

1 **ROBUST CLASSIFICATION UNDER ℓ_0 ATTACK FOR THE**
2 **GAUSSIAN MIXTURE MODEL**

3 PAYAM DELGOSHA*, HAMED HASSANI†, AND RAMTIN PEDARSANI‡

4 **Abstract.** It is well-known that machine learning models are vulnerable to small but cleverly-
5 designed adversarial perturbations that can cause misclassification. While there has been major
6 progress in designing attacks and defenses for various adversarial settings, many fundamental and
7 theoretical problems are yet to be resolved. In this paper, we consider classification in the presence
8 of ℓ_0 -bounded adversarial perturbations, a.k.a. sparse attacks. This setting is significantly different
9 from other ℓ_p -adversarial settings, with $p \geq 1$, as the ℓ_0 -ball is non-convex and highly non-smooth.
10 Under the assumption that data is distributed according to the Gaussian mixture model, our goal
11 is to characterize the optimal robust classifier and the corresponding robust classification error as
12 well as a variety of trade-offs between robustness, accuracy, and the adversary’s budget. To this end,
13 we develop a novel classification algorithm called `FilTrun` that has two main modules: Filtration
14 and Truncation. The key idea of our method is to first filter out the *non-robust* coordinates of the
15 input and then apply a carefully-designed truncated inner product for classification. By analyzing
16 the performance of `FilTrun`, we derive an upper bound on the optimal robust classification error.
17 We further find a lower bound by designing a specific adversarial strategy that enables us to derive
18 the corresponding robust classifier and its achieved error. For the case that the covariance matrix of
19 the Gaussian mixtures is diagonal, we show that as the input’s dimension gets large, the upper and
20 lower bounds converge; i.e. we characterize the asymptotically-optimal robust classifier. Throughout,
21 we discuss several examples that illustrate interesting behaviors such as the existence of a *phase*
22 *transition* for adversary’s budget determining whether the effect of adversarial perturbation can be
23 fully neutralized or not.

24 **1. Introduction.** Machine learning has been widely used in a variety of appli-
25 cations including image recognition, virtual assistants, autonomous driving, many of
26 which are safety-critical. Adversarial attacks to machine learning models in the form
27 of a small perturbation added to the input have been shown to be effective in causing
28 classification errors [4, 33, 10, 5, 17]. Formally, the adversary aims to perturb the
29 data in a small ℓ_p -neighborhood so that the perturbed data is “close” to the original
30 data (e.g. imperceptible perturbation in the case of an image) and misclassification
31 occurs. There have been a variety of attacks and defenses proposed in the literature
32 which mostly focus on ℓ_2 or ℓ_∞ bounded perturbations [2, 19, 35]. The state-of-the-art
33 empirical defense against adversarial attacks is iterative training with adversarial ex-
34 amples [18]. While adversarial training can improve robustness, it is shown that there
35 is a fundamental tradeoff between robustness and test accuracy, and such defenses
36 typically lack good generalization performance [34, 32, 26, 1, 36, 13].

37 The focus of this paper is different from such prior work as we consider the problem
38 of robust classification under ℓ_0 -bounded attacks. In this setting, given a pre-specified
39 budget k , the adversary can choose up to k coordinates and arbitrarily change the
40 value of the input at those coordinates. In other words, the adversary can change the
41 input within the so-called ℓ_0 -ball of radius k . In contrast with ℓ_p -balls ($p \geq 1$), the
42 ℓ_0 -ball is non-convex and highly non-smooth. Moreover, the ℓ_0 ball contains inherent
43 discrete (combinatorial) structures that can be exploited by both the learner and the
44 adversary. As a result, the ℓ_0 -adversarial setting bears several fundamental challenges
45 that are absent in other adversarial settings commonly studied in the literature and

*Department of Computer Science, University of Illinois at Urbana Champaign, IL,
delgosha@illinois.edu

†Department of Electrical and Systems Engineering, University of Pennsylvania, Philadelphia,
PA, hassani@seas.upenn.edu

‡Department of Electrical and Computer Engineering, University of California, Santa Barbara,
Santa Barbara, CA, ramtin@ece.ucsb.edu

46 most techniques from prior work do not readily apply in the ℓ_0 setting. Complicating
47 matters further, it can be shown that any piece-wise linear classifier, e.g. a feed-
48 forward deep neural network with ReLU activations, completely fails in the ℓ_0 setting
49 [31]. These all point to the fact that new methodologies are required in the ℓ_0 setting.

50 The ℓ_0 -adversarial setting involves sparse attacks that perturb only a small por-
51 tion of the input signal. This has a variety of applications including natural language
52 processing [14], malware detection [11], and physical attacks in object detection [16].
53 Prior work on ℓ_0 adversarial attacks can be divided into two categories of white-
54 box attacks that are gradient-based, e.g. [5, 22, 21], and black-box attacks based on
55 zeroth-order optimization, e.g. [29, 7]. Defense strategies against ℓ_0 -bounded attacks
56 have also been proposed, e.g. defenses based on randomized ablation [15] and de-
57 fensive distillation [23]. Moreover, [31] develops a simple mathematical framework to
58 show the existence of targeted adversarial examples with ℓ_0 -bounded perturbation in
59 arbitrarily deep neural networks.

60 Despite this interesting recent progress and practical relevance, many fundamen-
61 tal theoretical questions in the ℓ_0 -setting have so far been unanswered: *What are the*
62 *key properties of a robust classifier (recall that all piece-wise linear classifiers fail)?*
63 *What is the optimal robust classifier in standard theoretical settings such the Gauss-*
64 *ian mixture model for data? Is there a trade-off between robustness and accuracy?*
65 *How does the (optimal) robust classification error behave as the adversary’s budget k*
66 *increases? Are there any phase transitions?*

67 We consider the problem of classification with ℓ_0 -adversarially perturbed inputs
68 under the assumption that data is distributed according to the Gaussian mixture
69 model. We formally introduce this setting in Section 2, and address the questions
70 above in the proceeding sections. In particular, instead of searching for the exact
71 form of the optimal robust classifier (which is intractable), we follow a design-based
72 approach: We introduce a novel algorithm for classification as well as strategies for
73 the adversary. We then precisely characterize the error performance of these method-
74 ologies, and consequently, analyse the optimal robust classification error, tradeoffs be-
75 tween robustness and accuracy, phase transitions, etc. We envision that our proposed
76 classification method introduces important modules and insights that are necessary to
77 obtain robustness against ℓ_0 -adversaries for general data distributions (and practical
78 datasets), going beyond the theoretical setting of this paper.

79 **Summary of Contributions.** The main contributions of this paper are as follows:

- 80 • We propose a new robust classification algorithm called **FilTrun** that is based
81 on two main modules: **Filtration** and **Truncation** (See Section 3.1.1 and Al-
82 gorithm 3.1 therein). The filtration module removes the *non-robust* coordi-
83 nates (features) from the input by zeroing out their values. The result is then
84 passed through the truncation module which returns a label by computing
85 a *truncated inner product* with a weight vector whose weights are optimized
86 according to the distribution of un-filtered (surviving) coordinates. The trun-
87 cation module is inspired by tools from robust statistics and guarantees that
88 major outlier values in the input vector, which are possibly caused by the
89 adversary, do not pass to affect the final decision. We highlight that the
90 proposed classifier is highly nonlinear. This is consistent with the simple
91 observation that any linear classifier fails to be robust in the presence of ℓ_0
92 attacks.
- 93 • We analytically derive the robust classification error of the proposed clas-
94 sifier. This in particular serves as an upper bound on the optimal robust

95 classification error (See Theorem 3.2 and Corollary 3.5).

- 96 • We introduce adversarial strategies which, given sufficient budget, perturb
97 the input in a way that the information about the true label is totally erased
98 within the adversarially modified coordinates. The key idea is to pick a
99 subset of the coordinates and to modify their distribution so that they become
100 independent from the true label. This leads to a lower bound for the optimal
101 robust error. (See Theorems 3.8 and 3.11).
- 102 • In the case of having a diagonal covariance matrix for the Gaussian mixtures,
103 we prove that our proposed algorithm FilTrun is indeed *asymptotically-*
104 *optimal*, i.e. as the input dimension d approaches infinity, the upper and
105 lower bounds converge to the same analytical expression (See Theorems 3.13
106 in Section 3.3.2). To the best of our knowledge, this is the first result that
107 establishes optimality for the robust classification error of any mathematical
108 model with ℓ_0 attack.
- 109 • We discuss our results through several example scenarios. In certain scenarios,
110 a phase transition is observed in the sense that for a threshold α_0 , when the
111 adversary’s budget is asymptotically below d^{α_0} , its effect can be completely
112 neutralized, while if the adversary’s budget is above d^{α_0} , no classifier can
113 do better than a naive classifier. In some other scenarios, no sharp phase
114 transition is existent, leading to a trade-off between robustness and accuracy.

2. Problem Formulation. We consider the binary Gaussian mixture model where the distribution for the data generation is specified by the label being $y \sim \text{Unif}\{\pm 1\}$ and $\mathbf{x} \sim \mathcal{N}(y\boldsymbol{\mu}, \Sigma)$, i.e. the Gaussian distribution with mean $y\boldsymbol{\mu}$ and covariance matrix Σ , where $\boldsymbol{\mu} \in \mathbb{R}^d$ and Σ is positive definite. Hereafter, we denote this distribution by $(\mathbf{x}, y) \sim \mathcal{D}$ and refer to y as the label and to \mathbf{x} as the input. Our results correspond to arbitrary choices of $\boldsymbol{\mu}$ and Σ , however, we consider as running example an important special case in which Σ is a diagonal matrix, i.e. the coordinates of \mathbf{x} are independent conditioned on y . Focusing on classification, we consider functions of the form $\mathcal{C} : \mathbb{R}^d \rightarrow \{-1, 1\}$ that predict the label from the input. As a metric for the discrepancy between the prediction of the classifier on the input \mathbf{x} and the true label y , we consider the 0-1 loss $\ell(\mathcal{C}; \mathbf{x}, y) = \mathbb{1}[\mathcal{C}(\mathbf{x}) \neq y]$. We consider classification in the presence of an adversary that perturbs the input \mathbf{x} within the ℓ_0 -ball of radius k :

$$\mathcal{B}_0(\mathbf{x}, k) := \{\mathbf{x}' \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{x}'\|_0 \leq k\},$$

115 where for $\mathbf{x} = (x_1, \dots, x_d)$ we define $\|\mathbf{x}\|_0 := \sum_{i=1}^d \mathbb{1}[x_i \neq 0]$. In other words, the
116 adversary can arbitrarily modify at most k coordinates of \mathbf{x} to obtain \mathbf{x}' , and feed the
117 new vector \mathbf{x}' to the classifier. We call k the *budget* of the adversary. In this setting,
118 the robust classification error of a classifier \mathcal{C} is defined to be the following:

$$(2.1) \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}(\mathcal{C}, k) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\max_{\mathbf{x}' \in \mathcal{B}_0(\mathbf{x}, k)} \ell(\mathcal{C}; \mathbf{x}', y) \right].$$

120 We aim to design classifiers with minimum robust classification error. Hence, we define
121 the *optimal robust classification error* by minimizing (2.1) over all possible classifiers:

$$(2.2) \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*(k) := \inf_{\mathcal{C}} \mathcal{L}_{\boldsymbol{\mu}, \Sigma}(\mathcal{C}, k).$$

123 Our goal in this paper is to precisely characterize $\mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*(k)$ parameterized by $\Sigma, \boldsymbol{\mu}$ and
124 in different regimes of the adversary’s budget k .

125 It is well known that in the absence of the adversary, i.e. when $k = 0$, the Bayes
 126 optimal classifier is the linear classifier $\mathcal{C}(\mathbf{x}) = \text{sgn}(\langle \Sigma^{-1} \boldsymbol{\mu}, \mathbf{x} \rangle)$ which achieves the
 127 *optimal standard error* of $\bar{\Phi}(\|\boldsymbol{\nu}\|_2)$ where $\boldsymbol{\nu} := \Sigma^{-1/2} \boldsymbol{\mu}$ and $\bar{\Phi}(x) := 1 - \Phi(x)$ denotes
 128 the complementary CDF of a standard normal distribution. In order to fix the base-
 129 line, specifically to have a meaningful asymptotic discussion, we may assume without
 130 loss of generality that

$$131 \quad (2.3) \quad \|\boldsymbol{\nu}\|_2 = \|\Sigma^{-1/2} \boldsymbol{\mu}\|_2 = 1.$$

132 Hence, the optimal standard error, which is a lower bound for (2.2), becomes $\bar{\Phi}(1)$.

133 To highlight some of the main challenges of the ℓ_0 -adversarial setting, we note
 134 that linear classifiers in general have been very successful in the Gaussian mixture
 135 setting. Apart from the fact that the Bayes-optimal classifier is linear (when there is
 136 no adversary), even when the adversarial corruptions are chosen in a ℓ_p -ball for $p \geq 1$
 137 it can be shown that the optimal robust classifiers in many cases are also linear (see
 138 [3, 9]). In contrast, in the presence of ℓ_0 -adversaries, it is not hard to show that *any*
 139 linear classifier completely fail. More precisely, when \mathcal{C} is linear and $k \geq 1$, we have
 140 $\mathcal{L}_{\boldsymbol{\mu}, \Sigma}(\mathcal{C}, k) = \frac{1}{2}$. Such failure of linear classifiers showcases, on the one hand, how
 141 powerful the adversary is, and on the other hand, the necessity of new methodologies
 142 in designing robust classifiers.

143 **Further Related Work.** For ℓ_p adversaries, $p \geq 1$, Gaussian mixture models have
 144 been the main setting used in prior work to investigate optimal rules, trade-offs, and
 145 various other phenomena for robust classification; See e.g. [28, 3, 9, 12, 27, 8, 25, 6,
 146 20, 24]. Further, [30] considers data to be uniformly distributed on the sphere or cube
 147 and shows the inevitability of adversarial examples in ℓ_p -settings, $p \geq 0$. In contrast,
 148 to the best of our knowledge, our work provides the first comprehensive study on the
 149 ℓ_0 -adversarial setting using the Gaussian mixture model.

150 **Notation.** Given two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, $\mathbf{x} \odot \mathbf{y} \in \mathbb{R}^d$ denotes the elementwise product
 151 of \mathbf{x} and \mathbf{y} , i.e. $(x_1 y_1, \dots, x_d y_d)$. Moreover, $\text{sort}(\mathbf{x})$ denotes the vector containing the
 152 elements in \mathbf{x} in descending order. For $a \in \mathbb{R}$, $\text{sgn}(a)$ returns the sign of a . We use
 153 $[d]$ to denote the set $\{1, \dots, d\}$ and $[i : j]$ denotes the set $\{i, i + 1, \dots, j\}$. Given a
 154 vector $\mathbf{x} \in \mathbb{R}^d$ and a subset $A \subseteq [d]$, $\mathbf{x}_A = (x_a : a \in A) \in \mathbb{R}^{|A|}$ denotes the subvector
 155 of \mathbf{x} consisting of the coordinates in A . Given a matrix Σ , its diagonal part, denoted
 156 by $\tilde{\Sigma}$, has the same diagonal entries as Σ and its other entries are 0. Given a matrix
 157 $A \in \mathbb{R}^{d \times d}$, $\|A\|_\infty$ denotes the operator norm of A induced by the vector ℓ_∞ norm, i.e.
 158 $\|A\|_\infty := \sup_{\mathbf{x} \neq \mathbf{0}} \|\mathbf{A}\mathbf{x}\|_\infty / \|\mathbf{x}\|_\infty = \max_{1 \leq i \leq d} \sum_{j=1}^d |A_{i,j}|$.

