mathcal{E}(\boldsymbol{\xi}) = \mathcal{V}(\boldsymbol{\xi}) - \mathbb{E}[\boldsymbol{\xi}].$$

We complete the picture by connecting regular measures of error and deviation in a manner that resembles the relation between regular measures of regret and risk; see Theorem 4.4.

**5.6 Theorem** (error-deviation). *For a regular measure of error  $\mathcal{E}$ , a regular measure of deviation  $\mathcal{D}$  is obtained by*

$$\mathcal{D}(\boldsymbol{\xi}) = \min_{\gamma \in \mathbb{R}} \mathcal{E}(\boldsymbol{\xi} - \gamma).$$

*For every regular measure of deviation  $\mathcal{D}$  there is a regular measure of error  $\mathcal{E}$ , not necessarily unique, that constructs  $\mathcal{D}$  through this minimization formula.*

*Moreover, if  $\mathcal{E}$  is paired with a regular measure of regret  $\mathcal{V}$  via Theorem 5.5, then*

$$\mathcal{S}(\boldsymbol{\xi}) = \operatorname{argmin}_{\gamma \in \mathbb{R}} \mathcal{E}(\boldsymbol{\xi} - \gamma) = \operatorname{argmin}_{\gamma \in \mathbb{R}} \{\gamma + \mathcal{V}(\boldsymbol{\xi} - \gamma)\} \quad (5.2)$$

*and this set is a nonempty compact interval as long as  $\mathcal{V}(\boldsymbol{\xi} - \gamma)$ , or equivalently  $\mathcal{E}(\boldsymbol{\xi} - \gamma)$ , is finite for some  $\gamma \in \mathbb{R}$ .*

**Proof.** This fact stems from [281], but here stated as it appears in [288, Theorem 8.21] which incorporates the refinements of [271].  $\square$

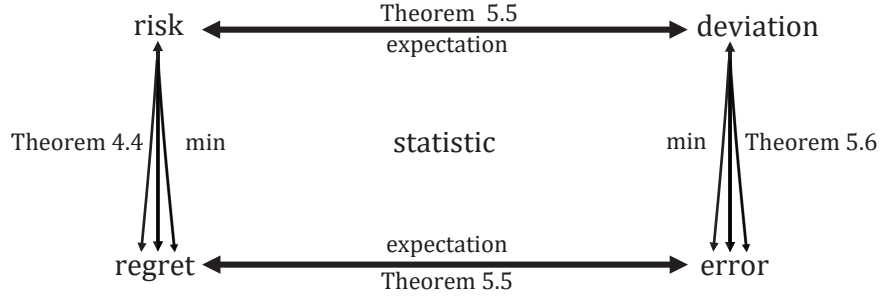


Figure 6: Relations between regular measures of risk, regret, deviation, and error that form a risk quadrangle.

The key insight of [276] (with technical refinements in [271]) as summarized in Theorems 4.4, 5.5, and 5.6 is that regular measures of risk, regret, deviation, and error can be connected in a *risk quadrangle*. Subsection 4.2 shows how passing to a regular measure of regret may simplify risk minimization. Similarly, a regular measure of deviation can be evaluated by first minimizing regret as in (5.2). One can construct new regular measures of risk and deviation by starting from either a regular measure of regret or a regular measure of error using the relations summarized in Figure 6. The “horizontal” connections in the figure are one-to-one as seen in Theorem 5.5, while the “vertical” connections are not unique in the sense that multiple regular measures of regret produce the same regular measure of risk, with a similar situation taking place when passing from error to deviation.

The regular measures of risk, regret, deviation, and error labeled (a) in Examples 4.3, 5.2, and 5.4 are connected in a risk quadrangle in the sense of Figure 6. Likewise, the regular measures of risk, regret, deviation, and error labeled (b,c,d) in Examples 4.3, 5.2, and 5.4 are connected, respectively, and then also form risk quadrangles. For many risk quadrangles, we refer to Section 7 and [276, 160]. Further connections with maximum entropy is developed in [123].

The relations between regular measures of error and deviation lead to a decomposition of the generalized regression problem (5.1) first identified in [281] for positively homogeneous functionals. Here, we state the fact as it appears in [271].

**5.7 Proposition** (decomposition). *For a regular measure of error  $\mathcal{E}$  and a regular measure of deviation  $\mathcal{D}$  paired by Theorem 5.6, one has for any  $\boldsymbol{\eta}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n \in \mathcal{L}^2$ ,*

$$(\gamma^*, c^*) \in \operatorname{argmin}_{\gamma, c} \mathcal{E}(\boldsymbol{\eta} - \gamma - \langle c, \boldsymbol{\xi} \rangle) \iff c^* \in \operatorname{argmin}_c \mathcal{D}(\boldsymbol{\eta} - \langle c, \boldsymbol{\xi} \rangle) \text{ and } \gamma^* \in \mathcal{S}(\boldsymbol{\eta} - \langle c^*, \boldsymbol{\xi} \rangle),$$

where  $\boldsymbol{\xi} = (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n)$  and  $\mathcal{S}$  is the statistic corresponding to  $\mathcal{E}$ .

The proposition shows that we can determine the coefficients  $(\gamma^*, c^*)$  in the generalized regression problem (5.1) using two steps: first, determine  $c^*$  by minimizing deviation and, second, fix  $\gamma^*$  by computing a statistic. From Theorem 5.5 we see that minimizing deviation is equivalent to minimizing risk modified by an expectation. Thus,  $c^*$  can often be computed by algorithms for minimizing risk, with minor adjustments.

Separate roles for the “slope”  $c$  and the “intercept”  $\gamma$  emerge from the proposition. The former is selected to minimize the “nonconstancy” of  $\boldsymbol{\eta} - \langle c, \boldsymbol{\xi} \rangle$ , while the latter translates the functional  $\boldsymbol{\xi} \mapsto \langle c^*, \boldsymbol{\xi} \rangle$  up or down to match the correct statistic.

The relative quality of a statistical model obtained by solving the generalized regression problem (5.1) can be assessed using the following concept proposed in [271]; see also precursors in [273]. For a regular measure of error  $\mathcal{E}$  paired with a regular measure of deviation  $\mathcal{D}$  via Theorem 5.6, the *coefficient of determination*, or *R-squared*, of the random vector  $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)$  relative to  $\boldsymbol{\eta}$  is given as<sup>12</sup>

$$R^2 = 1 - \frac{\inf_{(\gamma, c) \in \mathbb{R}^{1+n}} \mathcal{E}(\boldsymbol{\eta} - \gamma - \langle c, \boldsymbol{\xi} \rangle)}{\mathcal{D}(\boldsymbol{\eta})}.$$

The definition extends a classical concept from least-squares regression, where  $\mathcal{D}$  is the variance and  $\mathcal{E}$  is the mean-squared error. It is immediate from the definition that if  $(\gamma^*, c^*)$  is a minimizer of the generalized regression problem and  $\mathcal{E}(\boldsymbol{\eta} - \gamma^* - \langle c^*, \boldsymbol{\xi} \rangle) = 0$ , then  $R^2 = 1$ . However, such vanishing error only takes place when the residual  $\boldsymbol{\eta} - \gamma^* - \langle c^*, \boldsymbol{\xi} \rangle = \mathbf{0}$ ; the model  $\gamma^* + \langle c^*, \boldsymbol{\xi} \rangle$  predicts  $\boldsymbol{\eta}$  perfectly. Under less ideal circumstances, one would assess the quality of  $(\gamma^*, c^*)$  by seeing how close  $R^2$  is to 1. In general,  $R^2 \geq 0$  and  $R^2 = 0$  when  $(\gamma, 0)$  is a minimizer in (5.1). In that case,  $\boldsymbol{\xi}$  provides “no information” about  $\boldsymbol{\eta}$ .

In Example 2.5, we obtain five statistical models of the form  $\gamma + c\boldsymbol{\xi}$ ; see the lines in Figure 2. The quality of the models can be assessed using  $R^2$ . Quantile regression with  $\alpha = 0.75$  (highest dashed line) produces a deviation of 7.97 and an error of 1.83 and thus  $R^2 = 0.77$ . For quantile regression with  $\alpha = 0.95$  (lowest dotted line), we obtain a deviation of 12.32 and an error of 3.29 and thus  $R^2 = 0.73$ .

The paper [281] brought to the forefront the possibilities of constructing conservative statistical models through the use of “nonstandard” measures of error and other adjustments. Later developments include theoretical refinements in [271], applications to reliability engineering in [140], and applications to naval architecture in [285, 35]. Support vector regression is the topic of [209].

### 5.3 Superquantile Regression

The formula (3.5) for superquantiles as minimizers of some functionals gives rise to *superquantile regression* [273] that supplements least-squares and quantile regression. It stems from the insight that the functionals  $\mathcal{R}_\alpha, \mathcal{V}_\alpha, \mathcal{D}_\alpha, \mathcal{E}_\alpha$  on  $\mathcal{L}^2$  given by

$$\begin{aligned} \mathcal{R}_\alpha(\boldsymbol{\xi}) &= \frac{1}{1-\alpha} \int_\alpha^1 \text{s-rsk}_\beta(\boldsymbol{\xi}) d\beta & \mathcal{D}_\alpha(\boldsymbol{\xi}) &= \mathcal{R}_\alpha(\boldsymbol{\xi}) - \mathbb{E}[\boldsymbol{\xi}] \\ \mathcal{V}_\alpha(\boldsymbol{\xi}) &= \frac{1}{1-\alpha} \int_0^1 \max\{0, \text{s-rsk}_\beta(\boldsymbol{\xi})\} d\beta & \mathcal{E}_\alpha(\boldsymbol{\xi}) &= \mathcal{V}_\alpha(\boldsymbol{\xi}) - \mathbb{E}[\boldsymbol{\xi}] \end{aligned}$$

for  $\alpha \in [0, 1)$  are regular measures of risk, regret, deviation, and error, respectively, and form a risk quadrangle [269], i.e., they are connected in the sense of Theorems 4.4, 5.5, and 5.6. The corresponding statistic is  $\mathcal{S}_\alpha(\boldsymbol{\xi}) = \{\text{s-rsk}_\alpha(\boldsymbol{\xi})\}$  so  $\gamma \mapsto \mathcal{E}_\alpha(\boldsymbol{\xi} - \gamma)$  has the  $\alpha$ -superquantile of  $\boldsymbol{\xi}$  as its unique minimizer.

<sup>12</sup>Here,  $\infty/\infty$  is interpreted as 1 and  $0/0$  as 0.

We observed that *no* error measure of the expectation kind achieves an  $\alpha$ -superquantile as its minimizer when  $\alpha \in (0, 1)$ ; see [115, Theorem 11]. Thus, one needs to consider more complicated measures of error with  $\mathcal{E}_\alpha$  being one possibility, others are derived in [160]. For an introduction to the area of *elicitation*, we refer to [176], with a machine learning perspective, and [315, 28, 352]; see also the recent paper [104].

Superquantile regression predicts a random variable  $\boldsymbol{\eta}$  from the random vector  $\boldsymbol{\xi} = (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n)$  using the statistical model  $\gamma^* + \langle c^*, \boldsymbol{\xi} \rangle$ , where  $(\gamma^*, c^*)$  is a minimizer of (5.1) with  $\mathcal{E} = \mathcal{E}_\alpha$  for some  $\alpha \in [0, 1)$ . In view of Proposition 5.7, it might be computationally more efficient to obtain  $(\gamma^*, c^*)$  in two steps: First, solve

$$c^* \in \operatorname{argmin}_{c \in \mathbb{R}^n} \left\{ \frac{1}{1-\alpha} \int_\alpha^1 \operatorname{s-rsk}_\beta(\boldsymbol{\eta} - \langle c, \boldsymbol{\xi} \rangle) d\beta - \mathbb{E}[\boldsymbol{\eta} - \langle c, \boldsymbol{\xi} \rangle] \right\}. \quad (5.3)$$

Second, set  $\gamma^* = \operatorname{s-rsk}_\alpha(\boldsymbol{\eta} - \langle c^*, \boldsymbol{\xi} \rangle)$ .

If the probability distribution of  $(\boldsymbol{\eta}, \boldsymbol{\xi})$  is finite with equal probability at each atom, then (5.3) is equivalent to a linear optimization problem. Let  $\{\eta_j \in \mathbb{R}, \xi^j \in \mathbb{R}^n, j = 1, \dots, s\}$  be the possible outcomes of  $(\boldsymbol{\eta}, \boldsymbol{\xi})$ . As shown in [273], a minimizer of (5.3) is equivalently obtained by solving

$$\begin{aligned} & \underset{c, u, v, \tau}{\operatorname{minimize}} \quad \frac{1}{1-\alpha} \sum_{i=s_\alpha}^{s-1} (\beta_i - \beta_{i-1}) u_i + \frac{1}{s(1-\alpha)} \sum_{i=s_\alpha}^{s-1} \sum_{j=1}^s a_i v_{ij} + \frac{\tau}{s(1-\alpha)} - \frac{1}{s} \sum_{j=1}^s (\eta_j - \langle c, \xi^j \rangle) \\ & \text{subject to} \quad \eta_j - \langle c, \xi^j \rangle - u_i \leq v_{ij}, \quad i = s_\alpha, \dots, s-1, \quad j = 1, \dots, s \\ & \quad \quad \quad 0 \leq v_{ij}, \quad i = s_\alpha, \dots, s-1, \quad j = 1, \dots, s \\ & \quad \quad \quad \eta_j - \langle c, \xi^j \rangle \leq \tau, \quad j = 1, \dots, s \\ & c \in \mathbb{R}^n, \quad u = (u_{s_\alpha}, \dots, u_{s-1}) \in \mathbb{R}^{s-s_\alpha}, \quad v = (v_{s_\alpha, 1}, \dots, v_{s-1, s}) \in \mathbb{R}^{(s-s_\alpha)s}, \quad \tau \in \mathbb{R}, \end{aligned}$$

where  $s_\alpha = \lceil s\alpha \rceil$  is the smallest integer no smaller than  $s\alpha$ ,  $\beta_{s_\alpha-1} = \alpha$ ,  $\beta_i = i/s$ , for  $i = s_\alpha, s_\alpha + 1, \dots, s$ , and  $a_i = \ln(1 - \beta_{i-1}) - \ln(1 - \beta_i)$  for  $i = s_\alpha, s_\alpha + 1, \dots, s-1$ . For a related formulation that makes explicit connection with mixed risk, we refer to [116]. Alternatively, we can solve (5.3) by approximating the integral using standard numerical integration techniques [273, 160]. As laid out in [288, Example 8.41], one can also solve (5.3) using the subgradient method.

Superquantile regression as define above is *different* from minimizing  $\operatorname{s-rsk}_\alpha(\boldsymbol{\eta} - \gamma - \langle c, \boldsymbol{\xi} \rangle)$ . However, this distinction is not always made in the literature; see for example [175, Example 2.2]. It is also different than attempts to estimate conditional superquantiles, i.e., superquantiles of random variables that depend on a vector of predictors; see [297, 41, 148] for kernel-based approaches and [238, 182, 55] for methods passing through (3.3). Still, we can view [55] as approaching superquantile regression via an approximation based on the mixed error measure in Example 5.2(d) with  $\mathcal{E}_i, i = 1, \dots, q$ , being Koenker-Bassett error measures at different probability levels  $\alpha_i$ . That reference also estimates conditional superquantiles using an appropriately shifted least-squares regression curve based on superquantiles of the resulting residuals.

## 6 Duality Theory

A fundamental fact from convex analysis is that a convex function, with minor exceptions, can be expressed as the supremum over a collection of affine functions. Since regular measures of risk, regret, error, and deviation are convex by definition and also satisfy the required technical conditions, every such functional has an alternative representation. This offers computational possibilities and help with interpretation. Most significantly, it makes connections with distributionally robust optimization as previewed in the discussion of Mr. Averse and Ms. Ambiguous (cf. Subsection 3.1 and (3.6)). We now examine such connections broadly, and go much beyond superquantiles.

### 6.1 Conjugacy

On the space  $\mathcal{L}^2$  of random variables, we adopt the inner product  $(\xi, \pi) \mapsto \mathbb{E}[\xi\pi]$  and this leads to a definition of conjugates.

