

# More Supervision, Less Computation: Statistical-Computational Tradeoffs in Weakly Supervised Learning

**Authored by:**

Han Liu  
Constantine Caramanis  
Zhaoran Wang  
Zhuoran Yang  
Xinyang Yi

## **Abstract**

We consider the weakly supervised binary classification problem where the labels are randomly flipped with probability  $1-\alpha$ . Although there exist numerous algorithms for this problem, it remains theoretically unexplored how the statistical accuracies and computational efficiency of these algorithms depend on the degree of supervision, which is quantified by  $\alpha$ . In this paper, we characterize the effect of  $\alpha$  by establishing the information-theoretic and computational boundaries, namely, the minimax-optimal statistical accuracy that can be achieved by all algorithms, and polynomial-time algorithms under an oracle computational model. For small  $\alpha$ , our result shows a gap between these two boundaries, which represents the computational price of achieving the information-theoretic boundary due to the lack of supervision. Interestingly, we also show that this gap narrows as  $\alpha$  increases. In other words, having more supervision, i.e., more correct labels, not only improves the optimal statistical accuracy as expected, but also enhances the computational efficiency for achieving such accuracy.

## **1 Paper Body**

Practical classification problems usually involve corrupted labels. Specifically, let  $\{(\mathbf{x}_i, z_i)\}_{i=1}^n$  be  $n$  independent data points, where  $\mathbf{x}_i \in \mathbb{R}^d$  is the covariate vector and  $z_i \in \{0, 1\}$  is the uncorrupted label. Instead of observing  $\{(\mathbf{x}_i, z_i)\}_{i=1}^n$ , we observe  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  in which  $y_i$  is the corrupted label. In detail, with probability  $(1 - \alpha)$ ,  $y_i$  is chosen uniformly at random over  $\{0, 1\}$ , and with probability  $\alpha$ ,  $y_i = z_i$ . Here  $\alpha \in [0, 1]$  quantifies the degree of

supervision: a larger  $\gamma$  indicates more supervision since we have more uncorrupted labels in this case. In this paper, we are particularly interested in the effect of  $\gamma$  on the statistical accuracy and computational efficiency for parameter estimation in this problem, particularly in the high dimensional settings where the dimension  $d$  is much larger than the sample size  $n$ . There exists a vast body of literature on binary classification problems with corrupted labels. In particular, the study of randomly perturbed labels dates back to [1] in the context of random classification noise model. See, e.g., [12, 20] for a survey. Also, classification problems with missing labels are also extensively studied in the context of semi-supervised or weakly supervised learning by [14, 17, 21], among others. Despite the extensive study on this problem, its information-theoretic and computational boundaries remain unexplored in terms of theory. In a nutshell, the information-theoretic boundary refers to the optimal statistical accuracy achievable by any algorithms, while the computational boundary refers to the optimal statistical accuracy achievable by the algorithms under a computational budget that has a polynomial dependence on the problem scale  $(d, n)$ . Moreover, it remains unclear how these two boundaries vary along with  $\gamma$ . One interesting question to ask is 29th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

how the degree of supervision affects the fundamental statistical and computational difficulties of this problem, especially in the high dimensional regime. In this paper, we sharply characterize both the information-theoretic and computational boundaries of the weakly supervised binary classification problems under the minimax framework. Specifically, we consider the Gaussian generative model where  $X|Z = z \sim N(\mu_z, \Sigma)$  and  $z \in \{0, 1\}$  is the true label. Suppose  $\{(x_i, z_i)\}_{i=1}^n$  are  $n$  independent samples of  $(X, Z)$ . We assume that  $\{y_i\}_{i=1}^n$  are generated from  $\{z_i\}_{i=1}^n$  in the aforementioned manner. We focus on the high dimensional regime, where  $d \gg n$  and  $\mu_1 - \mu_0$  is  $s$ -sparse, i.e.,  $\mu_1 - \mu_0$  has  $s$  nonzero entries. We are interested in estimating  $\mu_1 - \mu_0$  from the observed samples  $\{(x_i, y_i)\}_{i=1}^n$ . By a standard reduction argument [24], the fundamental limits of this estimation task are captured by a hypothesis testing problem, namely,  $H_0 : \mu_1 - \mu_0 = 0$  versus  $H_1 : \mu_1 - \mu_0$  is  $s$ -sparse and  $(\mu_1 - \mu_0)^T (\mu_1 - \mu_0) := \gamma n \neq 0$ ,

(1.1)

where  $\gamma n$  denotes the signal strength that scales with  $n$ . Consequently, we focus on studying the fundamental limits of  $\gamma n$  for solving this hypothesis testing problem.  $\gamma n$  is

Efficient

$\gamma n \geq s^2 \log d$ ,  $\gamma n = o(s^2 \log d)$ ,  $\gamma n \geq s \log d$ ,  $\gamma n = o(s \log d)$

$\gamma n = ?$

?

$\gamma n \geq s^2 \log d$ ,  $\gamma n \geq 2n$

?

Intractable Impossible 0

?

$\gamma n = o$

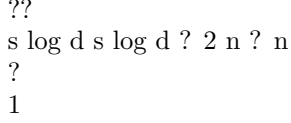


Figure 1: Computational-statistical phase transitions for weakly supervised binary classification. Here  $\gamma$  denotes the degree of supervision, i.e., the label is corrupted to be uniformly random with probability  $1 - \gamma$ , and  $n$  is the signal strength, which is defined in (1.1). Here  $a \wedge b$  denotes  $\min\{a, b\}$ .