159 **3. Main Results.** In this section, we state our main results that include (i) the
 160 proposed algorithm and its performance analysis that serves as an upper bound on
 161 the optimal robust classification error (Section 3.1), (ii) lower bound on the optimal
 162 robust classification error (Section 3.2), and (iii) discussion on the optimality of the
 163 proposed algorithm (Section 3.3). Throughout, we illustrate our theoretical results
 164 and their ramifications via several examples.

165 **3.1. Upper Bound on the Optimal Robust Classification Error: Algo-
 166 rithm Description and Theoretical Guarantees.** In Section 3.1.1, we introduce
 167 `FilTrun`, our proposed robust classification algorithm, and in Section 3.1.2, we ana-
 168 lyze its performance.

169 **3.1.1. Algorithm Description.** We describe our proposed algorithm `FilTrun`,
 170 a robust classifier which is based on two main modules: Truncation and Filtration. We

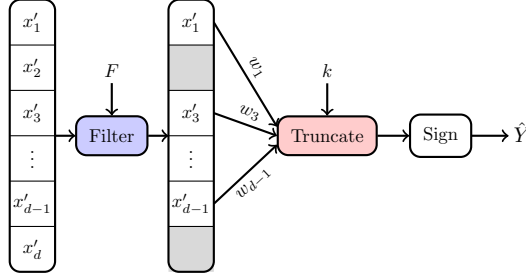


Fig. 1: Schematic of FilTrun.

171 first introduce each of these modules and then proceed with describing the classifier.
 172 **Truncation.** Given vectors $\mathbf{w}, \mathbf{x} \in \mathbb{R}^d$ and an integer $0 \leq k < d/2$, we define the
 173 k -truncated inner product of \mathbf{w} and \mathbf{x} as the summation of the element-wise product
 174 of \mathbf{w} and \mathbf{x} after removing the top and bottom k elements, and denote it by $\langle \mathbf{w}, \mathbf{x} \rangle_k$.
 175 More precisely, let $\mathbf{z} := \mathbf{w} \odot \mathbf{x} \in \mathbb{R}^d$ be the element-wise product of \mathbf{w} and \mathbf{x}
 176 and let $\mathbf{s} = (s_1, \dots, s_d) = \text{sort}(\mathbf{z})$ be obtained by sorting coordinates of \mathbf{z} in descending
 177 order. We then define

$$178 \quad (3.1) \quad \langle \mathbf{w}, \mathbf{x} \rangle_k := \sum_{i=k+1}^{d-k} s_i.$$

179 Note that when $k = 0$, this reduces to the normal inner product $\langle \mathbf{w}, \mathbf{x} \rangle$. Trunca-
 180 tion is a natural method to remove “outliers” which might exist in the data due to
 181 an adversary modifying some coordinates. Therefore, we expect the truncated inner
 182 product to be robust against ℓ_0 perturbations. The following lemma formalizes this.
 183 The proof of Lemma 3.1 is given in Appendix A.

184 **LEMMA 3.1.** *Given $\mathbf{x}, \mathbf{x}', \mathbf{w} \in \mathbb{R}^d$, for integer k satisfying $\|\mathbf{x} - \mathbf{x}'\|_0 \leq k < d/2$,*
 185 *we have*

$$186 \quad |\langle \mathbf{w}, \mathbf{x}' \rangle_k - \langle \mathbf{w}, \mathbf{x} \rangle| \leq 8k \|\mathbf{w} \odot \mathbf{x}\|_\infty.$$

187 In the context of our problem, this lemma suggests that if the budget of the
 188 adversary is at most k , we can bound the difference between the k -truncated inner
 189 product between \mathbf{w} and the adversarially modified sample \mathbf{x}' and the (non-truncated)
 190 inner product between \mathbf{w} and the original sample \mathbf{x} . Recall that in the absence of
 191 the adversary, the optimal Bayes classifier is a linear classifier of the form $\text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle)$
 192 with $\mathbf{w} = \Sigma^{-1} \boldsymbol{\mu}$. Hence, motivated by Lemma 3.1, one can argue that $\text{sgn}(\langle \mathbf{w}, \mathbf{x}' \rangle_k)$
 193 would be robust against ℓ_0 adversarial attacks with budget at most k assuming we
 194 can appropriately control the bound of Lemma 3.1. However, this is not enough—it
 195 turns out that in certain cases, we need to *filter out* some of the input coordinates
 196 and perform the truncation on the remaining coordinates, which we call the *surviving*
 197 coordinates.

198 **Filtration** refers to discarding some of the coordinates of the input. Intuitively,
 199 these coordinates are the *non-robust* features which do more harm than good when
 200 the input is adversarially corrupted. More precisely, given a fixed and nonempty
 201 subset of coordinates $F \subseteq [d]$, we define the classifier $\mathcal{C}_F^{(k)}$ as follows:

202 (3.2)
$$\mathcal{C}_F^{(k)}(\mathbf{x}') := \text{sgn}(\langle \mathbf{w}(F), \mathbf{x}'_F \rangle_k),$$

203 where

204
$$\mathbf{w}(F) := \Sigma_F^{-1} \boldsymbol{\mu}_F,$$

205 and

206 (3.3)
$$\Sigma_F = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [(\mathbf{x}_F - \boldsymbol{\mu}_F)(\mathbf{x}_F - \boldsymbol{\mu}_F)^T | y = 1]$$

207 is the covariance matrix of \mathbf{x}_F conditioned on y , which is essentially the submatrix of
 208 Σ corresponding to the elements in F . Note that $\mathbf{w}(F)$ is the optimal Bayes classifier
 209 of y given \mathbf{x}_F in the absence of the adversary. It is easy to see that when Σ is diagonal,
 210 $\mathbf{w}(F) = \mathbf{w}_F$, but this might not hold in general.

211 Algorithm 3.1 and Figure 1 illustrate the classification procedure `FilTrun` given
 212 in (3.2). So far we have not explained how the set F is chosen and the algorithm
 213 works with any such set given as an input. Later we discuss how the set F is chosen
 214 (see Remarks 3.4 and 3.15).

Algorithm 3.1 `FilTrun`

Input:

- k : adversary's ℓ_0 budget
- $\boldsymbol{\mu}, \Sigma$: parameters of the Gaussian distribution
- F : the set of surviving coordinates
- \mathbf{x}' : the corrupted input

Output:

- $\mathcal{C}_F^{(k)}(\mathbf{x}')$
 - 1: **function** `FILTRUN`($k, \boldsymbol{\mu}, \Sigma, F, \mathbf{x}'$)
 - 2: **Filtering:** Construct $\boldsymbol{\mu}_F, \Sigma_F$ and \mathbf{x}'_F corresponding to the coordinates in F
 - 3: Compute $\mathbf{w}(F) \leftarrow \Sigma_F^{-1} \boldsymbol{\mu}_F$
 - 4: **Truncation:** Compute $\langle \mathbf{w}(F), \mathbf{x}'_F \rangle_k$
 - 5: **Return** $\text{sgn}(\langle \mathbf{w}(F), \mathbf{x}'_F \rangle_k)$
 - 6: **end function**
-

215 **3.1.2. Upper bound on the robust classification error of `FilTrun`.** The
 216 theorem 3.2 below states an upper bound for the robust error associated with the clas-
 217 sification algorithm `FilTrun` introduced in Section 3.1.1. In particular, this yields an
 218 upper bound on the optimal robust classification error. The proof of Theorem 3.2 is
 219 given in Appendix B.

220 **THEOREM 3.2.** *Assume that $\boldsymbol{\mu}, \Sigma$ are given such that (2.3) holds. For a given*
 221 *nonempty $F \subseteq [d]$ and $0 \leq k < d/2$, we have*

222 (3.4)
$$\mathcal{L}_{\boldsymbol{\mu}, \Sigma}(\mathcal{C}_F^{(k)}, k) \leq \frac{1}{\sqrt{2 \log d}} + \bar{\Phi} \left(\|\boldsymbol{\nu}(F)\|_2 - \frac{16k\sqrt{2 \log d} \|\tilde{\Sigma}_F^{1/2} \Sigma_F^{-1/2}\|_\infty \|\boldsymbol{\nu}(F)\|_\infty}{\|\boldsymbol{\nu}(F)\|_2} \right),$$

223 where Σ_F is defined in (3.3), $\tilde{\Sigma}_F$ is the diagonal part of Σ_F , and

224
$$\boldsymbol{\nu}(F) := \Sigma_F^{-1/2} \boldsymbol{\mu}_F.$$

225 As a consequence, we obtain
 226 (3.5)

$$226 \quad \mathcal{L}_{\mu, \Sigma}^*(k) \leq \frac{1}{\sqrt{2 \log d}} + \min_{F \subseteq [d]} \bar{\Phi} \left(\|\nu(F)\|_2 - \frac{16k\sqrt{2 \log d} \|\tilde{\Sigma}_F^{1/2} \Sigma_F^{-1/2}\|_\infty \|\nu(F)\|_\infty}{\|\nu(F)\|_2} \right).$$

227 *Remark 3.3.* Recall from Section 3.1.1 that F is the set of coordinates used for
 228 classification (i.e. the information in the coordinates F^c is discarded). Therefore, we
 229 essentially work with \mathbf{x}_F as an input. If the adversary is not present, the optimal clas-
 230 sification error is achieved via the Bayes linear classifier which has error $\bar{\Phi}(\|\nu(F)\|_2)$.
 231 However, due to the existence of an adversary, we need to perform truncation which
 232 influences the error through the second term inside the argument of $\bar{\Phi}$ in (3.4).

233 *Remark 3.4.* The bound in Theorem 3.2 can be used as a guide to choose the
 234 set of surviving coordinates F . More precisely, we can choose F which minimizes the
 235 right hand side in (3.5). Later, in Section 3.3, we discuss a simpler mechanism for
 236 choosing F when the covariance matrix Σ is diagonal (see Remark 3.15 therein).

237 Here, we outline the proof of Theorem 3.2. Due to the symmetry, we only
 238 need to analyze the classification error when $y = 1$. In this case, an error oc-
 239 curs only when there exists some $\mathbf{x}' \in \mathcal{B}_0(\mathbf{x}, k)$ such that $\langle \mathbf{w}(F), \mathbf{x}'_F \rangle_k \leq 0$. But
 240 since $\|\mathbf{x}'_F - \mathbf{x}_F\|_0 \leq \|\mathbf{x}' - \mathbf{x}\|_0 \leq k$, Lemma 3.1 implies that for such \mathbf{x}' , we have
 241 $|\langle \mathbf{w}(F), \mathbf{x}'_F \rangle_k - \langle \mathbf{w}(F), \mathbf{x}_F \rangle| \leq 8k \|\mathbf{w}(F) \odot \mathbf{x}_F\|_\infty$. Therefore, the robust classifica-
 242 tion error is upper bounded by $\mathbb{P}(\langle \mathbf{w}(F), \mathbf{x}_F \rangle \leq 8k \|\mathbf{w}(F) \odot \mathbf{x}_F\|_\infty)$. But the random
 243 variable $\langle \mathbf{w}(F), \mathbf{x}_F \rangle$ is Gaussian with a known distribution, and the proof follows by
 244 bounding $\|\mathbf{w}(F) \odot \mathbf{x}_F\|_\infty$. See Appendix B for details.

245 When the covariance matrix Σ is diagonal, Σ_F is also diagonal and $\tilde{\Sigma}_F^{1/2} \Sigma_F^{-1/2} = I$.
 246 Moreover, $\nu(F) = \nu_F$ where $\nu := \Sigma^{-1/2} \mu$. This yields the following corollary of
 247 Theorem 3.2.

248 **COROLLARY 3.5.** Assume that μ, Σ are given such that (2.3) holds and Σ is di-
 249 agonal. Then, for nonempty $F \subseteq [d]$ we have

$$250 \quad \mathcal{L}_{\mu, \Sigma}(\mathcal{C}_F^{(k)}, k) \leq \frac{1}{\sqrt{2 \log d}} + \bar{\Phi} \left(\|\nu_F\|_2 - \frac{16k\sqrt{2 \log d} \|\nu_F\|_\infty}{\|\nu_F\|_2} \right),$$

251 and in particular

$$252 \quad \mathcal{L}_{\mu, \Sigma}^*(k) \leq \frac{1}{\sqrt{2 \log d}} + \min_{F \subseteq [d]} \bar{\Phi} \left(\|\nu_F\|_2 - \frac{16k\sqrt{2 \log d} \|\nu_F\|_\infty}{\|\nu_F\|_2} \right).$$

253 Now we discuss the above bounds via two examples, which we use as running
 254 examples to discuss our results in the subsequent sections as well. In the following,
 255 $I_d \in \mathbb{R}^{d \times d}$ and $\mathbf{1}_d \in \mathbb{R}^d$ denote the $d \times d$ identity matrix and the all-ones vector of
 256 size d , respectively.

257 **EXAMPLE 3.6.** Let $\Sigma = I_d$ and $\mu = \frac{1}{\sqrt{d}} \mathbf{1}_d$. In the absence of the adversary,
 258 the optimal Bayes classification error is $\bar{\Phi}(1)$. Moreover, simplifying the bounds in
 259 Corollary 3.5, we get

$$260 \quad \mathcal{L}_{\mu, \Sigma}(\mathcal{C}_F^{(k)}, k) \leq \frac{1}{\sqrt{2 \log d}} + \bar{\Phi} \left(\sqrt{\frac{|F|}{d}} - \frac{16k\sqrt{2 \log d}}{\sqrt{|F|}} \right).$$

261 This is minimized when $F = [d]$, resulting in

$$262 \quad \mathcal{L}_{\mu, \Sigma}^*(k) \leq \frac{1}{\sqrt{2 \log d}} + \bar{\Phi} \left(1 - \frac{16k\sqrt{2 \log d}}{\sqrt{d}} \right).$$

263 Note that if $k = o(\sqrt{d/\log d})$, the upper bound is approximately $\bar{\Phi}(1)$ which is the
 264 optimal classification error in the absence of the adversary. This means that for
 265 $k = o(\sqrt{d/\log d})$, the effect of the adversary can be completely neutralized. We will
 266 show a lower bound for this example later in Section 3.2 (see Example 3.9 therein)
 267 which shows that when $k \geq \sqrt{d} \log d$, no classifier can do asymptotically better than
 268 a naive classifier. This establishes a phase transition at $k = \sqrt{d}$ up to logarithmic
 269 terms.

270 **EXAMPLE 3.7.** Let $\Sigma = I_d$ and $\mu = (d^{-\frac{1}{3}}, cd^{-\frac{1}{2}}, cd^{-\frac{1}{2}}, \dots, cd^{-\frac{1}{2}})$ where c is cho-
 271 sen such that $\|\mu\|_2 = 1$, resulting in an optimal standard error of $\bar{\Phi}(1)$ in the absence
 272 of the adversary. It turns out that the set F that optimizes the bound in Corollary 3.5
 273 is the set $[2 : d]$, i.e. we need to discard the first coordinate. In addition to this, we
 274 can see that if the classifier does not discard the first coordinate, it can neutralize
 275 adversarial attacks with budget of at most $d^{\frac{1}{3}-\epsilon}$, while discarding the first coordinate
 276 makes the classifier immune to adversarial budgets up to $d^{\frac{1}{2}-\epsilon}$. In fact, although the
 277 first coordinate is more informative compared to the other coordinates, due to this
 278 very same reason it is more susceptible to adversarial attacks, and it can do more
 279 harm than good when the input is adversarially corrupted. This example highlights
 280 the importance of the filtration phase.

