

---

# Robustness and scalability under heavy tails, without strong convexity

---

Matthew J. Holland

Institute of Scientific and Industrial Research  
Osaka University

## Abstract

Real-world data is laden with outlying values. The challenge for machine learning is that the learner typically has no prior knowledge of whether the feedback it receives (losses, gradients, etc.) will be heavy-tailed or not. In this work, we study a simple, cost-efficient algorithmic strategy that can be leveraged when both losses and gradients can be heavy-tailed. The core technique introduces a simple robust validation sub-routine, which is used to boost the confidence of inexpensive gradient-based sub-processes. Compared with recent robust gradient descent methods from the literature, dimension dependence (both risk bounds and cost) is substantially improved, without relying upon strong convexity or expensive per-step robustification. We also empirically show that the proposed procedure cannot simply be replaced with naive cross-validation.

## 1 INTRODUCTION

Uncertainty is inherent in real world physical and social systems. This implies that machine learning methods, driven by data generated from these systems, are inherently uncertain. Coupled with our lack of knowledge regarding the mechanisms underlying these systems, it is not possible to provide exact, certain statements regarding algorithm performance. This uncertainty is manifested clearly in the “risk minimization” formulation of learning problems: over some set of candidates  $\mathcal{W} \subseteq \mathbb{R}^d$ , the *risk* is  $R_{\mathbf{P}}(w) := \mathbf{E}_{\mathbf{P}} L(w; Z)$ , namely the expected value of a loss  $L : \mathcal{W} \times \mathcal{Z} \rightarrow \mathbb{R}_+$  evaluated at  $w$ , where  $Z \sim \mathbf{P}$  denotes our random data, taking values in a set  $\mathcal{Z}$ . In the context of machine

learning (by empirical risk minimization), this notion was popularized in the statistical community by the work of Vapnik (1982), and in the computer science community by the work of Haussler (1992). A learning algorithm will have access to  $n$  data points sampled from  $\mathbf{P}$ , denoted  $Z_1, \dots, Z_n$ . Write  $(Z_1, \dots, Z_n) \mapsto \hat{w}_n$  to denote the output of an arbitrary learning algorithm. The usual starting point for analyzing algorithm performance is the *estimation error*  $R_{\mathbf{P}}(\hat{w}_n) - R_{\mathbf{P}}^*$ , where  $R_{\mathbf{P}}^* := \inf\{R_{\mathbf{P}}(w) : w \in \mathcal{W}\}$ , or more precisely, the distribution of this error. Since we never know much about the underlying data-generating process, typically all we can assume is that  $\mathbf{P}$  belongs to some class  $\mathcal{P}$  of probability measures on  $\mathcal{Z}$ , and typical guarantees are given in the form of  $\mathbf{P}\{R_{\mathbf{P}}(\hat{w}_n) - R_{\mathbf{P}}^* > \varepsilon(n, \delta, \mathbf{P}, \mathcal{W})\} \leq \delta$ , over a class  $\mathbf{P} \in \mathcal{P}$ . In words, the procedure yielding  $\hat{w}_n$  will obtain  $\varepsilon$ -good performance with  $(1 - \delta)$ -high confidence over the draw of the sample. Citing Holland (2021), in order to have meaningful performance guarantees, the following properties are important.

1. **Transparency:** can we actually compute the output  $\hat{w}_n$  that we study in theory?
2. **Strength:** what form do bounds on  $\varepsilon(n, \delta, \mathbf{P}, \mathcal{W})$  take? How rich is the class  $\mathcal{P}$ ?
3. **Scalability:** how do computational costs scale with the above-mentioned factors?

Balancing these three points is critical to developing guarantees for *algorithms that will actually be used in practice*.

**Our problem setting** This work considers the setup of *potentially heavy-tailed* data, and in contrast with the recent work of Holland (2021), we only assume a *convex* loss, rather than strong convexity. All the learner can know is that for some  $m < \infty$ ,

$$\mathcal{P} \subseteq \left\{ \mathbf{P} : \sup_{w \in \mathcal{W}} \mathbf{E}_{\mathbf{P}} |L(w; Z)|^m < \infty \right\}, \quad (1)$$

where typically  $m = 2$ . Thus, it is unknown whether the losses (or partial derivatives, etc.) are congenial in

a sub-Gaussian sense (where (1) holds for all  $m$ ), or heavy-tailed in the sense that all higher-order moments could be infinite or undefined. The goal then comes down to obtaining the strongest possible guarantees for a tractable learning algorithm, given (1). We next review the related technical literature, and give an overview of our contributions.

## 2 CONTEXT AND CONTRIBUTIONS

**Challenges without strong convexity** When one is lucky enough to have a  $\mu$ -strongly convex risk  $R_P$ , using a very simple basic idea, a wide range of *distance-based* algorithmic strategies are available (Minsker, 2015; Hsu and Sabato, 2016; Holland, 2021). For example, say we construct  $k$  candidates  $\hat{w}^{(1)}, \dots, \hat{w}^{(k)}$ , and we know that with high probability, a majority of the candidates are  $\varepsilon$ -good in terms of the risk  $R_P$ . Since  $R_P$  is unknown, we can never know which candidates are the  $\varepsilon$ -good ones. However, this barrier can be circumvented by utilizing the fact that  $\mu$ -strong convexity of  $R_P$  implies that any  $\varepsilon$ -good candidate must be at least  $\sqrt{2\varepsilon/\mu}$ -close to  $w^*$ , the minimizer of  $R_P$  on  $\mathcal{W}$ . It follows that on the “good event” in which the majority of candidates are  $\varepsilon$ -good, it is sufficient to simply “follow the majority.” This can be done in various ways, but in the end all such procedures comes down to computing and comparing distances  $\|w - \hat{w}^{(j)}\|$  for all  $j \in [k]$ . This can be done without knowing which of the  $\hat{w}^{(j)}$  are  $\varepsilon$ -good, which makes the problem tractable.