Our main results are illustrated in Figure 1. Specifically, we identify the impossible, intractable, and efficient regimes for the statistical-computational phase transitions under certain regularity conditions. (i) For  $n = o[s \log d/n \wedge (1/2 \wedge s \log d/n)]$ , any algorithm is asymptotically powerless in solving the hypothesis testing problem. (ii) For  $n = \Theta[s \log d/n \wedge (1/2 \wedge s \log d/n)]$  and  $n = o[s^2/n \wedge (1/2 \wedge s \log d/n)]$ , any tractable algorithm that has a polynomial oracle complexity under an extension of the statistical query model [18] is asymptotically powerless. We will rigorously define the computational model in §2. (iii) For  $n = \Theta[s^2/n \wedge (1/2 \wedge s \log d/n)]$ , there is an efficient algorithm with a polynomial oracle complexity that is asymptotically powerful in solving the testing problem. Here  $s \log d/n \wedge (1/2 \wedge s \log d/n)$  gives the information-theoretic boundary, while  $s^2/n \wedge (1/2 \wedge s \log d/n)$  gives the computational boundary. Moreover, by a reduction from the estimation problem to the testing problem, these boundaries for testing imply the ones for estimating  $\gamma \wedge 1$  as well. Consequently, there exists a significant gap between the computational and information-theoretic boundaries for small  $\gamma$ . In other word, to achieve the information-theoretic boundary, one has to pay the price of intractable computation. As  $\gamma$  tends to one, this gap between computational and information-theoretic boundaries narrows and eventually vanishes. This indicates that, having more supervision not only improves the statistical accuracy, as shown by the decay of information-theoretic boundary in Figure 1, but more importantly, enhances the computational efficiency by reducing the computational price for attaining information-theoretic optimality. This phenomenon of more supervision, less computation is observed for the first time in this paper. 1.1

#### More Related Work, Our Contribution, and Notation

Besides the aforementioned literature on weakly supervised learning and label corruption, our work is also connected to a recent line of work on statistical-computational tradeoffs [25, 8, 13, 15, 19, 26–28]. In comparison, we quantify the statistical-computational tradeoffs for weakly supervised learning for the first time. Furthermore, our results are built on an oracle computational model

in [8] that slightly extends the statistical query model [18], and hence do not hinge on unproven conjectures on computational hardness like planted clique. Compared with our work, [8] focuses on the computational hardness of learning heterogeneous models, whereas we consider the interplay between supervision and statistical-computational tradeoffs. A similar computational model is used in [27] to study structural normal mean model and principal component analysis,

which exhibit different statistical-computational phase transitions. In addition, our work is related to sparse linear discriminant analysis and two-sample testing of sparse means, which correspond to our special cases of  $\eta = 1$  and  $\eta = 0$ , respectively. See, e.g., [7, 23] for details. In contrast with their results, our results capture the effects of  $\eta$  on statistical and computational tradeoffs. In summary, the contribution of our work is two-fold: (i) We characterize the computational and statistical boundaries of the weakly supervised binary classification problem for the first time. Compared with existing results for other models, our results do not rely on unproven conjectures. (ii) Based on our theoretical characterization, we propose the “more supervision, less computation” phenomenon, which is observed for the first time. Notation. We denote the  $\ell_2$ -divergence between two distributions  $P, Q$  by  $D_2(P, Q)$ . For two nonnegative sequences  $a_n, b_n$  indexed by  $n$ , we use  $a_n = o(b_n)$  as a shorthand for  $\lim_{n \rightarrow \infty} a_n / b_n = 0$ . We say  $a_n = \Theta(b_n)$  if  $a_n / b_n \asymp c$  for some absolute constant  $c > 0$  when  $n$  is sufficiently large. We use  $a \vee b$  and  $a \wedge b$  to denote  $\max\{a, b\}$  and  $\min\{a, b\}$ , respectively. For any positive integer  $k$ , we denote  $\{1, 2, \dots, k\}$  by  $[k]$ . For  $v \in \mathbb{R}^d$ , we denote by  $\|v\|_p$  the  $\ell_p$ -norm of  $v$ . In addition, we denote the operator norm of a matrix  $A$  by  $\|A\|_2$ .

## 2 Background

In this section, we formally define the statistical model for weakly supervised binary classification. Then we follow it with the statistical query model that connects computational complexity and statistical optimality. 2.1

### Problem Setup

Consider the following Gaussian generative model for binary classification. For a random vector  $X \in \mathbb{R}^d$  and a binary random variable  $Z \in \{0, 1\}$ , we assume  $X|Z=0 \sim N(\mu_0, \Sigma)$ ,

$$X|Z=1 \sim N(\mu_1, \Sigma), \quad (2.1)$$

where  $P(Z=0) = P(Z=1) = 1/2$ . Under this model, the optimal classifier by Bayes rule corresponds to the Fisher’s linear discriminative analysis (LDA) classifier. In this paper, we focus on the noisy label setting where true label  $Z$  is replaced by a uniformly random label in  $\{0, 1\}$  with probability  $1 - \eta$ . Hence,  $\eta$  characterizes the degree of supervision in the model. In specific, if  $\eta = 0$ , we observe the true label  $Z$ , thus the problem belongs to supervised learning. Whereas if  $\eta = 1$ , the observed label is completely random, which contains no information of the model in (2.1). This setting is thus equivalent to learning a Gaussian mixture model, which is an unsupervised problem. In the general setting with noisy labels, we denote the observed label by  $Y$ , which is linked to the true label  $Z$  via  $P(Y=Z) = (1 + \eta)/2$ ,  $P(Y=1-Z) = (1 - \eta)/2$ . (2.2)