**6.1 Definition** (conjugate). *For  $\mathcal{F} : \mathcal{L}^2 \rightarrow [-\infty, \infty]$ , the functional  $\mathcal{F}^* : \mathcal{L}^2 \rightarrow [-\infty, \infty]$  defined by*

$$\mathcal{F}^*(\pi) = \sup \{ \mathbb{E}[\xi\pi] - \mathcal{F}(\xi) \mid \xi \in \mathcal{L}^2 \}$$

*is the conjugate of  $\mathcal{F}$ .*

A main motivation for restricting the attention to square-integrable random variables is that a conjugate is then defined on  $\mathcal{L}^2$  as well. In contrast, if  $\xi$  were only integrable, then  $\pi$  would have to be bounded for  $\mathbb{E}[\xi\pi]$  to be finite; see [294] for such extensions. A central result from convex analysis is the Fenchel-Moreau theorem; see, e.g., [265, Theorem 5].

**6.2 Theorem** (Fenchel-Moreau). *For a lsc convex functional  $\mathcal{F} : \mathcal{L}^2 \rightarrow (-\infty, \infty]$  with<sup>13</sup>  $\mathcal{F} \not\equiv \infty$ , the conjugate  $\mathcal{F}^*$  is also lsc and convex with  $\mathcal{F}^* \not\equiv \infty$ ,  $\mathcal{F}^*(\xi) > -\infty$  for all  $\xi \in \mathcal{L}^2$ , and  $(\mathcal{F}^*)^* = \mathcal{F}$ . Thus,*

$$\mathcal{F}(\xi) = \sup \{ \mathbb{E}[\xi\pi] - \mathcal{F}^*(\pi) \mid \pi \in \mathcal{L}^2 \}.$$

Since regular measures of risk, regret, deviation, and error are lsc and convex by definition and also finite at some point, the Fenchel-Moreau theorem 6.2 applies and furnishes alternative expressions. This observation is most useful if we can characterize the corresponding conjugates. While this is possible in many situations, we concentrate on the attractive case of positively homogeneous functionals; see [265, 27] for comprehensive treatments. A summary of properties appears in [271, Section 2].

We recall that the *domain* of  $\mathcal{F} : \mathcal{L}^2 \rightarrow [-\infty, \infty]$  is the set  $\text{dom } \mathcal{F} = \{ \xi \in \mathcal{L}^2 \mid \mathcal{F}(\xi) < \infty \}$ . With this notation, we have the following direct consequence of the Fenchel-Moreau theorem 6.2.

**6.3 Proposition** (conjugacy under positive homogeneity). *For a positively homogeneous, lsc, and convex functional  $\mathcal{F} : \mathcal{L}^2 \rightarrow (-\infty, \infty]$  with  $\mathcal{F} \not\equiv \infty$ , one has*

$$\mathcal{F}(\xi) = \sup \{ \mathbb{E}[\xi\pi] \mid \pi \in \text{dom } \mathcal{F}^* \}.$$

---

<sup>13</sup>The notation  $\mathcal{F} \not\equiv \infty$  simply indicates that the functional with  $\mathcal{F}(\xi) = \infty$  for all  $\xi \in \mathcal{L}^2$  is ruled out.

A nonempty, closed, and convex set  $\mathcal{C} \subset \mathcal{L}^2$  is the domain of  $\mathcal{F}^*$  for some positively homogeneous, lsc, and convex  $\mathcal{F} : \mathcal{L}^2 \rightarrow (-\infty, \infty]$  with  $\mathcal{F} \not\equiv \infty$  and then

$$\mathcal{C} = \{\boldsymbol{\pi} \in \mathcal{L}^2 \mid \mathbb{E}[\boldsymbol{\xi}\boldsymbol{\pi}] \leq \mathcal{F}(\boldsymbol{\xi}) \quad \forall \boldsymbol{\xi} \in \mathcal{L}^2\}.$$

The proposition asserts that positively homogeneous, lsc, and convex functionals are fully characterized by a subset of  $\mathcal{L}^2$ .

The restriction to positively homogeneous functionals still captures many important situations. For example, if  $\mathcal{V}$  is a positively homogeneous regular measure of regret, then the regular measure of risk  $\mathcal{R}$  constructed by Theorem 4.4 is also positively homogeneous. If  $\mathcal{E}$  is a positively homogeneous regular measure of error, then the regular measure of deviation  $\mathcal{D}$  constructed by Theorem 5.6 is also positively homogeneous. These facts follow immediately from the theorems. Since  $\boldsymbol{\xi} \mapsto \mathbb{E}[\max\{0, \boldsymbol{\xi}\}]$  is positively homogeneous,  $\mathcal{V}$  in Example 4.3(a) is positively homogeneous and then also s-rsk $_{\alpha}$ . The measures of regret and risk in Example 4.3(b,c) are also positively homogeneous. For Example 4.3(d),  $\mathcal{V}$  is positively homogeneous when  $\mathcal{V}_1, \dots, \mathcal{V}_q$  are positively homogeneous and then  $\mathcal{R}$  also has this property.

## 6.2 Risk Envelopes and Dual Algorithms

With the development of risk measures and related concepts, it was quickly realized that dual expressions were available via the Fenchel-Moreau theorem 6.2; see [67, 68, 98, 278]. In particular, this insight led to the dual formula (3.6) for superquantiles in [67]. While these papers concentrate on bounded random variables or random variables in  $\mathcal{L}^2$ , subsequent efforts [242, 294, 295] address more general spaces of random variables and also furnish many examples. We follow the development in [276] and concentrate on random variables in  $\mathcal{L}^2$  for simplicity.

As a direct application of Proposition 6.3, we obtain that for a positively homogeneous regular measure of risk  $\mathcal{R}$ , one has

$$\mathcal{R}(\boldsymbol{\xi}) = \sup \{\mathbb{E}[\boldsymbol{\xi}\boldsymbol{\pi}] \mid \boldsymbol{\pi} \in \Pi\}, \tag{6.1}$$

where, following the terminology of [278],  $\Pi = \text{dom } \mathcal{R}^*$  is the *risk envelope* of  $\mathcal{R}$ . Thus, a measure of risk of this kind is fully characterized by the domain of its conjugate. Convex analysis furnishes details about such domains. The following fact is taken from [288, Proposition 8.30].

**6.4 Proposition** (properties of risk envelopes). *The risk envelope  $\Pi$  of a positively homogeneous regular measure of risk  $\mathcal{R}$  is closed and convex and, for every  $\boldsymbol{\pi} \in \Pi$ , one has  $\mathbb{E}[\boldsymbol{\pi}] = 1$ . Moreover, it has the expression*

$$\Pi = \{\boldsymbol{\pi} \in \mathcal{L}^2 \mid \mathbb{E}[\boldsymbol{\xi}\boldsymbol{\pi}] \leq \mathcal{R}(\boldsymbol{\xi}) \quad \forall \boldsymbol{\xi} \in \mathcal{L}^2\}.$$

*A nonempty, closed, and convex set  $\mathcal{C} \subset \mathcal{L}^2$  is the risk envelope of some positively homogeneous regular measure of risk on  $\mathcal{L}^2$  provided that  $\mathcal{C}$  also satisfies:*

$$\mathbb{E}[\boldsymbol{\pi}] = 1 \quad \forall \boldsymbol{\pi} \in \mathcal{C};$$

*for each nonconstant  $\boldsymbol{\xi} \in \mathcal{L}^2$ , there exists  $\boldsymbol{\pi} \in \mathcal{C}$  such that  $\mathbb{E}[\boldsymbol{\xi}\boldsymbol{\pi}] > \mathbb{E}[\boldsymbol{\xi}]$ .*



**6.5 Example** (risk envelopes). *The risk envelopes corresponding to the regular measures of risk in Example 4.3 are as follows; see, e.g., [288, Example 8.31] for a justification:*

(a) Superquantile risk envelope:

$$\Pi = \{ \boldsymbol{\pi} \mid 0 \leq \boldsymbol{\pi}(\omega) \leq 1/(1 - \alpha) \quad \forall \omega, \quad \mathbb{E}[\boldsymbol{\pi}] = \mathbf{1} \}. \quad (6.2)$$

(b) Worst-case risk envelope:

$$\Pi = \{ \boldsymbol{\pi} \mid \boldsymbol{\pi}(\omega) \geq 0 \quad \forall \omega, \quad \mathbb{E}[\boldsymbol{\pi}] = \mathbf{1} \}.$$

(c) Mean-plus-standard-deviation risk envelope<sup>14</sup>:

$$\Pi = \{ \mathbf{1} + \lambda \boldsymbol{\pi} \mid \mathbb{E}[\boldsymbol{\pi}^2] \leq \mathbf{1}, \quad \mathbb{E}[\boldsymbol{\pi}] = \mathbf{0} \}.$$

(d) Mixed risk envelope:

$$\Pi = \left\{ \sum_{i=1}^q \lambda_i \boldsymbol{\pi}_i \mid \boldsymbol{\pi}_i \in \Pi_i \right\},$$

where  $\Pi_i$  is the risk envelope of the positively homogeneous regular measure of risk  $\mathcal{R}_i$ .

The dual representation of risk measures in (6.1) offers several computational possibilities. In particular, it brings forward an expectation, which can subsequently be approximated using Monte Carlo sampling or other techniques.

**Dual Algorithms for Risk Minimization.** The discussion in Subsection 3.3 about dual algorithms for superquantile minimization hints to possibilities that derive from (6.1). Concretely, let us consider a quantity of interest represented by  $f : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$ , a constraint set  $X \subset \mathbb{R}^n$ , a positively homogeneous regular measure of risk  $\mathcal{R}$ , and a finitely distributed random vector  $\boldsymbol{\xi}$  with support  $\Xi \subset \mathbb{R}^m$  of cardinality  $s$  and probabilities  $\{p_\xi > 0, \xi \in \Xi\}$ . This leads to the problem

$$\underset{x \in X}{\text{minimize}} \quad \mathcal{R}(f(\boldsymbol{\xi}, x)) = \sup \{ \mathbb{E}[f(\boldsymbol{\xi}, x)\boldsymbol{\pi}] \mid \boldsymbol{\pi} \in \Pi \} = \sup \left\{ \sum_{\xi \in \Xi} p_\xi f(\xi, x) q_\xi \mid q \in \mathcal{Q} \right\}, \quad (6.3)$$

where  $\Pi$  is the risk envelope of  $\mathcal{R}$ . The second equality recognizes that the underlying probability space  $(\Omega, \mathcal{A}, \mathbb{P})$  defining  $\mathcal{L}^2$  can in this case have  $\Omega = \Xi$  and  $\mathbb{P}(\{\xi\}) = p_\xi$ . Thus, all random variables  $\boldsymbol{\pi} \in \mathcal{L}^2$  are represented by  $s$ -dimensional vectors of the form  $q = (q_\xi, \xi \in \Xi) \in \mathbb{R}^s$  and the maximization over  $\Pi$  in (6.3) is equivalent to a maximization over a subset  $\mathcal{Q}$  of  $\mathbb{R}^s$ . Specifically,  $\mathcal{Q} = \{q \in \mathbb{R}^s \mid \exists \boldsymbol{\pi} \in \Pi \text{ such that } q_\xi = \boldsymbol{\pi}(\xi), \xi \in \Xi\}$ . In the case of superquantiles, (6.2) gives  $\Pi$  and this yields

$$\mathcal{Q} = \left\{ q \in \mathbb{R}^s \mid 0 \leq q_\xi \leq \frac{1}{1 - \alpha}, \xi \in \Xi, \quad \sum_{\xi \in \Xi} p_\xi q_\xi = 1 \right\}. \quad (6.4)$$

By setting  $\bar{p}_\xi = p_\xi q_\xi$ , we reconcile the present development with the formula (3.6).

---

<sup>14</sup>The constant random variable with value 1 is denoted by  $\mathbf{1}$ .

Regardless of the circumstances, the risk-minimization problem (6.3) is equivalently stated as a minsup problem for which there are many algorithmic approaches. These include subgradient-type methods which apply, especially, if  $f(\xi, \cdot)$  is convex for all  $\xi \in \Xi$  (see, e.g., [288, Section 2.I]), the outer approximation algorithm (see, e.g., [288, Section 6.C] and the discussion in Subsection 3.3), and “primal-dual” algorithms (see, e.g., [167]). While these approaches provide important stepping stones, there are bound to be implementation challenges depending on the specific risk measure under consideration. In the absence of a finite distribution for  $\xi$ , the first equality in (6.3) still holds but it becomes harder to compute a maximizer, which may not even exist, causing additional challenges; see our discussion in Subsection 6.4.

For recent derivations of risk envelope formulas, we refer to [15, 316], which focus on regular measures of deviation and the use of set operations, and [272], which concentrates on mixed risk measures and second-order superquantiles; see also [302, 99]. Dual expressions for measures of risk and related quantities remain a crucial stepping stone in several contexts. The papers [124, 88] utilize them in sensitivity analysis and [165] in derivation of optimality condition; see also [280]. Dual expressions are key in developing algorithms, including in the computationally challenging areas of PDE-constrained optimization [164, 166, 167] and statistical learning [61, 185, 173, 175].

### 6.3 Connections with Distributionally Robust Optimization

It has long been recognized that a probability distribution adopted in an application is just a *model* of the uncertainty associated with unsettled parameters; see [82] for the first recorded work in the area of stochastic programming but the perspective can be traced back to the early studies of games as summarized in [330]. This leads to *distributionally robust optimization problems* of the form

$$\underset{x \in X}{\text{minimize}} \sup_{P \in \mathcal{P}} \int f(\xi, x) dP(\xi), \quad (6.5)$$

where  $f : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$  represents a quantity of interest,  $X$  is a constraint set, and  $\mathcal{P}$  is an *ambiguity set* of candidate distributions  $P$  on  $\mathbb{R}^m$ . (Here,  $P$  is a probability distribution on  $\mathbb{R}^m$ , which clashes slightly with the earlier use of  $P$  for a cumulative distribution function.) Distributionally robust optimization problems may be perceived to “optimally” address distributional ambiguity [327]. Key developments include [337, 304, 351]; see also the recent reviews [257, 260, 191]. A main focus in this area is the development of tractable reformulations of the minsup problem (6.5) as a minimization problem. We see next that such tractable reformulations are immediately available for certain  $\mathcal{P}$ .

The following fact is well known from convex analysis; see, e.g., [276].

**6.6 Proposition** (conjugacy under monotonicity). *For a lsc convex functional  $\mathcal{F} : \mathcal{L}^2 \rightarrow (-\infty, \infty]$  with  $\mathcal{F} \neq \infty$ , one has that  $\mathcal{F}$  is monotone if and only if  $\pi \geq 0$  for all  $\pi \in \text{dom } \mathcal{F}^*$ .*

The proposition together with Example 6.5 confirm that superquantile risk and worst-case risk measures are monotone, while mean-plus-standard-deviation risk measures are not.

Combining Propositions 6.4 and 6.6, we conclude that a risk envelope  $\Pi$  of a positively homogeneous, monotone, regular measure of risk contains exclusively nonnegative random variables with expectation one. This has the consequence that we can interpret the expression  $\mathbb{E}[\boldsymbol{\xi}\boldsymbol{\pi}] = \int \boldsymbol{\xi}(\omega)\boldsymbol{\pi}(\omega) d\mathbb{P}(\omega)$  as the expectation of  $\boldsymbol{\xi}$  under a different probability distribution  $\bar{\mathbb{P}}$ . From this perspective,  $\boldsymbol{\pi}$  is the density (Radon-Nikodym derivative)  $d\bar{\mathbb{P}}/d\mathbb{P}$ . Consequently, a monotone, positively homogeneous, regular measure of risk can be expressed by (6.1) using a risk envelope  $\Pi$  or, equivalently, as the supremum of  $\int \boldsymbol{\xi}(\omega)d\bar{\mathbb{P}}(\omega)$  over an ambiguity set with probability distributions  $\bar{\mathbb{P}}$ , which brings us back to (6.5). We refer to [294] for a formal treatment and to [276, Section 6] for a brief introduction.

This insight implies that *every* monotone, positively homogeneous, regular measure of risk inherently addresses ambiguity leading to some distributionally robust problem akin to (6.5). Conversely, a distributionally robust problem (6.5) corresponds to a risk-minimization problem provided that the densities defined by  $\mathcal{P}$  form a set  $\Pi$  satisfying the properties in the second half of Proposition 6.4, which implicitly requires that the distributions in  $\mathcal{P}$  are absolutely continuous with respect to the underlying “base” distribution. The risk-minimization problem is, in turn, equivalently to a regret minimization problem; cf. Subsection 4.2. The risk-minimization and the regret-minimization problems might be tractable alternatives to the original distributionally robust problem (6.5).