Unfortunately,  $\mu$ -strong convexity is a luxury that is often unavailable. In particular for high-dimensional settings, it is common for the strong convexity parameter  $\mu$  to shrink rapidly as  $d$  grows, making  $1/\mu$ -dependent error bounds vacuous (Bach and Moulines, 2014). Algorithmically, if strong convexity cannot be guaranteed, then the distance-based strategy just described will fail, since for any particular minimizer  $w^*$ , it is perfectly plausible to have a  $\varepsilon$ -good candidate which is arbitrarily far from  $w^*$ . Even when we assume  $\lambda_1$ -smoothness of the risk, all we can say is that  $\varepsilon$ -badness implies  $\sqrt{2\varepsilon/\lambda_1}$ -farness from all minimizers; the converse need not hold. The traditional approach: if from sample  $\mathbf{Z}_n$  we obtain independent candidates  $\hat{w}^{(1)}, \dots, \hat{w}^{(k)}$ , and we have a second sample  $\mathbf{Z}'_n = (Z'_1, \dots, Z'_n)$  available for “validation,” then classical procedures return

$$\hat{w}_n = \hat{w}^{(\star)}, \text{ where}$$

$$\star \in \arg \min \left\{ \frac{1}{n} \sum_{i=1}^n L(\hat{w}^{(j)}; Z'_i) : j = 1, \dots, k \right\}. \quad (2)$$

This technique of confidence boosting for bounded losses is well-known; see (Kearns and Vazirani, 1994,

Ch. 4.2) or (Shalev-Shwartz et al., 2010, Thm. 26). Under exp-concave distributions, Mehta (2016) also recently made use of this technique. Problems arise, however, when the losses can be potentially heavy-tailed. The quality of the validated final candidate is only as good as the precision of the risk estimate, and the empirical risk is well-known to be sub-optimal under potentially heavy-tailed data (Devroye et al., 2016).

**Limitations of existing procedures** To begin, empirical risk minimization (ERM) is the cornerstone of classical learning theory, which studies the statistical properties of any minimizer of the empirical risk, i.e., the sample mean of the losses. Concrete implementations of ERM just require minimizing a finite sum, and thus are computationally quite congenial, and scale well, taken at face value. However, formal guarantees for ERM-based procedures are limited; the empirical mean is known to be sensitive to outliers, and this sensitivity appears in weak formal guarantees. Concretely, under potentially heavy-tailed losses, the empirical mean is sub-optimal in that it cannot guarantee sub-Gaussian deviation bounds. Put roughly, it cannot guarantee error bounds better than those which scale as  $\Omega(1/\delta)$ ; see Catoni (2012) and Devroye et al. (2016) for more details. Furthermore, since ERM leaves the method of implementation completely abstract, this leaves open a large conceptual gap. Feldman (2017) showed lucidly how there exist both “good” and “bad” ERM solutions; the problem with transparency is that we can never know whether any particular ERM candidate is one of the good ones or not. In contrast, starting with seminal work by Brownlees et al. (2015), a recent line of work has led to new statistical learning procedures to address the weak guarantees and lack of robustness of ERM. The basic idea is simply to minimize a different estimator of the risk, for example median-of-means estimators (Minsker, 2018) or M-estimators (Brownlees et al., 2015). Under weak moment bounds like (1), their minimizers enjoy  $\mathcal{O}(1/\sqrt{n})$  rates with  $\mathcal{O}(\log(\delta^{-1}))$  dependence on the confidence. This provides a significant improvement in terms of the strength of guarantees over ERM, but unfortunately the issue of transparency remains. Like ERM, the algorithmic side of the problem is left abstract here, and in general may even be a much more difficult computational task. As such, the gap between formal guarantees and the guarantees that hold for any given output of the algorithm may be even more severe than in the case of ERM.

Furthermore, several new families of algorithms have been designed in the past few years to tackle the potentially heavy-tailed setting using a tractable procedure. Such algorithms may naturally be called *robust gradient descent* (RGD), since the core update takes the

same form as steepest-descent applied directly to the risk, i.e.,  $\hat{w}_{t+1} \leftarrow \hat{w}_t - \alpha \nabla R_{\mathbb{P}}(w_t)$ , except that  $\nabla R_{\mathbb{P}}$  is replaced with an estimator  $\hat{G}_n(w) \approx \nabla R_{\mathbb{P}}(w)$  which has deviations with near-optimal confidence intervals under potentially heavy-tailed data (Chen et al., 2017a; Prasad et al., 2018; Lecué et al., 2018; Holland, 2019; Holland and Ikeda, 2019a,b). Under strong convexity, all these proposals enjoy excess risk bounds with optimal dependence on  $n$  and  $1/\delta$  under potentially heavy-tailed data, and can be implemented as-is. Unfortunately, instead of a simple one-dimensional robust mean estimate as in Brownlees et al. (2015), all RGD methods rely on sub-routines that work in  $d$ -dimensions. This makes the procedures much more expensive computationally for “big” learning tasks, and leads to an undesirable dependence on the ambient dimension  $d$  in the statistical guarantees as well. Moreover, when strong convexity is not available, the propagation of statistical error over time for RGD methods becomes much worse, leading to bounds that are extremely sensitive to mis-specified smoothness parameters, step sizes, and total number of iterations.

**Our contributions** Considering the above limitations of existing techniques, notably the lack of scalability and weak formal guarantees available for RGD methods under weak convexity, here we study a different algorithmic approach to the problem. Our approach has equal generality, and the hope is to achieve as-good or better dependence on  $n$ ,  $d$ , and  $1/\delta$  under potentially heavy-tailed data (losses and/or partial derivatives), without strong convexity, and in provably less time for larger problems. The main technique that we investigate is a natural robustification of classical confidence boosting (2), applied to traditional stochastic gradient descent routines run in parallel, though we note that the basic argument can be easily generalized to other optimization strategies (e.g., accelerated methods, adaptive methods, quasi-Newton techniques, etc.). Our main contributions:

- We study a general-purpose robust learning procedure (Algorithm 1), obtaining sharp risk bounds (Theorem 1) that improve on the poor dependence of RGD methods on the dimension and number of iterations under weak convexity (see Table 1).
- The archetype given in Algorithm 1 is concrete, easy to implement as-is, and trivial to run in parallel. All else equal, for high-dimensional learning tasks we can expect to obtain a result as good or better than existing serial RGD methods in far less time.
- Empirically, we study the robustness of our proposed procedure and relevant competitors to vari-

ous perturbations in the experimental conditions, simulating a lack of prior knowledge about noise levels and convexity. The proposed procedure can easily match benchmark RGD methods in less time, over a variety of test settings. We also verify that a naive cross-validation heuristic does not achieve the same level of performance.