We consider the hypothesis testing problem of detecting whether  $\mu_0 = \mu_1$  given  $n$  i.i.d. samples  $\{y_i, x_i\}_{i=1}^n$  of  $(Y, X)$ , namely  $H_0 : \mu_0 = \mu_1$  versus  $H_1 : \mu_0 \neq \mu_1$ .

$$(2.3)$$

We focus on the high dimensional and sparse regime, where  $d \gg n$  and  $\mu_0 - \mu_1$  is  $s$ -sparse, i.e.,  $\mu_0 - \mu_1 \in B_0(s)$ , where  $B_0(s) := \{\mu \in \mathbb{R}^d : \|\mu\|_0 \leq s\}$ .

s}. Throughout this paper, use the sample size  $n$  to drive the asymptotics. We introduce a shorthand notation  $\theta := (\theta_0, \theta_1, \theta, \theta)$  to represent the parameters of the aforementioned model. Let  $P_\theta$  be the joint distribution of  $(Y, X)$  under our statistical model with parameter  $\theta$ , and  $P_n^\theta$  be the product distribution of  $n$  i.i.d. samples accordingly. We denote the parameter spaces of the null and alternative hypotheses by  $G_0$  and  $G_1$  respectively. For any test function  $\phi : \{(y_i, x_i)\}_{i=1}^n \rightarrow \{0, 1\}$ , the classical testing risk is defined as the summation of 3

$$\text{type-I and type-II errors, namely } R_n(\phi; G_0, G_1) := \sup_{\theta \in G_1} P_n^\theta(\phi = 1) + \sup_{\theta \in G_0} P_n^\theta(\phi = 0). \quad (2.3)$$

The minimax risk is defined as the smallest testing risk of all possible test functions, that is,  $R_n^*(G_0, G_1) := \inf_{\phi} R_n(\phi; G_0, G_1)$ , where the infimum is taken over all measurable test functions. Intuitively, the separation between two Gaussian components under  $H_1$  and the covariance matrix  $\Sigma$  together determine the hardness of detection. To characterize such dependence, we define the signal-to-noise ratio (SNR) as  $\text{SNR}(\theta) := (\theta_0 - \theta_1)^2 / \text{tr}(\Sigma(\theta_0 + \theta_1))$ . For any nonnegative sequence  $\{\gamma_n\}_{n=1}^\infty$ , let  $G_1(\gamma_n) := \{\theta : \text{SNR}(\theta) \geq \gamma_n\}$  be a sequence of alternative parameter spaces with minimum separation  $\gamma_n$ . The following minimax rate characterizes the information-theoretic limits of the detection problem. Definition 2.1 (Minimax rate). We say a sequence  $\{\gamma_n\}_{n=1}^\infty$  is a minimax rate if for any sequence  $\{\gamma_n\}_{n=1}^\infty$  satisfying  $\gamma_n = o(\gamma_n^*)$ , we have  $\liminf_{n \rightarrow \infty} R_n^*(G_0, G_1(\gamma_n)) = 1$ ; for any sequence  $\{\gamma_n\}_{n=1}^\infty$  satisfying  $\gamma_n = \Omega(\gamma_n^*)$ , we have  $\liminf_{n \rightarrow \infty} R_n^*(G_0, G_1(\gamma_n)) = 0$ . The minimax rate in Definition 2.1 characterizes the statistical difficulty of the testing problem. However, it fails to shed light on the computational efficiency of possible testing algorithms. The reason is that this concept does not make any computational restriction on the test functions. The minimax risk in (2.4) might be attained only by test functions that have exponential computational complexities. This limitation of Definition 2.1 motivates us to study statistical limits under computational constraints. 2.2

### Computational Model

Statistical query models [8, 11, 18, 27] capture computational complexity by characterizing the total number of rounds an algorithm interacts with data. In this paper, we consider the following statistical query model, which admits bounded query functions but allows the responses of query functions to be unbounded. Definition 2.2 (Statistical query model). In the statistical query model, an algorithm  $A$  is allowed to query an oracle  $T$  rounds, but not to access data  $\{(y_i, x_i)\}_{i=1}^n$  directly. At each round,  $A$  queries the oracle  $r$  with a query function  $q \in \mathcal{Q}_A$ , in which  $\mathcal{Q}_A \subseteq \{q : \{0, 1\}^n \rightarrow [-M, M]\}$  denotes the query space of  $A$ . The oracle  $r$  outputs a realization of a random variable  $Z_q \in \mathbb{R}$  satisfying  $\mathbb{E}[Z_q] = \mathbb{E}[q(Y, X)]$  and  $\text{var}(Z_q) \leq \sigma^2$ , where  $q \in \mathcal{Q}_A$ .  $\gamma_q = [\text{tr}(\Sigma(q)) + \log(1/\delta)] \leq M/n$

$$\delta \leq 2[\text{tr}(\Sigma(q)) + \log(1/\delta)] \leq (M^2 + \mathbb{E}[q(Y, X)]^2) n. \quad (2.5)$$

Here  $\delta \in (0, 1)$  is the tolerance parameter and  $\delta \in [0, 1)$  is the tail probability.