In summary, the introductory discussion in Subsection 3.1 about Mr. Averse, a risk-averse decision maker relying on superquantiles, and Ms. Ambiguity, which is risk-neutral but uncertain about the underlying probability distribution, extends to all monotone, positively homogeneous, regular measures of risk. These risk measures can therefore be used to model *both* risk-averseness and distributional ambiguity. Regardless of the initial motivation, the duality between risk-averseness and distributional ambiguity allows us to switch between the two perspectives, adopting the one that is computationally attractive or affords other advantages.

Further connections appear in [117], which shows that an ambiguity set centered at an empirical distribution makes the distributionally robust problem (6.5) in some sense close to that of minimizing a *mean-plus-standard-deviation* risk measure; see also [118] and references therein. Recent reviews of distributionally robust models in engineering and portfolio optimization include [146, 243]. Adversarial learning [204] is also of the form (6.5) with a specific ambiguity set that only shifts the support.

## 6.4 Risk Identifiers and Subgradients

While solving distributionally robust optimization problems such as (6.5) or utilizing dual formulas in risk minimization, it might be important to identify which probability distribution  $P \in \mathcal{P}$  or which random variable  $\boldsymbol{\pi} \in \Pi$ , if any, attains the maximum. These issues are also intimately tied to subgradients of functionals. Convex analysis provides key insights; see, e.g., [265].

**6.7 Definition** (subgradients of functional). *For a convex functional  $\mathcal{F} : \mathcal{L}^2 \rightarrow [-\infty, \infty]$  and a point  $\boldsymbol{\xi}_0 \in \mathcal{L}^2$  at which  $\mathcal{F}$  is finite,  $\boldsymbol{\pi} \in \mathcal{L}^2$  is a subgradient of  $\mathcal{F}$  at  $\boldsymbol{\xi}_0$  when*

$$\mathcal{F}(\boldsymbol{\xi}) \geq \mathcal{F}(\boldsymbol{\xi}_0) + \mathbb{E}[\boldsymbol{\pi}(\boldsymbol{\xi} - \boldsymbol{\xi}_0)] \quad \forall \boldsymbol{\xi}, \boldsymbol{\xi}_0 \in \mathcal{L}^2.$$

*The set of all subgradients of  $\mathcal{F}$  at  $\boldsymbol{\xi}_0$  is denoted by  $\partial\mathcal{F}(\boldsymbol{\xi}_0)$ .*

The *interior* of  $\mathcal{C} \subset \mathcal{L}^2$ , denoted by  $\text{int } \mathcal{C}$ , consists of every  $\xi \in \mathcal{C}$  for which there exists  $\rho > 0$  such that  $\{\xi_0 \mid \|\xi_0 - \xi\|_{\mathcal{L}^2} \leq \rho\} \subset \mathcal{C}$ . In particular,  $\text{int}(\text{dom } \mathcal{F}) = \mathcal{L}^2$  when  $\mathcal{F}$  is real-valued. From [294, Proposition 3.1], we know that a monotone convex functional  $\mathcal{F} : \mathcal{L}^2 \rightarrow (-\infty, \infty]$  has at least one subgradient at every point in  $\text{int}(\text{dom } \mathcal{F})$ . The following classical fact from convex analysis provides a means to calculate subgradients; see, for example, [288, Proposition 8.36] for a short proof based on the Fenchel-Moreau theorem 6.2. The claim about nonemptiness follows from [265, Corollary 8B and Theorem 11].

**6.8 Proposition** (subgradients from conjugates). *For a lsc convex functional  $\mathcal{F} : \mathcal{L}^2 \rightarrow (-\infty, \infty]$  and a point  $\xi$  at which  $\mathcal{F}$  is finite, one has*

$$\partial\mathcal{F}(\xi) = \text{argmax} \{ \mathbb{E}[\xi\pi] - \mathcal{F}^*(\pi) \mid \pi \in \mathcal{L}^2 \}.$$

*This subset of  $\mathcal{L}^2$  is nonempty provided that  $\xi \in \text{int}(\text{dom } \mathcal{F})$ .*

For a positively homogeneous regular measure of risk, we say that  $\hat{\pi} \in \text{argmax}\{\mathbb{E}[\xi\pi] \mid \pi \in \Pi\}$  is a *risk identifier* of  $\mathcal{R}$  at  $\xi$ . Thus, the maximum in (6.1) is indeed attained as long as  $\xi \in \text{int}(\text{dom } \mathcal{R})$  and  $\mathcal{R}$  is positively homogeneous and regular. The term risk identifier was coined in [280] but the quantity was known much earlier simply as a subgradient of  $\mathcal{R}$ . A more general treatment of existence of risk identifiers appears in [158], with a particular focus on distributionally robust optimization. Expressions for risk identifiers in the case of mixed measures of risk appear in [272].

We now have tools to compute subgradients of functions involving risk measures and related functionals as already recognized by [294, 280]. For concreteness, we limit the focus to compositions with linear functions; see, e.g., [272, Section 4] for more general cases.

Given a positively homogeneous regular measure of risk  $\mathcal{R}$  and  $\eta, \xi_1, \dots, \xi_n \in \mathcal{L}^2$ , we consider the problem

$$\underset{x \in X \subset \mathbb{R}^n}{\text{minimize}} \quad \varphi(x) = \mathcal{R}(\eta - \langle \xi, x \rangle), \tag{6.6}$$

where  $\xi = (\xi_1, \dots, \xi_n)$ . For example,  $\eta - \langle \xi, x \rangle$  might represent the (random) shortfall of “production”  $\langle \xi, x \rangle$  relative to “demand”  $\eta$ . We also consider the generalized regression problem (5.1), which in view of Proposition 5.7 effectively reduces to

$$\underset{c \in C \subset \mathbb{R}^n}{\text{minimize}} \quad \psi(c) = \mathcal{D}(\eta - \langle c, \xi \rangle), \tag{6.7}$$

where  $\mathcal{D}$  is the regular measure of deviation that pairs (in the sense of Theorem 5.6) with the regular measure of error of interest in (5.1). There are accessible subgradient formulas for these objective functions, which allow us to bring in subgradient methods, cutting-plane methods, and related algorithms. The following proposition is taken from [288, Proposition 8.38]; see [294, 280] for similar formulas.

**6.9 Proposition** (subgradients in risk and deviation minimization). *For a real-valued, positively homogeneous, and regular measure of risk  $\mathcal{R}$ , with risk envelope  $\Pi$ , the measure of deviation  $\mathcal{D}$  paired with*

$\mathcal{R}$  in Theorem 5.5, and random variables  $\boldsymbol{\eta}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n \in \mathcal{L}^2$ , with  $\boldsymbol{\xi} = (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n)$ , consider  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$  in (6.6) and  $\psi : \mathbb{R}^n \rightarrow \mathbb{R}$  in (6.7). Then,

$$\begin{aligned}\partial\varphi(x) &= \left\{ -\mathbb{E}[\boldsymbol{\xi}\hat{\boldsymbol{\pi}}] \mid \hat{\boldsymbol{\pi}} \in \operatorname{argmax} \left\{ \mathbb{E}[(\boldsymbol{\eta} - \langle \boldsymbol{\xi}, x \rangle)\boldsymbol{\pi}] \mid \boldsymbol{\pi} \in \Pi \right\} \right\} \quad \forall x \in \mathbb{R}^n \\ \partial\psi(c) &= \left\{ \mathbb{E}[\boldsymbol{\xi}] - \mathbb{E}[\boldsymbol{\xi}\hat{\boldsymbol{\pi}}] \mid \hat{\boldsymbol{\pi}} \in \operatorname{argmax} \left\{ \mathbb{E}[(\boldsymbol{\eta} - \langle c, \boldsymbol{\xi} \rangle)\boldsymbol{\pi}] \mid \boldsymbol{\pi} \in \Pi \right\} \right\} \quad \forall c \in \mathbb{R}^n.\end{aligned}$$

In the case of superquantiles, we obtain the following example.

**6.10 Example** (subgradients for superquantile functions). For  $\alpha \in (0, 1)$  and finitely distributed random variables  $\boldsymbol{\eta}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n$ , with values  $\{\eta^i, \xi_1^i, \dots, \xi_n^i, i = 1, \dots, s\}$  and probabilities  $\{p_i > 0, i = 1, \dots, s\}$ , consider the function  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$  given by

$$\varphi(x) = \text{s-rsk}_\alpha(\boldsymbol{\eta} - \langle \boldsymbol{\xi}, x \rangle),$$

where  $\boldsymbol{\xi} = (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n)$ . Let  $\xi^i = (\xi_1^i, \dots, \xi_n^i)$ . A subgradient of  $\varphi$  at  $x$  is then given by

$$-\sum_{i=1}^s p_i \xi^i \hat{q}_i, \quad \text{where } \hat{q}_i = \begin{cases} \frac{1}{1-\alpha} & \text{if } i \in \mathbb{I}_+(x) \\ \frac{\sigma-\alpha}{(1-\alpha)\tau} & \text{if } i \in \mathbb{I}_0(x) \\ 0 & \text{otherwise,} \end{cases}$$

with  $\sigma = 1 - \sum_{i \in \mathbb{I}_+(x)} p_i$ ,  $\tau = \sum_{i \in \mathbb{I}_0(x)} p_i$ , and, denoting the  $\alpha$ -quantile of  $\boldsymbol{\eta} - \langle \boldsymbol{\xi}, x \rangle$  by  $Q(\alpha)$ , also

$$\mathbb{I}_+(x) = \{i \mid \eta^i - \langle \xi^i, x \rangle > Q(\alpha)\} \quad \text{and} \quad \mathbb{I}_0(x) = \{i \mid \eta^i - \langle \xi^i, x \rangle = Q(\alpha)\}.$$

**Detail.** Using  $\mathcal{Q}$  in (6.4) and the arguments around that equation, we seek

$$\hat{q} \in \operatorname{argmax}_{q \in \mathcal{Q}} \sum_{i=1}^s p_i (\eta^i - \langle \xi^i, x \rangle) q_i,$$

which is essentially available explicitly. The asserted subgradient follows by Proposition 6.9.

The computational work required to obtain a subgradient is of order  $O(s \ln s)$  as it essentially requires sorting the numbers  $\{\eta^i - \langle \xi^i, x \rangle, i = 1, \dots, s\}$ . We also obtain the alternative formula:  $\text{s-rsk}_\alpha(\boldsymbol{\eta} - \langle \boldsymbol{\xi}, x \rangle) = \sum_{i=1}^s p_i (\eta^i - \langle \xi^i, x \rangle) \hat{q}_i$ .  $\square$

## 7 Additional Examples

There is a vast landscape of risk measures beyond the instances listed in Example 4.3; [276] and [307, Chapter 6] review many others. This section discusses a few additional risk measures and provides recipes for constructing even more. We also touch on other measures of regret, deviation, and error.

We refer to [249] for a discussion of how to compare risk measures quantitatively, and this might help in the process of choosing a suitable one.

**7.1 Example** (mean-plus-upper-semideviation). For  $\lambda \in (0, \infty)$ , the measures of risk, regret, deviation, and error given by

$$\begin{aligned}\mathcal{R}(\xi) &= \mathbb{E}[\xi] + \lambda \sqrt{\mathbb{E}[\max\{0, \xi - \mathbb{E}[\xi]\}^2]} & \mathcal{D}(\xi) &= \lambda \sqrt{\mathbb{E}[\max\{0, \xi - \mathbb{E}[\xi]\}^2]} \\ \mathcal{V}(\xi) &= \lambda \sqrt{\mathbb{E}[\max\{0, \xi - \mathbb{E}[\xi]\}^2]} + \max\{0, 2\mathbb{E}[\xi]\} & \mathcal{E}(\xi) &= \lambda \sqrt{\mathbb{E}[\max\{0, \xi - \mathbb{E}[\xi]\}^2]} + |\mathbb{E}[\xi]|\end{aligned}$$

are regular and form a risk quadrangle in the sense of Figure 6. Moreover, the functionals are all positively homogenous. The corresponding statistic is  $\mathcal{S}(\xi) = \mathbb{E}[\xi]$ . The risk envelope of  $\mathcal{R}$ , furnishing the alternative expression (6.1), is given by

$$\Pi = \{1 + \lambda\pi - \lambda\mathbb{E}[\pi] \mid \pi(\omega) \geq 0 \forall \omega, \mathbb{E}[\pi^2] \leq 1\}.$$

We refer to  $\mathcal{R}$  as a mean-plus-upper-semideviation risk measure.

**Detail.** Most of these properties can be found in [307, Section 6.3] or follow immediately from the various definitions. The construction of measures of regret and error is motivated by the approach taken in the proofs of Theorems 8.9 and 8.21 in [288]. Extensions beyond square-integrable random variables, to  $p$ -integrable random variables, are discussed in [307, Section 6.3]. That reference also covers when the resulting risk measures are monotone.

Markowitz already in [211] studied semideviations. Later efforts include [230, 231, 232], with a focus on connections with stochastic dominance, and the algorithmic developments in [143, 144].  $\square$

**7.2 Example** (entropic risk). The measures of risk, regret, deviation, and error given by

$$\begin{aligned}\mathcal{R}(\xi) &= \ln \mathbb{E}[\exp(\xi)] & \mathcal{D}(\xi) &= \ln \mathbb{E}[\exp(\xi - \mathbb{E}[\xi])] \\ \mathcal{V}(\xi) &= \mathbb{E}[\exp(\xi) - 1] & \mathcal{E}(\xi) &= \mathbb{E}[\exp(\xi) - \xi - 1]\end{aligned}$$

are regular and form a risk quadrangle in the sense of Figure 6. The corresponding statistic is  $\mathcal{S}(\xi) = \ln \mathbb{E}[\exp(\xi)]$ . While  $\mathcal{R}$  is monotone, it is not positively homogeneous and therefore lacks an expression of the form (6.1). Still, the Fenchel-Moreau theorem 6.2 applies and yields the alternative expression

$$\mathcal{R}(\xi) = \sup \{ \mathbb{E}[\xi\pi] - \mathbb{E}[\pi \ln \pi] \mid \pi(\omega) \geq 0 \forall \omega, \mathbb{E}[\pi] = 1 \},$$

where  $0 \ln 0$  is defined as 0. We refer to  $\mathcal{R}(\xi)$  as entropic risk.

**Detail.** These facts can be deduced from [307, Section 6.3] and [276]. We note that  $\mathcal{R}$  may not be real-valued unless, for example, the underlying probability space is finite; see [307, Section 6.3] for a detailed discussion. The measure of regret  $\mathcal{V}$  stems from a “normalized” exponential disutility function and thus has a long history within expected utility theory. We refer to [172] for extensions.  $\square$

Examples 4.3(d), 5.2(d), 5.4(d), and the supporting discussion show how we can construct new regular measures of risk, regret, error, and deviation by combining existing ones. Subsection 4.3 relies on the same principle, but with the narrower focus on using superquantiles as the “base” measures of

risk from which others emerge. Simple scaling may also be used. If  $\mathcal{R}_0, \mathcal{V}_0, \mathcal{D}_0, \mathcal{E}_0$  are regular measures of risk, regret, deviation, and error forming a risk quadrangle in the sense of Figure 6, with statistic  $\mathcal{S}_0$ , and  $\lambda \in (0, \infty)$ , then the functionals  $\mathcal{R}, \mathcal{V}, \mathcal{D}, \mathcal{E}$  given by

$$\begin{aligned}\mathcal{R}(\boldsymbol{\xi}) &= \lambda \mathcal{R}_0(\lambda^{-1} \boldsymbol{\xi}) & \mathcal{D}(\boldsymbol{\xi}) &= \lambda \mathcal{D}_0(\lambda^{-1} \boldsymbol{\xi}) \\ \mathcal{V}(\boldsymbol{\xi}) &= \lambda \mathcal{V}_0(\lambda^{-1} \boldsymbol{\xi}) & \mathcal{E}(\boldsymbol{\xi}) &= \lambda \mathcal{E}_0(\lambda^{-1} \boldsymbol{\xi})\end{aligned}$$

are regular and form a risk quadrangle, with statistic  $\mathcal{S}(\boldsymbol{\xi}) = \lambda \mathcal{S}_0(\lambda^{-1} \boldsymbol{\xi})$ . Here, monotonicity is preserved by this scaling process. Naturally, the scaling is only interesting when a functional is *not* positively homogeneous such as in Example 7.2.