281 **3.2. Lower Bound on Optimal Robust Classification Error: Strategies**
 282 **for the Adversary.** In this section, we provide a lower bound on the optimal robust
 283 classification error. This is accomplished by introducing an attack strategy for the
 284 adversary, and showing that given such a fixed attack, no classifier can achieve better
 285 than the lower bound that we introduce. The strategy is best understood when
 286 the covariance matrix is diagonal. Therefore, we first assume that Σ is diagonal and
 287 denote the diagonal elements of Σ by $\sigma_1^2, \dots, \sigma_d^2$. We later use our strategy for diagonal
 288 covariance matrices to get a general lower bound for arbitrary Σ (see Theorem 3.11
 289 at the end of this section).

290 Assume that the adversary observes realizations $(\mathbf{x}, y) \sim \mathcal{D}$ generated from the
 291 Gaussian mixture model with parameters μ, Σ , where Σ is diagonal. A randomized
 292 strategy for the adversary with budget k is identified by a probability distribution
 293 which upon observing such realizations (\mathbf{x}, y) , generates a random vector \mathbf{x}' that
 294 satisfies $\mathbb{P}(\|\mathbf{x}' - \mathbf{x}\|_0 \leq k \mid \mathbf{x}, y) = 1$. The goal of the adversary is to design this
 295 randomized strategy in a way that the corrupted vector \mathbf{x}' bears very little information
 296 (or even no information) about the label y . In this way, the loss in (2.2) will be
 297 maximized. Before rigorously defining our proposed strategy for the adversary, we
 298 illustrated its main idea when $d = 1$ in Figure 2.

299 Recall that $\nu = \Sigma^{-1/2} \mu$. Since Σ is diagonal, $\nu_i = \mu_i / \sigma_i$. We will fix a set of
 300 coordinates $A \subseteq [d]$ and a specific value for the budget $k(A) = \|\nu_A\|_1 \log d$. We in-
 301 troduce a randomized strategy for the adversary with the following properties: (i) it
 302 can change up to $k(A)$ coordinates of the input; and (ii) all the changed coordinates
 303 belong to A , i.e. the coordinates in A^c are left untouched. We denote this adversarial
 304 strategy by $\text{Adv}(A)$. Given $A \subset [d]$, having observed (\mathbf{x}, y) , $\text{Adv}(A)$ follows the pro-
 305 cedure explained below. Let $\mathbf{Z} = (Z_1, \dots, Z_d) \in \mathbb{R}^d$ be a random vector that $\text{Adv}(A)$
 306 constructs using the true input \mathbf{x} . First of all, recall that $\text{Adv}(A)$ does not touch the

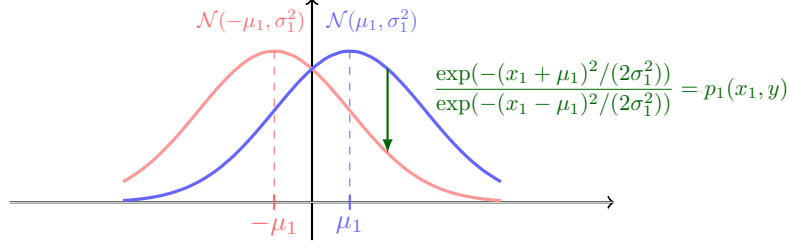


Fig. 2: The idea behind our proposed strategy for the adversary when $d = 1$. Assume $\mu_1 > 0$ and the adversary observes a realization (x_1, y) such that $y = 1$, meaning that x_1 is a realization of $\mathcal{N}(\mu_1, \sigma_1^2)$ (i.e. the blue curve). If $x_1 \leq 0$, the adversary leaves it unchanged, i.e. $x'_1 = x_1$. On the other hand, if $x_1 > 0$, we compute the ratio between the two densities (which is precisely $p_1(x_1, y)$ shown in the figure), and with probability $p_1(x_1, y)$ we pick x'_1 from an arbitrary distribution (e.g. $\text{Unif}[-1, 1]$). When $y = -1$, we follow a similar procedure, but reversed. It is easy to see that by doing so, the distribution of x'_1 is the same when $y = 1$ and $y = -1$, hence x'_1 bears no information about y .

307 coordinates that are not in A , i.e. for $i \in A^c$ we let $Z_i = x_i$. For each $i \in A$, the
 308 adversary's act is simple: it either leaves the value unchanged, i.e. $Z_i = x_i$, or it
 309 erases the value, i.e. $Z_i \sim \text{Unif}[-1, 1]$ —a completely random value between -1 and
 310 $+1$. This binary decision is encoded through a Bernoulli random variable I_i taking
 311 value 0 with probability $p_i(x_i, y)$ and value 1 otherwise. Here $p_i(x_i, y)$ is defined as

$$312 \quad p_i(x_i, y) := \begin{cases} \frac{\exp(-(x_i + y\mu_i)^2 / 2\sigma_i^2)}{\exp(-(x_i - y\mu_i)^2 / 2\sigma_i^2)} & \text{if } \text{sgn}(x_i) = \text{sgn}(y\mu_i) \\ 0 & \text{otherwise} \end{cases}$$

313 Note that the condition $\text{sgn}(x_i) = \text{sgn}(y\mu_i)$ ensures that $p_i(x_i, y) \leq 1$. In summary,
 314 for each $i \in A$, $\text{Adv}(A)$ lets

$$315 \quad (3.6) \quad Z_i = x_i \times I_i + \text{Unif}[-1, 1] \times (1 - I_i),$$

316 where $I_i = \text{Bernoulli}(1 - p_i(x_i, y_i))$, and the random variables I_i are generated com-
 317 pletely independently w.r.t. all the other variables. It is easy to see that the following
 318 holds for the conditional density of \mathbf{Z}_A given y

$$319 \quad (3.7) \quad \begin{aligned} f_{\mathbf{Z}_A|y}(z_A|1) &= f_{\mathbf{Z}_A|y}(z_A|-1) \\ &= \prod_{i \in A} \left[\frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(|z_i| + |\mu_i|)^2}{2\sigma_i^2}\right) + \frac{\alpha_i}{2} \mathbb{1}[z_i \in [-1, 1]] \right], \end{aligned}$$

320 where for $i \in A$

$$321 \quad \alpha_i := \mathbb{P}(I_i = 1|y = 1) = \mathbb{P}(I_i = 1|y = -1) = \int_0^\infty [1 - p_i(t, 1)] f_{x_i|y}(t|1) dt.$$

322 In other words, α_i is the probability of changing coordinate i . Finally, $\text{Adv}(A)$ checks
 323 if the vectors \mathbf{Z} and \mathbf{x} differ within the budget constraint $k(A) := \|\nu_A\|_1 \log d$. Define
 324 \mathbf{x}' as follows:

$$325 \quad (3.8) \quad \mathbf{x}' := \begin{cases} \mathbf{Z} & \text{if } \sum_{i \in A} I_i \leq \|\nu_A\|_1 \log d \\ \mathbf{x} & \text{o.t.w.} \end{cases}$$

326 It can be shown that with high probability, \mathbf{Z} is indeed within the specified budget
 327 and $\mathbf{x}' = \mathbf{Z}$. From this definition, it is evident that with probability one we have

$$328 \quad (3.9) \quad \|\mathbf{x}' - \mathbf{x}\|_0 \leq \|\boldsymbol{\nu}_A\|_1 \log d,$$

329 and hence $\text{Adv}(A)$ is a randomized adversarial strategy that only changes the coordi-
 330 nates in A and has budget $k(A) = \|\boldsymbol{\nu}_A\|_1 \log d$. Now we use this adversarial strategy
 331 to show the following result. The proof of Theorem 3.8 is given in Appendix C.

332 **THEOREM 3.8.** *Assume that the covariance matrix Σ is diagonal and let $\boldsymbol{\nu} =$
 333 $\Sigma^{-1/2}\boldsymbol{\mu}$. Then for any subset $A \subseteq [d]$, we have*

$$334 \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*(\|\boldsymbol{\nu}_A\|_1 \log d) \geq \bar{\Phi}(\|\boldsymbol{\nu}_{A^c}\|_2) - \frac{1}{\log d}.$$

335 The main idea behind this result and the above adversarial strategy is that due
 336 to (3.7), \mathbf{Z}_A is independent from y and since the coordinates of the input are inde-
 337 pendent from each other, and since with high probability $\mathbf{x}' = \mathbf{Z}$, the coordinates in
 338 A have no useful information for the classifier. Hence, the classifier can do no better
 339 than the optimal Bayes classifier for the remaining coordinates in A^c , which results
 340 in a classification error of $\bar{\Phi}(\|\boldsymbol{\nu}_{A^c}\|_2)$.

341 We now apply the bound of Thm 3.8 to Examples 3.6, 3.7 that we discussed in
 342 Section 3.1.2.

343 **EXAMPLE 3.9.** *Assume that $\boldsymbol{\mu}$ and Σ are as in Example 3.6. Applying the bound
 344 in Theorem 3.8, we get*

$$345 \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*\left(\frac{|A|}{\sqrt{d}} \log d\right) \geq \bar{\Phi}\left(\sqrt{1 - \frac{|A|}{d}}\right) - \frac{1}{\log d}.$$

346 Therefore, setting $A = [d]$, we obtain a lower bound of almost $\bar{\Phi}(0) = 1/2$ for adver-
 347 sarial budget $\sqrt{d} \log d$. In other words, if the adversarial budget is more than $\sqrt{d} \log d$,
 348 asymptotically no classifier can do better than a random guess. This together with the
 349 discussion in Example 3.6 establishes a phase transition around \sqrt{d} (modulo logarith-
 350 mic terms).

351 **EXAMPLE 3.10.** *Assume that $\boldsymbol{\mu}$ and Σ are as in Example 3.7. Applying the bound
 352 of Theorem 3.8 with $A = [d]$, we obtain $\mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*(k) \geq \bar{\Phi}(0) - 1/\log d \approx 1/2$ where
 353 $k = (d^{-\frac{1}{3}} + c(d-1)/\sqrt{d}) \log d \approx \sqrt{d} \log d$. Hence, comparing this to Example 3.7, we
 354 find similar to Example 3.9 above that a phase transition occurs around adversarial
 355 budget \sqrt{d} up to logarithmic terms.*

356 Now we state our general lower bound which holds for an arbitrary covariance
 357 matrix. This is Theorem 3.11 below, whose proof is provided in Appendix D. Given $\boldsymbol{\mu}$
 358 and Σ , we define the $d \times d$ matrix R where the i, j entry in R is $R_{i,j} = \Sigma_{i,j} / \sqrt{\Sigma_{i,i} \Sigma_{j,j}}$.
 359 In other words, $R_{i,j}$ is the correlation coefficient between the i th and the j th coordi-
 360 nates in our Gaussian noise. Equivalently, with $\tilde{\Sigma}$ being the diagonal part of Σ , we
 361 may write

$$362 \quad (3.10) \quad R := \tilde{\Sigma}^{-\frac{1}{2}} \Sigma \tilde{\Sigma}^{-\frac{1}{2}}.$$

363 It is evident that since Σ is assumed to be positive definite, R is also positive definite.
 364 Furthermore, we define $\mathbf{u} = (u_1, \dots, u_d)$ where

$$365 \quad (3.11) \quad u_i = \frac{\mu_i}{\sqrt{\Sigma_{i,i}}} \quad 1 \leq i \leq d.$$

366 THEOREM 3.11. With \mathbf{u} and R defined as in (3.10) and (3.11) respectively, for
 367 all $A \subseteq [d]$, we have

$$368 \quad \mathcal{L}_{\mu, \Sigma}^* \left(\frac{1}{\sqrt{\zeta_{\min}}} \|\mathbf{u}_A\|_1 \log d \right) \geq \bar{\Phi}(\|\mathbf{u}_{A^c}\|_2) - \frac{1}{\log d},$$

369 where $\zeta_{\min} > 0$ denotes the minimum eigenvalue of R .

370 *Remark 3.12.* Note that when Σ is diagonal, we have $R = I_d$, $\zeta_{\min} = 1$, and $\mathbf{u} =$
 371 $\boldsymbol{\nu} = \Sigma^{-1/2} \boldsymbol{\mu}$. Therefore, the bound in Theorem 3.11 reduces to that of Theorem 3.8.

372 **3.3. Optimality of FilTrun in the diagonal regime.** We have already seen
 373 for our two running examples that up to logarithmic terms, our lower and upper
 374 bounds match (Examples 3.6 and 3.7 for upper bound, and their matching lower
 375 bounds in Examples 3.9 and 3.10, respectively). First, in Section 3.3.1, we show that
 376 our lower and upper bounds indeed match up to logarithmic terms in the *diagonal*
 377 *regime*, i.e. when the covariance matrix is diagonal. Then, in Section This in particular
 378 implies that our robust classification algorithm **FilTrun** is optimal in this regime.

379 **3.3.1. Comparing the Bounds.** In Theorem 3.13 below, in the diago-
 380 nal regime we compare our upper bound of Corollary 3.5 and our lower bound of
 381 Theorem 3.8. Proof of Theorem 3.13 is given in Appendix E. Recall that $\boldsymbol{\nu} := \Sigma^{-1/2} \boldsymbol{\mu}$
 382 and we assume (2.3) holds. When Σ is diagonal and its diagonal entries are $\sigma_1^2, \dots, \sigma_d^2$,
 383 we have $\nu_i = \mu_i / \sigma_i$. Without loss of generality, we may assume that the coordinates
 384 of $\boldsymbol{\nu}$ are decreasingly ordered such that

$$385 \quad (3.12) \quad |\nu_1| \geq |\nu_2| \geq \dots \geq |\nu_d|.$$

386 Given $c \in [0, 1]$, we define

$$387 \quad (3.13) \quad \lambda_c := \min\{\lambda : \|\boldsymbol{\nu}_{[1:\lambda]}\|_2 \geq c\}.$$

388 THEOREM 3.13. If Σ is diagonal and the coordinates in $\boldsymbol{\nu}$ are sorted as in (3.12),
 389 then:

390 1. For $0 \leq c < 1$, we have

$$391 \quad \mathcal{L}_{\mu, \Sigma}^* \left(\frac{\|\boldsymbol{\nu}_{[1:\lambda_c]}\|_1}{\log d} \right) \leq \frac{1}{\sqrt{2} \log d} + \bar{\Phi} \left(\sqrt{1 - c^2} - \frac{16\sqrt{2}}{\sqrt{1 - c^2} \sqrt{\log d}} \right).$$

392 2. For $0 < c \leq 1$, we have

$$393 \quad \mathcal{L}_{\mu, \Sigma}^* (\|\boldsymbol{\nu}_{[1:\lambda_c]}\|_1 \log d) \geq \bar{\Phi}(\sqrt{1 - c^2}) - \frac{1}{\log d}.$$

394 *Remark 3.14.* Roughly speaking, Theorem 3.13 says that up to logarithmic terms,
 395 we have

$$396 \quad \mathcal{L}_{\mu, \Sigma}^* (\|\boldsymbol{\nu}_{[1:\lambda_c]}\|_1) \approx \bar{\Phi}(\sqrt{1 - c^2}).$$

397 Recall from our previous discussion that we are interested in studying adversarial
 398 budgets scaling as d^α , which justifies neglecting the multiplicative logarithmic terms.
 399 Furthermore, following the proof of Theorem 3.13, the upper bound in the first part
 400 is obtained by our robust classifier by setting $F = \{\lambda_c, \dots, d\}$. Roughly speaking, the
 401 classifier discards the coordinates in $\boldsymbol{\nu}$ which constitute fraction c of the ℓ_2 norm of $\boldsymbol{\nu}$,

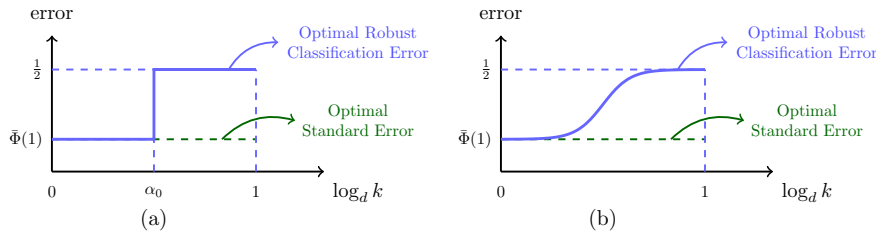


Fig. 3: Asymptotic behavior in the diagonal regime: Illustration of scenarios with (a) a phase transition, and (b) no phase transition

402 and performs a truncated inner product classification on the remaining coordinates.
 403 But the ℓ_2 norm of the remaining coordinates is roughly $\sqrt{1-c^2}$, and the effect
 404 of truncation is vanishing as long as the adversarial power is below $\|\nu_{[1:\lambda_c]}\|_1$ by a
 405 logarithmic factor. Note that although the top coordinates in ν are relatively more
 406 important in terms of the classification power, due to the same reason, they are more
 407 susceptible to adversarial attack.