The quantity  $\varphi(QA) \geq 0$  in  $\varphi_q$  measures the capacity of  $QA$  in logarithmic scale, e.g., for countable  $QA$ ,  $\varphi(QA) = \log(\#QA)$ . The number  $T$  is defined as the oracle complexity. We denote by  $R[\varphi, n, T, \varphi(QA)]$  the set of oracles satisfying (2.5), and by  $A(T)$  the family of algorithms that queries an oracle no more than  $T$  rounds. This version of statistical query model is used in [8], and reduces to the VSTAT model proposed in [9,11] by the transformation  $q_\varphi(y, x) = q(y, x)/(2M) + 1/2$  for any  $q \in QA$ . The computational model in Definition 2.2 enables us to handle query functions that are bounded by an unknown and fixed number  $M$ . Note that by incorporating the tail probability  $\varphi$ , the response  $Z_q$  is allowed to be unbounded. To understand the intuition behind Definition 2.2, we remark that (2.5) resembles the Bernstein's inequality for bounded random variables [25]: 
$$P\left(\sum_{i=1}^n q(Y_i, X_i) \geq E[q(Y, X)] + t\right) \leq 2 \exp\left(-\frac{t^2}{2\text{Var}[q(Y, X)] + Mt}\right). \quad (2.6)$$

We first replace  $\text{Var}[q(Y, X)]$  by its upper bound  $M^2 \varphi\{E[q(Y, X)]\}^2$ , which is tight when  $q$  takes  $n$  values in  $\{[-M, M]\}$ . Then inequality (2.5) is obtained by replacing  $\sum_{i=1}^n q(Y_i, X_i)$  in (2.6) by  $Z_q$  and then bounding the suprema over the query space  $QA$ . In the definition of  $\varphi_q$  in (2.5), we

incorporate the effect of uniform concentration over the query space  $QA$  by adding the quantity  $\varphi(QA)$ , which measures the capacity of  $QA$ . In addition, under the Definition 2.2, the algorithm  $A$  does not interact directly with data. Such a restriction characterizes the fact that in statistical problems, the effectiveness of an algorithm only depends on the global statistical properties, not the information of individual data points. For instance, algorithms that only rely on the convergence of the empirical distribution to the population distribution are contained in the statistical query model; whereas algorithms that hinge on the first data point  $(y_1, x_1)$  is not allowed. This restriction captures a vast family of algorithms in statistics and machine learning, including applying gradient method to maximize likelihood function, matrix factorization algorithms, expectation-maximization algorithms, and sampling algorithms [9]. Based on the statistical query model, we study the minimax risk under oracle complexity constraints. For the testing problem (2.3), let  $A(T_n)$  be a class of testing algorithms under the statistical query model with query complexity no more than  $T_n$ , with  $\{T_n\}_{n \geq 1}$  being a sequence of positive integers depending on the sample size  $n$ . For any  $A \in A(T_n)$  and any oracle  $r \in R[\varphi, n, T_n, \varphi(QA)]$  that responds to  $A$ , let  $H(A, r)$  be the set of test functions that deterministically depend on  $A$ 's queries to the oracle  $r$  and the corresponding responses. We use  $P_\varphi$  to denote the distribution of the random variables returned by oracle  $r$  when the model parameter is  $\varphi$ . For a general hypothesis testing problem, namely,  $H_0: \varphi \in G_0$  versus  $H_1: \varphi \in G_1$ , the minimax testing risk with respect to an algorithm  $A$  and a statistical oracle  $r \in R[\varphi, n, T_n, \varphi(QA)]$  is defined as 
$$R_n(G_0, G_1; A, r) := \inf \sup P_\varphi(\varphi \in G_1) + \sup P_\varphi(\varphi \in G_0). \quad (2.7)$$
 
$$\varphi \in G_0$$

$\varphi \in G_1$

Compared with the classical minimax risk in (2.4), the new notion in (2.7) incorporates the computational budgets via oracle complexity. In specific, we only consider the test functions obtained by an algorithm with at most  $T_n$  queries

to a statistical oracle. If  $T_n$  is a polynomial of the dimensionality  $d$ , (2.7) characterizes the statistical optimality of computational efficient algorithms. This motivates us to define the computationally tractable minimax rate, which contrasts with Definition 2.1. Definition 2.3 (Computationally tractable minimax rate). Let  $G_1(\gamma_n) := \{\gamma : \gamma(\gamma) \leq \gamma_n\}$  be a sequence of model spaces with minimum separation  $\gamma_n$ , where  $\gamma(\gamma)$  is the SNR. A sequence  $\{\gamma_n\}_{n=1}^\infty$  is called a computationally tractable minimax rate if: For any sequence  $\{\gamma_n\}_{n=1}^\infty$  satisfying  $\gamma_n = o(\gamma^*)$ , any constant  $\epsilon > 0$ , and any  $A \in A(d, \epsilon)$ , there exists an oracle  $r \in R[\gamma, n, T_n, \gamma(QA)]$  such that  $\lim_{n \rightarrow \infty} R_n[G_0, G_1(\gamma_n); A, r] = 1$ ; For any sequence  $\{\gamma_n\}_{n=1}^\infty$  satisfying  $\gamma_n = \gamma^*$ , there exist a constant  $\epsilon > 0$  and an algorithm  $A \in A(d, \epsilon)$  such that, for any  $r \in R[\gamma, n, T_n, \gamma(QA)]$ , we have  $\lim_{n \rightarrow \infty} R_n[G_0, G_1(\gamma_n); A, r] = 0$ .