Adding a weighted functional to the expected value is also meaningful; see, for example, [102] for usage in machine learning. Again starting with  $\mathcal{R}_0, \mathcal{V}_0, \mathcal{D}_0, \mathcal{E}_0$ , all regular and forming a risk quadrangle, and  $\lambda \in (0, \infty)$ , we obtain that the functionals given by

$$\begin{aligned}\mathcal{R}(\boldsymbol{\xi}) &= (1 - \lambda) \mathbb{E}[\boldsymbol{\xi}] + \lambda \mathcal{R}_0(\boldsymbol{\xi}) & \mathcal{D}(\boldsymbol{\xi}) &= \lambda \mathcal{D}_0(\boldsymbol{\xi}) \\ \mathcal{V}(\boldsymbol{\xi}) &= (1 - \lambda) \mathbb{E}[\boldsymbol{\xi}] + \lambda \mathcal{V}_0(\boldsymbol{\xi}) & \mathcal{E}(\boldsymbol{\xi}) &= \lambda \mathcal{E}_0(\boldsymbol{\xi})\end{aligned}$$

are regular and form a risk quadrangle, with statistic  $\mathcal{S}(\boldsymbol{\xi}) = \mathcal{S}_0(\boldsymbol{\xi})$ . Positive homogeneity is preserved by this process regardless of  $\lambda \in (0, \infty)$  and the same holds for monotonicity in the case of  $\mathcal{D}$  and  $\mathcal{E}$ . Monotonicity of  $\mathcal{R}$  and  $\mathcal{V}$  is preserved when  $\lambda \leq 1$ . The properties of scaling and expectation-mixing follow straightforwardly from the various definitions; see [276] for a recording of these facts.

For additional means to construct new functionals from existing ones, we refer to [102, 103] and [193], which combines risk measures using epi-sums (inf-convolutions) as motivated by risk sharing. The paper [331] studies signed Choquet integrals and their convexity and [197] utilizes distributional transforms.

**7.3 Example** (entropic value-at-risk). *For  $\alpha \in [0, 1)$ , the measure of risk given by*

$$\mathcal{R}_\alpha(\boldsymbol{\xi}) = \inf_{\gamma > 0} \gamma^{-1} \ln \left( \frac{\mathbb{E}[\exp(\gamma \boldsymbol{\xi})]}{1 - \alpha} \right)$$

*is convex, positively homogeneous, monotone, and also satisfies the constancy property. The risk envelope of  $\mathcal{R}_\alpha$  furnishing the alternative expression (6.1) is given by*

$$\Pi = \{ \boldsymbol{\pi} \mid \mathbb{E}[\boldsymbol{\pi} \ln \boldsymbol{\pi}] \leq -\ln(1 - \alpha), \mathbb{E}[\boldsymbol{\pi}] = 1, \boldsymbol{\pi}(\omega) \geq 0 \forall \omega \},$$

*where again  $0 \ln 0$  is defined as 0. Thus, this risk measure yields the worst-case expected value with respect to probability distributions that has Kullback-Leibler divergence from  $\mathbb{P}$  of at most  $-\ln(1 - \alpha)$ . In particular,  $\mathcal{R}_0(\boldsymbol{\xi}) = \mathbb{E}[\boldsymbol{\xi}]$  and  $\lim_{\alpha \nearrow 1} \mathcal{R}_\alpha(\boldsymbol{\xi}) = \sup \boldsymbol{\xi}$ . The  $\alpha$ -superquantile of  $\boldsymbol{\xi}$  never exceeds  $\mathcal{R}_\alpha(\boldsymbol{\xi})$ . We refer to  $\mathcal{R}_\alpha(\boldsymbol{\xi})$  as the entropic value-at-risk (at level  $\alpha$ ).*

**Detail.** These properties are given in [4, 6, 5], which pioneered entropic value-at-risk. With the involvement of the moment-generating function  $\gamma \mapsto \mathbb{E}[\exp(\gamma \boldsymbol{\xi})]$  in the definition of  $\mathcal{R}_\alpha$ , it is natural

that these claims are restricted to random variables for which this function is finite; see [6, 5] for details. An advantage of entropic value-at-risk is the ease by which it can be computed under various independence assumptions.

For far-reaching extensions leveraging more general divergences in the construction of risk envelopes, we refer to [250], where, for example, Kusuoka representations appear as well. Efforts in this direction already emerged in [6].  $\square$

For additional measures of risk, we refer to [43], which considers a broad class of star-shaped risk measures that extend beyond the convex ones. Connections between risk measures and scoring rules emerge from [310]. The paper [184] discusses risk measures motivated by cumulative prospect theory and [194] examines tail-focused generalizations of quantiles and superquantiles. Further constructions involving divergences can be found in [100, 78].

## 8 Computational Tools

The number and capabilities of computational tools for risk minimization are rapidly expanding. Standard packages for optimization address many problems, especially after reformulations to common formats such as linear and nonlinear programs; see for example Subsections 3.2 and 5.3. There are also specialized packages that address risk minimization directly.

*Portfolio Safeguard* [13] allows users to specify measures of risk, regret, deviation, and error using a high-level language. It also has built-in custom algorithms that address large-scale problems and those involving nonsmoothness, for example caused by the formula (3.2). It also has a tailored format for superquantile regression (referred to as CVaR regression in the documentation). However, being a commercial product, the details about some of the algorithms are unavailable. Portfolio Safeguard has Matlab and R interfaces.

The Python package *SPQR* [313] implements first-order algorithms for superquantile minimization from [173]; see also the discussion in Subsection 3.3. Another Python package *SQwash* [314] provides connections with PyTorch to solve superquantile problems involving neural networks. It also implements dual smoothing algorithms from [175, 253] as surveyed in Subsection 3.3. These two packages are especially tailored to statistical estimation problems.

There are several packages motivated by financial applications, an area where superquantiles are known by the alternative names conditional value-at-risk, tail value-at-risk, and expected shortfall. To solve risk-minimization problems, the Python package *Riskfolio-Lib* [263] leverages the optimization package *CvxPy* and closely integrates with pandas data structures. *AzaPy* [23] is another Python package that handles mixed superquantile risk measures. *Cvar-Portfolio* [62] is a portfolio optimization package in Python that addresses high-dimensional problems.

Beyond Python, we find packages in R [256], C++ [317], and Julia [87, 59, 264], with the latter reference furnishing implementation of an exceptionally wide range of risk measures. The Rapid Optimization Library [283], developed in C++, focuses on large-scale engineering applications and has capabilities to address problems with uncertainty including risk minimization.



## 9 Extensions and Challenges

Since decision making under uncertainty arise in nearly all human activity, this survey is unable to review all aspects in detail. We have ignored the fact that some optimization problems involve infinite-dimensional  $x$  such as in PDE-constrained optimization; see, e.g., [163]. While this introduces numerous challenges, it does *not* change the definition of risk measures or their desirable properties. We still apply a risk measure to  $f(\boldsymbol{\xi}, x)$ , which is a random variable for any (infinite-dimensional)  $x$  as long as  $f(\cdot, x)$  is measurable, just as before. In fact, many articles cited in the survey deal with such general cases. This section mentions two other domains: measures of reliability and multi-stage problems. We end with a discussion of research opportunities.

### 9.1 Measures of Reliability

The risk of a random variable has the same unit as the random variable itself; it is a number that represents (conservatively) the unknown future value. Definition 2.7 reflects this through its treatment of constant random variables. In reliability engineering and operations research, it is common to consider other ways of quantifying uncertainty and ensuring safety. Here, we propose terminology for these alternative ways as a counterpart to Definition 2.7.

**9.1 Definition** (reliability measure). *A measure of reliability  $\mathcal{P}$  assigns to a random variable  $\boldsymbol{\xi}$  a number  $\mathcal{P}(\boldsymbol{\xi}) \in [0, 1]$  as a quantification of its reliability, with this number being either 0 or 1 if  $\boldsymbol{\xi}$  is a constant random variable.*

A measure of reliability is a functional on a space of random variables, with (hopefully) desirable properties just as measures of risk, regret, deviation, and error. Since  $\mathcal{P}(\boldsymbol{\xi}) \in [0, 1]$ , it can be interpreted as a probability with the following prominent examples.

**Probability of exceedance.** For  $\tau \in \mathbb{R}$ , the choice  $\mathcal{P}_\tau(\boldsymbol{\xi}) = \text{prob}\{\boldsymbol{\xi} > \tau\}$  utilizes the *probability of exceedance* of the threshold  $\tau$ . This measure of reliability is in one-to-one correspondence with the quantile risk measure  $\mathcal{R}_\alpha(\boldsymbol{\xi}) = Q(\alpha)$ , where  $Q(\alpha)$  is the  $\alpha$ -quantile of  $\boldsymbol{\xi}$ ,  $\alpha \in (0, 1)$ . Specifically,  $\mathcal{R}_\alpha(\boldsymbol{\xi}) \leq \tau$  if and only if  $\mathcal{P}_\tau(\boldsymbol{\xi}) \leq 1 - \alpha$ . If the threshold  $\tau = 0$ , then  $\mathcal{P}_\tau(\boldsymbol{\xi})$  is called the *failure probability* of  $\boldsymbol{\xi}$ ; see, e.g., [270]. Probabilities of exceedance give rise to chance constraints as seen in [288, Examples 2.57 and 3.20].

**Buffered probability of exceedance.** For an integrable random variable  $\boldsymbol{\xi}$  with superquantiles  $\bar{Q}(\alpha)$ ,  $\alpha \in [0, 1]$ , the *buffered probability* of exceeding  $\tau \in \mathbb{R}$  is

$$\text{b-prob}\{\boldsymbol{\xi} > \tau\} = \begin{cases} 0 & \text{if } \tau \geq \bar{Q}(1) \\ 1 - \bar{Q}^{-1}(\tau) & \text{if } \bar{Q}(0) < \tau < \bar{Q}(1) \\ 1 & \text{otherwise.} \end{cases}$$

Here,  $\bar{Q}^{-1}(\tau)$  is the solution to the equation  $\tau = \bar{Q}(\alpha)$ , which is unique when  $\bar{Q}(0) < \tau < \bar{Q}(1)$  because  $\bar{Q}$  is a continuous function on  $(0, 1)$  and also increasing on  $(0, \text{prob}\{\boldsymbol{\xi} < \bar{Q}(1)\})$ ; see [269,

Theorem 2] and [207, Proposition A.1]. For each  $\tau$ , this defines the measure of reliability  $\mathcal{P}_\tau(\boldsymbol{\xi}) = \text{b-prob}\{\boldsymbol{\xi} > \tau\}$  first studied in [268] for the case  $\tau = 0$  under the name *buffered failure probability* and later extended to address natural hazards [66], classification problems [224, 226], and cash-flow matching [300]; see [207, 205] for theoretical foundations and [277] for a decomposition method to solve resulting optimization problems. In addition to [268], other application papers in reliability engineering include [215, 355, 39, 40]. Sensitivity analysis of buffered probabilities appears in [350, 286]. The paper [159] defines higher-order buffered probabilities. The measure of reliability  $\mathcal{P}_\tau$  is in one-to-one correspondence with the  $\alpha$ -superquantile measure of risk: For integrable  $\boldsymbol{\xi}$ ,  $\alpha \in (0, 1]$ , and  $\tau \in \mathbb{R}$ , one has  $\text{s-rsk}_\alpha(\boldsymbol{\xi}) \leq \tau$  if and only if  $\mathcal{P}_\tau(\boldsymbol{\xi}) \leq 1 - \alpha$ .

## 9.2 Dynamic and Multi-Stage Optimization

Many real-world problems involve intricate and gradual revelation of the “true” values of uncertain parameters, usually intertwined with decisions, which even might be continuously adjusted. Modeling of such situations lead to optimal control problems, Markov decision processes, and multi-stage stochastic programming. Decisions now need to account for the uncertainty associated with stochastic processes as well as the opportunities for later recourse actions. This complicates notions of risk-averseness and how to model it in a meaningful and computationally tractable manner.

Again, the initial efforts were motivated by financial applications [21] and dealt with cadlag processes [49, 50]; see also [51]. The papers [262, 293] sever connections with the “static” case and develop dynamic and conditional risk measures from an axiomatic point of view. Around the same time, [84] defines polyhedral risk measures for multi-stage problems; see also [127]. Informally, a conditional risk measure at time  $t$  assigns to a sequence of random variables (representing future “costs”) a random variable modeling the amount we would be willing to pay at time  $t$  to avoid facing the future costs. A collection of such conditional risk measures, one for each time period, amounts to a dynamic measure of risk. Time consistency of a dynamic measure of risk is then the property that it will deem, at the present time, a future sequence of random costs  $\boldsymbol{\xi}^1$  no worse than another sequence  $\boldsymbol{\xi}^2$  whenever it deems  $\boldsymbol{\xi}^1$  at least as good as  $\boldsymbol{\xi}^2$  at a future time  $t$  and  $\boldsymbol{\xi}^1$  is identical to  $\boldsymbol{\xi}^2$  between now and  $t$ . We refer to the tutorials [290, 305], the monograph [246], and the recent papers [79, 92] for many more details. Connections with distributionally robust optimization appear in [252]. The paper [139] discusses the difference between quantifying the risk using an end state as compared to quantifying it using a “composite” risk incurred at each stage in a decision process. Risk measures in the context of Markov decision processes are pioneered in [289, 291], which furnish new thinking about time consistency, dynamic programming equations, and a value iteration method. Later efforts include [54, 301, 342, 33].

The papers [107, 53] address reinforcement learning, while [214, 251, 334] study optimal control and [308] examines dual dynamic programming from a risk-averse perspective.

## 9.3 Challenges and Open Problems

Despite 25 years of rapid development, the area of risk measures and related concepts remains relevant with many opportunities and open problems.

A wealth of opportunities emerge in connection with *autonomous systems* such as driverless cars and delivery vehicles; see the recent review [334]. The resulting problems are often “dynamic” with the need for a combination of learning and decision making, but with conservativeness and reliability as a main focus.

As seen above, there have been many efforts to utilize superquantiles in *statistical learning* but fewer to leverage other risk measures; see [102, 103] for notable exceptions. In view of Section 5, it appears that measures of error and deviation should be central to this area, in fact more than measures of risk. Besides [345], there seem to be few attempts at utilizing measures of risk for *adversarial training* and similar approaches to “robust” learning. Since adversarial training corresponds to the choice of a worst-case risk measure, it only represents one of many possibilities for promoting conservativeness and robustness.

There are surprisingly few studies of *Bayesian formulations* of risk-averse optimization problems under uncertainty. A recent effort is [341], which applies a risk measure to a posterior distribution, but this can be extended in several directions and presents alternatives to distributionally robust formulations. Situations with *decision-dependent probability measures* and *partially observable Markov decision processes* are likewise fertile ground; see [76] for emerging ideas.

As discussed in Subsection 6.3, risk measures may arise in response to ambiguity about the “true” probability distribution. A source of such ambiguity could be a quantity of interest  $f(\boldsymbol{\xi}, x)$  for which we know only the marginal distributions of the random vector  $\boldsymbol{\xi} = (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_m)$ . The paper [112] examines the situation in the context of distributionally robust optimization. It would be useful to develop *risk envelopes* in concrete situations with such “missing information.”

A solution  $(\gamma^*, c^*)$  of the generalized regression problem (5.1) produces a “best” approximation  $\gamma^* + \langle c^*, \boldsymbol{\xi} \rangle$  of the random variable  $\boldsymbol{\eta}$  in the sense of a regular measure of error  $\mathcal{E}$ . Still,  $\gamma^* + \langle c^*, \boldsymbol{\xi} \rangle$  may not approximate well the conditional statistic  $\mathcal{S}(\boldsymbol{\eta}_\xi)$  for all  $\xi$ . Here,  $\mathcal{S}$  is the statistic corresponding to  $\mathcal{E}$  in the sense of Definition 5.1 and  $\boldsymbol{\eta}_\xi$  is a random variable with the distribution of  $\boldsymbol{\eta}$  conditional on  $\boldsymbol{\xi}$  having the value  $\xi$ . Theorem 5.1 of [271] addresses such *statistic tracking* for cases with additive noise and certain risk measures. However, we would like to understand this issue in broader settings.

The construction of *surrogates* (cf. Example 2.5) and supporting *designs of experiments* may need to be carried out conservatively. As in [35], we often seek surrogates that predict response quantities from more easily available parameters and models. With our orientation toward response quantities such as “cost” and “damage,” where high values are undesirable, we tend to be more concerned about underestimating the response than overestimating. Design of experiments may need to be carry out with these factors in mind, while also considering risk-averseness relative to the outcomes of the experiments. For example, with few experiments left, the common expected improvement criterion for selecting the next design point might be replaced with a risk-averse alternative. We refer to [162] for initial efforts in this direction, but numerous opportunities remain with the possibility of improving how complex physical systems are assessed and designed. Ideas from distributionally robust Bayesian optimization [154] may provide guidance for efforts in this direction.