408 *Remark 3.15.* In view of Theorem 3.13 and Remark 3.14, we can introduce the
 409 following mechanism for choosing the surviving set F for the adversary given adver-
 410 sarial power k . Let $r(k) = \min\{r : \|\nu_{[1:r]}\|_1 \geq k \log d\}$ and set $F = [r(k) : d]$. Then
 411 the classifier $\mathcal{C}_F^{(k)}$ achieves the optimal robust classification error of almost $\bar{\Phi}(\sqrt{1-c^2})$
 412 where $c = \|\nu_{[1:r(k)]}\|_2$.

413 **3.3.2. Asymptotic Analysis, Phase Transitions, and Trade-offs.** In this
 414 section, we perform a thorough analysis when the adversarial budget scales as d^α
 415 using our results in the diagonal regime. Here, we describe the main messages. (i)
 416 We show that our bounds asymptotically match in the diagonal regime and **FilTrun**
 417 is indeed optimal. (ii) Through the asymptotic analysis, we observe that in some
 418 scenarios, a sharp phase transition on the optimal robust error occurs as we increase
 419 $\alpha := \log_d k$ (See Figure 3-(a)). We have already given examples of such scenarios (e.g.
 420 Example 3.6). In such cases, below the transition, i.e. when $\alpha < \alpha_0$, the optimal
 421 robust error is the same as the optimal standard error. And when we are above
 422 the transition, i.e. when $\alpha > \alpha_0$, any classifier becomes useless as the robust error
 423 becomes $\frac{1}{2}$. As a result, asymptotically speaking, there exists no tradeoff between
 424 robustness and standard accuracy in scenarios where there is a sharp transition.

425 However, there are other scenarios where instead of a sharp phase transition, in
 426 the asymptotic regime, the optimal robust error continuously increases as a function
 427 of adversary's budget (see Figure 3-(b)). In such scenarios, there exists a non-trivial
 428 tradeoff between robustness and standard accuracy. I.e. to achieve optimal robust
 429 error it is necessary to filter many informative coordinates which hurts the standard
 430 accuracy. See Example 3.21 below.

431 In order to perform an asymptotic analysis, we assume that the dimension of
 432 the space, d , goes to infinity. More precisely, we assume that we have a sequence
 433 $(\mu^{(d)}, \Sigma^{(d)})$ where for each d , $\mu^{(d)} \in \mathbb{R}^d$ and $\Sigma^{(d)}$ is a diagonal covariance matrix with
 434 nonzero diagonal entries. We define

$$435 \quad \nu^{(d)} := (\Sigma^{(d)})^{-1/2} \mu^{(d)}.$$

436 As usual, as in (2.3), in order to keep the optimal classification error in the absence

437 of the adversary fixed, we assume that

$$438 \quad (3.14) \quad \|\boldsymbol{\nu}^{(d)}\|_2 = 1 \quad \forall d.$$

439 Furthermore, without loss of generality, we assume that the coordinates in $\boldsymbol{\nu}$ are
440 sorted in a descending order with respect to their magnitude, i.e.

$$441 \quad (3.15) \quad |\nu_1^{(d)}| \geq |\nu_2^{(d)}| \geq \dots \geq |\nu_d^{(d)}| \quad \forall d.$$

442 To simplify the notation, we use $\mathcal{L}_d^*(\cdot)$ as a shorthand for $\mathcal{L}_{\boldsymbol{\mu}^{(d)}, \Sigma^{(d)}}^*(\cdot)$. We are mainly
443 interested in studying the asymptotic behavior of $\mathcal{L}_d^*(k_d)$ when k_d is a sequence of
444 adversarial budgets so that k_d behaves like d^α . Motivated by Theorem 3.13, it is
445 natural to define

$$446 \quad (3.16) \quad \lambda_c^{(d)} := \min\{\lambda : \|\boldsymbol{\nu}_{[1:\lambda]}^{(d)}\|_2 \geq c\} \quad \text{for } 0 < c \leq 1.$$

447 Furthermore, for $0 < c \leq 1$, we define

$$448 \quad (3.17) \quad \Psi_d(c) := \log_d \|\boldsymbol{\nu}_{[1:\lambda_c^{(d)}]}^{(d)}\|_1.$$

449 Note that since $c > 0$, $\lambda_c^{(d)} \geq 1$ and $\|\boldsymbol{\nu}_{[1:\lambda_c^{(d)}]}^{(d)}\|_1 > 0$. Therefore, $\Psi_d(c)$ is well-defined.

450 Furthermore, it is easy to verify the following properties for the function $\Psi_d(\cdot)$:

451 LEMMA 3.16. $\Psi_d(\cdot)$ is nonincreasing and $\Psi_d(c) \in [-1/2, 1/2]$ for all $c \in (0, 1]$.

452 *Proof.* Note that

$$453 \quad \Psi_d(c) = \log_d \|\boldsymbol{\nu}_{[1:\lambda_c^{(d)}]}^{(d)}\|_1 \leq \log_d \|\boldsymbol{\nu}^{(d)}\|_1 \leq \log_d(\sqrt{d}\|\boldsymbol{\nu}^{(d)}\|_2) = \log_d \sqrt{d} = \frac{1}{2}.$$

454 On the other hand, note that for $c > 0$, we have $\lambda_c^{(d)} \geq 1$ and $\Psi_d(c) \geq \log_d |\nu_1^{(d)}| =$
455 $\log_d \|\boldsymbol{\nu}\|_\infty$. Furthermore, we have $1 = \|\boldsymbol{\nu}^{(d)}\|_2^2 \leq d\|\boldsymbol{\nu}^{(d)}\|_\infty^2$ which implies that
456 $\|\boldsymbol{\nu}^{(d)}\|_\infty \geq 1/\sqrt{d}$. Consequently, $\Psi_d(c) \geq \log_d 1/\sqrt{d} = -1/2$. This completes the
457 proof. \square

458 Roughly speaking, Theorem 3.13 implies that if k_d behaves like $d^{\Psi_d(c)}$, then
459 $\mathcal{L}^*(k_d) \approx \Phi(\sqrt{1-c^2})$. In order to transform this into a formal asymptotic argument,
460 we assume that for all $c \in (0, 1]$, the sequence $\Psi_d(c)$ is convergent, and we define
461 $\Psi_\infty(c) := \lim_{d \rightarrow \infty} \Psi_d(c)$ as the limit. Since $\Psi_d(\cdot)$ is nondecreasing, if the pointwise
462 limit $\Psi_\infty(\cdot)$ exists, it is also nondecreasing and we may define

$$463 \quad \Psi_\infty(0) := \lim_{c \downarrow 0} \Psi_d(c).$$

464 Additionally, we can show the following lemma.

465 LEMMA 3.17. If $\Psi_\infty(\cdot)$ exists as above, then $\Psi_\infty(c) \in [0, 1/2]$ for all $c \in [0, 1]$.

466 *Proof.* For all $c > 0$ and all d , we have

$$467 \quad \|\boldsymbol{\nu}_{[1:\lambda_c^{(d)}]}^{(d)}\|_1 \geq \|\boldsymbol{\nu}_{[1:\lambda_c^{(d)}]}^{(d)}\|_2^2 \geq c^2.$$

468 Therefore

$$469 \quad \Psi_\infty(c) = \lim_{d \rightarrow \infty} \Psi_d(c) = \lim_{d \rightarrow \infty} \log_d \|\boldsymbol{\nu}_{[1:\lambda_c^{(d)}]}^{(d)}\|_1 \geq \liminf_{d \rightarrow \infty} 2 \log_d c = 0.$$

470 Sending c to zero we also realize that $\Psi_\infty(0) \geq 0$. \square

471 Given these, we can formalize the following asymptotic behavior for the optimal
 472 robust classification error. The proof of Theorem 3.18 below is given in Appendix F.

473 **THEOREM 3.18.** *If $\Psi_d(\cdot)$ converges pointwise to a nondecreasing function $\Psi_\infty : [0, 1] \rightarrow [0, 1/2]$ as above, then the following hold for all $c \in [0, 1]$:*

- 475 1. *If $\limsup_{d \rightarrow \infty} \log_d k_d < \Psi_\infty(c)$, then $\limsup_{d \rightarrow \infty} \mathcal{L}_d^*(k_d) \leq \bar{\Phi}(\sqrt{1 - c^2})$.*
 476 2. *If $\liminf_{d \rightarrow \infty} \log_d k_d > \Psi_\infty(c)$, then $\liminf_{d \rightarrow \infty} \mathcal{L}_d^*(k_d) \geq \bar{\Phi}(\sqrt{1 - c^2})$.*

477 It is sometimes more convenient to state the above theorem in terms of the pseudo
 478 inverse of the function $\Psi_\infty(\cdot)$ defined as follows. For $\alpha \in [0, 1]$, we define

479 (3.18)
$$\Psi_\infty^{-1}(\alpha) := \inf\{\bar{\Phi}(\sqrt{1 - c^2}) : \Psi_\infty(c) \geq \alpha\} \wedge \frac{1}{2}.$$

480 Note that since $\Psi_\infty(c) \leq 1/2$ for all $c \in [0, 1]$, we have

481
$$\Psi_\infty^{-1}(\alpha) = \frac{1}{2} \quad \forall c > \frac{1}{2}.$$

482 With this, we can restate Theorem 3.18 as follows.

483 **COROLLARY 3.19.** *In the setup of Theorem 3.18, for $\alpha \in [0, 1]$ we have*

- 484 1. *If $\limsup \log_d k_d < \alpha$ then $\limsup \mathcal{L}_d^*(k_d) \leq \Psi_\infty^{-1}(\alpha)$.*
 485 2. *If $\liminf \log_d k_d > \alpha$ then $\liminf \mathcal{L}_d^*(k_d) \geq \Psi_\infty^{-1}(\alpha)$.*

486 We now discuss this asymptotic result through some examples.

487 **EXAMPLE 3.20.** *Let $\boldsymbol{\mu}^{(d)}$ and $\Sigma^{(d)}$ be as in Example 3.6, i.e. $\Sigma^{(d)} = I_d$ and $\boldsymbol{\mu}^{(d)} =$
 488 $\frac{1}{\sqrt{d}}\mathbf{1}_d$. Therefore, we have*

489
$$\boldsymbol{\nu}^{(d)} = (\Sigma^{(d)})^{-\frac{1}{2}} \boldsymbol{\mu}^{(d)} = \left(\frac{1}{\sqrt{d}}, \frac{1}{\sqrt{d}}, \dots, \frac{1}{\sqrt{d}} \right).$$

490 Using (3.16), we have $\lambda_c^{(d)} = \lfloor dc^2 \rfloor$ and

491
$$\Psi_d(c) = \log_d \|\boldsymbol{\nu}^{(d)}\|_{[1:\lambda_c^{(d)}]} = \log_d \frac{\lfloor dc^2 \rfloor}{\sqrt{d}} = \frac{1}{2} + o(1).$$

492 Therefore, sending $d \rightarrow \infty$, we realize that

493
$$\Psi_\infty(c) = \frac{1}{2} \quad \forall c \in [0, 1].$$

494 Moreover, using (3.18), we get

495
$$\Psi_\infty^{-1}(\alpha) = \begin{cases} \bar{\Phi}(1) & \alpha \leq \frac{1}{2} \\ \frac{1}{2} & \alpha > \frac{1}{2}. \end{cases}$$

496 Figure 4 illustrates $\Psi_\infty(\cdot)$ and $\Psi_\infty^{-1}(\cdot)$ for this example. Therefore, employing Corol-
 497 lary 3.19, we realize that

- 498 1. *If $\limsup \log_d k_d < 1/2$ then $\limsup \mathcal{L}_d^*(k_d) \leq \bar{\Phi}(1)$*
 499 2. *If $\liminf \log_d k_d > 1/2$ then $\mathcal{L}^*(k_d) \geq 1/2$.*

500 In other words, we observe a phase transition around \sqrt{d} in the sense that if the
 501 adversary's budget is asymptotically below \sqrt{d} , the classifier can achieve the robust
 502 classification error $\bar{\Phi}(1)$, i.e. as if there is no adversary, while if the adversary's budget
 503 is asymptotically above \sqrt{d} , no classifier can achieve a robust classification error better
 504 than that of a trivial classifier. This is consistent with the previous observations in
 505 this case, i.e. Examples 3.6 and 3.9.

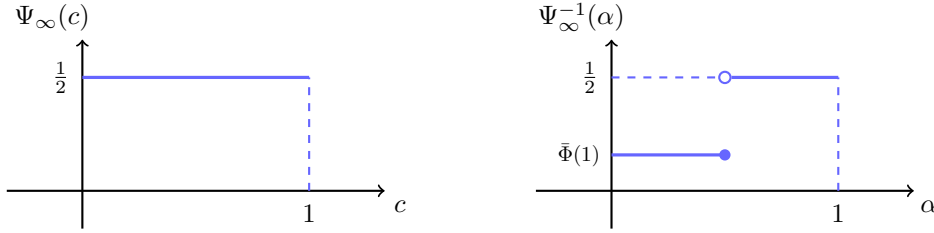


Fig. 4: $\Psi_\infty(\cdot)$ and $\Psi_\infty^{-1}(\cdot)$ for Example 3.20. This observe a phase transition at \sqrt{d} where below this threshold, adversary's effect can completely be neutralized, while above this threshold, the classifier can only achieve the trivial bound.

506 It is interesting to observe that not always we have a phase transition as in the
 507 above example. Below we discuss an example in which we have no phase transition,
 508 and the asymptotic robust classification error gradually increases as a function of the
 509 adversary's budget.