3

### Main Results

Throughout this paper, we assume that the covariance matrix  $\Sigma$  in (2.1) is known. Specifically, for some positive definite  $\Sigma \in \mathbb{R}^{d \times d}$ , the parameter spaces of the null and alternative hypotheses are defined as  $G_0(\gamma) := \{\gamma = (\gamma_0, \gamma_1, \gamma_2, \gamma_3) : \gamma_0 \in \mathbb{R}^d\}$ ,

(3.1) d

$G_1(\gamma; \gamma_n) := \{\gamma = (\gamma_0, \gamma_1, \gamma_2, \gamma_3) : \gamma_0 \in \mathbb{R}^d, \gamma_1 \in \mathbb{R}^d, \gamma_2 \in B_0(s), \gamma_3 \in \gamma_n\}$ .

Accordingly, the testing problem of detecting whether  $\gamma_0 = \gamma_1$  is to distinguish  $H_0 : \gamma \in G_0(\gamma)$  versus  $H_1 : \gamma \in G_1(\gamma; \gamma_n)$ .

(3.2) (3.3)

In §3.1, we present the minimax rate of the detection problem from an information-theoretic perspective. In §3.2, under the statistical query model introduced in §2.2, we provide a computational lower bound and a nearly matching upper bound that is achieved by an efficient testing algorithm. 3.1

### Information-theoretic Limits

Now we turn to characterize the minimax rate given in Definition 2.1. For parameter spaces (3.1) and (3.2) with known  $\Sigma$ , we show that in highly sparse setting where  $s = o(d)$ , we have  $\gamma_n^* = s \log d/n \asymp (1/2) s \log d/n$ , (3.4) 5

To prove (3.4), we first present a lower bound which shows that the hypothesis testing problem in (3.3) is impossible if  $\gamma_n = o(\gamma_n^*)$ . Theorem 3.1. For the hypothesis testing problem in (3.3) with known  $\Sigma$ , we assume that there exists a small constant  $\epsilon > 0$  such that  $s = o(d^{1/2-\epsilon})$ . Let  $\gamma_n^*$  be defined in (3.4). For any sequence  $\{\gamma_n\}_{n=1}^\infty$  such that  $\gamma_n = o(\gamma_n^*)$ , any hypothesis test is asymptotically powerless, namely,  $\limsup_{n \rightarrow \infty} R_n[G_0(\gamma), G_1(\gamma; \gamma_n)] = 1$ .

n??

?

By Theorem 3.1, we observe a phase transition in the necessary SNR for powerful detection when  $\gamma$  decreases from one to zero. Starting with rate  $s \log d/n$  in the supervised setting where  $\gamma = 1$ , the required SNR gradually increases as label qualities decrease. Finally, when  $\gamma$  reaches zero, which corresponds to the unsupervised setting, powerful detection requires the SNR to be  $\gamma \asymp (s \log d/n)$ . It is worth noting that when  $\gamma = (s \log d/n)^{1/4}$ , we still have  $(n^3 s$

$\log d)^{1/4}$  uncorrupted labels. However, our lower bound (along with the upper bound shown in Theorem 3.2) indicates that the information contained in these uncorrupted labels are buried in the noise, and cannot essentially improve the detection quality compared with the unsupervised setting. Next we establish a matching upper bound for the detection problem in (3.3). We denote the condition number of the covariance matrix  $\Sigma$  by  $\kappa$ , i.e.,  $\kappa := \lambda_{\max}(\Sigma)/\lambda_{\min}(\Sigma)$ , where  $\lambda_{\max}(\Sigma)$  and  $\lambda_{\min}(\Sigma)$  are the largest and smallest eigenvalues of  $\Sigma$ , respectively. Note that marginally  $Y$  is uniformly distributed over  $\{0, 1\}$ . For ease of presentation, we assume that the sample size is  $2n$  and each class contains exactly  $n$  data points. Note that we can always discard some samples in the larger class to make the sample sizes of both classes to be equal. Due to the law of large numbers, this trick will not affect the analysis of sample complexity in the sense of order wise.  $d$  Given  $2n$  i.i.d. samples  $\{(y_i, x_i)\}_{i=1}^{2n}$  of  $(Y, X) \in \{0, 1\} \times \mathbb{R}^d$ , we define  $w_i = (x_i^2 - x_{i+1}^2)/2$ , for all  $i \in [n]$ .

(0) In addition, we split the dataset  $\{(y_i, x_i)\}_{i=1}^{2n}$  into two disjoint parts  $\{(0, x_i)\}_{i=1}^n$  and  $\{(1, x_i)\}_{i=1}^n$ .

(1)

$u_i = x_i$

(0)

$y_i = x_i$ , for all  $i \in [n]$ .