Risk-based formulations also arise for *games and equilibrium problems* as exemplified by [261, 201, 234, 189, 95]. The consideration of risk in these settings is complicated by the two sources of uncertainty

for an agent: lack of knowledge about the action of the other agents and inherent randomness in the environment. Moreover, each agent might have their own tolerance for risk, which raises the question whether there is a central planner, utilizing *some* measure of risk, that achieves a distributed solution. Since games and equilibrium problems tend to be large scale, nonsmooth, and nonconvex, there are also numerous computational challenges in this area.

While extending back to [268], *buffered probabilities* have not been studied in the context of surrogate models. The recent sensitivity results in [286] may provide a useful stepping stone toward surrogates for buffered probabilities of exceedance. While a chance constraint naturally extends from a single random variable not exceeding a threshold to a random vector taking values in a multi-dimensional set, a buffered probability constraint is presently limited to a single random variable. With the introduction of *measures of reliability*, we hope to spur the development of many alternatives to failure probabilities and buffered probabilities as a parallel to the numerous measures of risk.

**Acknowledgement.** The author is thankful for valuable input and encouragement from Amir Ahmadi-Javid, Harbir Antil, Christian Fröhlich, Drew Kouri, Boris Kramer, and Ruodu Wang and for the insights provided by two reviewers. This work is supported in part by ONR (Mathematical and Resource Optimization), ONR (Science of Autonomy), and AFOSR (Mathematical Optimization).

## References

- [1] A. E. Abbas and A. H. Cadenbach. On the use of utility theory in engineering design. *IEEE Systems J.*, 12(2):1129–1139, 2018.
- [2] C. Acerbi. Spectral measures of risk: a coherent representation of subjective risk aversion. *J. Banking and Finance*, 26(7):1505–1518, 2002.
- [3] C. Acerbi and D. Tasche. On the coherence of expected shortfall. *J. Banking and Finance*, 26(7):1487–1503, 2002.
- [4] A. Ahmadi-Javid. An information-theoretic approach to constructing coherent risk measures. In *2011 IEEE International Symposium on Information Theory Proceedings*, pages 2125–2127, 2011.
- [5] A. Ahmadi-Javid. Addendum to: Entropic value-at-risk: A new coherent risk measure. *J. Optimization Theory and Applications*, 155(3):1124–1128, 2012.
- [6] A. Ahmadi-Javid. Entropic value-at-risk: A new coherent risk measure. *J. Optimization Theory and Applications*, 155(3):1105–1123, 2012.
- [7] S. Ahmed, U. Cakmak, and A. Shapiro. Coherent risk measures in inventory problems. *European J. Operational Research*, 182(1):226–238, 2007.
- [8] P. Akella, A. Dixit, M. Ahmadi, J. W. Burdick, and A. D. Ames. Sample-based bounds for coherent risk measures: Applications to policy synthesis and verification. *Preprint arXiv:2204.09833*, 2022.

- [9] J. C. Alais, P. Carpentier, and M. De Lara. Multi-usage hydropower single dam management: chance-constrained optimization and stochastic viability. *Energy Systems*, 8(1):7–30, 2017.
- [10] S. Alexander, T.F. Coleman, and Y. Li. Minimizing CVaR and VaR for a portfolio of derivatives. *J. Banking & Finance*, 30(2):583–605, 2006.
- [11] A. M. A. Alghamdi, P. Chen, and M. Karamehmedovic. Optimal design of photonic nanojets under uncertainty. *Preprint arXiv:2209.02454*, 2022.
- [12] T. S. Alotaibi, L. Dalla Valle, and M. J. Craven. The worst case GARCH-copula CVaR approach for portfolio optimisation: Evidence from financial markets. *J. Risk and Financial Management*, 15(10):482, 2022.
- [13] American Optimal Decisions, Inc. [www.aorda.com](http://www.aorda.com), accessed, November 10, 2022.
- [14] E. Anderson, H. Xu, and D. Zhang. Varying confidence levels for CVaR risk measures and minimax limits. *Mathematical Programming*, 180:327–370, 2020.
- [15] M. Ang, J. Sun, and Q. Yao. On the dual representation of coherent risk measures. *Annals of Operations Research*, 262:29–46, 2018.
- [16] H. Antil, S. Dolgov, and A. Onwunta. Ttrisk: Tensor train decomposition algorithm for risk averse optimization. *Numerical Linear Algebra with Applications*, 30(3):e2481, 2023.
- [17] H. Antil, D. P. Kouri, and J. Pfeifferer. Risk-averse control of fractional diffusion with uncertain exponent. *SIAM J. Control and Optimization*, 59(2):1161–1187, 2021.
- [18] S. Arpon, T. Homem-de-Mello, and B. Pagnoncelli. Scenario reduction for risk-averse stochastic programs. *Preprint optimization-online.org*, 2018.
- [19] P. Artzner, F. Delbaen, J.-M. Eber, and D. Heath. Thinking coherently. *Risk*, 10:68–71, 1997.
- [20] P. Artzner, F. Delbaen, J.-M. Eber, and D. Heath. Coherent measures of risk. *Mathematical Finance*, 9(3):203–228, 1999.
- [21] P. Artzner, F. Delbaen, J. M. Eber, D. Heath, and H. Ku. Coherent multiperiod risk measurement. Preprint, ETH, 2002.
- [22] C. Audet and W. Hare. *Derivative-free and blackbox optimization*. Springer, Cham, 2017.
- [23] AzaPy. <https://github.com/Mircea-MMXXI/azapy>, accessed, November 10, 2022.
- [24] F. Bagheri, H. Dagdougui, and M. Gendreau. Stochastic optimization and scenario generation for peak load shaving in smart district microgrid: sizing and operation. *Energy and Buildings*, 275:112426, 2022.

- [25] H. G. Basova, R. T. Rockafellar, and J. O. Royset. A computational study of the buffered failure probability in reliability-based design optimization. In *Proceedings of the 11th International Conference on Application of Statistics and Probability in Civil Engineering, Zurich, Switzerland*, 2011.
- [26] S. P. Bath and L. A. Prashanth. Concentration of risk measures: A Wasserstein distance approach. In *Advances in Neural Information Processing Systems*, 2019.
- [27] H. H. Bauschke and P. L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer, 2011.
- [28] F. Bellini and V. Bignozzi. On elicitable risk measures. *Quantitative Finance*, 15(5):725–733, 2015.
- [29] A. Ben-Tal, L. El Ghaoui, and A. Nemirovski. *Robust Optimization*. Princeton University Press, 2009.
- [30] A. Ben-Tal and M. Teboulle. Expected utility, penalty functions and duality in stochastic non-linear programming. *Management Science*, 32:1445–1466, 1986.
- [31] A. Ben-Tal and M. Teboulle. An old-new concept of convex risk measures: the optimal certainty equivalent. *Mathematical Finance*, 17:449–476, 2007.
- [32] D. Bertsimas, D. B. Brown, and C. Caramanis. Theory and applications of robust optimization. *SIAM Review*, 53(3):464–501, 2011.
- [33] T. R. Bielecki, I. Cialenco, and A. Ruszczyński. Risk filtering and risk-averse control of Markovian systems subject to model uncertainty. *Mathematical Methods of Operations Research*, to appear, 2023.
- [34] J. R. Birge and F. Louveaux. *Introduction to Stochastic Programming*. Springer Series in Operations Research and Financial Engineering. Springer, 2. edition, 2011.
- [35] L. Bonfiglio and J. O. Royset. Multidisciplinary risk-adaptive set-based design of supercavitating hydrofoils. *AIAA J.*, 57(8):3360–3378, 2019.
- [36] S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities. A Nonasymptotic Theory of Independence*. Oxford University Press, 2013.
- [37] D. B. Brown. Large deviations bounds for estimating conditional value-at-risk. *Operations Research Letters*, 35:722–730, 2007.
- [38] J. V. Burke and T. Hoheisel. Epi-convergence properties of smoothing by infimal convolution. *Set-Valued and Variational Analysis*, 25:1–23, 2017.

- [39] J.-E. Byun, W. de Oliveira, and J. O. Royset. S-BORM: Reliability-based optimization of general systems using buffered optimization and reliability method. *Reliability Engineering & System Safety*, 236:109314, 2023.
- [40] J.-E. Byun and J. O. Royset. Data-driven optimization of reliability using buffered failure probability. *Structural Safety*, 98:102232, 2022.
- [41] Z. Cai and X. Wang. Nonparametric estimation of conditional VaR and expected shortfall. *Journal of Econometrics*, 147(1):120–130, 2008.
- [42] A. M. Campbell, M. Gendreau, and B. W. Thomas. The orienteering problem with stochastic travel and service times. *Annals of Operations Research*, 186(1):61–81, 2011.
- [43] E. Castagnoli, G. Cattelan, F. Maccheroni, C. Tebaldi, and R. Wang. Star-shaped risk measures. *Operations Research*, 70(5):2637–2654, 2022.
- [44] A. Chaudhuri, M. Norton, and B. Kramer. Risk-based design optimization via probability of failure, conditional value-at-risk, and buffered probability of failure. In *Proceeding of AIAA Scitech 2020 Forum*, 2020.
- [45] A. Chaudhuri, B. Peherstorfer, and K. Willcox. Multifidelity cross-entropy estimation of conditional value-at-risk for risk-averse design optimization. In *Proceeding of AIAA Scitech 2020 Forum*, 2129, 2020.
- [46] P. Chen, M. R. Haberman, and O. Ghattas. Optimal design of acoustic metamaterial cloaks under uncertainty. *J. Computational Physics*, 431:110114, 2021.
- [47] X. Chen. Smoothing methods for nonsmooth, nonconvex minimization. *Mathematical Programming*, 134:71–99, 2012.
- [48] X. Chen and K.-K. Kim. Efficient var and cvar measurement via stochastic kriging. *INFORMS J. Computing*, 28(4):629–644, 2016.
- [49] P. Cheridito, F. Delbaen, and M. Kupper. Coherent and convex monetary risk measures for bounded cadlag processes. *Stochastic Processes and their Applications*, 112(1):1–22, 2004.
- [50] P. Cheridito, F. Delbaen, and M. Kupper. Coherent and convex monetary risk measures for unbounded cadlag processes. *Finance and Stochastics*, 9(3):369–387, 2005.
- [51] P. Cheridito, F. Delbaen, and M. Kupper. Dynamic monetary risk measures for bounded discrete-time processes. *Electronic J. Probability*, 11:57–106, 2006.
- [52] P. Cheridito and T. Li. Dual characterization of properties of risk measures on Orlicz hearts. *Mathematics and Financial Economics*, 2(1):29–55, 2008.

- [53] Y. Chow, M. Ghavamzadeh, L. Janson, and M. Pavone. Risk-constrained reinforcement learning with percentile risk criteria. *J. Machine Learning Research*, 18:1–51, 2018.
- [54] Y. Chow, A. Tamar, S. Mannor, and M. Pavone. Risk-sensitive and robust decision-making: a CVaR optimization approach. *Preprint arXiv:1506.02188*, 2015.
- [55] S. Y. Chun, A. Shapiro, and S. Uryasev. Conditional value-at-risk and average value-at-risk: Estimation and asymptotics. *Operations Research*, 60(4):739–756, 2012.
- [56] M. Claus, V. Krätschmer, and R. Schultz. Weak continuity of risk functionals with applications to stochastic programming. *SIAM J. Optimization*, 27(1):91–109, 2017.
- [57] R. Cominetti and A. Torrico. Additive consistency of risk measures and its application to risk-averse routing in networks. *Mathematics of Operations Research*, 41(4):1510–1521, 2022.
- [58] C. W. Commander, P. M. Pardalos, V. Ryabchenko, S. Uryasev, and G. Zrazhevsky. The wireless network jamming problem. *J. Combinatorial Optimization*, 14:481–498, 2007.
- [59] ConditionalValueAtRisk. <https://github.com/jaantollander/ConditionalValueAtRisk>, accessed, November 10, 2022.
- [60] S. Cuomo, V. S. Di Cola, F. Giampaolo, G. Rozza, M. Raissi, and F. Piccialli. Scientific machine learning through physics-informed neural networks: Where we are and what’s next. *J. Scientific Computing*, 92:88, 2022.
- [61] S. Curi, K. Y. Levy, S. Jegelka, and A. Krause. Adaptive sampling for stochastic risk-averse learning. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020*, 2020.
- [62] CVaR-Portfolio. <https://github.com/jaydu1/CVaR-Portfolio>, accessed, November 10, 2022.
- [63] D. Davis and D. Drusvyatskiy. Stochastic model-based minimization of weakly convex functions. *SIAM J. Optimization*, 29(1):207–239, 2018.
- [64] D. Davis, D. Drusvyatskiy, Y. T. Lee, S. Padmanabhan, and G. Ye. A gradient sampling method with complexity guarantees for Lipschitz functions in high and low dimensions. In *In 36th Conference on Neural Information Processing Systems (NeurIPS) 2022*, 2022.
- [65] D. Davis, D. Drusvyatskiy, K. J. MacPhee, and C. Paquette. Subgradient methods for sharp weakly convex functions. *J. Optimization Theory and Applications*, 179:962–982, 2018.
- [66] J. R. Davis and S. Uryasev. Analysis of tropical storm damage using buffered probability of exceedance. *Natural Hazards*, 83(1):465–483, 2016.
- [67] F. Delbaen. Draft: Coherent risk measures. Lecture Notes, Pisa, 2000.



- [68] F. Delbaen. Coherent risk measures on general probability spaces. In K. Sandmann and P. J. Schönbucher, editors, *Advances in Finance and Stochastics*, pages 1–37. Springer, Berlin, 2002.
- [69] F. Delbaen. Law of large numbers for risk measures. *Preprint arXiv:2109.10612*, 2021.
- [70] F. Delbaen and K. Owari. Convex functions on dual Orlicz spaces. *Positivity*, 23(5):1051–1064, 2019.
- [71] D. Dentcheva, Y. Lin, and S. Penev. Stability and sample-based approximations of composite stochastic optimization problems. *Operations Research*, to appear, 2022.
- [72] D. Dentcheva and G. Martinez. Two-stage stochastic optimization problems with stochastic ordering constraints on the recourse. *European J. Operational Research*, 219:1–8, 2012.
- [73] D. Dentcheva, S. Penev, and A. Ruszczyński. Statistical estimation of composite risk functionals and risk optimization problems. *Annals of the Institute of Statistical Mathematics*, 69:737–760, 2017.
- [74] D. Dentcheva and A. Ruszczyński. Optimization with stochastic dominance constraints. *SIAM J. Optimization*, 14(2):548–566, 2003.
- [75] D. Dentcheva and A. Ruszczyński. Common mathematical foundations of expected utility and dual utility theories. *SIAM J. Optimization*, 23(1):381–405, 2013.
- [76] D. Dentcheva and A. Ruszczyński. Risk forms: representation, disintegration, and application to partially observable two-stage systems. *Mathematical Programming*, 181:297–317, 2020.
- [77] M. Denuit, J. Dhaene, M. Goovaerts, R. Kaas, and R. Laeven. Risk measurement with equivalent utility principles. *Statistics & Decisions*, 24:1–25, 2006.
- [78] P. Dommel and A. Pichler. Convex risk measures based on divergence. *Preprint arXiv:2003.07648*, 2020.
- [79] P. Dommel and A. Pichler. Foundations of multistage stochastic programming. *Preprint arXiv:2102.07464*, 2021.
- [80] J. C. Duchi and H. Namkoong. Learning models with uniform performance via distributionally robust optimization. *The Annals of Statistics*, 49(3):1378–1406, 2021.
- [81] J. C. Duchi and F. Ruan. Stochastic methods for composite and weakly convex optimization problems. *SIAM J. Optimization*, 28(4):3229–3259, 2018.
- [82] J. Dupačová. On minimax solutions of stochastic linear programming problems. *Časopis pro Pěstování Matematiky*, 91(4):423–430, 1966.
- [83] T. Dylan, G. Frederic, and Y. J. Yuan. Risk-averse action selection using extreme value theory estimates of the CVaR. *Preprint arXiv:1912.01718*, 2019.