510 **EXAMPLE 3.21.** Let $\Sigma = I_d$. Assume that $d = 2^n - 1$ for some integer n and
 511 define

$$512 \quad \boldsymbol{\mu}^{(d)} = \left(\frac{\sqrt{1/n}}{1}, \frac{\sqrt{1/n}}{\sqrt{2}}, \frac{\sqrt{1/n}}{\sqrt{2}}, \dots, \frac{\sqrt{1/n}}{\sqrt{d/2}}, \dots, \frac{\sqrt{1/n}}{\sqrt{d/2}} \right).$$

513 More precisely, we split the unit ℓ_2 norm of $\boldsymbol{\mu}^{(d)}$ into n blocks, where the first block
 514 is the first coordinate, the second block is the second two coordinate, the i th block
 515 constitutes of 2^i coordinates, and the final block is the last $d/2$ coordinates. Moreover,
 516 the power is uniformly distributed within each block. It is easy to see that for $c =$
 517 $\sqrt{m/n}$ for $1 \leq m \leq n$, we have $\lambda_c^{(d)} = 2^m - 1$ and

$$518 \quad \Psi_d(c) = \Psi_d \left(\sqrt{\frac{m}{n}} \right) = \log_d \left(\sqrt{\frac{1}{n}} \frac{\sqrt{2^m - 1}}{\sqrt{2} - 1} \right) = \frac{c^2}{2} + o(1).$$

519 Therefore, $\Psi_d(\cdot)$ converges pointwise to $\Psi_\infty(\cdot)$ such that $\Psi_\infty(c) = c^2/2$ for $0 \leq c \leq 1$.
 520 Thereby, we have

$$521 \quad \Psi_\infty^{-1}(\alpha) = \begin{cases} \bar{\Phi}(1 - 2\alpha) & 0 \leq \alpha \leq 1/2 \\ \frac{1}{2} & 1/2 < \alpha \leq 1. \end{cases}$$

522 Figure 5 illustrates $\Psi_\infty(\cdot)$ and $\Psi_\infty^{-1}(\cdot)$ in this examples. As we can see, unlike Exam-
 523 ple 3.20, we do not have a phase transition here. In fact, the asymptotic optimal robust
 524 classification error continuously increases as a function of adversarial ℓ_0 budget.

525 **4. Conclusion.** In this paper, we studied the binary Gaussian mixture model
 526 under ℓ_0 attack. We developed a novel nonlinear classifier called **FilTrun** that first
 527 cleverly selects the robust coordinates of the input and then classifies based on a trun-
 528 cated inner product operation. Analyzing the performance of our proposed method,
 529 we derived an upper bound on optimal robust classification error. We further derived
 530 a lower bound on this, and showed the efficacy of **FilTrun**: when the covariance
 531 matrix of Gaussian mixtures is diagonal, **FilTrun** is asymptotically optimal.

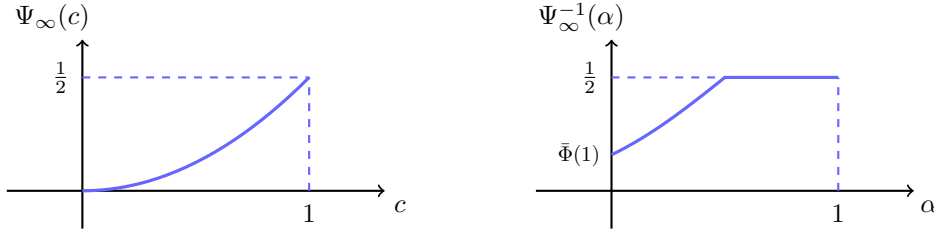


Fig. 5: $\Psi_\infty(\cdot)$ and $\Psi_\infty^{-1}(\cdot)$ for Examples 3.21. Unlike Example 3.20, we do not have a phase transition here and the asymptotic optimal robust classification error continuously increases as a function of the adversarial ℓ_0 budget.

532 There are many directions to be pursued. Deriving a tighter lower bound and
 533 resolving the optimality gap for the case of non-diagonal covariance matrices remains
 534 open. Applying the key ideas of FilTrun, filtration and truncation, to a more com-
 535 plicated setting (e.g. neural networks) can be of great importance from a practical
 536 viewpoint. A crucial message of this paper is to emphasize the importance of non-
 537 linear operations such as truncation for designing defense against ℓ_0 attacks. Finally,
 538 analyzing robust classification error with ℓ_0 attacks for more complex stylized mod-
 539 els such as multi-class Gaussian mixtures, two-layer neural networks, neural tangent
 540 kernel models, etc. is a promising future direction.

541

REFERENCES

- 542 [1] A. A. AL MAKDAH, V. KATEWA, AND F. PASQUALETTI, *A fundamental performance limitation*
 543 *for adversarial classification*, IEEE Control Systems Letters, 4 (2019), pp. 169–174.
- 544 [2] A. ATHALYE, N. CARLINI, AND D. A. WAGNER, *Obfuscated gradients give a false sense of*
 545 *security: Circumventing defenses to adversarial examples*, in Proceedings of the 35th Inter-
 546 national Conference on Machine Learning, ICML, Stockholm, Sweden, July 10-15, 2018,
 547 pp. 274–283, <http://proceedings.mlr.press/v80/athalye18a.html>.
- 548 [3] A. N. BHAGOJI, D. CULLINA, AND P. MITTAL, *Lower bounds on adversarial robustness from*
 549 *optimal transport*, in Advances in Neural Information Processing Systems, 8-14 Decem-
 550 ber 2019, Vancouver, BC, Canada, 2019, pp. 7496–7508, [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/8968-lower-bounds-on-adversarial-robustness-from-optimal-transport)
 551 [8968-lower-bounds-on-adversarial-robustness-from-optimal-transport](http://papers.nips.cc/paper/8968-lower-bounds-on-adversarial-robustness-from-optimal-transport).
- 552 [4] B. BIGGIO, I. CORONA, D. MAIORCA, B. NELSON, N. ŠRNDIĆ, P. LASKOV, G. GIACINTO, AND
 553 F. ROLI, *Evasion attacks against machine learning at test time*, in Joint European confer-
 554 ence on machine learning and knowledge discovery in databases, Springer, 2013, pp. 387–
 555 402.
- 556 [5] N. CARLINI AND D. A. WAGNER, *Towards evaluating the robustness of neural networks*, in
 557 2017 IEEE Symposium on Security and Privacy, San Jose, CA, USA, May 22-26., 2017,
 558 pp. 39–57, <https://doi.org/10.1109/SP.2017.49>, <https://doi.org/10.1109/SP.2017.49>.
- 559 [6] L. CHEN, Y. MIN, M. ZHANG, AND A. KARBASI, *More data can expand the generalization gap*
 560 *between adversarially robust and standard models*, in International Conference on Machine
 561 Learning, PMLR, 2020, pp. 1670–1680.
- 562 [7] F. CROCE, M. ANDRIUSHCHENKO, N. D. SINGH, N. FLAMMARION, AND M. HEIN, *Sparse-rs: a*
 563 *versatile framework for query-efficient sparse black-box adversarial attacks*, arXiv preprint
 564 arXiv:2006.12834, (2020).
- 565 [8] C. DAN, Y. WEI, AND P. RAVIKUMAR, *Sharp statistical guarantees for adversarially robust*
 566 *gaussian classification*, in International Conference on Machine Learning, PMLR, 2020,
 567 pp. 2345–2355.
- 568 [9] E. DOBRIBAN, H. HASSANI, D. HONG, AND A. ROBEY, *Provable tradeoffs in adversarially robust*
 569 *classification*, arXiv preprint arXiv:2006.05161, (2020).
- 570 [10] I. J. GOODFELLOW, J. SHLENS, AND C. SZEGEDY, *Explaining and harnessing adversarial exam-*
 571 *ples*, arXiv preprint arXiv:1412.6572, (2014).

- 572 [11] K. GROSSE, N. PAPERNOT, P. MANOHARAN, M. BACKES, AND P. MCDANIEL, *Adversarial*
573 *perturbations against deep neural networks for malware classification*, arXiv preprint
574 arXiv:1606.04435, (2016).
- 575 [12] J. HAYES, *Provable trade-offs between private \mathcal{E} robust machine learning*, arXiv preprint
576 arXiv:2006.04622, (2020).
- 577 [13] A. JAVANMARD, M. SOLTANOLKOTABI, AND H. HASSANI, *Precise tradeoffs in adversarial training*
578 *for linear regression*, in Conference on Learning Theory, PMLR, 2020, pp. 2034–2078.
- 579 [14] D. JIN, Z. JIN, J. T. ZHOU, AND P. SZOLOVITS, *Is bert really robust? natural language attack*
580 *on text classification and entailment*, arXiv preprint arXiv:1907.11932, 2 (2019).
- 581 [15] A. LEVINE AND S. FEIZI, *Robustness certificates for sparse adversarial attacks by randomized*
582 *ablation.*, in AAAI, 2020, pp. 4585–4593.
- 583 [16] J. LI, F. SCHMIDT, AND Z. KOLTER, *Adversarial camera stickers: A physical camera-based*
584 *attack on deep learning systems*, in International Conference on Machine Learning, PMLR,
585 2019, pp. 3896–3904.
- 586 [17] A. MADRY, A. MAKELOV, L. SCHMIDT, D. TSIPRAS, AND A. VLADU, *Towards deep learning*
587 *models resistant to adversarial attacks*, arXiv preprint arXiv:1706.06083, (2017).
- 588 [18] A. MADRY, A. MAKELOV, L. SCHMIDT, D. TSIPRAS, AND A. VLADU, *Towards deep learning*
589 *models resistant to adversarial attacks*, in 6th International Conference on Learning Repre-
590 *sentations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference*
591 *Track Proceedings, OpenReview.net, 2018, <https://openreview.net/forum?id=rJzIBfZAb>.*
- 592 [19] Z. MARZI, S. GOPALAKRISHNAN, U. MADHOW, AND R. PEDARSANI, *Sparsity-based defense*
593 *against adversarial attacks on linear classifiers*, in 2018 IEEE International Sympos-
594 *ium on Information Theory, ISIT, Vail, CO, USA, June 17-22, 2018, pp. 31–35, <https://doi.org/10.1109/ISIT.2018.8437638>, <https://doi.org/10.1109/ISIT.2018.8437638>.*
- 595 [20] Y. MIN, L. CHEN, AND A. KARBASI, *The curious case of adversarially robust models: More*
596 *data can help, double descend, or hurt generalization*, arXiv preprint arXiv:2002.11080,
597 (2020).
- 598 [21] A. MODAS, S.-M. MOOSAVI-DEZFOOLI, AND P. FROSSARD, *Sparsefool: a few pixels make a big*
600 *difference*, in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
601 *Recognition, 2019, pp. 9087–9096.*
- 602 [22] N. PAPERNOT, P. MCDANIEL, S. JHA, M. FREDRIKSON, Z. B. CELIK, AND A. SWAMI, *The*
603 *limitations of deep learning in adversarial settings*, in 2016 IEEE European symposium on
604 *security and privacy (EuroS&P)*, IEEE, 2016, pp. 372–387.
- 605 [23] N. PAPERNOT, P. MCDANIEL, X. WU, S. JHA, AND A. SWAMI, *Distillation as a defense to ad-*
606 *versarial perturbations against deep neural networks*, in 2016 IEEE Symposium on Security
607 *and Privacy (SP)*, IEEE, 2016, pp. 582–597.
- 608 [24] B. PURANIK, U. MADHOW, AND R. PEDARSANI, *Adversarially robust classification based on glrt*,
609 arXiv preprint arXiv:2011.07835, (2020).
- 610 [25] M. S. PYDI AND V. JOG, *Adversarial risk via optimal transport and optimal couplings*, in
611 *International Conference on Machine Learning, PMLR, 2020, pp. 7814–7823.*
- 612 [26] A. RAGHUNATHAN, S. M. XIE, F. YANG, J. C. DUCHI, AND P. LIANG, *Adversarial training can*
613 *hurt generalization*, arXiv preprint arXiv:1906.06032, (2019).
- 614 [27] E. RICHARDSON AND Y. WEISS, *A bayes-optimal view on adversarial examples*, arXiv preprint
615 arXiv:2002.08859, (2020).
- 616 [28] L. SCHMIDT, S. SANTURKAR, D. TSIPRAS, K. TALWAR, AND A. MADRY, *Adversarially robust*
617 *generalization requires more data*, in Advances in Neural Information Processing Systems,
618 2018, pp. 5014–5026.
- 619 [29] L. SCHOTT, J. RAUBER, M. BETHGE, AND W. BRENDEL, *Towards the first adversarially robust*
620 *neural network model on mnist*, arXiv preprint arXiv:1805.09190, (2018).
- 621 [30] A. SHAFABI, W. R. HUANG, C. STUDER, S. FEIZI, AND T. GOLDSTEIN, *Are adversarial examples*
622 *inevitable?*, arXiv preprint arXiv:1809.02104, (2018).
- 623 [31] A. SHAMIR, I. SAFRAN, E. RONEN, AND O. DUNKELMAN, *A simple explanation for the existence*
624 *of adversarial examples with small hamming distance*, arXiv preprint arXiv:1901.10861,
625 (2019).
- 626 [32] D. SU, H. ZHANG, H. CHEN, J. YI, P.-Y. CHEN, AND Y. GAO, *Is robustness the cost of*
627 *accuracy?—a comprehensive study on the robustness of 18 deep image classification models*,
628 in Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 631–
629 648.
- 630 [33] C. SZEGEDY, W. ZAREMBA, I. SUTSKEVER, J. BRUNA, D. ERHAN, I. J. GOODFELLOW, AND
631 R. FERGUS, *Intriguing properties of neural networks*, in International Conference on Learning
632 *Representations, 2014, Banff, AB, Canada, April 14-16, 2014, [http://arxiv.org/abs/](http://arxiv.org/abs/1312.6199)*
633 *1312.6199.*

- 634 [34] D. TSIPRAS, S. SANTURKAR, L. ENGSTROM, A. TURNER, AND A. MADRY, *Robustness may be*
635 *at odds with accuracy*, in International Conference on Learning Representations, no. 2019,
636 2019.
- 637 [35] E. WONG AND J. Z. KOLTER, *Provable defenses against adversarial examples via the convex*
638 *outer adversarial polytope*, in Proceedings of the 35th International Conference on Machine
639 Learning, ICML, Stockholm, Sweden, July 10-15, 2018, [http://proceedings.mlr.press/v80/](http://proceedings.mlr.press/v80/wong18a.html)
640 [wong18a.html](http://proceedings.mlr.press/v80/wong18a.html).
- 641 [36] H. ZHANG, Y. YU, J. JIAO, E. XING, L. EL GHAOUI, AND M. JORDAN, *Theoretically princi-*
642 *pled trade-off between robustness and accuracy*, in International Conference on Machine
643 Learning, PMLR, 2019, pp. 7472–7482.

644 **Appendix A. Proof of Lemma 3.1.**

645 In this section, we prove Lemma 3.1. First we need to define some notations and
646 discuss some lemmas.

647 Given $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$, we define the sample average of \mathbf{x} as $\text{Mean}(\mathbf{x}) :=$
648 $\sum_{i=1}^d x_i/d$. Moreover, we define *truncated sum* $\text{TSum}_k(\mathbf{x})$ for $k < n/2$ as follows. Let
649 $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ be the set of sorted values in \mathbf{x} . We define

$$650 \quad \text{TSum}_k(\mathbf{x}) := \sum_{i=k+1}^{d-k} x_{(i)},$$

651 which is the truncated sum of the elements in \mathbf{x} after removing the top and bottom k
652 values. For instance, $\text{TSum}_1(1, 1, 2, 3, 4, 5) = 1 + 2 + 3 + 4 = 10$. Moreover, we define
653 the truncated mean of \mathbf{x} as follows:

$$654 \quad \text{TMean}_k(\mathbf{x}) := \frac{\text{TSum}_k(S)}{d - 2k}.$$

655 Note that when $k = 0$, the above quantities reduce to the sum and the sample average,
656 respectively. It is straightforward to see that

$$657 \quad (\text{A.1}) \quad \left| \text{TSum}_k(\mathbf{x}) - \sum_{i=1}^n x_i \right| \leq 2kM \quad \text{given } |x_i| \leq M \quad \forall 1 \leq i \leq n.$$

658 **LEMMA A.1.** *Assume that $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$ and $\mathbf{x}' = (x'_1, \dots, x'_d) \in \mathbb{R}^d$
659 are given such that \mathbf{x}' is identical to \mathbf{x} in all but at most $k < d/2$ coordinates, i.e.
660 $\|\mathbf{x} - \mathbf{x}'\|_0 \leq k$. Moreover, assume that for some $M < \infty$, we have $|x_i| \leq M$ for all
661 $1 \leq i \leq d$. Then, if $x'_{(1)} \leq x'_{(2)} \leq \dots \leq x'_{(d)}$ are the sorted coordinates in \mathbf{x}' , we have*

$$662 \quad |x'_{(i)}| \leq M \quad \forall k+1 \leq i \leq d-k.$$

663 Essentially, what Lemma A.1 states is that if we modify at most k coordinates in
664 a vector whose elements are bounded by M , in the resulting vector, after truncating
665 the top and bottom k coordinates, all the surviving values are also bounded by M .