(3.5)

(1) and  $\{(1, x_i)\}_{i=1}^n$ ,

(3.6)

We note that computing sample differences in (3.5) and (3.6) is critical for our problem because we focus on detecting the difference between  $\Sigma_0$  and  $\Sigma_1$ , and computing differences can avoid estimating  $\mathbb{E}P(X)$  that might be dense. For any integer  $s \leq [d]$ , we define  $B_2(s) := B_0(s) \cap S^{d-1}$  as the set of  $s$ -sparse vectors on the unit sphere in  $\mathbb{R}^d$ . With  $\{w_i\}_{i=1}^n$  and  $\{u_i\}_{i=1}^n$ , we introduce two test functions  $\phi_1, \phi_2 : \mathbb{R}^d \rightarrow \mathbb{R}$  as  $\phi_1(v) = \frac{1}{n} \sum_{i=1}^n (v \cdot w_i)^2$ ,  $\phi_2(v) = \frac{1}{n} \sum_{i=1}^n (v \cdot u_i)^2$ , (3.7)  $\phi_1(v) = \frac{1}{n} \sum_{i=1}^n (v \cdot w_i)^2$ ,  $\phi_2(v) = \frac{1}{n} \sum_{i=1}^n (v \cdot u_i)^2$ , (3.8)  $\phi_1(v) = \frac{1}{n} \sum_{i=1}^n (v \cdot w_i)^2$ ,  $\phi_2(v) = \frac{1}{n} \sum_{i=1}^n (v \cdot u_i)^2$  where  $\gamma_1, \gamma_2 \in (0, 1]$  are algorithmic parameters that will be specified later. To provide some intuitions, we consider the case where  $\Sigma = I$ . Test function  $\phi_1$  seeks a sparse direction that explains the most variance of  $w_i$ . Therefore, such a test is closely related to the sparse principal component detection problem [3]. Test function  $\phi_2$  simply selects the coordinate of  $n$   $u_i$  that has the largest magnitude and compares it with  $\gamma_2$ . This test is closely related to detecting sparse normal mean in high dimensions [16]. Based on these two ingredients, we construct our final testing function  $\phi$  as  $\phi = \phi_1 \vee \phi_2$ , i.e., if any of  $\phi_1$  and  $\phi_2$  is true, then  $\phi$  rejects the null. The following theorem establishes a sufficient condition for test function  $\phi$  to be asymptotically powerful. Theorem 3.2. Consider the testing problem (3.3) where  $\Sigma$  is known and has condition number  $\kappa$ . For test functions  $\phi_1$  and  $\phi_2$  defined in (3.7) and (3.8) with parameters  $\gamma_1$  and  $\gamma_2$  given by  $\gamma_1 = \frac{s \log(ed/s)}{n}$ ,  $\gamma_2 = 8 \log d/n$ .

We define the ultimate test function as  $\phi = \phi_1 \vee \phi_2$ . We assume that  $s \leq C$



$\geq d$  for some absolute constant  $C_s$  and  $n \geq 64 \lceil s \log(ed/s) \rceil$ . Then if  $\beta_n \geq C \lceil s \log(ed/s)/n \rceil^{1/2} \lceil s \log d/n \rceil$ , (3.9) holds

where  $C$  is an absolute constant, then test function  $\hat{\beta}$  is asymptotically powerful. In specific, we have  $\sup_{\beta \in \mathcal{B}_2(s)} P_n(\beta) \leq 1 + o(1) + o(1/n)$

$$\sup_{\beta \in \mathcal{B}_2(s)} P_n(\beta) \leq 1 + o(1) + o(1/n) \quad (3.10)$$

Theorem 3.2 provides a non-asymptotic guarantee. When  $n$  goes to infinity, (3.10) implies that the test function  $\hat{\beta}$  is asymptotically powerful. When  $s = o(d)$  and  $\beta$  is a constant, (3.9) yields  $\beta_n = \lceil s \log d/n \rceil^{1/2} \lceil s \log d/n \rceil$ , which matches the lower bound given in Theorem 3.1. Thus we conclude that  $\beta_n$  defined in (3.4) is the minimax rate of testing problem in (3.3). We remark that when  $s = \beta(d)$ ,  $\beta = 1$ , i.e., the standard (low-dimensional) setting of two sample testing, the bound provided in (3.9) is sub-optimal as [22] shows that SNR rate  $d/n$  is sufficient for asymptotically powerful detection when  $n = \beta(d)$ . It is thus worth noting that we focus on the highly sparse setting  $s = o(d)$  and provided sharp minimax rate for  $\beta$  in this regime. In the definition of  $\beta_1$  in (3.7), we search over the set  $\mathcal{B}_2(s)$ . Since  $\mathcal{B}_2(s)$  contains  $d$  distinct sets of supports, computing  $\beta_1$  requires exponential running time. 3.2

#### Computational Limits

In this section, we characterize the computationally tractable minimax rate  $\beta_n$  given in Definition 2.3. Moreover, we focus on the setting where  $\beta$  is known a priori and the parameter spaces for the null and alternative hypotheses are defined in (3.1) and (3.2), respectively. The main result is that, in highly sparse setting where  $s = o(d)$ , we have  $\beta_n = s^2/n \lceil s \log d/n \rceil^{1/2}$ . (3.11)

We first present the lower bound in the next result.

**Theorem 3.3.** For the testing problem in (3.3) with  $\beta$  known a priori, we make the same assumptions as in Theorem 3.1. For any sequence  $\{\beta_n\}_{n=1}^\infty$  such that  $\beta_n \geq s^2/n \lceil s \log d/n \rceil^{1/2}$ , (3.12) holds

where  $\beta_n$  is defined in (3.4), any computationally tractable test is asymptotically powerless under the statistical query model. That is, for any constant  $\epsilon > 0$  and any  $A \in \mathcal{A}(d)$ , there exists an oracle  $\mathcal{R}[\beta, n, T_n, \beta(QA)]$  such that  $\lim_{n \rightarrow \infty} \mathbb{E} [G_0(\beta), G_1(\beta, \beta_n); A, \mathcal{R}] = 1$ .