- [84] A. Eichhorn and W. Römisch. Polyhedral risk measures in stochastic programming. *SIAM J. Optimization*, 16:69–95, 2005.
- [85] P. Embrechts. Extreme value theory as a risk management tool. *North American Actuarial J.*, 3(2):30–41, 1999.
- [86] P. Embrechts, A. Schied, and R. Wang. Robustness in the optimization of risk measures. *Operations Research*, 70(1):95–110, 2022.
- [87] EMP. <https://github.com/xhub/EMP.jl>, accessed, November 10, 2022.
- [88] O. G. Ernst, A. Pichler, and B. Sprungh. Wasserstein sensitivity of risk and uncertainty propagation. *SIAM/ASA J. Uncertainty Quantification*, 10(3):915–948, 2022.
- [89] C. I. Fabian. Handling CVaR objectives and constraints in two-stage stochastic models. *European J. Operational Research*, 191(3):888–911, 2008.
- [90] C. I. Fabian, C. Wolf, A. Koberstein, and L. Suhl. Risk-averse optimization in two-stage stochastic models: Computational aspects and a study. *SIAM J. Optimization*, 25(1):28–52, 2015.
- [91] J. Fairbrother, A. Turner, and S. W. Wallace. Problem-driven scenario generation: an analytical approach for stochastic programs with tail risk measure. *Mathematical Programming*, 191:141–182, 2022.
- [92] J. Fan and A. Ruszczyński. Process-based risk measures and risk-averse control of discrete-time systems. *Mathematical Programming*, 191:113–140, 2022.
- [93] Y. Fan, S. Lyu, Y. Ying, and B.-G. Hu. Learning with average top-k loss. *Preprint arXiv:1705.08826*, 2017.
- [94] E. Fernandez, Y. Hinojosa, J. Puerto, and F. Saldanha-da-Gama. New algorithmic framework for conditional value at risk: Application to stochastic fixed-charge transportation. *European J. Operational Research*, 277(1):215–226, 2019.
- [95] M. Ferris and A. Philpott. Dynamic risk equilibrium. *Operations Research*, 70(3):1933–1952, 2022.
- [96] C. Filippi, G. Guastaroba, and M. G. Speranza. Conditional value-at-risk beyond finance: a survey. *International Transactions in Operational Research*, 27(3):1277–1319, 2020.
- [97] H. Föllmer and A. Schied. Convex measures of risk and trading constraints. *Finance and Stochastics*, 6(4):429–447, 2002.
- [98] H. Föllmer and A. Schied. Robust preferences and convex measures of risk. In K. Sandmann and P. J. Schönbucher, editors, *Advances in Finance and Stochastics*, page 39–56. Springer, Berlin, 2002.

- [99] H. Föllmer and A. Schied. *Stochastic Finance: An Introduction in Discrete Time*. de Gruyter, 4. edition, 2016.
- [100] H. Föllmer and S. Weber. The axiomatic approach to risk measures for capital determination. *Annual Review of Financial Economics*, 7(1):301–337, 2015.
- [101] D. Friedman, R. M. Isaac, D. James, and S. Sunder. *Risky Curves*. Routledge, 2014.
- [102] C. Fröhlich and R. C. Williamson. Risk measures and upper probabilities: Coherence and stratification. *Preprint arXiv:2206.03183*, 2022.
- [103] C. Fröhlich and R. C. Williamson. Tailoring to the tails: Risk measures for fine-grained tail sensitivity. *Transactions on Machine Learning Research*, 1, 2023.
- [104] R. Frongillo and I. A. Kash. Elicitation complexity of statistical properties. *Biometrika*, 108(4):857–879, 2021.
- [105] H. Gangammanavar and S. Sen. Two-scale stochastic optimization for controlling distributed storage devices. *IEEE Transactions on Smart Grid*, 9(4):2691–2702, 2016.
- [106] F. Gao and S. Wang. Asymptotic behavior of the empirical conditional value-at-risk. *Insurance: Mathematics and Economics*, 49(3):345–352, 2011.
- [107] J. Garc and F. Fern. A comprehensive survey on safe reinforcement learning. *J. Machine Learning Research*, 16(1):1437–1480, 2015.
- [108] S. Garreis, T. M. Surowiec, and M. Ulbrich. An interior-point approach for solving risk-averse PDE-constrained optimization problems with coherent risk measures. *SIAM J. Optimization*, 31(1):1–29, 2021.
- [109] G. Garrigos and R. M. Gower. Handbook of convergence theorems for (stochastic) gradient methods. *Preprint arXiv:2301.11235*, 2023.
- [110] B. Geihe, M. Lenz, M. Rumpf, and R. Schultz. Risk averse elastic shape optimization with parametrized fine scale geometry. *Mathematical Programming*, 141(1):383–403, 2013.
- [111] A. Ghasemi, H. Jamshidi Monfared, A. Loni, and M. Marzband. CVaR-based retail electricity pricing in day-ahead scheduling of microgrids. *Energy*, 227:120529, 2021.
- [112] M. Ghossoub, J. Hall, and D. Saunders. Maximum spectral measures of risk with given risk factor marginal distributions. *Mathematics of Operations Research*, 48(2):1158–1182, 2022.
- [113] M. B. Giles and A.-L. Haji-Ali. Multilevel nested simulation for efficient risk estimation. *SIAM/ASA J. Uncertainty Quantification*, 7(2):497–525, 2019.
- [114] P. W. Glynn, L. Fan, M. C. Fu, J.-Q. Hu, and Y. Peng. Technical note—central limit theorems for estimated functions at estimated points. *Operations Research*, 68(5):1557–1563, 2020.

- [115] T. Gneiting. Making and evaluating point forecasts. *J. American Statistical Association*, 106(494):746–762, 2011.
- [116] A. Golodnikov, V. Kuzmenko, and S. Uryasev. CVaR regression based on the relation between CVaR and mixed-quantile quadrangles. *J. Risk and Financial Management*, 12(3):107, 2019.
- [117] J. Gotoh, M. J. Kim, and A. E. B. Lim. Robust empirical optimization is almost the same as mean-variance optimization. *Operations Research Letters*, 46(4):448–452, 2018.
- [118] J. Gotoh, M. J. Kim, and A. E. B. Lim. Calibration of distributionally robust empirical optimization models. *Operations Research*, 69(5):1630–1650, 2021.
- [119] J. Gotoh and Y. Takano. Newsvendor solutions via conditional value-at-risk minimization. *European J. Operational Research*, 179:80–96, 2007.
- [120] J. Gotoh, A. Takeda, and R. Yamamoto. Interaction between financial risk measures and machine learning methods. *Computational Management Science*, 11(4):365–402, 2014.
- [121] J. Gotoh and S. Uryasev. Two pairs of families of polyhedral norms versus  $\ell_p$ -norms: proximity and applications in optimization. *Mathematical Programming*, 156(1):391–431, 2016.
- [122] J. Gotoh and S. Uryasev. Support vector machines based on convex risk functions and general norms. *Annals of Operations Research*, 249(1):301–328, 2017.
- [123] B. Grechuk, A. Molyboha, and M. Zabaranin. Maximum entropy principle with general deviation measures. *Mathematics of Operations Research*, 34(2):445–467, 2009.
- [124] B. Grechuk and M. Zabaranin. Sensitivity analysis in applications with deviation, risk, regret, and error measures. *SIAM J. Optimization*, 27(4):2481–2507, 2017.
- [125] B. Grechuk and M. Zabaranin. Regression analysis: likelihood, error and entropy. *Mathematical Programming*, 174(1):145–166, 2019.
- [126] V. Guigues, V. Krätschmer, and A. Shapiro. A central limit theorem and hypotheses testing for risk-averse stochastic programs. *SIAM J. Optimization*, 28(2):1337–1366, 2018.
- [127] V. Guigues and W. Römisch. Sampling-based decomposition methods for multistage stochastic programs based on extended polyhedral risk measures. *SIAM J. Optimization*, 22:286–312, 2012.
- [128] S. Guo and H. Xu. Robust spectral risk optimization when the subjective risk aversion is ambiguous: a moment-type approach. *Mathematical Programming*, 194:305–340, 2022.
- [129] G. A. Hanasusanto, D. Kuhn, S. W. Wallace, and S. Zymler. Distributionally robust multi-item newsvendor problems with multimodal demand distributions. *Mathematical Programming*, 152(1):1–32, 2015.

- [130] T. B. Hashimoto, M. Srivastava, H. Namkoong, and P. Liang. Fairness without demographics in repeated loss minimization. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*. PMLR, 2018.
- [131] M. Heinkenschloss, B. Kramer, and T. Takhtaganov. Adaptive reduced-order model construction for conditional value-at-risk estimation. *SIAM/ASA J. Uncertainty Quantification*, 8(2):668–692, 2020.
- [132] M. Heinkenschloss, B. Kramer, T. Takhtaganov, and K. Willcox. Conditional-value-at-risk estimation via reduced-order models. *SIAM/ASA J. Uncertainty Quantification*, 6(4):1395–1423, 2018.
- [133] R. Hemmati, H. Saboori, and S. Saboori. Stochastic risk-averse coordinated scheduling of grid integrated energy storage units in transmission constrained wind-thermal systems within a conditional value-at-risk framework. *Energy*, 113:762–775, 2016.
- [134] C. Hess and R. Seri. Generic consistency for approximate stochastic programming and statistical problems. *SIAM J. Optimization*, 29(1):290–317, 2019.
- [135] N. Ho-Nguyen and S. J. Wright. Adversarial classification via distributional robustness with Wasserstein ambiguity. *Mathematical Programming*, 198:1411–1447, 2023.
- [136] M. J. Holland and E. M. Haress. Spectral risk-based learning using unbounded losses. *Preprint arXiv:2105.04816*, 2021.
- [137] L. J. Hong, Z. Hu, and G. Liu. Monte Carlo methods for value-at-risk and conditional value-at-risk. *ACM Transactions on Modeling and Computer Simulation*, 24(4):1–37, 2014.
- [138] P. Huber. *Robust Statistics*. Wiley, 1981.
- [139] D. A Iancu, M. Petrik, and D. Subramanian. Tight approximations of dynamic risk measures. *Mathematics of Operations Research*, 40(3):655–682, 2015.
- [140] J. D. Jakeman, D. P. Kouri, and J. G. Huerta. Surrogate modeling for efficiently, accurately and conservatively estimating measures of risk. *Reliability Engineering & System Safety*, 221:108280, 2022.
- [141] X. Jiang, R. Bai, S. W. Wallace, G. Kendall, and D. Landa-Silva. Soft clustering-based scenario bundling for a progressive hedging heuristic in stochastic service network design. *Computers & Operations Research*, 128:105182, 2021.
- [142] D. Kahneman and A. Tversky. Prospect theory: An analysis of decision under risk. *Econometrica*, 47(2):263–291, 1979.
- [143] D. S. Kalogierias and W. B. Powell. Recursive optimization of convex risk measures: Mean-semideviation models. *Preprint arXiv:1804.00636*, 2018.

- [144] D. S. Kalogierias and W. B. Powell. Zeroth-order stochastic compositional algorithms for risk-aware learning. *SIAM J. Optimization*, 32(2):386–416, 2022.
- [145] M. G. Kapteyn, D. J. Knezevic, and K. Willcox. Toward predictive digital twins via component-based reduced-order models and interpretable machine learning. In *AIAA Scitech 2020 Forum*, page 0418, 2020.
- [146] M. G. Kapteyn, K. E. Willcox, and A. B. Philpott. Distributionally robust optimization for engineering design under uncertainty. *International J. Numerical Methods in Engineering*, 120(7):835–859, 2019.
- [147] G. E. Karniadakis, I. G. Kevrekidis, L. Lu, P. Perdikaris, S. Wang, and L. Yang. Physics informed machine learning. *Nature Reviews Physics*, 3:422–440, 2022.
- [148] K. Kato. Weighted Nadaraya-Watson estimation of conditional expected shortfall. *Journal of Financial Econometrics*, 10(2):265–291, 2012.
- [149] K. Kawaguchi and H. Lu. Ordered SGD: A new stochastic optimization framework for empirical risk minimization. *Preprint arXiv:1907.04371*, 2019.
- [150] A. Khayyer, A. Vinel, and J. J. Kennedy. Efficient global estimation of conditional-value-at-risk through stochastic kriging and extreme value theory. *Preprint arXiv:2403.19018*, 2024.
- [151] J. Khim, L. Leqi, A. Prasad, and P. Ravikumar. Uniform convergence of rank-weighted learning. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 5254–5263. PMLR, 13–18 Jul 2020.
- [152] R. Kiesel, R. Ruhlicke, G. Stahl, and J. Zheng. The Wasserstein metric and robustness in risk management. *Risks*, 4(3), 2016.
- [153] M. Kijima and M. Ohnishi. Mean-risk analysis of risk aversion and wealth effects on optimal portfolios with multiple investment opportunities. *Annals of Operations Research*, 45:147–163, 1993.
- [154] J. Kirschner, I. Bogunovic, S. Jegelka, and A. Krause. Distributionally robust Bayesian optimization. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 2174–2184. PMLR, 26–28 Aug 2020.
- [155] R. K. Kolla, P. L.A., S. P. Bhat, and K. Jagannathan. Concentration bounds for empirical conditional value-at-risk: The unbounded case. *Operations Research Letters*, 47(1):16–20, 2019.
- [156] P. Kolvenbach, O. Lass, and S. Ulbrich. An approach for robust PDE-constrained optimization with application to shape optimization of electrical engines and of dynamic elastic structures under uncertainty. *Optimization Engineering*, 19:697–731, 2018.

- [157] H. Konno and H. Shirakawa. Equilibrium relations in the mean-absolute deviation capital market. *Asia-Pacific Financial Markets*, 1:21–35, 1994.
- [158] D. P. Kouri. A measure approximation for distributionally robust PDE-constrained optimization problems. *SIAM J. Numerical Analysis*, 5(6):3147–3172, 2017.
- [159] D. P. Kouri. Higher-moment buffered probability. *Optimization Letters*, 13(6):1223–1237, 2019.
- [160] D. P. Kouri. Spectral risk measures: the risk quadrangle and optimal approximation. *Mathematical Programming*, 174(1):525–552, 2019.
- [161] D. P. Kouri. PDE-constrained optimization under uncertainty. Sand2021-10992pe, Sandia National Laboratories, 2021.
- [162] D. P. Kouri, J. D. Jakeman, and J. G. Huerta. Risk-adapted optimal experimental design. *SIAM/ASA J. Uncertainty Quantification*, 10(2):687–716, 2022.
- [163] D. P. Kouri and A. Shapiro. Optimization of PDEs with uncertain inputs. In H. Antil, D. P. Kouri, M.-D. Lacasse, and D. Ridzal, editors, *Frontiers in PDE-Constrained Optimization*, The IMA Volumes in Mathematics and its Applications, pages 41–81. Springer, Cham, 2018.
- [164] D. P. Kouri and T. M. Surowiec. Risk-averse PDE-constrained optimization using the conditional value-at-risk. *SIAM J. Optimization*, 26(1):365–396, 2016.
- [165] D. P. Kouri and T. M. Surowiec. Existence and optimality conditions for risk-averse PDE-constrained optimization. *SIAM/ASA J. Uncertainty Quantification*, 6(2):787–815, 2018.
- [166] D. P. Kouri and T. M. Surowiec. Epi-regularization of risk measures. *Mathematics of Operations Research*, 45(2):774–795, 2020.
- [167] D. P. Kouri and T. M. Surowiec. A primal-dual algorithm for risk minimization. *Mathematical Programming*, 193:337–363, 2022.
- [168] P. Krokmal, J. Palmquist, and S. Uryasev. Portfolio optimization with conditional value-at-risk objective and constraints. *J. Risk*, 4:43–68, 2002.
- [169] P. Krokmal, M. Zabaranin, and S. Uryasev. Modeling and optimization of risk. *Surveys in Operations Research and Management Science*, 16(2):49–66, 2011.
- [170] A. Künzi-Bay and J. Mayer. Computational aspects of minimizing conditional value-at-risk. *Computational Management Science*, 3:3–27, 2006.
- [171] S. Kusuoka. On law-invariant coherent risk measures. In S. Kusuoka and T. Maruyama, editors, *Advances in Mathematical Economics, Volume 3*, pages 83–95. Springer, 2001.
- [172] R. J. A. Laeven and M. Stajda. Entropy coherent and entropy convex measures of risk. *Mathematics of Operations Research*, 438(2):265–293, 2013.