## 3 THEORETICAL ANALYSIS

### 3.1 Preliminaries

**Notation** First we establish some basic notation, and organize numerous technical assumptions in one place for ease of reference. For any positive integer  $k$ , write  $[k] := \{1, \dots, k\}$ . For any index  $\mathcal{I} \subseteq [n]$ , write  $\mathbf{Z}_{\mathcal{I}} := (Z_i)_{i \in \mathcal{I}}$ , defined analogously for independent copy  $\mathbf{Z}'_{\mathcal{I}}$ . To keep the notation simple, in the special case of  $\mathcal{I} = [n]$ , we write  $\mathbf{Z}_n := \mathbf{Z}_{[n]} = (Z_1, \dots, Z_n)$ . We shall use  $\mathbf{P}$  as a generic symbol to denote computing probability; in most cases this will be the product measure induced by the sample  $\mathbf{Z}_n$  or  $\mathbf{Z}'_n$ . For any function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , denote by  $\partial f(u)$  the sub-differential of  $f$  evaluated at  $u$ . Variance of the loss is denoted by  $\sigma_{\mathbb{P}}^2(w) := \text{var}_{\mathbb{P}} L(w; Z) = \mathbf{E}_{\mathbb{P}}(L(w; Z) - R_{\mathbb{P}}(w))^2$  for each  $w \in \mathcal{W}$ . Denote by  $I\{\text{event}\}$  the indicator function that returns a value of 1 when **event** is true, and 0 otherwise.

**Running assumptions** The two key running assumptions that we make are related to independence and convexity. First, we assume that all the observed data are independent, i.e., the random variables  $Z_i$  and  $Z'_i$  taken over all  $i \in [n]$  are independent copies of  $Z \sim \mathbb{P}$ . Second, for each  $z \in \mathcal{Z}$ , we assume the map  $w \mapsto L(w; z)$  is a real-valued convex function over  $\mathbb{R}^d$ , and that the parameter set  $\mathcal{W} \subseteq \mathbb{R}^d$  is non-empty, convex, and compact, with diameter  $\Delta := \sup\{\|u - v\| : u, v \in \mathcal{W}\}$ . All results derived in the next sub-section will be for an arbitrary choice of  $\mathbb{P} \in \mathcal{P}$ , where  $\mathcal{P}$  satisfies (1) with  $m = 2$ . We say a function  $f$  is  $\lambda_1$ -smooth if its gradient is  $\lambda_1$ -Lipschitz continuous. Finally, to make formal statements technically simpler, we assume that  $R_{\mathbb{P}}(\cdot)$  achieves its minimum on the interior of  $\mathcal{W}$ .

### 3.2 Error bounds when both losses and gradients can be heavy-tailed

Recalling the challenges described in section 2, we consider a straightforward robustification of the classical validation-based approach using robust mean estimators in a sub-routine. The full procedure is summarized in Algorithm 1. First we outline the key points of the procedure, then give a general-purpose excess risk bound in Theorem 1.

Method	Error	Cost
RV-SGDave	$\mathcal{O}\left(\sqrt{\frac{\log(\delta^{-1})}{n}}(\sigma_{L,P} + \sigma_{G,P})\right) + \mathcal{O}\left(\frac{\lambda_1 \log(\delta^{-1})}{n}\right)$	$\mathcal{O}(dn \log(\delta^{-1}))$
RGD (Chen et al., 2017b)	$\mathcal{O}\left((1 + \lambda_1 \alpha)^T \sqrt{\frac{d(\sigma_{G,P}^2 \log(d\delta^{-1}) + \log(n))}{nT}}\right)$	$\mathcal{O}(Tdn \log(\delta^{-1}))$

Table 1: High-probability error bounds and computational cost estimates for RV-SGDave (Algorithm 1), compared with modern RGD methods, *without* assuming strong convexity. Error denotes confidence intervals for  $R_P(\hat{w}_n) - R_P^*$  with  $\hat{w}$  being the output of each procedure after  $T$  steps (noting RV-SGDave has  $T = n$  by definition). Detailed explanation is in supplementary materials.

**Algorithm 1** Divide-and-conquer with robust validation; RV-SGDave  $[\mathbf{Z}_n, \mathbf{Z}'_n, \hat{w}_0; k]$ .

**inputs:** samples  $\mathbf{Z}_n$  and  $\mathbf{Z}'_n$ , initial value  $\hat{w}_0 \in \mathcal{W}$ , parameter  $1 \leq k \leq n$ .

Split  $\bigcup_{j=1}^k \mathcal{I}_j = [n]$ , with  $|\mathcal{I}_j| \geq \lfloor n/k \rfloor$ , and  $\mathcal{I}_j \cap \mathcal{I}_l = \emptyset$  when  $j \neq l$ .

For each  $j \in [k]$ , set  $\bar{w}^{(j)}$  to the mean of the sequence  $\text{SGD}[\hat{w}_0; \mathcal{Z}_{\mathcal{I}_j}, \mathcal{W}]$ .

Compute  $\star = \arg \min_{j \in [k]} \text{Valid}[\bar{w}^{(j)}; \mathbf{Z}'_n]$ .

**return:**  $\hat{w}_{\text{RV}} = \bar{w}^{(\star)}$ .

**Core procedure** Viewed at a high level, Algorithm 1 is comprised of three extremely simple steps: partition, train, and validate. For our purposes, the key to improving on traditional ERM-style boosting techniques is to ensure the validation step is done with sufficient precision, even when the losses can be heavy-tailed. To achieve this, we shall require that there exist a constant  $c > 0$  which does not depend on the distribution  $\mathbf{P}$ , such that for any choice of confidence level  $\delta \in (0, 1)$  and large enough  $n$ , the sub-routine **Valid** satisfies

$$|\text{Valid}[w; \mathbf{Z}'_n] - R_P(w)| \leq c \sqrt{\frac{(1 + \log(\delta^{-1}))\sigma_P^2(w)}{n}} \quad (3)$$

with probability no less than  $1 - \delta$ . Recall that we are denoting  $\sigma_P^2(w) := \text{var}_P L(w; Z)$ , thus the only requirement on the class of data distributions is finite variance, readily allowing for both heavy-tailed losses and gradients. The training step can be done in any number of ways; for concreteness and clarity of the results, we elect to use a simple stochastic gradient descent sub-process. Unpacking the notation from Algorithm 1, the basic update used is traditional projected (sub-