We the lower bound in (3.12) differs from  $\beta_n$  in (3.11) by a logarithmic term when  $\beta$  remark 2 that  $\beta \geq 1/n \lceil s \log d/n \rceil^{1/2}$ . We expect this gap to be eliminated by more delicate analysis under the statistical query model. Now putting Theorems 3.1 and 3.3 together, we describe the more supervision, less computation phenomenon as follows. (i) When  $0 \leq \beta \leq (\log^2 d/n)^{1/4}$ , the computational lower bound implies that the uncorrupted labels are unable to improve the quality of computationally tractable detection compared with the unsupervised setting. In addition, in this region, the gap between  $\beta_n$  and  $\beta$  remains the same. (ii) When  $(\log^2 d/n)^{1/4} \leq \beta \leq (s \log d/n)^{1/4}$ , the information-theoretic lower bound shows that the uncorrupted labels cannot improve the quality of detection compared with unsupervised setting. However, more uncorrupted labels improve the statistical performances of hypothesis tests

that are computationally tractable by shrinking the gap between  $\beta_n$  and  $\beta_n^*$ .  
(iii) When  $(s \log d/n)^{1/4} \leq 1$ , having more uncorrupted labels improves both statistical optimality and the computational efficiency. In specific, in this case, the gap between  $\beta_n$  and  $\beta_n^*$  vanishes and we have  $\beta_n = \beta_n^* = 1/2 \pm s \log d/n$ . Now we derive a nearly matching upper bound under the statistical query model, which establishes the computationally tractable minimax rate together with Theorem 3.3. We construct a computationally efficient testing procedure that combines two test functions which yields the two parts in  $\beta_n^*$  respectively. Similar to  $\beta_n$  defined in (3.7), the first test function discards the information of labels, which works for the purely unsupervised setting where  $\gamma = 0$ . For  $j \in [d]$ , we denote by  $\beta_j$  the  $j$ -th diagonal element of  $\beta$ . Under the statistical query model, we consider the 2d query functions  $q_j(y, x) := x_j / \beta_j - 1 \{ -x_j / \beta_j - R \log d \}$ ,  $q_j(y, x) := (x_j^2 / \beta_j - 1) - 1 \{ -x_j / \beta_j - R \log d \}$ , for all  $j \in [d]$ , (3.14)

where  $R \geq 0$  is an absolute constant. Here we apply truncation to the query functions to obtain bounded queries, which is specified by the statistical query model in Definition 2.2. We denote by  $z_{qj}$  and  $z_{q^2j}$  the realizations of the random variables output by the statistical oracle for query functions  $q_j$  and  $q_j^2$ , respectively. As for the second test function, similar to (3.8), we consider  $q \cdot v(y, x) = (2y - 1) \cdot v \cdot \text{diag}(\beta)^{1/2} x - v \cdot \text{diag}(\beta)^{1/2} x - R \log d$  (3.15)

for all  $v \in B_2(1)$ . We denote by  $Z_{qv}$  the output of the statistical oracle corresponding to query function  $q \cdot v$ . With these 4d query functions, we introduce test functions  $\beta_1 := 1 \sup (z_{q^2j} - z_{qj}^2) / C$ ,  $\beta_2 := 1 \sup z_{qv} - 2$ , (3.16)

where  $\beta_1$  and  $\beta_2$  are positive parameters that will be specified later and  $C$  is an absolute constant. Theorem 3.4. For the test functions  $\beta_1$  and  $\beta_2$  defined in (3.16), we define the ultimate test function as  $\beta = \beta_1 \wedge \beta_2$ . We set  $\beta_1 = R^2 \log d \pm \log(4d/\beta)/n$ ,  $\beta_2 = R \log d \pm \log(4d/\beta)/n$ , (3.17) where  $\beta = o(1)$ . For the hypothesis testing problem in (3.3), we further assume that  $\beta_0 \wedge \beta_1 \geq C_0$  for some constant  $C_0 \geq 0$ . Under the assumption that  $\beta_0 \wedge \beta_1 \geq C_0$ ,  $\sup (\beta_{0,j} - \beta_{1,j})^2 / \beta_j = 1/2 \pm \log^2 d \pm \log(d/\beta)/n \pm \log d \pm \log(d/\beta)/n$ , (3.18)

the risk of  $\beta$  satisfies that  $R_n(\beta) = \sup_{\beta_0 \in \mathcal{G}_0} P(\beta = 1) + \sup_{\beta_1 \in \mathcal{G}_1} P(\beta = 0) \leq 5$ . Here we denote by  $\beta_{0,j}$  and  $\beta_{1,j}$  the  $j$ -th entry of  $\beta_0$  and  $\beta_1$ , respectively. If we set the tail probability of the statistical query model to be  $\beta = 1/d$ , (3.18) shows that  $\beta$  is asymptotically powerful if  $\sup_{j \in [d]} (\beta_{0,j} - \beta_{1,j})^2 / \beta_j = [(1/2 \pm \log^2 d/n) \pm (\log^3 d/n)^{1/2}]$ . When the energy of  $\beta_0 \wedge \beta_1$  is spread over its support,  $\beta_0 \wedge \beta_1 \geq C_0$  and  $\beta_0 \wedge \beta_1 \geq 2/s$  are close. Under the assumption that the condition number  $\kappa$  of  $\beta$  is a constant, (3.18) is implied by  $\beta_n \leq (s^2 \log^3 d/n)^{1/2} \pm (1/2 \pm s \log^3 d/n)$ .