666 *Proof of Lemma A.1.* Let i_1, \dots, i_l for $l \leq k$ be the coordinates where \mathbf{x}' differs
667 from \mathbf{x} , i.e. $x_{i_j} \neq x'_{i_j}$ for $1 \leq j \leq l$. Note that if $|x'_{i_j}| > M$ for any of $1 \leq j \leq l$,
668 then x'_{i_j} will definitely fall into the top or bottom k coordinates in the sorted list
669 $x'_{(1)} \leq \dots \leq x'_{(d)}$, since all the $d-l \geq d-k$ remaining coordinates in \mathbf{x}' are bounded by
670 M . This means that all the surviving coordinates $x'_{(k+1)}, \dots, x'_{(d-k)}$ after truncating
671 top and bottom k coordinates in \mathbf{x}' are indeed bounded by M which completes the
672 proof. \square

673 **LEMMA A.2.** *Assume that $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$ is given such that $|x_i| \leq M$ for
674 all $1 \leq i \leq d$. Also, assume that $\mathbf{x}' = (x'_1, \dots, x'_d) \in \mathbb{R}^d$ is identical to \mathbf{x} in all but at
675 most k coordinates, i.e. $\|\mathbf{x} - \mathbf{x}'\|_0 \leq k$. Then, we have*

$$676 \quad |\text{TSum}_k(\mathbf{x}) - \text{TSum}_k(\mathbf{x}')| \leq 6kM.$$

677 *Proof.* Let $x_{\sigma(1)} \leq \dots \leq x_{\sigma(d)}$ and $x'_{\sigma'(1)} \leq \dots \leq x'_{\sigma'(d)}$ be the sorted elements in

678 \mathbf{x} and \mathbf{x}' with permutations σ and σ' , respectively. Following the definition, we have

$$\begin{aligned}
679 \quad \text{TSum}_k(\mathbf{x}) &= \sum_{i=k+1}^{d-k} x_{\sigma(i)} = \sum_{i:\sigma^{-1}(i) \in \{k+1, \dots, d-k\}} x_i \\
680 &= \sum_{i=1}^d \mathbb{1}[\sigma^{-1}(i) \in \{k+1, \dots, d-k\}] x_i. \\
681
\end{aligned}$$

682 Similarly, we have

$$683 \quad \text{TSum}_k(\mathbf{x}') = \sum_{i=1}^d \mathbb{1}[\sigma'^{-1}(i) \in \{k+1, \dots, d-k\}] x'_i.$$

To simplify the notation, for $1 \leq i \leq d$, we define

$$y_i := \mathbb{1}[\sigma^{-1}(i) \in \{k+1, \dots, d-k\}] x_i,$$

and

$$y'_i := \mathbb{1}[\sigma'^{-1}(i) \in \{k+1, \dots, d-k\}] x'_i.$$

684 Moreover, let

$$\begin{aligned}
685 \quad A_1 &:= \{1 \leq i \leq d : \sigma^{-1}(i) \in \{k+1, \dots, d-k\} \text{ and } \sigma'^{-1}(i) \notin \{k+1, \dots, d-k\}\} \\
686 \quad A_2 &:= \{1 \leq i \leq d : \sigma^{-1}(i) \notin \{k+1, \dots, d-k\} \text{ and } \sigma'^{-1}(i) \in \{k+1, \dots, d-k\}\} \\
687 \quad A_3 &:= \{1 \leq i \leq d : \sigma^{-1}(i) \in \{k+1, \dots, d-k\} \\
688 &\quad \text{and } \sigma'^{-1}(i) \in \{k+1, \dots, d-k\} \text{ and } x_i \neq x'_i\} \\
689 \quad A &:= A_1 \cup A_2 \cup A_3.
\end{aligned}$$

691 Note that if $i \notin A$, either $\sigma^{-1}(i) \notin \{k+1, \dots, d-k\}$ and $\sigma'^{-1}(i) \notin \{k+1, \dots, d-k\}$,
692 in which case $y_i = y'_i = 0$; or $\sigma^{-1}(i) \in \{k+1, \dots, d-k\}$, $\sigma'^{-1}(i) \in \{k+1, \dots, d-k\}$,
693 and $x_i = x'_i$, in which case $y_i = y'_i = x_i = x'_i$. This means that $y_i = y'_i$ for $i \notin A$ and

$$\begin{aligned}
694 \quad (\text{A.2}) \quad |\text{TSum}_k(\mathbf{x}) - \text{TSum}_k(\mathbf{x}')| &\leq \sum_{i \in A} |y_i - y'_i| \\
&\leq \sum_{i \in A_1} |y_i - y'_i| + \sum_{i \in A_2} |y_i - y'_i| + \sum_{i \in A_3} |y_i - y'_i|.
\end{aligned}$$

695 Note that for $i \in A_1$, we have $y'_i = 0$ and $y_i = x_i$, implying $|y_i - y'_i| = |x_i| \leq M$. On
696 the other hand, for $i \in A_2$, $y_i = 0$ and $y'_i = x'_i$. But since $\sigma'^{-1}(i) \in \{k+1, \dots, d-k\}$,
697 using Lemma A.1, we have $|y_i - y'_i| = |x'_i| \leq M$. Moreover, for $i \in A_3$, we have
698 $y_i = x_i$ and $y'_i = x'_i$. Also, from Lemma A.1, we have $|x'_i| \leq M$. Thereby, $|y_i - y'_i| \leq$
699 $|x_i| + |x'_i| \leq 2M$. Putting all these together, we get

$$700 \quad (\text{A.3}) \quad \sum_{i \in A_1} |y_i - y'_i| + \sum_{i \in A_2} |y_i - y'_i| + \sum_{i \in A_3} |y_i - y'_i| \leq M|A_1| + M|A_2| + 2M|A_3|.$$

701 Observe that

$$702 \quad (\text{A.4}) \quad |A_1| \leq |\{1 \leq i \leq d : \sigma'^{-1}(i) \notin \{k+1, \dots, d-k\}\}| = 2k.$$

703 Similarly,

$$704 \quad (\text{A.5}) \quad |A_2| \leq 2k.$$

705 On the other hand,

$$706 \quad (\text{A.6}) \quad |A_3| \leq |\{1 \leq i \leq d : x_i \neq x'_i\}| \leq k.$$

707 Using (A.4), (A.5), and (A.6) back into (A.3) and comparing with (A.2), we realize
708 that

$$709 \quad |\text{TSum}_k(\mathbf{x}) - \text{TSum}_k(\mathbf{x}')| \leq 6kM,$$

710 which completes the proof. \square

711 The following is a direct consequence of Lemma A.2.

712 **COROLLARY A.3.** *Given $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ and integer k satisfying $\|\mathbf{x} - \mathbf{x}'\|_0 \leq k < d/2$,*
713 *we have*

$$714 \quad |\text{TSum}_k(\mathbf{x}) - \text{TSum}_k(\mathbf{x}')| \leq 6k \min\{\|\mathbf{x}\|_\infty, \|\mathbf{x}'\|_\infty\}.$$

715 We are now ready to give the proof of Lemma 3.1:

716 *Proof of Lemma 3.1.* We have

$$\begin{aligned} 717 \quad |\langle \mathbf{w}, \mathbf{x}' \rangle_k - \langle \mathbf{w}, \mathbf{x} \rangle| &\leq |\langle \mathbf{w}, \mathbf{x}' \rangle_k - \langle \mathbf{w}, \mathbf{x} \rangle_k| + |\langle \mathbf{w}, \mathbf{x} \rangle_k - \langle \mathbf{w}, \mathbf{x} \rangle| \\ 718 &\leq |\langle \mathbf{w}, \mathbf{x}' \rangle_k - \langle \mathbf{w}, \mathbf{x} \rangle_k| + 2k\|\mathbf{w} \odot \mathbf{x}\|_\infty \\ 719 &= |\text{TSum}_k(\mathbf{w} \odot \mathbf{x}') - \text{TSum}_k(\mathbf{w} \odot \mathbf{x})| + 2k\|\mathbf{w} \odot \mathbf{x}\|_\infty \\ 720 &\stackrel{(a)}{\leq} 6k\|\mathbf{w} \odot \mathbf{x}\|_\infty + 2k\|\mathbf{w} \odot \mathbf{x}\|_\infty \\ 721 &= 8k\|\mathbf{w} \odot \mathbf{x}\|_\infty, \end{aligned}$$

723 where in step (a) we have used $\|\mathbf{w} \odot \mathbf{x}' - \mathbf{w} \odot \mathbf{x}\|_0 \leq \|\mathbf{x}' - \mathbf{x}\|_0 \leq k$ together with
724 Corollary A.3. This completes the proof. \square

725 **Appendix B. Proof of the Upper Bound (Theorem 3.2).**

726 Given $\mathbf{x} \in \mathbb{R}^d$ and $y \in \{\pm 1\}$, define

$$727 \quad \ell^{(k)}(\mathcal{C}_F^{(k)}; \mathbf{x}, y) := \max_{\mathbf{x}' \in \mathcal{B}_0(\mathbf{x}, k)} \ell(\mathcal{C}_F^{(k)}; \mathbf{x}', y).$$

728 We have

$$\begin{aligned} 729 \quad \ell^{(k)}(\mathcal{C}_F^{(k)}; \mathbf{x}, 1) &= \mathbb{1} [\exists \mathbf{x}' \in \mathcal{B}_0(\mathbf{x}, k) : \mathcal{C}_F^{(k)}(\mathbf{x}') \neq 1] \\ 730 &= \mathbb{1} [\exists \mathbf{x}' \in \mathcal{B}_0(\mathbf{x}, k) : \langle \mathbf{w}(F), \mathbf{x}'_F \rangle_k \leq 0] \end{aligned}$$

732 Using Lemma 3.1, for \mathbf{x}' such that $\|\mathbf{x}' - \mathbf{x}\|_0 \leq 0$, since $\|\mathbf{x}'_F - \mathbf{x}_F\|_0 \leq \|\mathbf{x}' - \mathbf{x}\|_0 \leq k$,
733 we have

$$734 \quad |\langle \mathbf{w}(F), \mathbf{x}'_F \rangle_k - \langle \mathbf{w}(F), \mathbf{x}_F \rangle| \leq 8k\|\mathbf{w}(F) \odot \mathbf{x}_F\|_\infty.$$

735 This means that

$$736 \quad \mathbb{1} [\exists \mathbf{x}' \in \mathcal{B}_0(\mathbf{x}, k) : \langle \mathbf{w}(F), \mathbf{x}'_F \rangle_k \leq 0] \leq \mathbb{1} [\langle \mathbf{w}(F), \mathbf{x}_F \rangle \leq 8k\|\mathbf{w}(F) \odot \mathbf{x}_F\|_\infty],$$

737 and
 (B.1)

$$738 \quad \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\ell^{(k)}(\mathcal{C}_F^{(k)}; \mathbf{x}, 1) | y = 1 \right] \leq \mathbb{P}(\langle \mathbf{w}(F), \mathbf{x}_F \rangle \leq 8k \|\mathbf{w}(F) \odot \mathbf{x}_F\|_\infty | y = 1).$$

739 Let Σ_F be as defined in (3.3) and let $\tilde{\Sigma}_F$ be the diagonal part of Σ_F . Note that since
 740 Σ is positive definite, $\tilde{\Sigma}_F$ is diagonal with positive diagonal entries. Hence, we may
 741 write
 (B.2)

$$742 \quad \|\mathbf{w}(F) \odot \mathbf{x}_F\|_\infty = \|(\tilde{\Sigma}^{1/2} \mathbf{w}(F)) \odot (\tilde{\Sigma}^{-1/2} \mathbf{x}_F)\|_\infty \leq \|\tilde{\Sigma}_F^{1/2} \mathbf{w}(F)\|_\infty \|\tilde{\Sigma}_F^{-1/2} \mathbf{x}_F\|_\infty.$$

743 Let σ_i^2 denote the i th diagonal coordinate of Σ . Fix $i \in F$ and note that conditioned
 744 on $y = 1$, we have $x_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$. On the other hand, with $\mathbf{a} := \tilde{\Sigma}_F^{-1/2} \mathbf{x}_F$, we have
 745 $a_i \sim \mathcal{N}(\sigma_i^{-1} \mu_i, 1)$. Note that $\bar{\Phi}(\sigma_i^{-1} \mu_i)$ is the optimal Bayes classification error of y
 746 given x_i only, which is indeed not smaller than the optimal Bayes classification error
 747 of y given the whole vector \mathbf{x} , which is in turn equal to $\bar{\Phi}(\|\boldsymbol{\nu}\|_2) = \bar{\Phi}(1)$. Since $\bar{\Phi}$ is
 748 decreasing, this implies $\sigma_i^{-1} \mu_i \leq 1$. Consequently, by union bound, we have

$$749 \quad \mathbb{P}\left(\|\tilde{\Sigma}_F^{-1/2} \mathbf{x}_F\|_\infty > 1 + \sqrt{2 \log d}\right) \leq \sum_{i \in F} \mathbb{P}\left(a_i - \sigma_i^{-1} \mu_i > \sqrt{2 \log d}\right)$$

$$750 \quad \leq d \bar{\Phi}(\sqrt{2 \log d})$$

$$751 \quad \leq d \frac{1}{\sqrt{2\pi} \sqrt{2 \log d}} e^{-\log d}$$

$$752 \quad \leq \frac{1}{\sqrt{2 \log d}}.$$

754 Thereby, we get

$$755 \quad (B.3) \quad \mathbb{P}\left(\|\tilde{\Sigma}_F^{-1/2} \mathbf{x}_F\|_\infty > 2\sqrt{2 \log d} | y = 1\right) \leq \frac{1}{\sqrt{2 \log d}}.$$

756 On the other hand, we have

$$757 \quad (B.4) \quad \|\tilde{\Sigma}_F^{1/2} \mathbf{w}(F)\|_\infty = \|\tilde{\Sigma}_F^{1/2} \Sigma_F^{-1/2} \boldsymbol{\nu}(F)\|_\infty \leq \|\tilde{\Sigma}_F^{1/2} \Sigma_F^{-1/2}\|_\infty \|\boldsymbol{\nu}(F)\|_\infty,$$