)gradient descent, with update denoted by

$$\text{SGD}[w; Z, \alpha, \mathcal{W}] := \Pi_{\mathcal{W}}(w - \alpha G(w; Z)). \quad (4)$$

Here  $\alpha \geq 0$  denotes a step-size parameter,  $\Pi_{\mathcal{W}}$  denotes projection to  $\mathcal{W}$  with respect to the  $\ell_2$  norm, and the standard assumption is that the random vector  $G(w; Z)$  satisfies  $\mathbf{E}_P G(w; Z) \in \partial R_P(w)$ , for each  $w \in \mathcal{W}$ . Then for arbitrary sequence  $(Z_1, \dots, Z_m)$ , we define

$$\begin{aligned} \text{SGD}[\hat{w}_0; (Z_1, \dots, Z_m), \mathcal{W}] := \\ \{ \text{SGD}[\hat{w}_t; Z_{t+1}, \alpha_t, \mathcal{W}] : t = 0, 1, \dots, m-1 \}, \end{aligned}$$

noting we have suppressed step-sizes from the notation for readability. Replacing the generic sequence  $(Z_1, \dots, Z_m)$  here with  $\mathbf{Z}_{\mathcal{I}_j}$  for each  $j \in [k]$  yields the iterate sequences used in Algorithm 1, with each  $\bar{w}^{(j)}$  being simply the arithmetic (vector) mean of the iterates. As all the data we work with are independent copies of  $Z \sim P$ , the order in which we take the indices from each  $\mathcal{I}_j$  does not matter. Under weak assumptions on the underlying loss distribution, the output  $\hat{w}_{\text{RV}}$  of this algorithm enjoys strong excess risk bounds, as the following theorem shows.

**Theorem 1.** *Let  $R_P$  be  $\lambda_1$ -smooth in the  $\ell_2$  norm,  $\mathbf{E}_P \|G(w; Z) - \nabla R_P(w)\|_2^2 \leq \sigma_{G,P}^2 < \infty$ , and  $\sigma_P^2(w) \leq \sigma_{L,P}^2 < \infty$  for all  $w \in \mathcal{W}$ . Run Algorithm 1 with sub-routine **Valid** satisfying (3), given a total sample size  $n \geq 2k$  split into  $\mathbf{Z}_{n/2}$  and  $\mathbf{Z}'_{n/2}$ , and SGD sub-processes using step sizes  $\alpha_t = 1/(\lambda_1 + (1/a))$ , where  $a = \Delta/\sqrt{n\sigma_{G,P}^2/2k}$ . If we set  $k = \lceil \log(2\lceil \log(\delta^{-1}) \rceil \delta^{-1}) \rceil$ , then for any confidence parameter  $0 < \delta \leq 1/3$ , we have*

$$\begin{aligned} R_P(\hat{w}_{\text{RV}}) - R_P^* \leq \\ 2c \sqrt{\frac{2(1 + \log(2\lceil \log(\delta^{-1}) \rceil \delta^{-1}))\sigma_{L,P}^2}{n}} \\ + 3 \left( \frac{k\Delta^2\lambda_1}{n} + \sqrt{\frac{2k\Delta^2\sigma_{G,P}^2}{n}} \right) \end{aligned}$$

with probability no less than  $1 - 3\delta$ .

---

**Algorithm 2** Catoni type M-estimate.
 

---

**inputs:** sample  $\{u_1, \dots, u_n\}$ , parameters  $\sigma > 0$  and  $0 < \delta < 1$ .

Set  $q^2 = \frac{2\sigma^2 \log(2\delta^{-1})}{n - 2\log(2\delta^{-1})}$  and  $s^2 = \frac{n(\sigma^2 + q^2)}{2\log(2\delta^{-1})}$ .

**return:**  $\arg \min_{\theta \in \mathbb{R}} \sum_{i=1}^n \rho\left(\frac{u_i - \theta}{s}\right)$ .

---

We shall look at concrete implementations of `Valid` in section 4, but as an initial example, setting `Valid` to be a properly scaled M-estimator (sub-routine given in Algorithm 2) satisfies (3) with  $c \leq 2$  whenever  $n \geq 4\log(\delta^{-1})$ . Additional examples and supporting lemmas are given in the supplementary materials.

**Proof sketch** Here we give an overview of the proof of Theorem 1. We have data sequences  $\mathbf{Z}_n$  and  $\mathbf{Z}'_n$ . The former is used to obtain independent candidates  $\bar{w}^{(1)}, \dots, \bar{w}^{(k)}$ , and the latter is used to select among these candidates. As mentioned earlier, distance-based strategies (Minsker, 2015; Hsu and Sabato, 2016) require that the *majority* of these candidates are  $\varepsilon$ -good, in order to ensure that points near the majority coincide with  $\varepsilon$ -good points. In our present setup, where strong convexity is not available, we are taking a very different approach. Now we only require that *at least one* of the candidates is  $\varepsilon$ -good. Making this explicit,

$$\mathcal{E}_1(\varepsilon; k) := \bigcup_{j=1}^k \left\{ R_{\mathbb{P}}(\bar{w}^{(j)}) - R_{\mathbb{P}}^* \leq \varepsilon \right\} \quad (5)$$

is our first event of interest. Note that *if* for each  $j \in [k]$  we have an upper bound  $\varepsilon_{\mathbb{P}}(\cdot)$  depending on the sample size for the sub-process outputs  $\bar{w}^{(j)}$  such that

$$\mathbf{E} \left[ R_{\mathbb{P}}(\bar{w}^{(j)}) - R_{\mathbb{P}}^* \right] \leq \varepsilon_{\mathbb{P}}(\lfloor n/k \rfloor),$$

where expectation is taken over the subset indexed by  $\mathcal{I}_j$ , then using Markov's inequality and taking a union bound, it follows that setting  $\varepsilon = e\varepsilon_{\mathbb{P}}$ , we have  $\mathbf{P} \mathcal{E}_1(e\varepsilon_{\mathbb{P}}; k) \geq 1 - e^{-k}$ . Asking for one  $\varepsilon$ -good candidate is a much weaker requirement than asking for the majority to be  $\varepsilon$ -good, but we must pay the price in a different form, as we require that `Valid` provide a good estimate of the true risk for *all* of the  $k$  candidates. In particular, writing  $b_{\mathbb{P}}(n, \delta)$  for a confidence interval to be specified shortly, this is the following event:

$$\mathcal{E}_2(\delta; k) := \bigcap_{j=1}^k \left\{ \left| \text{Valid}[\bar{w}^{(j)}; \mathbf{Z}'_n] - R_{\mathbb{P}}(\bar{w}^{(j)}) \right| \leq b_{\mathbb{P}}(n, \delta) \right\}.$$