- [173] Y. Laguel, J. Malick, and Z. Harchaoui. First-order optimization for superquantile-based supervised learning. In *30th IEEE International Workshop on Machine Learning for Signal Processing, MLSP 2020, Espoo, Finland, September 21-24, 2020*, pages 1–6, 2020.
- [174] Y. Laguel, K. Pillutla, J. Malick, and Z. Harchaoui. A superquantile approach to federated learning with heterogeneous devices. In *Proceedings of 55th Annual Conference on Information Sciences and Systems (CISS), Baltimore, MD, USA, March 24-26, 2021*, pages 1–6, 2021.
- [175] Y. Laguel, K. Pillutla, J. Malick, and Z. Harchaoui. Superquantiles at work: Machine learning applications and efficient subgradient computation. *Set-Valued and Variational Analysis*, 29:967–996, 2021.
- [176] N. Lambert, D. M. Pennock, and Y. Shoham. Eliciting properties of probability distributions. In *Proceedings of the 9th ACM Conference on Electronic Commerce*, pages 129–138, 2008.
- [177] G. Lan. *First-Order and Stochastic Optimization Methods for Machine Learning*. Springer, 2020.
- [178] G. Lan, A. Nemirovski, and A. Shapiro. Validation analysis of mirror descent stochastic approximation method. *Mathematical Programming*, 134:425–458, 2012.
- [179] D. Lee and B. Kramer. Bi-fidelity conditional value-at-risk estimation by dimensionally decomposed generalized polynomial chaos expansion. *Structural and Multidisciplinary Optimization*, 66:33, 2023.
- [180] D. Lee and B. Kramer. Multifidelity conditional value-at-risk estimation by dimensionally decomposed generalized polynomial chaos-kriging. *Reliability Engineering & System Safety*, 235:109208, 2023.
- [181] J. Lee, S. Park, and J. Shin. Learning bounds for risk-sensitive learning. In *In 34th Conference on Neural Information Processing Systems (NeurIPS) 2020*, 2020.
- [182] S. Leorato, F. Peracchi, and A. V. Tanase. Asymptotically efficient estimation of the conditional expected shortfall. *Computational Statistics & Data Analysis*, 56(4):768–784, 2012.
- [183] L. Leqi, A. Huang, Z. C. Lipton, and K. Azizzadenesheli. Supervised learning with general risk functionals. In *Proceedings of the 39th International Conference on Machine Learning (ICML 2022)*, PMLR 162, 2022.
- [184] L. Leqi, A. Prasad, and P. Ravikumar. On human-aligned risk minimization. In *Proceedings of NeurIPS*, 2019.
- [185] D. Levy, Y. Carmon, J. C. Duchi, and A. Sidford. Large-scale methods for distributionally robust optimization. In *Proceedings of NeurIPS*, 2020.
- [186] A. S. Lewis and S. J. Wright. A proximal method for composite minimization. *Mathematical Programming*, 158:501–546, 2016.



- [187] J. Y.-M. Li. Technical note—closed-form solutions for worst-case law invariant risk measures with application to robust portfolio optimization. *Operations Research*, 66(6):1533–1541, 2018.
- [188] T. Li, A. Beirami, M. Sanjabi, and V. Smith. On tilted losses in machine learning: Theory and applications. *J. Machine Learning Research*, 24:1–79, 2023.
- [189] T. Lianas, E. Nikolova, and N. E. Stier-Moses. Risk-averse selfish routing. *Mathematics of Operations Research*, 44(1):38–57, 2019.
- [190] C. Lim, H. D. Sherali, and S. Uryasev. Portfolio optimization by minimizing conditional value-at-risk via nondifferentiable optimization. *Computational Optimization and Applications*, 46(3):391–415, 2010.
- [191] F. Lin, X. Fang, and Z. Gao. Distributionally robust optimization: A review on theory and applications. *Numerical Algebra, Control and Optimization*, 12(1):159–212, 2022.
- [192] T. Lin, Z. Zheng, and M. I. Jordan. Gradient-free methods for deterministic and stochastic nonsmooth nonconvex optimization. In *In 36th Conference on Neural Information Processing Systems (NeurIPS) 2022*, 2022.
- [193] F. Liu, T. Mao, R. Wang, and L. Wei. Inf-convolution, optimal allocations, and model uncertainty for tail risk measures. *Mathematics of Operations Research*, 47(3):2494–2519, 2022.
- [194] F. Liu and R. Wang. A theory for measures of tail risk. *Mathematics of Operations Research*, 46(3):1109–1128, 2021.
- [195] J. Liu, Y. Cui, and J.-S. Pang. Solving nonsmooth and nonconvex compound stochastic programs with applications to risk measure minimization. *Mathematics of Operations Research*, 47(4):3051–3083, 2022.
- [196] J. Liu and J.-S. Pang. Risk-based robust statistical learning by stochastic difference-of-convex value-function optimization. *Operations Research*, 71(2):397–414, 2022.
- [197] P. Liu, A. Schied, and R. Wang. Distributional transforms, probability distortions, and their applications. *Mathematics of Operations Research*, 46(4):1490–1512, 2021.
- [198] M. O. Lorenz. Methods of measuring concentration of wealth. *Publications of American Statistical Association*, 9(70):209–219, 1905.
- [199] M. Lu, J. G. Shanthikumar, and Z.-J. M. Shen. Technical note—operational statistics: Properties and the risk-averse case. *Naval Research Logistics*, 62(3):206–214, 2015.
- [200] M. Lu and Z.-J. M. Shen. A review of robust operations management under model uncertainty. *Production and Operations Management*, 30(6):1927–1943, 2021.

- [201] J. P. Luna, C. Sagastizabal, and M. Solodov. An approximation scheme for a class of risk-averse stochastic equilibrium problems. *Mathematical Programming*, 157:451–481, 2016.
- [202] Z. Luo. Nonparametric kernel estimation of CVaR under  $\alpha$ -mixing sequences. *Statistical Papers*, 61(2):615–643, 2020.
- [203] Z. Luo and S. Ou. The almost sure convergence rate of the estimator of optimized certainty equivalent risk measure under  $\alpha$ -mixing sequences. *Communications in Statistics - Theory and Methods*, 46(16):8166–8177, 2017.
- [204] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- [205] A. Mafusalov, A. Shapiro, and S. Uryasev. Estimation and asymptotics for buffered probability of exceedance. *European J. Operational Research*, 270(3):826–836, 2018.
- [206] A. Mafusalov and S. Uryasev. CVaR (superquantile) norm: Stochastic case. *European J. Operational Research*, 249(1):200–208, 2016.
- [207] A. Mafusalov and S. Uryasev. Buffered probability of exceedance: Mathematical properties and optimization. *SIAM J. Optimization*, 28(2):1077–1103, 2018.
- [208] M. Mahsuli. *Probabilistic Models, Methods, and Software for Evaluating Risk to Civil Infrastructure*. Ph.d. thesis, University of British Columbia, 2012.
- [209] A. Malandii and S. Uryasev. Support vector regression: Risk quadrangle framework. *Preprint arXiv:2212.09178*, 2022.
- [210] H. M. Markowitz. Portfolio selection. *J. Finance*, 7(1):77–91, 1952.
- [211] H. M. Markowitz. *Portfolio selection, efficient diversification of investments*. Wiley, New York, 1959.
- [212] A. J. McNeil, R. Frey, and P. Embrechts. *Quantitative Risk Management: Concepts, Techniques and Tools*. Princeton University Press, 2015.
- [213] Z. Mhammedi, B. Guedj, and R. C. Williamson. PAC-Bayesian bound for the conditional value at risk. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 17919–17930. Curran Associates, Inc., 2020.
- [214] C. W. Miller and I. Yang. Optimal control of conditional value-at-risk in continuous time. *SIAM J. Control and Optimization*, 55(2):856–884, 2017.

- [215] R. Minguez, E. Castillo, and J. L. Lara. Iterative scenario reduction technique to solve reliability based optimization problems using the buffered failure probability. In G. Deodatis, editor, *Proceedings of ICOSSAR*, 2013.
- [216] P. Mohajerin Esfahani and D. Kuhn. Data-driven distributionally robust optimization using the wasserstein metric: Performance guarantees and tractable reformulations. *Mathematical Programming*, 171(1):115–166, 2018.
- [217] T. Morimura, M. Sugiyama, H. Kashima, H. Hachiya, and T. Tanaka. Nonparametric return distribution approximation for reinforcement learning. In J. Fürnkranz and T. Joachims, editors, *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, page 799–806. Omnipress, 2010.
- [218] S. Nasini, M. Labbe, and L. Brotcorne. Multi-market portfolio optimization with conditional value at risk. *European J. Operational Research*, 300(1):350–365, 2022.
- [219] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM J. Optimization*, 19(4):1574–1609, 2009.
- [220] A. Nemirovski and D. Yudin. On Cezari’s convergence of the steepest descent method for approximating saddle point of convex-concave functions. *Soviet Math. Doklady*, 19(2), 1978.
- [221] A. Nemirovski and D. Yudin. *Problem Complexity and Method Efficiency in Optimization*. Wiley, 1983.
- [222] E. Newman, L. Ruthotto, J. Hart, and B. van Bloemen Waanders. Train like a (var) pro: efficient training of neural networks with variable projection. *SIAM J. Mathematics of Data Science*, 3(4):1041–1066, 2021.
- [223] M. Norton, V. Khokhlov, and S. Uryasev. Calculating CVaR and bPOE for common probability distributions with application to portfolio optimization and density estimation. *Annals of Operations Research*, 299:1281–1315, 2021.
- [224] M. Norton, A. Mafusalov, and S. Uryasev. Soft margin support vector classification as buffered probability minimization. *J. Machine Learning Research*, 18:1–43, 2017.
- [225] M. Norton and J. O. Royset. Diametrical risk minimization: Theory and computations. *Machine Learning*, 112:2933–2951, 2023.
- [226] M. Norton and S. Uryasev. Maximization of AUC and buffered AUC in binary classification. *Mathematical Programming B*, 174:575–612, 2019.
- [227] N. Noyan. Risk-averse two-stage stochastic programming with an application to disaster management. *Computers & Operations Research*, 39:541–559, 2012.

- [228] N. Noyan, M. Merakli, and S. Küçükyavuz. Two-stage stochastic programming under multivariate risk constraints with an application to humanitarian relief network design. *Mathematical Programming*, 191:7–45, 2022.
- [229] N. Noyan and G. Rudolf. Kusuoka representations of coherent risk measures in general probability spaces. *Annals of Operations Research*, 229:591–605, 2015.
- [230] W. Ogryczak and A. Ruszczyński. From stochastic dominance to mean-risk models: semideviations and risk measures. *European J. Operational Research*, 116:33–50, 1999.
- [231] W. Ogryczak and A. Ruszczyński. On consistency of stochastic dominance and mean-semideviation models. *Mathematical Programming*, 89:217–232, 2001.
- [232] W. Ogryczak and A. Ruszczyński. Dual stochastic dominance and related mean-risk models. *SIAM J. Optimization*, 13:60–78, 2002.
- [233] S. O. Ottesen and A. Tomsgard. A stochastic model for scheduling energy flexibility in buildings. *Energy*, 88:364–376, 2015.
- [234] J.-S. Pang, S. Sen, and U. V. Shanbhag. Two-stage non-cooperative games with risk-averse players. *Mathematical Programming*, 165:235–290, 2017.
- [235] K. Pavlikov and S. Uryasev. CVaR norm and applications in optimization. *Optimization Letters*, 8(7):1999–2020, 2014.
- [236] K. Pavlikov and S. Uryasev. CVaR distance between univariate probability distributions and approximation problems. *Annals of Operations Research*, 262(1):67–88, 2018.
- [237] B. Peherstorfer, K. Willcox, and M. Gunzburger. Survey of multifidelity methods in uncertainty propagation, inference, and optimization. *SIAM Review*, 60(3):550–591, 2018.
- [238] F. Peracchi and A. V. Tanase. On estimating the conditional expected shortfall. *Applied Stochastic Models in Business and Industry*, 24:471–493, 2008.
- [239] P. Perdikaris, D. Venturi, J. O. Royset, and G. E. Karniadakis. Multi-fidelity modelling via recursive co-kriging and Gaussian-Markov random fields. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 471(2179):20150018, 2015.
- [240] G. Ch. Pflug. Some remarks on the value-at-risk and the conditional value-at-risk. In S. Uryasev, editor, *Probabilistic Constrained Optimization, Nonconvex Optimization and Its Applications*, pages 272–281. Springer, Boston, 2000.
- [241] G. Ch. Pflug. On distortion functionals. *Statistics and Decisions*, 24:45–60, 2006.
- [242] G. Ch. Pflug. Subdifferential representation of risk measures. *Mathematical Programming*, 108:339–354, 2006.

- [243] G. Ch. Pflug and M. Pohl. A review on ambiguity in stochastic portfolio optimization. *Set-Valued and Variational Analysis*, 26(4):733–757, 2018.
- [244] G. Ch. Pflug and W. Römisch. *Modeling, measuring and managing risk*. World Scientific, 2007.
- [245] G. Ch. Pflug and N. Wozabal. Asymptotic distribution of law-invariant risk functionals. *Finance and Stochastics*, 14:397–418, 2010.
- [246] G.Ch. Pflug and A. Pichler. *Multistage Stochastic Optimization*. Springer, 2014.
- [247] A. Pichler. Evaluations of risk measures for different probability measures. *SIAM J. Optimization*, 23(1):530–551, 2013.
- [248] A. Pichler. The natural Banach space for version independent risk measures. *Insurance: Mathematics and Economics*, 53(2):405–415, 2013.
- [249] A. Pichler. A quantitative comparison of risk measures. *Annals of Operations Research*, 254(1):251–275, 2017.
- [250] A. Pichler and R. Schlotter. Entropy based risk measures. *European J. Operational Research*, 285(1):223–236, 2020.
- [251] A. Pichler and R. Schlotter. Risk-averse optimal control in continuous time by nesting risk measures. *Mathematics of Operations Research*, 48(3):1657–1678, 2022.
- [252] A. Pichler and A. Shapiro. Mathematical foundations of distributionally robust multistage optimization. *SIAM J. Optimization*, 31(4):3044–3067, 2021.
- [253] K. Pillutla, Y. Laguel, J. Malick, and Z. Harchaoui. Superquantile-based learning: A direct approach using gradient-based optimization. *J. Signal Processing Systems*, 94:161–177, 2022.
- [254] K. Pillutla, Y. Laguel, J. Malick, and Z. Harchaoui. Federated learning with superquantile aggregation for heterogeneous data. *Machine Learning*, to appear, 2023.
- [255] E. Polak. *Optimization. Algorithms and Consistent Approximations*, volume 124 of *Applied Mathematical Sciences*. Springer, 1997.
- [256] PortfolioAllocation. <https://github.com/EulersNumber/PortfolioAllocation>, accessed, November 10, 2022.
- [257] K. Postek, D. den Hertog, and B. Melenberg. Computationally tractable counterparts of distributionally robust constraints on risk measures. *SIAM Review*, 58(4):603–650, 2016.
- [258] L. A. Prashanth, K. Jagannathan, and R. K. Kolla. Concentration bounds for CVaR estimation: The cases of light-tailed and heavy-tailed distributions. In *Proceedings of the 37th International Conference on Machine Learning*. PMLR, 2020.