758 where $\|\tilde{\Sigma}_F^{1/2} \Sigma_F^{-1/2}\|_\infty$ denotes the operator norm of $\tilde{\Sigma}_F^{1/2} \Sigma_F^{-1/2}$ induced by the vector
 759 ℓ_∞ norm. Using (B.2), (B.3), and (B.4) back into (B.1) and simplifying, we get

$$760 \quad \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\ell^{(k)}(\mathcal{C}_F^{(k)}; \mathbf{x}, 1) | y = 1 \right]$$

$$761 \quad \leq \frac{1}{\sqrt{2 \log d}} + \mathbb{P}\left(\langle \mathbf{w}(F), \mathbf{x}_F \rangle \leq 16k \sqrt{2 \log d} \|\tilde{\Sigma}_F^{1/2} \Sigma_F^{-1/2}\|_\infty \|\boldsymbol{\nu}(F)\|_\infty | y = 1\right)$$

763 It is easy to see that conditioned on $y = 1$, $\langle \mathbf{w}(F), \mathbf{x}_F \rangle \sim \mathcal{N}(\|\boldsymbol{\nu}(F)\|_2^2, \|\boldsymbol{\nu}(F)\|_2^2)$.
 764 Using this in the above bound, we get

$$765 \quad \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\ell^{(k)}(\mathcal{C}_F^{(k)}; \mathbf{x}, 1) | y = 1 \right]$$

$$766 \quad \leq \frac{1}{\sqrt{2 \log d}} + \bar{\Phi}\left(\|\boldsymbol{\nu}(F)\|_2 - \frac{16k \sqrt{2 \log d} \|\tilde{\Sigma}_F^{1/2} \Sigma_F^{-1/2}\|_\infty \|\boldsymbol{\nu}(F)\|_\infty}{\|\boldsymbol{\nu}(F)\|_2}\right).$$

768 Due to the symmetry, we have the same bound conditioned on $y = -1$ which yields
769 the desired result.

770 **Appendix C. Lower Bound in the Diagonal Regime (Theorem 3.8).**

771 Before giving the proof of Theorem 3.8, we need the following lemma.

772 LEMMA C.1. *For any random adversarial strategy with budget k which has a den-*
773 *sity function $f_{\mathbf{x}'|x,y}$, we have*

$$774 \mathcal{L}_{\mu,\Sigma}^*(k) \geq \frac{1}{2} \mathbb{P}(f_{\mathbf{x}'|y}(\mathbf{x}'|1) = f_{\mathbf{x}'|y}(\mathbf{x}'|-1)) + \mathbb{P}\left(f_{\mathbf{x}'|y}(\mathbf{x}'|-1) > f_{\mathbf{x}'|y}(\mathbf{x}'|1) \middle| y = 1\right),$$

775 *Proof.* Note that the right hand side is indeed the Bayes optimal error associated
776 with the MAP estimator assuming that the classifier knows adversary's strategy. Since
777 the classifier does not know the adversary's strategy in general, the right hand side is
778 indeed a lower bound on the optimal robust classification error. \square

779 Now we are ready to prove Theorem 3.8.

780 *Proof of Theorem 3.8.* Note that when A is empty, there is no adversarial modifi-
781 cation and the standard Bayes analysis implies that $\mathcal{L}_{\mu,\Sigma}^*(0) = \bar{\Phi}(\|\nu\|_2) = \bar{\Phi}(\|\nu_{A^c}\|_2)$
782 and the desired bound holds. Hence, we may assume that A is nonempty for the rest
783 of the proof.

784 Note that due to (3.9), the randomized strategy $\text{Adv}(A)$ is valid for the adversary
785 given the budget $\|\nu\|_1 \log d$. Thereby we may use Lemma C.1 with $\text{Adv}(A)$ to bound
786 $\mathcal{L}_{\mu,\Sigma}^*(\|\nu_A\|_1 \log d)$ from below. Before that, we show that with high probability under
787 the above randomized strategy for the adversary, recalling the definition of random
788 variables I_i for $i \in A$ from (3.6), we have $\sum_{i \in A} I_i \leq \|\nu_A\|_1 \log d$ and hence $\mathbf{x}' = \mathbf{Z}$.
789 It is easy to see that for each i , $\mathbb{P}(I_i = 1|y = 1) = \mathbb{P}(I_i = 1|y = -1)$; therefore,

$$\begin{aligned} 790 \mathbb{P}(I_i = 1) &= \mathbb{P}(I_i = 1|y = \text{sgn}(\mu_i)) \\ 791 &= \int_0^\infty [1 - p_i(t, \text{sgn}(\mu_i))] f_{x_i|y}(t|\text{sgn}(\mu_i)) dt \\ 792 &= \int_0^\infty \left[1 - \frac{\exp(-(t + |\mu_i|)^2/2\sigma_i^2)}{\exp(-(t - |\mu_i|)^2/2\sigma_i^2)}\right] \exp(-(t - |\mu_i|)^2/2\sigma_i^2) dt \\ 793 &= 1 - \bar{\Phi}(|\nu_i|) \\ 794 &= \text{Erf}(|\nu_i|/\sqrt{2}) \\ 795 &\leq \left(\sqrt{\frac{2}{\pi}}|\nu_i|\right) \wedge 1. \end{aligned}$$

797 Hence, we have

$$798 \mathbb{P}(I_i = 1) = \mathbb{P}(I_i = 1|y = 1) = \mathbb{P}(I_i = 1|y = -1) \leq \left(\sqrt{\frac{2}{\pi}}|\nu_i|\right) \wedge 1.$$

799 Therefore, using Markov's inequality, if I is the indicator of the event $\sum_{i \in A} I_i >$
800 $\|\nu_A\|_1 \log d$, we have

$$801 \text{(C.1)} \quad \mathbb{P}(I = 1) = \mathbb{P}(I = 1|y = 1) = \mathbb{P}(I = 1|y = -1) \leq \frac{\sqrt{2/\pi} \sum_{i \in A} |\nu_i|}{\|\nu_A\|_1 \log d} \leq \frac{1}{\log d}.$$

802 Now, we bound $\mathcal{L}_{\mu,\Sigma}^*(\|\nu_A\|_1 \log d)$ from below in the following two cases.

803 Case 1: $A = [d]$. In this case, using Lemma C.1, we have

$$\begin{aligned}
804 \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*(\|\boldsymbol{\nu}_A\|_1 \log d) &\geq \frac{1}{2} \mathbb{P}(f_{\mathbf{x}'|y}(\mathbf{x}'|1) = f_{\mathbf{x}'|y}(\mathbf{x}'|-1)) \\
805 &\stackrel{(a)}{=} \frac{1}{2} \mathbb{P}(f_{\mathbf{x}'|y}(\mathbf{x}'|1) = f_{\mathbf{x}'|y}(\mathbf{x}'|-1) | y = 1) \\
806 &\geq \frac{1}{2} \mathbb{P}(f_{\mathbf{x}'|y}(\mathbf{x}'|1) = f_{\mathbf{x}'|y}(\mathbf{x}'|-1), I = 0 | y = 1) \\
807 &\stackrel{(b)}{=} \frac{1}{2} \mathbb{P}(f_{\mathbf{Z}|y}(\mathbf{Z}|1) = f_{\mathbf{Z}|y}(\mathbf{Z}|-1) | y = -1) \\
808 &\geq \frac{1}{2} \mathbb{P}(f_{\mathbf{Z}|y}(\mathbf{Z}|1) = f_{\mathbf{Z}|y}(\mathbf{Z}|-1) | y = 1) - \frac{1}{2} \mathbb{P}(I = 1 | y = 1) \\
809 &\stackrel{(c)}{\geq} \frac{1}{2} - \frac{1}{2 \log d}, \\
810
\end{aligned}$$

811 where (a) uses the symmetry, (b) uses the fact that when $I = 0$, by definition we have
812 $\mathbf{x}' = \mathbf{Z}$, and (c) uses (3.7) and (C.1).

813 Case 2: $A \subsetneq [d]$. Using Lemma C.1, we have

$$\begin{aligned}
&\mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*(\|\boldsymbol{\nu}_A\|_1 \log d) \geq \mathbb{P}(f_{\mathbf{x}'|y}(\mathbf{x}'|-1) > f_{\mathbf{x}'|y}(\mathbf{x}'|1) | y = 1) \\
&\geq \mathbb{P}(f_{\mathbf{x}'|y}(\mathbf{x}'|-1) > f_{\mathbf{x}'|y}(\mathbf{x}'|1), I = 0 | y = 1) \\
814 \quad (C.2) \quad &\stackrel{(a)}{=} \mathbb{P}(f_{\mathbf{Z}|y}(\mathbf{Z}|-1) > f_{\mathbf{Z}|y}(\mathbf{Z}|1), I = 0 | y = 1) \\
&\geq \mathbb{P}(f_{\mathbf{Z}|y}(\mathbf{Z}|-1) > f_{\mathbf{Z}|y}(\mathbf{Z}|1) | y = 1) - \mathbb{P}(I = 1 | y = 1) \\
&\stackrel{(b)}{\geq} \mathbb{P}(f_{\mathbf{Z}|y}(\mathbf{Z}|-1) > f_{\mathbf{Z}|y}(\mathbf{Z}|1) | y = 1) - \frac{1}{\log d}
\end{aligned}$$

815 where (a) uses the fact that by definition, when $I = 0$, we have $\mathbf{x}' = \mathbf{Z}$, and (b)
816 uses (C.1). Note that since Z_i are conditionally independent given y , we have

$$817 \quad f_{\mathbf{Z}|y}(\mathbf{Z}|y) = f_{\mathbf{Z}_A|y}(\mathbf{Z}_A|y) f_{\mathbf{Z}_{A^c}|y}(\mathbf{Z}_{A^c}|y).$$

818 But from (3.7), we have $f_{\mathbf{Z}_A|y}(\mathbf{Z}_A|1) = f_{\mathbf{Z}_A|y}(\mathbf{Z}_A|-1)$ with probability one. Using
819 this in (C.2), we get

$$\begin{aligned}
820 \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*(\|\boldsymbol{\nu}_A\|_1 \log d) &\geq \mathbb{P}(f_{\mathbf{Z}_{A^c}|y}(\mathbf{Z}_{A^c}|-1) > f_{\mathbf{Z}_{A^c}|y}(\mathbf{Z}_{A^c}|1) | y = 1) - \frac{1}{\log d} \\
821 &= \bar{\Phi}(\|\boldsymbol{\nu}_{A^c}\|_2) - \frac{1}{\log d}. \\
822
\end{aligned}$$

823 We may combine the two cases following the convention that when $A = [d]$,
824 $A^c = \emptyset$ and $\|\boldsymbol{\nu}_{A^c}\|_2 = 0$. This completes the proof. \square

825 Appendix D. Proof of the General Lower Bound (Theorem 3.11).

826 In this section, we prove Theorem 3.11 by providing a general lower bound for
827 the optimal robust classification error which relaxes the diagonal assumption for the
828 covariance matrix. Our strategy is to approximate the covariance matrix by a diagonal
829 matrix and use our lower bound of Theorem 3.8. It turns out that the optimal robust
830 classification error is monotone with respect to the positive definite ordering of the
831 covariance matrix. Lemma D.1 below formalizes this. Intuitively speaking, the reason
832 is that more noise makes the classification more difficult, resulting in an increase in
833 the optimal robust classification error.

834 LEMMA D.1. Assume that $\boldsymbol{\mu} \in \mathbb{R}^d$ and Σ_1 and Σ_2 are two positive definite co-
835 variance matrices such that $\Sigma_1 \preceq \Sigma_2$. Then for $0 \leq k \leq d$ we have

$$836 \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma_1}^*(k) \leq \mathcal{L}_{\boldsymbol{\mu}, \Sigma_2}^*(k).$$

837 *Proof.* Let $y \sim \text{Unif}(\pm 1)$, $\mathbf{x}_1 \sim \mathcal{N}(y\boldsymbol{\mu}, \Sigma_1)$ and $\mathbf{x}_2 \sim \mathcal{N}(y\boldsymbol{\mu}, \Sigma_2)$. Since $\Sigma_1 \preceq \Sigma_2$,
838 we may write $\Sigma_2 = \Sigma_1 + A$ such that $A \succeq 0$. In addition to this, we may couple $\mathbf{x}_1, \mathbf{x}_2$
839 on the same probability space as $\mathbf{x}_2 = \mathbf{x}_1 + \mathbf{Z}$ where $\mathbf{Z} \sim \mathcal{N}(0, A)$ is independent
840 from all other variables. Now, fix a classifier $\mathcal{C}_2 : \mathbb{R}^d \rightarrow \{\pm 1\}$ and note that

$$\begin{aligned} \mathcal{L}_{\boldsymbol{\mu}, \Sigma_2}(\mathcal{C}_2, k) &= \mathbb{P}(\exists \mathbf{x}' \in \mathcal{B}_0(\mathbf{x}_2, k) : \mathcal{C}_2(\mathbf{x}') \neq y) \\ &= \mathbb{P}(\exists \mathbf{x}' \in \mathcal{B}_0(\mathbf{x}_1 + \mathbf{Z}, k) : \mathcal{C}_2(\mathbf{x}') \neq y) \\ 841 \quad (\text{D.1}) \quad &= \mathbb{P}(\exists \mathbf{x}'' \in \mathcal{B}_0(\mathbf{x}_1, k) : \mathcal{C}_2(\mathbf{x}'' + \mathbf{Z}) \neq y) \\ &\geq \inf_{\tilde{\mathcal{C}}_2: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \{\pm 1\}} \mathbb{P}(\exists \mathbf{x}'' \in \mathcal{B}_0(\mathbf{x}_1, k) : \tilde{\mathcal{C}}_2(\mathbf{x}'', \mathbf{Z}) \neq y) \end{aligned}$$

842 Now, fix $\tilde{\mathcal{C}}_2 : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \{\pm 1\}$ and note that using the independence of Z , we may
843 write

$$\begin{aligned} &\mathbb{P}(\exists \mathbf{x}'' \in \mathcal{B}_0(\mathbf{x}_1, k) : \tilde{\mathcal{C}}_2(\mathbf{x}'', \mathbf{Z}) \neq y) \\ 844 \quad (\text{D.2}) \quad &= \mathbb{E} \left[\mathbb{E} \left[\mathbb{1} \left[\exists \mathbf{x}'' \in \mathcal{B}_0(\mathbf{x}_1, k) : \tilde{\mathcal{C}}_2(\mathbf{x}'', \mathbf{Z}) \neq y \right] \middle| \mathbf{Z} \right] \right] \\ &= \int \mathbb{P}(\exists \mathbf{x}'' \in \mathcal{B}_0(\mathbf{x}_1, k) : \tilde{\mathcal{C}}_2(\mathbf{x}_1, \mathbf{z}) \neq y) f_{\mathbf{Z}}(\mathbf{z}) d\mathbf{z} \end{aligned}$$

845 But for $\mathbf{z} \in \mathbb{R}^d$, if we let $\tilde{\mathcal{C}}_{2, \mathbf{z}}(\mathbf{x}) := \tilde{\mathcal{C}}_2(\mathbf{x}, \mathbf{z})$, we get

$$\begin{aligned} 846 \quad \mathbb{P}(\exists \mathbf{x}'' \in \mathcal{B}_0(\mathbf{x}_1, k) : \tilde{\mathcal{C}}_2(\mathbf{x}_1, \mathbf{z}) \neq y) &= \mathbb{P}(\exists \mathbf{x}'' \in \mathcal{B}_0(\mathbf{x}_1, k) : \tilde{\mathcal{C}}_{2, \mathbf{z}}(\mathbf{x}_1) \neq y) \\ 847 \quad &\geq \inf_{\mathcal{C}_1: \mathbb{R}^d \rightarrow \{\pm 1\}} \mathbb{P}(\exists \mathbf{x}'' \in \mathcal{B}_0(\mathbf{x}_1, k) : \tilde{\mathcal{C}}_1(\mathbf{x}_1) \neq y) \\ 848 \quad &= \mathcal{L}_{\boldsymbol{\mu}, \Sigma_1}^*(k). \end{aligned}$$