Intuitively, while we only require that at least one of the  $k$  candidates be good, we must reliably know which is *best* at the available precision, which requires paying the price of the intersection defining  $\mathcal{E}_2(\delta; k)$ . Recalling the requirement (3), if we condition on  $\mathbf{Z}_n$ , the candidates  $\bar{w}^{(1)}, \dots, \bar{w}^{(k)}$  become non-random elements of  $\mathcal{W}$ , which means that setting  $b_{\mathbb{P}}(n, \delta) = c\sqrt{(1 + \log(\delta^{-1}))\sigma_{L, \mathbb{P}}^2/n}$ , a union bound gives us  $\mathbf{P}(\mathcal{E}_2(\delta; k); \mathbf{Z}_n) \geq 1 - k\delta$ . This inequality holds as-is for any realization of  $\mathbf{Z}_n$ , so we can thus integrate to obtain

$$\mathbf{P} \mathcal{E}_2(\delta; k) = \int \mathbf{P}(\mathcal{E}_2(\delta; k); \mathbf{Z}_n) \mathbf{P}(d\mathbf{Z}_n) \geq 1 - k\delta.$$

The good event of interest then has probability

$$\mathbf{P} \left[ \mathcal{E}_1 \left( e\varepsilon_{\mathbb{P}} \left( \lfloor \frac{n}{k} \rfloor \right); k \right) \cap \mathcal{E}_2(\delta; k) \right] \geq 1 - e^{-k} - k\delta.$$

On this good event, we know that there does exist an  $\varepsilon$ -good candidate, even though we can never know which it is; call it  $\bar{w}_{\text{LUCK}} \in \{\bar{w}^{(1)}, \dots, \bar{w}^{(k)}\}$ . Furthermore, even though this candidate is unknown, since we have  $b_{\mathbb{P}}(n, \delta)$ -good risk estimates for all  $k$  candidates, the choice of  $\bar{w}^{(\star)}$ , with  $\star = \arg \min_{j \in [k]} \text{Valid}[\bar{w}^{(j)}; \mathbf{Z}'_n]$ , cannot be much worse. More precisely, we have

$$\begin{aligned} R_{\mathbb{P}}(\bar{w}^{(\star)}) - R_{\mathbb{P}}^* &= R_{\mathbb{P}}(\bar{w}^{(\star)}) - \text{Valid}[\bar{w}^{(\star)}] + \text{Valid}[\bar{w}^{(\star)}] - R_{\mathbb{P}}^* \\ &\leq R_{\mathbb{P}}(\bar{w}^{(\star)}) - \text{Valid}[\bar{w}^{(\star)}] + \text{Valid}[\bar{w}_{\text{LUCK}}] - R_{\mathbb{P}}^* \\ &= \left[ R_{\mathbb{P}}(\bar{w}^{(\star)}) - \text{Valid}[\bar{w}^{(\star)}] \right] + \\ &\quad \left[ \text{Valid}[\bar{w}_{\text{LUCK}}] - R_{\mathbb{P}}(\bar{w}_{\text{LUCK}}) \right] + \left[ R_{\mathbb{P}}(\bar{w}_{\text{LUCK}}) - R_{\mathbb{P}}^* \right] \\ &\leq 2b_{\mathbb{P}}(n, \delta) + e\varepsilon_{\mathbb{P}}(\lfloor n/k \rfloor). \end{aligned}$$

We have effectively proved the following lemma.

**Lemma 2** (Boosting the confidence under potentially heavy tails). *Assume we have a learning algorithm `Learn` such that for  $n \geq 1$  and  $\delta \in (0, 1)$ , we have*

$$\mathbf{P} \left\{ R_{\mathbb{P}}(\text{Learn}[\mathbf{Z}_n]) - R_{\mathbb{P}}^* > \frac{\varepsilon_{\mathbb{P}}(n)}{\delta} \right\} \leq \delta.$$

*Splitting the data  $\mathbf{Z}_n$  using sub-indices  $\mathcal{I}_1, \dots, \mathcal{I}_k$ , if we set*

$$\star = \arg \min_{j \in [k]} \text{Valid}[\text{Learn}[\mathbf{Z}_{\mathcal{I}_j}]; \mathbf{Z}'_n],$$

*then when `Valid` satisfies (3), it follows that for any  $\delta \in (0, 1)$ , we have*

$$\begin{aligned} R_{\mathbb{P}}(\text{Learn}[\mathbf{Z}_{\mathcal{I}_{\star}}]) - R_{\mathbb{P}}^* &\leq \\ &\sup_{w \in \mathcal{W}} 2c\sqrt{\frac{(1 + \log(\delta^{-1}))\sigma_{\mathbb{P}}^2(w)}{n}} + e\varepsilon_{\mathbb{P}}\left(\left\lfloor \frac{n}{k} \right\rfloor\right) \end{aligned}$$