- [259] M. Qi, H.-Y. Mak, and Z.-J. M. Shen. Data-driven research in retail operations—a review. *Naval Research Logistics*, 67(8):595–616, 2020.
- [260] H. Rahimian and S. Mehrotra. Frameworks and results in distributionally robust optimization. *Open J. Mathematical Optimization*, 3(4):1–85, 2022.
- [261] D. Ralph and Y. Smeers. Risk trading and endogenous probabilities in investment equilibria. *SIAM J. Optimization*, 25(4):2589–2611, 2015.
- [262] F. Riedel. Dynamic coherent risk measures. *Stochastic Processes and their Applications*, 112:185–200, 2004.
- [263] Riskfolio-Lib. <https://github.com/dcajasn/Riskfolio-Lib>, accessed, November 10, 2022.
- [264] RiskMeasures. <https://github.com/rubsc/RiskMeasures.jl>, accessed, November 10, 2022.
- [265] R. T. Rockafellar. *Conjugate Duality and Optimization*. SIAM, 1974.
- [266] R. T. Rockafellar. Coherent approaches to risk in optimization under uncertainty. In *Tutorials in Operations Research: OR Tools and Applications: Glimpses of Future Technologies*, Tutorials in Operations Research, pages 38–61. INFORMS, Cantonville, MD, 2007.
- [267] R. T. Rockafellar. Risk and utility in the duality framework of convex analysis. In D. H. Bailey et al., editor, *From Analysis to Visualization, a Celebration of the Life and Legacy of Jonathan M. Borwein*, pages 21–42. Springer, 2020.
- [268] R. T. Rockafellar and J. O. Royset. On buffered failure probability in design and optimization of structures. *Reliability Engineering & System Safety*, 95:499–510, 2010.
- [269] R. T. Rockafellar and J. O. Royset. Random variables, monotone relations, and convex analysis. *Mathematical Programming B*, 148(1):297–331, 2014.
- [270] R. T. Rockafellar and J. O. Royset. Engineering decisions under risk-averseness. *ASCE-ASME J. Risk and Uncertainty in Engineering Systems, Part A: Civil Engineering*, 1(2):04015003, 2015.
- [271] R. T. Rockafellar and J. O. Royset. Measures of residual risk with connections to regression, risk tracking, surrogate models, and ambiguity. *SIAM J. Optimization*, 25(2):1179–1208, 2015.
- [272] R. T. Rockafellar and J. O. Royset. Superquantile/CVaR risk measures: Second-order theory. *Annals of Operations Research*, 262:3–29, 2018.
- [273] R. T. Rockafellar, J. O. Royset, and S. I. Miranda. Superquantile regression with applications to buffered reliability, uncertainty quantification, and conditional value-at-risk. *European J. Operational Research*, 234(1):140–154, 2014.
- [274] R. T. Rockafellar and S. Uryasev. Optimization of conditional Value-at-Risk. *J. Risk*, 2:493–517, 2000.

- [275] R. T. Rockafellar and S. Uryasev. Conditional Value-at-Risk for general loss distributions. *J. Banking & Finance*, 26(7):1443–1471, 2002.
- [276] R. T. Rockafellar and S. Uryasev. The fundamental risk quadrangle in risk management, optimization and statistical estimation. *Surveys in Operations Research and Management Science*, 18:33–53, 2013.
- [277] R. T. Rockafellar and S. Uryasev. Minimizing buffered probability of exceedance by progressive hedging. *Mathematical Programming*, 181(2):453–472, 2020.
- [278] R. T. Rockafellar, S. Uryasev, and M. Zabarankin. Deviation measures in risk analysis and optimization. Technical Report 2002-7, Department of Industrial and Systems Engineering, University of Florida, 2002.
- [279] R. T. Rockafellar, S. Uryasev, and M. Zabarankin. Generalized deviations in risk analysis. *Finance and Stochastics*, 10:51–74, 2006.
- [280] R. T. Rockafellar, S. Uryasev, and M. Zabarankin. Optimality conditions in portfolio analysis with general deviation measures. *Mathematical Programming B*, 108(2):515–540, 2006.
- [281] R. T. Rockafellar, S. Uryasev, and M. Zabarankin. Risk tuning with generalized linear regression. *Mathematics of Operations Research*, 33(3):712–729, 2008.
- [282] R. T. Rockafellar and R. J-B Wets. *Variational Analysis*, volume 317 of *Grundlehren der Mathematischen Wissenschaft*. Springer, 3rd printing-2009 edition, 1998.
- [283] ROL. <https://www.sandia.gov/ccr/software/rapid-optimization-library-rol/>, accessed, November 10, 2022.
- [284] J. O. Royset. Consistent approximations in composite optimization. *Mathematical Programming*, 201:339–372, 2023.
- [285] J. O. Royset, L. Bonfiglio, G. Vernengo, and S. Brizzolara. Risk-adaptive set-based design and applications to shaping a hydrofoil. *ASME J. Mechanical Design*, 139(10):1014031–1014038, 2017.
- [286] J. O. Royset and J.-E. Byun. Gradients and subgradients of buffered failure probability. *Operations Research Letters*, 49(6):868–873, 2021.
- [287] J. O. Royset and R. J-B Wets. Variational theory for optimization under stochastic ambiguity. *SIAM J. Optimization*, 27(2):1118–1149, 2017.
- [288] J. O. Royset and R. J-B Wets. *An Optimization Primer*. Springer, 2021.
- [289] A. Ruszczyński. Risk-averse dynamic programming for Markov decision processes. *Mathematical Programming*, 125:235–261, 2010.

- [290] A. Ruszczyński. Advances in risk-averse optimization. In H. Topaloglu, editor, *Tutorials in Operations Research: Theory Driven by Influential Applications*, Tutorials in Operations Research, pages 168–190. INFORMS, Cantonville, MD, 2013.
- [291] A. Ruszczyński. Erratum to: Risk-averse dynamic programming for Markov decision processes. *Mathematical Programming*, 145:601–604, 2014.
- [292] A. Ruszczyński. Convergence of a stochastic subgradient method with averaging for nonsmooth nonconvex constrained optimization. *Optimization Letters*, 14:1615–1625, 2020.
- [293] A. Ruszczyński and A. Shapiro. Conditional risk mappings. *Mathematics of Operations Research*, 31(3):544–561, 2006.
- [294] A. Ruszczyński and A. Shapiro. Optimization of convex risk functions. *Mathematics of Operations Research*, 31:433–452, 2006.
- [295] A. Ruszczyński and A. Shapiro. Optimization of risk measures. In F. Dabbene G. Calafiore, editor, *Probabilistic and Randomized Methods for Design Under Uncertainty*, pages 119–157. Springer, London, 2006.
- [296] S. Sarykalin, G. Serraino, and S. Uryasev. Value-at-Risk vs. Conditional Value-at-Risk in risk management and optimization. In Z.-L. Chen and S. Raghavan, editors, *Tutorials in Operations Research: State-of-the-Art Decision-Making Tools in the Information-Intensive Age*, Tutorials in Operations Research, pages 270–294. INFORMS, Cantonville, MD, 2008.
- [297] O. Scaillet. Nonparametric estimation of conditional expected shortfall. *Insurance and Risk Management Journal*, 72:639–660, 2005.
- [298] R. Schultz and S. Tiedemann. Conditional value-at-risk in stochastic programs with mixed-integer recourse. *Mathematical programming*, 105(2):365–386, 2006.
- [299] S. Seguin, S. E. Fleten, P. Cote, A. Pichler, and C. Audet. Stochastic short-term hydropower planning with inflow scenario trees. *European J. Operational Research*, 259(3):1156–1168, 2017.
- [300] D. Shang, V. Kuzmenko, and S. Uryasev. Cash flow matching with risks controlled by buffered probability of exceedance and conditional value-at-risk. *Annals of Operations Research*, 260(1):501–514, 2018.
- [301] S. Shao, A. Gupta, and W. B. Haskell. Robustness to modeling errors in risk-sensitive Markov decision problems with Markov risk measures. *Preprint arXiv:2209.12937*, 2022.
- [302] A. Shapiro. Consistency of sample estimates of risk averse stochastic programs. *J. Applied Probability*, 50:533–541, 2013.
- [303] A. Shapiro. On Kusuoka representation of law invariant risk measures. *Mathematics of Operations Research*, 38(1):142–152, 2013.



- [304] A. Shapiro. Distributionally robust stochastic programming. *SIAM J. Optimization*, 27(4):2258–2275, 2017.
- [305] A. Shapiro. Tutorial on risk neutral, distributionally robust and risk averse multistage stochastic programming. *European J. Operational Research*, 188(1):1–13, 2021.
- [306] A. Shapiro, D. Dentcheva, and A. Ruszczyński. *Lectures on Stochastic Programming: Modeling and Theory*. SIAM, 1. edition, 2009.
- [307] A. Shapiro, D. Dentcheva, and A. Ruszczyński. *Lectures on Stochastic Programming: Modeling and Theory*. SIAM, 3. edition, 2021.
- [308] A. Shapiro, W. Tekaya, J. P. da Costa, and M. P. Soares. Risk neutral and risk averse stochastic dual dynamic programming method. *European J. Operational Research*, 224(2):375–391, 2013.
- [309] R. Singh, Q. Zhang, and Y. Chen. Improving robustness via risk averse distributional reinforcement learning. In A. Bayen et al., editor, *In Learning for Dynamics and Control, Proceeding of Machine Learning Research*, volume 120, page 958–968, 2020.
- [310] Z. J. Smith and J. E. Bickel. Weighted scoring rules and convex risk measures. *Operations Research*, 70(6):3371–3385, 2022.
- [311] T. Soma and Y. Yoshida. Statistical learning with conditional value at risk. *Preprint arXiv:2002.05826*, 2020.
- [312] M. G. Speranza. Linear programming models for portfolio optimization. *Finance*, 14:1437–1446, 1993.
- [313] SPQR. [yassine-laguel.github.io/spqr/](https://yassine-laguel.github.io/spqr/), accessed, November 10, 2022.
- [314] SQwash. <https://krishnap25.github.io/sqwash/>, accessed, November 10, 2022.
- [315] I. Steinwart, C. Pasin, R. C. Williamson, and S. Zhang. Elicitation and identification of properties. In *Annual Conference Computational Learning Theory*, 2014.
- [316] J. Sun, X. Yang, Q. Yao, and M. Zhang. Risk minimization, regret minimization and progressive hedging algorithms. *Mathematical Programming*, 181:509–530, 2020.
- [317] TailRisk. <https://github.com/open-risk/tailRisk>, accessed, November 10, 2022.
- [318] A. Takeda and M. Sugiyama.  $\nu$ -support vector machine as conditional value-at-risk minimization. In *Proceedings of the 25th International Conference on Machine Learning (ICML 2008)*, 2008.
- [319] M. D. Teter, J. O. Royset, and A. M. Newman. Modeling uncertainty of expert elicitation for use in risk-based optimization. *Annals of Operations Research*, 280(1-2):189–210, 2019.

- [320] A. Thelen, X. Zhang, O. Fink, Y. Lu, S. Ghosh, B. D. Youn, M. D. Todd, S. Mahadevan, C. Hu, and Z. Hu. A comprehensive review of digital twin – part 1: Modeling and twinning enabling technologies. *Preprint arXiv:2208.14197*, 2022.
- [321] A. Thelen, X. Zhang, O. Fink, Y. Lu, S. Ghosh, B. D. Youn, M. D. Todd, S. Mahadevan, C. Hu, and Z. Hu. A comprehensive review of digital twin – part 2: Roles of uncertainty quantification and optimization, a battery digital twin, and perspectives. *Preprint arXiv:2208.12904*, 2022.
- [322] P. Thomas and E. Learned-Miller. Concentration inequalities for conditional value at risk. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 6225–6233. PMLR, 09–15 Jun 2019.
- [323] I. Toumazis and C. Kwon. Worst-case conditional value-at-risk minimization for hazardous materials transportation. *Transportation Science*, 50:1174–1187, 2016.
- [324] A. Trindade, S. Uryasev, A. Shapiro, and G. Zrazhevsky. Financial prediction with constrained tail risk. *Journal of Banking and Finance*, 31(11):3524–3538, 2007.
- [325] N. A. Urpi, S. Curi, and A. Krause. Risk-averse offline reinforcement learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- [326] E. R. van Beesten and W. Romeijnders. Convex approximations for two-stage mixed-integer mean-risk recourse models with conditional value-at-risk. *Mathematical Programming*, 181:473–507, 2020.
- [327] B. P. G. van Parys, P. Mohajerin Esfahani, and D. Kuhn. From data to decisions: Distributionally robust optimization is optimal. *Management Science*, 6(6):3387–3402, 2021.
- [328] F. A. C. Viana, C. Gogu, and T. Goel. Surrogate modeling: tricks that endured the test of time and some recent developments. *Structural and Multidisciplinary Optimization*, 64:2881–2908, 2021.
- [329] N. Vijayan and L. A. Prashanth. Risk-sensitive reinforcement learning via distortion risk measures. *Preprint arXiv:2107.04422*, 2021.
- [330] J. von Neumann and O. Morgenstern. *Theory of games and economic behavior*. Princeton University Press, Princeton, NJ, 1944.
- [331] R. Wang, Y. Wei, and G. E. Willmot. Characterization, robustness, and aggregation of signed Choquet integrals. *Mathematics of Operations Research*, 45(3):993–1015, 2020.
- [332] R. Wang and R. Zitikis. An axiomatic foundation for the expected shortfall. *Management Science*, 67(3):1413–1429, 2020.

- [333] W. Wang and H. Xu. Robust spectral risk optimization when information on risk spectrum is incomplete. *SIAM J. Optimization*, 30(4):3198–3229, 2020.
- [334] Y. Wang and M. P. Chapman. Risk-averse autonomous systems: A brief history and recent developments from the perspective of optimal control. *Artificial Intelligence*, 311:103743, 2022.
- [335] Y. Wang and F. Gao. Deviation inequalities for an estimator of the conditional value-at-risk. *Operations Research Letters*, 38(3):236–239, 2010.
- [336] R. J-B Wets. Stochastic programs with fixed recourse: The equivalent deterministic program. *SIAM Review*, 16(3):309–339, 1974.
- [337] W. Wiesemann, D. Kuhn, and M. Sim. Distributionally robust convex optimization. *Operations Research*, 62(6):1358–1376, 2014.
- [338] J. Willard, X. Jia, S. Xu, M. Steinbach, and V. Kumar. Integrating scientific knowledge with machine learning for engineering and environmental systems. *ACM Computing Surveys (CSUR)*, 2022.
- [339] R. C. Williamson and A. K. Menon. Fairness risk measures. In *Proceedings of the 36th International Conference on Machine Learning, Long Beach, California*, 2019.
- [340] D. Wozabal and N. Wozabal. Asymptotic consistency of risk functionals. *J. Nonparametric Statistics*, 21(8):977–990, 2009.
- [341] D. Wu, H. Zhu, and E. Zhou. A Bayesian risk approach to data-driven stochastic optimization: Formulations and asymptotics. *SIAM J. Optimization*, 28(2):1588–1612, 2018.
- [342] L. Xia, L. Zhang, and P. W. Glynn. Risk-sensitive Markov decision processes with long-run CVaR criterion. *Productions and Operations Management*, to appear, 2023.
- [343] A. Xuan, X. Shen, Q. Guo, and H. Sun. A conditional value-at-risk based planning model for integrated energy system with energy storage and renewables. *Applied Energy*, 294:116971, 2021.
- [344] M. E. Yaari. The dual theory of choice under risk. *Econometrica*, 55(1):95–115, 1987.
- [345] D. Yang. Robust machine learning using superquantiles. Master’s thesis, Naval Postgraduate School, 2021.
- [346] M. Zabaranin and S. Uryasev. *Statistical Decision Problems. Selected Concepts and Portfolio Safeguard Case Studies*. Springer, New York, NY, 2014.
- [347] M. B. Zafar, I. Valera, M. G. Rogriguez, and K. P. Gummadi. Fairness constraints: Mechanisms for fair classification. In *Artificial Intelligence and Statistics*, pages 962–970. PMLR, 2017.

- [348] D. Zhang, S. W. Wallace, Z. Guo, Y. Dong, and M. Kaut. On scenario construction for stochastic shortest path problems in real road networks. *Transportation Research Part E: Logistics and Transportation Review*, 152:102410, 2021.
- [349] J. Zhang, H. Lin, S. Jegelka, S. Sra, and A. Jadbabaie. Complexity of finding stationary points of nonconvex nonsmooth functions. In *International Conference on Machine Learning*, pages 11173–11182, 2020.
- [350] T. Zhang, S. Uryasev, and Y. Guan. Derivatives and subderivatives of buffered probability of exceedance. *Operations Research Letters*, 47:130–132, 2019.
- [351] J. Zhen, D. Kuhn, and W. Wiesemann. Mathematical foundations of robust and distributionally robust optimization. *Preprint arXiv:2105.00760*, 2021.
- [352] J. Ziegel. Coherence and elicibility. *Mathematical Finance*, 26:901–918, 2016.
- [353] Z. Zou, D. P. Kouri, and W. Aquino. An adaptive local reduced basis method for solving PDEs with uncertain inputs and evaluating risk. *Computer Methods in Applied Mechanics and Engineering*, 345:302–322, 2019.
- [354] Z. Zou, D. P. Kouri, and W. Aquino. A locally adapted reduced-basis method for solving risk-averse PDE-constrained optimization problems. *SIAM/ASA Journal on Uncertainty Quantification*, 10(4):1629–1651, 2022.
- [355] G. M. Zrazhevsky, A. N. Golodnikov, S. P. Uryasev, and A. G. Zrazhevsky. Application of buffered probability of exceedance in reliability optimization problems. *Cybernetics and Systems Analysis*, 56(3):476–484, 2020.