850 Comparing this with (D.1) and (D.2), we realize that $\mathcal{L}_{\boldsymbol{\mu}, \Sigma_2}(\mathcal{C}_2, k) \geq \mathcal{L}_{\boldsymbol{\mu}, \Sigma_1}^*(k)$. Since
851 this holds for arbitrary \mathcal{C}_2 , optimizing for \mathcal{C}_2 yields the desired result. \square

852 Note that since Σ is positive definite, we have $\Sigma \succeq \alpha I_d$ where $\alpha > 0$ is the
853 minimum eigenvalue of Σ . Therefore, we may use Lemma D.1 together with the lower
854 bound of Theorem 3.8 for $\mathcal{L}_{\boldsymbol{\mu}, \alpha I_d}^*(\cdot)$ to obtain a lower bound for $\mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*(\cdot)$. However,
855 it turns out that it is more efficient in some scenarios to first normalize the diagonal
856 entries of the covariance matrix. More precisely, define the $d \times d$ matrix R where the i, j
857 entry in R is $R_{i,j} = \Sigma_{i,j} / \sqrt{\Sigma_{ii} \Sigma_{jj}}$. In other words, $R_{i,j}$ is the correlation coefficient
858 between the i th and the j th coordinates in our Gaussian noise. Equivalently, with $\tilde{\Sigma}$
859 being the diagonal part of Σ , we may write

$$860 \quad (\text{D.3}) \quad R := \tilde{\Sigma}^{-\frac{1}{2}} \Sigma \tilde{\Sigma}^{-\frac{1}{2}}.$$

861 It is evident that since Σ is assumed to be positive definite, R is also positive definite.
862 In fact, R is the covariance matrix of the normalized random vector \mathbf{x}' such that
863 $x'_i = x_i / \sqrt{\Sigma_{i,i}}$ where $\mathbf{x} \sim \mathcal{N}(y\boldsymbol{\mu}, \Sigma)$. Also, all the diagonal entries in R are equal

864 to 1, and when Σ is diagonal, $R = I_d$ is the identity matrix. Furthermore, we define
 865 $\mathbf{u} = (u_1, \dots, u_d)$ where

$$866 \quad (\text{D.4}) \quad u_i = \frac{\mu_i}{\sqrt{\Sigma_{i,i}}} \quad 1 \leq i \leq d.$$

867 In fact, with \mathbf{x}' being the normalized of \mathbf{x} as above, we have $\mathbf{u} = \mathbb{E}[\mathbf{x}'|y=1]$. In
 868 Lemma D.2, we show that such coordinate-wise normalization does not affect the
 869 optimal robust classification error. The main reason for this is that any coordinate-
 870 wise product of a vector by positive values does not change the ℓ_0 norm. This property
 871 is unique to the combinatorial ℓ_0 norm, and indeed does not hold for ℓ_p norms for
 872 $p \geq 1$.

873 LEMMA D.2. Given a vector $\mathbf{a} \in \mathbb{R}^d$ with strictly positive entries, if we define
 874 $\boldsymbol{\mu}' \in \mathbb{R}^d$ and $\Sigma' \in \mathbb{R}^{d \times d}$ as $\mu'_i = a_i \mu_i$ and $\Sigma'_{i,j} = a_i a_j \Sigma_{i,j}$, then we have

$$875 \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*(k) = \mathcal{L}_{\boldsymbol{\mu}', \Sigma'}^*(k) \quad \forall 0 \leq k \leq d.$$

876 In particular, with \mathbf{u} and R defined above, we have

$$877 \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*(k) = \mathcal{L}_{\mathbf{u}, R}^*(k) \quad \forall 0 \leq k \leq d.$$

878 *Proof.* Pick $\epsilon > 0$ together with a classifier \mathcal{C} such that

$$879 \quad (\text{D.5}) \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*(k) \geq \mathcal{L}_{\boldsymbol{\mu}, \Sigma}(\mathcal{C}, k) - \epsilon.$$

880 Let $\mathbf{x} \sim \mathcal{N}(y\boldsymbol{\mu}, \Sigma)$, i.e. $(\mathbf{x}, y) \sim \mathcal{D}$, and define $\mathbf{x}' := \mathbf{a} \odot \mathbf{x}$. Note that $\mathbf{x}' \sim \mathcal{N}(y\boldsymbol{\mu}', \Sigma')$.
 881 Let \mathcal{D}' denote the joint distribution of (\mathbf{x}', Y) . Recall that by definition $\mathcal{L}_{\boldsymbol{\mu}, \Sigma}(\mathcal{C}, k) =$
 882 $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\max_{\mathbf{x}' \in \mathcal{B}_0(\mathbf{x}, k)} \ell(\mathcal{C}; \mathbf{x}', y)]$. Note that $\mathbf{x}' \in \mathcal{B}_0(\mathbf{x}, k)$ iff $\|\mathbf{x}' - \mathbf{x}\|_0 \leq k$. Since
 883 all the entries in \mathbf{a} are nonzero, this is equivalent to $\|\mathbf{a} \odot \mathbf{x}' - \mathbf{a} \odot \mathbf{x}\|_0 \leq k$ which is in
 884 turn equivalent to $\mathbf{a} \odot \mathbf{x}' \in \mathcal{B}_0(\mathbf{a} \odot \mathbf{x}, k)$. Therefore, if \mathbf{a}^{-1} denotes the elementwise
 885 inverse of \mathbf{a} , we may write

$$886 \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}(\mathcal{C}, k) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\max_{\mathbf{x}'' \in \mathcal{B}_0(\mathbf{a} \odot \mathbf{x}, k)} \ell(\mathcal{C}; \mathbf{a}^{-1} \odot \mathbf{x}'', y) \right].$$

887 Let \mathcal{C}' be the classifier defined that $\mathcal{C}'(\mathbf{x}) := \mathcal{C}(\mathbf{a} \odot \mathbf{x})$. With this, we can rewrite the
 888 above as

$$\begin{aligned} 889 \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}(\mathcal{C}, k) &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\max_{\mathbf{x}'' \in \mathcal{B}_0(\mathbf{a} \odot \mathbf{x}, k)} \ell(\mathcal{C}'; \mathbf{x}'', y) \right] \\ 890 &= \mathbb{E}_{(\mathbf{x}', y) \sim \mathcal{D}'} \left[\max_{\mathbf{x}'' \in \mathcal{B}_0(\mathbf{x}', k)} \ell(\mathcal{C}'; \mathbf{x}'', y) \right] \\ 891 &= \mathcal{L}_{\boldsymbol{\mu}', \Sigma'}(\mathcal{C}', k) \\ 892 &\geq \mathcal{L}_{\boldsymbol{\mu}', \Sigma'}^*(k). \end{aligned}$$

894 Comparing this with (D.5) and sending to zero, we realize that $\mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*(k) \geq \mathcal{L}_{\boldsymbol{\mu}', \Sigma'}^*(k)$.
 895 Changing the order of $(\boldsymbol{\mu}, \Sigma)$ and $(\boldsymbol{\mu}', \Sigma')$ and replacing \mathbf{a} with \mathbf{a}^{-1} yields the other
 896 direction and completes the proof. \square

897 Using the above tools, we are now ready to prove Theorem 3.11.

898 *Proof of Theorem 3.11.* Note that since Σ is positive definite, R is also positive
 899 definite and $\zeta_{\min} > 0$. Moreover, we have $R \succeq \zeta_{\min} I_d$. Therefore, using Lemmas D.1
 900 and D.2 above, we realize that for all k , we have

$$901 \quad (\text{D.6}) \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*(k) = \mathcal{L}_{\boldsymbol{u}, R}^*(k) \geq \mathcal{L}_{\boldsymbol{u}, \zeta_{\min} I_d}^*(k).$$

902 Since $\zeta_{\min} I_d$ is diagonal, we may use our lower bound of Theorem 3.8 with $\boldsymbol{\nu} =$
 903 $(\zeta_{\min} I_d)^{-1/2} \boldsymbol{u} = \boldsymbol{u} / \sqrt{\zeta_{\min}}$ to obtain the following bound which holds for all $A \subseteq [d]$

$$904 \quad \mathcal{L}_{\boldsymbol{u}, \zeta_{\min} I_d}^* \left(\frac{1}{\sqrt{\zeta_{\min}}} \|\boldsymbol{u}_A\|_1 \log d \right) \geq \bar{\Phi}(\|\boldsymbol{u}_{A^c}\|_2) - \frac{1}{\log d}.$$

905 The proof is complete by comparing this with (D.6). \square

906 **Appendix E. Proof of Theorem 3.13.**

907 We use the bound in Corollary 3.5 with $F = [\lambda_c : d]$, which simplifies into the
 908 following with $k = \|\boldsymbol{\nu}_{[1:\lambda_c]}\|_1 / \log d$:

$$909 \quad (\text{E.1}) \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}^* \left(\frac{\|\boldsymbol{\nu}_{[1:\lambda_c]}\|_1}{\log d} \right) \leq \frac{1}{\sqrt{2 \log d}} + \bar{\Phi} \left(\|\boldsymbol{\nu}_{[\lambda_c:d]}\|_2 - \frac{\|\boldsymbol{\nu}_{[1:\lambda_c]}\|_1 \|\boldsymbol{\nu}_{[\lambda_c:d]}\|_{\infty}}{\|\boldsymbol{\nu}_{[\lambda_c:d]}\|_2} \frac{16\sqrt{2}}{\sqrt{\log d}} \right).$$

910 Note that we have

$$911 \quad (\text{E.2}) \quad \|\boldsymbol{\nu}_{[\lambda_c:d]}\|_2^2 = 1 - \|\boldsymbol{\nu}_{[1:\lambda_c-1]}\|_2^2 \geq 1 - c^2.$$

912 On the other hand,

$$913 \quad (\text{E.3}) \quad \begin{aligned} \|\boldsymbol{\nu}_{[1:\lambda_c]}\|_1 \|\boldsymbol{\nu}_{[\lambda_c:d]}\|_{\infty} &= \|\boldsymbol{\nu}_{[1:\lambda_c]}\|_1 |\nu_{\lambda_c}| \\ &\leq \|\boldsymbol{\nu}_{[1:\lambda_c]}\|_2^2 \\ &\leq \|\boldsymbol{\nu}\|_2^2 \\ &= 1 \end{aligned}$$

914 Substituting (E.2) and (E.3) back into (E.1), we get

$$915 \quad (\text{E.4}) \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}^* \left(\frac{\|\boldsymbol{\nu}_{[1:\lambda_c]}\|_1}{\log d} \right) \leq \frac{1}{\sqrt{2 \log d}} + \bar{\Phi} \left(\sqrt{1 - c^2} - \frac{16\sqrt{2}}{\sqrt{1 - c^2} \sqrt{\log d}} \right)$$

916 Furthermore, with $A = [1 : \lambda_c]$, the bound in Theorem 3.8 implies that

$$917 \quad (\text{E.5}) \quad \mathcal{L}_{\boldsymbol{\mu}, \Sigma}^*(\|\boldsymbol{\nu}_{[1:\lambda_c]}\|_1 \log d) \geq \bar{\Phi}(\sqrt{1 - c^2}) - \frac{1}{\log d}.$$

918 This completes the proof.

919 **Appendix F. Proof of Theorem 3.18.**

920 Note that since $\Psi_d(\cdot)$ is nondecreasing for all d , if $\Psi_{\infty}(c) = \lim \Psi_d(c)$ exists,
 921 $\Psi_{\infty}(\cdot)$ is indeed nondecreasing and $\Psi_{\infty}(0)$ is well-defined.

922 Part 1 First we assume that $c \in (0, 1)$. Since $\Psi_{\infty}(c) = \lim \Psi_d(c)$ and
 923 $\log \log d / \log d \rightarrow 0$, $\limsup \log_d k_d < \Psi_{\infty}(c)$ implies that for d large enough, we have

$$924 \quad \log_d k_d < \Psi_d(c) - \frac{\log \log d}{\log d}.$$

925 Thereby,

$$926 \quad \log_d k_d < \log_d \|\boldsymbol{\nu}_{[1:\lambda_c^{(d)}]}^{(d)}\|_1 - \frac{\log \log d}{\log d} = \log_d \frac{\|\boldsymbol{\nu}_{[1:\lambda_c^{(d)}]}^{(d)}\|_1}{\log d}.$$

927 Hence, Theorem 3.13 implies that

$$928 \quad \mathcal{L}_d^*(k_d) \leq \mathcal{L}_d^* \left(\frac{\|\boldsymbol{\nu}_{[1:\lambda_c^{(d)}]}^{(d)}\|_1}{\log d} \right) \leq \frac{1}{\sqrt{2} \log d} + \bar{\Phi} \left(\sqrt{1-c^2} - \frac{16\sqrt{2}}{\sqrt{1-c^2} \sqrt{\log d}} \right).$$

929 Sending d to infinity, we get $\limsup \mathcal{L}_d^*(k_d) \leq \bar{\Phi}(\sqrt{1-c^2})$. Next, we consider $c =$
 930 0 . Note that since $\Psi_\infty(\cdot)$ is nondecreasing, $\limsup \log_d k_d < \Psi_\infty(0)$ implies that
 931 $\limsup \log_d k_d < \Psi_\infty(c)$ for all $c > 0$. Consequently, the above bound implies
 932 that $\limsup \mathcal{L}^*(k_d) \leq \bar{\Phi}(\sqrt{1-c^2})$ for all $c > 0$. Sending c to zero, we realize that
 933 $\limsup \mathcal{L}^*(k_d) \leq \bar{\Phi}(0)$. Finally, for $c = 1$, note that the classifier that always outputs 1
 934 has misclassification error at most $1/2$. This implies that irrespective of the sequence
 935 k_d , we always have $\limsup \mathcal{L}_d^*(k_d) \leq 1/2 = \bar{\Phi}(\sqrt{1-1^2})$ and the bound automatically
 936 holds for $c = 1$.

937 Part 2 First we assume that $c \in (0, 1]$. Similar to the first part, $\liminf \log_d k_d >$
 938 $\Psi_\infty(c)$ implies that for d large enough, we have

$$939 \quad \log_d k_d > \Psi_d(c) + \frac{\log \log d}{\log d},$$

940 and

$$941 \quad \log_d k_d > \log_d \|\boldsymbol{\nu}_{[1:\lambda_c^{(d)}]}^{(d)}\|_1 + \frac{\log \log d}{\log d} = \log_d(\log d \|\boldsymbol{\nu}_{[1:\lambda_c^{(d)}]}^{(d)}\|_1).$$

942 Hence, Theorem 3.13 implies that

$$943 \quad \mathcal{L}_d^*(k_d) \geq \mathcal{L}_d^*(\log d \|\boldsymbol{\nu}_{[1:\lambda_c^{(d)}]}^{(d)}\|_1) \geq \bar{\Phi}(\sqrt{1-c^2}) - \frac{1}{\log d}.$$

944 Sending $d \rightarrow \infty$, we get $\liminf \mathcal{L}_d^*(k_d) \geq \bar{\Phi}(\sqrt{1-c^2})$. For the case $c = 0$, note
 945 that irrespective of the sequence k_d , we always have $\mathcal{L}_d^*(k_d) \geq \mathcal{L}_d^*(0) = \bar{\Phi}(\sqrt{1-0^2})$.
 946 Thereby, the result for $c = 0$ automatically holds.