*with probability no less than  $1 - k\delta - e^{-k}$ .*

Note that Lemma 2 here makes no direct requirements on the underlying loss or risk, beyond the need for a variance bound, which appears as  $\sigma_{\mathbb{P}}^2(w) \leq \sigma_{L,\mathbb{P}}^2 < \infty$  in the statement of Theorem 1. Indeed, convexity does not even make an appearance. This is in stark contrast with distance-based confidence boosting methods, which rely upon the strong convexity of the risk (Minsker, 2015; Hsu and Sabato, 2016). As such, so long as we can validate in the sense of (3), then Lemma 2 gives us a general-purpose tool from which we can construct algorithms with competitive risk bounds under potentially heavy-tailed data. This can be done for many practical procedures, and the only main step that remains is to clean up the good event probability and specify  $k$  to achieve the properties stated in Theorem 1. Letting  $\delta$  be that given in the theorem statement, first set  $\delta_0 = \delta / (2 \lceil \log(\delta^{-1}) \rceil) < \delta$ . Next, set the number of subsets to be

$$k = \lceil \log(1/\delta_0) \rceil = \lceil \log(2 \lceil \log(\delta^{-1}) \rceil \delta^{-1}) \rceil,$$

and note that with this setting of  $k$  and  $\delta_0$ , we have that

$$\begin{aligned} 1 - k\delta_0 &= 1 - \lceil \log(2 \lceil \log(\delta^{-1}) \rceil \delta^{-1}) \rceil \left( \frac{\delta}{2 \lceil \log(\delta^{-1}) \rceil} \right) \\ &\geq 1 - \left( \frac{\lceil \log(2) \rceil}{\lceil \log(\delta^{-1}) \rceil} + \frac{\lceil \log(\log(\delta^{-1})) \rceil}{\lceil \log(\delta^{-1}) \rceil} + 1 \right) \frac{\delta}{2} \\ &\geq 1 - \left( \frac{3}{2} \right) \delta \\ &\geq 1 - 2\delta. \end{aligned}$$

The inequalities follow readily via the fact that for arbitrary  $c_1, c_2 \geq 0$  we have  $\lceil c_1 + c_2 \rceil \leq \lceil c_1 \rceil + \lceil c_2 \rceil$ , and that  $\lceil \log(2) \rceil / \lceil \log(\delta^{-1}) \rceil \leq 1$  for all  $\delta \leq 1/2$ . As for the exponential term, note that

$$e^{-k} = \exp(-\lceil \log(\delta_0^{-1}) \rceil) \leq \exp(-\log(\delta_0^{-1})) = \delta_0 < \delta.$$

It thus immediately follows that the desired good event holds with probability no less than

$$1 - e^{-k} - k\delta_0 \geq 1 - \delta - 2\delta = 1 - 3\delta.$$

Tying together these basic facts readily allows us to prove Theorem 1 (detailed proof in the supplementary materials).

## 4 EMPIRICAL ANALYSIS

To study how the theoretical insights obtained in the previous section play out in practice, we carried out a series of tightly controlled numerical tests. The basic experimental design strategy that we employ is to calibrate all the methods (learning algorithms) of interest to achieve good performance under a particular

learning setup, and then we systematically modify characteristics of the learning tasks, leaving the methods fixed, to observe how performance changes in both an absolute and relative sense. Viewed from a high level, the main points we address can be categorized as follows:

- (E1) How do error trajectories of baseline methods change via robust validation?
- (E2) How does relative performance change in high dimensions without strong convexity?
- (E3) How do actual computation times compare as  $n$  and/or  $d$  grow?
- (E4) Can robust validation be replaced by cross-validation?

**Experimental setup** We essentially follow the standard “noisy convex minimization” tests used in the literature to test the robustness of RGD methods (Holland and Ikeda, 2019a). Complete details of the experimental setup are provided in the supplementary materials.<sup>1</sup> Put simply, we provide the learner with random losses of the form  $L(w; Z) = (\langle w - w^*, X \rangle + E)^2/2$ , where  $w^* \in \mathbb{R}^d$  is a pre-defined vector unknown to the learner,  $X$  is a  $d$ -dimensional random vector,  $E$  is zero-mean random noise, and  $X$  and  $E$  are independent of each other. This approach is advantageous in that we can compute the resulting risk  $R_{\mathbb{P}}(w) = \mathbf{E}_{\mathbb{P}} L(w; Z)$  exactly, and by modifying the distribution  $\mathbb{P}$ , we can ensure that even while allowing heavy-tailed losses/gradients, we still satisfy the key technical assumptions of Theorem 1, namely  $\lambda_1$ -smooth  $R_{\mathbb{P}}$  and gradients with  $\sigma_{G,\mathbb{P}}$ -bounded variance, plus with the  $\mu$ -strong convexity of  $R_{\mathbb{P}}$  is at our control, we can construct many flat directions, and observe behaviour as  $\mu \downarrow 0$ .

With respect to the different methods being studied, we use a mixture of classical baselines and modern alternatives to compare with our Algorithm 1 based on SGD with and without averaging, denoted RV-SGD and RV-SGDave respectively. The sub-routine Valid is carried out using a Catoni-type M-estimator (Catoni, 2012). For baselines, we do empirical risk minimization using batch gradient descent (denoted ERM-GD) and stochastic gradient descent, both with and without averaging (denoted SGD and SGD-Ave). Several representative robust gradient descent methods discussed in section 2 are implemented here, including RGD via median-of-means (Chen et al., 2017b; Prasad et al., 2018) (denoted RGD-MoM), median-of-means minimization by gradient descent (Lecué et al., 2018) (denoted

<sup>1</sup>Software repository:

<https://github.com/feedbackward/sgd-robust>

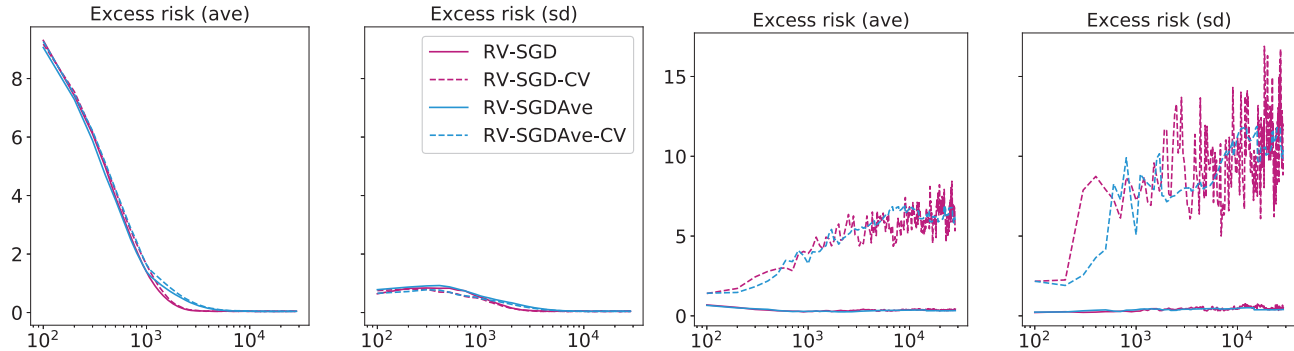


Figure 1: The negative impact of trying to modify Algorithm 1 to use a cross validation heuristic. Left: Normal noise. Right: log-Normal noise.

RGD-Lec), and RGD via M-estimation (Holland and Ikeda, 2019a) (denoted RGD-M). Finally, we also study a cross validation heuristic (marked by `-CV` suffix), where instead of splitting the sample for validation, we do  $k$ -fold cross validation using the full training data set. Essentially, the candidates  $\bar{w}^{(j)}$  are computed just as in Algorithm 1, except without the split into  $Z_{n/2}$  and  $Z'_{n/2}$ , and `Valid` is used to evaluate each candidate using the held-out data for each  $j \in [k]$ . The final candidate is the one for which `Valid` on the held-out data returns the smallest value. This gives the sub-processes double the data for training, but sacrifices the independence of the data used for validation. Detailed settings of each method, including random initialization, are given in the supplementary materials.

**Discussion of results** Representative results are given in Figures 2–1. Starting with proof-of-concept test results given in Figure 2, we see how even very noisy sub-processes can be ironed out easily using the simple robust validation sub-routine included in Algorithm 1, and that even running the algorithm for much longer than a single pass over the data, risk which is comparable to benchmark RGD methods can be realized at a much smaller cost, with comparable variance across trials, and that this holds under both sub-Gaussian and heavy-tailed data, without any modifications to the procedure being run. A particularly lucid improvement in the cost-performance tradeoff is evident in Figure 3, since near-identical performance can be achieved at a small fraction of the computational cost. Note that under Normal noise, running Algorithm 1 for just a single pass leaves room for improvement performance-wise, but as we saw in the low-dimension case, in practice this can be remedied by taking additional passes over the data. Finally, regarding the question of whether or not Algorithm 1 can be replaced with a naive  $k$ -fold cross-validation heuristic, the answer is clear (Figure 1): while the results are comparable under well-behaved

data (the Normal noise case here), when heavy tails are a possibility (e.g., the log-Normal case), the naive cross-validation method fails to get even near the performance of Algorithm 1.

## 5 FUTURE DIRECTIONS

The main take-away from this initial study is that even without strong convexity, under potentially heavy-tailed losses and/or gradients, there exists a computationally efficient procedure which improves upon the formal performance guarantees of modern robust GD techniques, is very competitive in practice, and scales much better to large tasks with many parameters. The basic archetype for such a procedure was illustrated using the concrete Algorithm 1, but naturally this can be extended in many directions, to deal with accelerated, adaptive, or variance-reducing sub-processes, under more general geometries and more challenging constraints applied to  $\mathcal{W}$ . In particular, if one considers a stochastic mirror descent type of generalization to the proposed algorithm, it would be interesting to compare the robust validation approach taken here with say the truncation-based approach studied recently by Juditsky et al. (2019), and how the performance of the respective methods changes under different constraints on prior knowledge of the underlying data-generating distribution.

## Acknowledgements

This work was partially supported by the JSPS KAKENHI Grant Number 19K20342, and the Kayamori Foundation of Information Science Advancement Grant Number K30-XXIII-518.

## References

Bach, F. and Moulines, E. (2014). Non-strongly-convex smooth stochastic approximation with convergence

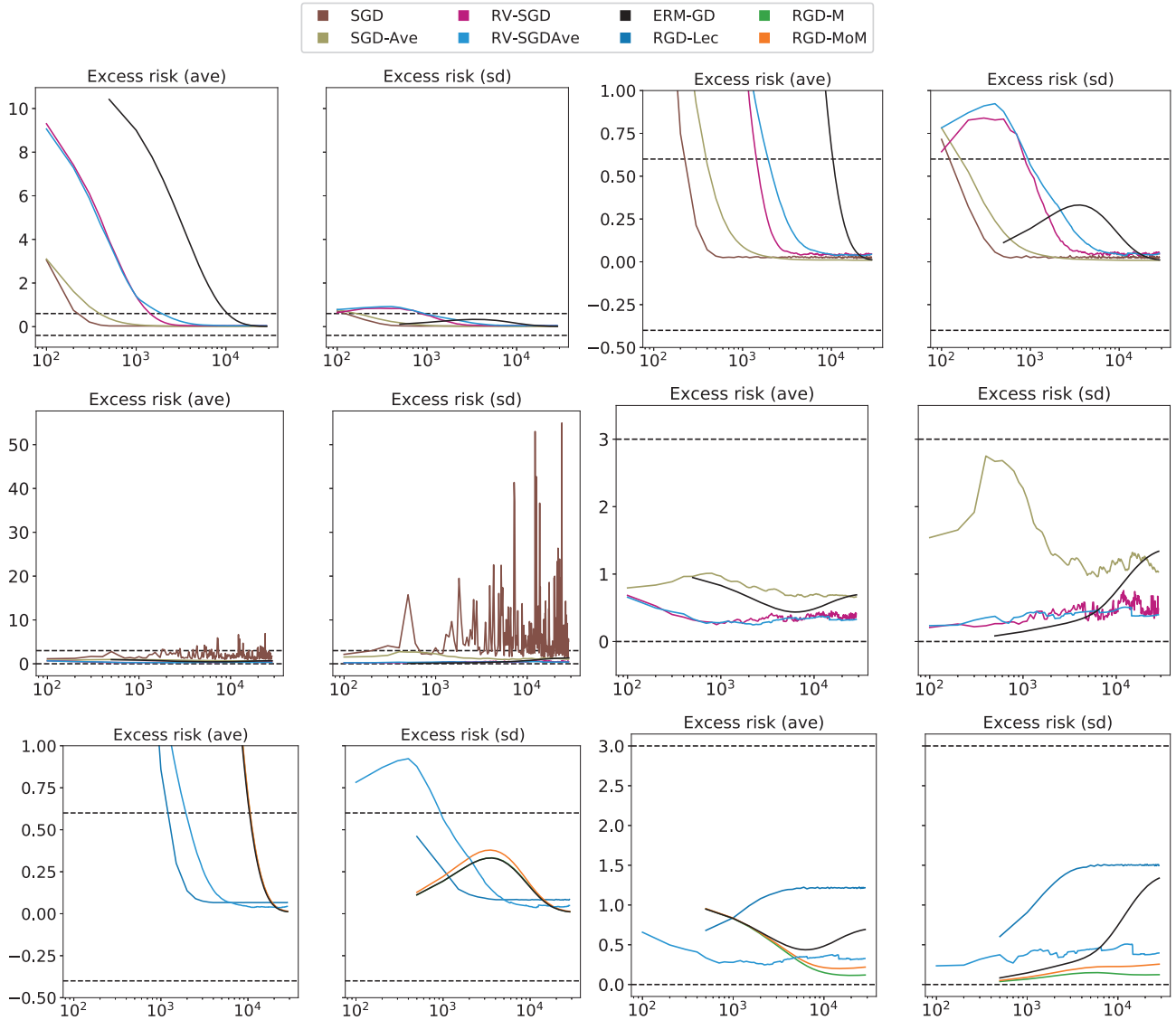


Figure 2: Excess risk statistics (ave and sd) as a function of cost in gradients (log scale, base 10). Top row: Normal noise. Middle row: log-Normal noise. Right-most plots show a zoomed-in view of left-most plots. Bottom row: comparison with RGD methods, Normal (left) and log-Normal (right).

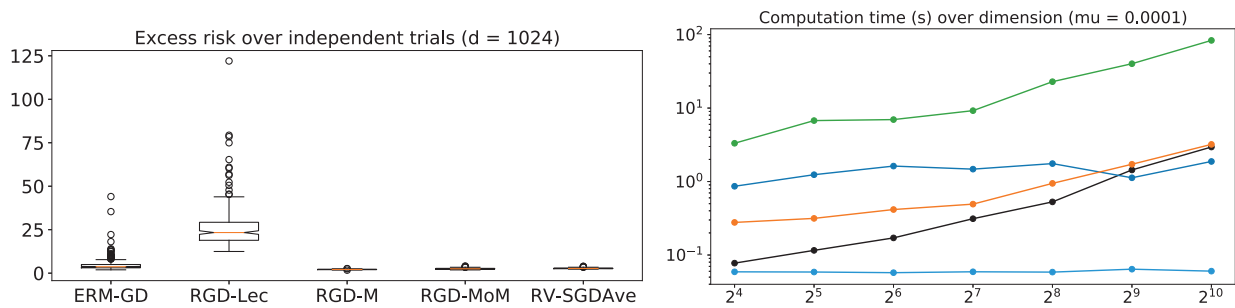


Figure 3: Performance under  $\mu = 0.0001$ . Left: box-plot of final excess risk for batch methods versus RV-SGDAve, with  $d = 1024$  under log-Normal noise. Right: median time cost over  $d$ .



- rate  $o(1/n)$ . In *Advances in Neural Information Processing Systems 26*, pages 773–781.
- Brownlees, C., Joly, E., and Lugosi, G. (2015). Empirical risk minimization for heavy-tailed losses. *Annals of Statistics*, 43(6):2507–2536.
- Catoni, O. (2012). Challenging the empirical mean and empirical variance: a deviation study. *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, 48(4):1148–1185.
- Chen, Y., Su, L., and Xu, J. (2017a). Distributed statistical machine learning in adversarial settings: Byzantine gradient descent. *arXiv preprint arXiv:1705.05491v2*.
- Chen, Y., Su, L., and Xu, J. (2017b). Distributed statistical machine learning in adversarial settings: Byzantine gradient descent. In *Proceedings of the ACM on Measurement and Analysis of Computing Systems*. ACM.
- Devroye, L., Lerasle, M., Lugosi, G., and Oliveira, R. I. (2016). Sub-gaussian mean estimators. *Annals of Statistics*, 44(6):2695–2725.
- Feldman, V. (2017). Generalization of ERM in stochastic convex optimization: The dimension strikes back. In *Advances in Neural Information Processing Systems 29 (NIPS 2016)*, pages 3576–3584.
- Hausler, D. (1992). Decision theoretic generalizations of the PAC model for neural net and other learning applications. *Information and Computation*, 100(1):78–150.
- Holland, M. J. (2019). Robust descent using smoothed multiplicative noise. In *22nd International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 89 of *Proceedings of Machine Learning Research*, pages 703–711.
- Holland, M. J. (2021). Scaling-up robust gradient descent techniques. In *35th AAAI Conference on Artificial Intelligence (AAAI 2021)*.
- Holland, M. J. and Ikeda, K. (2019a). Better generalization with less data using robust gradient descent. In *36th International Conference on Machine Learning (ICML)*, volume 97 of *Proceedings of Machine Learning Research*.
- Holland, M. J. and Ikeda, K. (2019b). Efficient learning with robust gradient descent. *Machine Learning*, 108(8):1523–1560.
- Hsu, D. and Sabato, S. (2016). Loss minimization and parameter estimation with heavy tails. *Journal of Machine Learning Research*, 17(18):1–40.
- Juditsky, A., Nazin, A., Nemirovsky, A., and Tsybakov, A. (2019). Algorithms of robust stochastic optimization based on mirror descent method. *arXiv preprint arXiv:1907.02707v1*.
- Kearns, M. J. and Vazirani, U. V. (1994). *An Introduction to Computational Learning Theory*. MIT Press.
- Lecué, G., Lerasle, M., and Mathieu, T. (2018). Robust classification via MOM minimization. *arXiv preprint arXiv:1808.03106v1*.
- Mehta, N. A. (2016). Fast rates with high probability in exp-concave statistical learning. *arXiv preprint arXiv:1605.01288*.
- Minsker, S. (2015). Geometric median and robust estimation in Banach spaces. *Bernoulli*, 21(4):2308–2335.
- Minsker, S. (2018). Uniform bounds for robust mean estimators. *arXiv preprint arXiv:1812.03523*.
- Prasad, A., Suggala, A. S., Balakrishnan, S., and Ravikumar, P. (2018). Robust estimation via robust gradient estimation. *arXiv preprint arXiv:1802.06485*.
- Shalev-Shwartz, S., Shamir, O., Srebro, N., and Sridharan, K. (2010). Learnability, stability and uniform convergence. *Journal of Machine Learning Research*, 11:2635–2670.
- Vapnik, V. (1982). *Estimation of Dependences Based on Empirical Data*. Springer Series in Statistics. Springer-Verlag.