Compared with Theorem 3.3, the above upper bound matches the computational lower bound up to a logarithmic factor and  $\beta_n$  is between  $s^2/n \pm (1/2 \pm s \log d/n)$  and  $(s^2 \log^3 d/n)^{1/2} \pm (1/2 \pm s \log^3 d/n)$ . Note that the

truncation on query functions in (3.13) and (3.14) yields an additional logarithmic term, which could be reduced to  $(s^2 \log d/n)^{1/2} \cdot (1/\sqrt{2} \cdot s \log d/n)$  using more delicate analysis. Moreover, the test function  $\phi_1$  is essentially based on a diagonal thresholding algorithm performed on the covariance matrix  $\Sigma$  of  $X$ . The work in [6] provides a more delicate analysis of this algorithm which establishes the  $s^2/n$  rate. Their algorithm can also be formulated into the statistical query model; we use the simpler version in (3.16) for ease of presentation. Therefore, with more sophisticated proof technique, it can be shown that  $s^2/n \cdot (1/\sqrt{2} \cdot s \log d/n)$  is the critical threshold for asymptotically powerful detection with computational efficiency. 3.3

### Implication for Estimation

Our aforementioned phase transition in the detection problems directly implies the statistical and computational trade-offs in the problem of estimation. We consider the problem of estimating the parameter  $\theta = \theta_0 + \theta_1$  of the binary classification model in (2.1) and (2.2), where  $\theta_1$  is  $s$ -sparse and  $\theta_0$  is known a priori. We assume that the signal to noise ratio is  $\phi(\theta) = \theta_1^T \theta_1 / \theta_0^T \theta_0 = o(n)$ . For any constant  $\epsilon > 0$  and any  $A = A(T)$  with  $T = O(d^\epsilon)$ , suppose we obtain an estimator  $\hat{\theta}$  of  $\theta$  by algorithm  $A$  under the statistical query model. If  $\hat{\theta}$  converges to  $\theta$  in the sense that  $\|\hat{\theta} - \theta\|_2 \leq C \sqrt{n^2 / \phi(\theta)}$ ,

$\|\hat{\theta} - \theta\|_2 \leq C \sqrt{n^2 / \phi(\theta)}$  then  $\|\hat{\theta} - \theta\|_2 = o(n)$ . Thus the test function  $\phi = 1\{\theta_1^T \theta_1 \geq \epsilon n\}$  is asymptotically powerful, which contradicts the computational lower bound in Theorem 3.3. Therefore, there exists a constant  $C$  such that  $\hat{\theta}$  constructed from polynomial number of queries. Acknowledgments

We would like to thank Vitaly Feldman for valuable discussions. 8

## 2 References

- [1] Angluin, D. and Laird, P. (1988). Learning from noisy examples. *Machine Learning*, 2 343-370.
- [2] Berthet, Q. and Rigollet, P. (2013). Computational lower bounds for sparse PCA. In *Conference on Learning Theory*.
- [3] Berthet, Q. and Rigollet, P. (2013). Optimal detection of sparse principal components in high dimension. *The Annals of Statistics*, 41 1780-1815.
- [4] Chandrasekaran, V. and Jordan, M. I. (2013). Computational and statistical tradeoffs via convex relaxation. *Proceedings of the National Academy of Sciences*, 110 1181-1190.
- [5] Chen, Y. and Xu, J. (2014). Statistical-computational tradeoffs in planted problems and submatrix localization with a growing number of clusters and submatrices. *arXiv preprint arXiv:1402.1267*.
- [6] Deshpande, Y. and Montanari, A. (2014). Sparse PCA via covariance thresholding. In *Advances in Neural Information Processing Systems*.
- [7] Fan, J., Feng, Y. and Tong, X. (2012). A road to classification in high dimensional space: The regularized optimal affine discriminant. *Journal of the Royal Statistical Society: Series B*, 74 745-771.
- [8] Fan, J., Liu, H., Wang, Z. and Yang, Z. (2016). Curse of heterogeneity: Computational barriers in sparse mixture models and phase re-

trieval. Manuscript. [9] Feldman, V., Grigorescu, E., Reyzin, L., Vempala, S. and Xiao, Y. (2013). Statistical algorithms and a lower bound for detecting planted cliques. In ACM Symposium on Theory of Computing. [10] Feldman, V., Guzman, C. and Vempala, S. (2015). Statistical query algorithms for stochastic convex optimization. arXiv preprint arXiv:1512.09170. [11] Feldman, V., Perkins, W. and Vempala, S. (2015). On the complexity of random satisfiability problems with planted solutions. In ACM Symposium on Theory of Computing. [12] Fr̄nay, B. and Verleysen, M. (2014). Classification in the presence of label noise: A survey. IEEE Transactions on Neural Networks and Learning Systems, 25 845–869. [13] Gao, C., Ma, Z. and Zhou, H. H. (2014). Sparse CCA: Adaptive estimation and computational barriers. arXiv preprint arXiv:1409.8565. [14] Garc̄a-Garc̄a, D. and Williamson, R. C. (2011). Degrees of supervision. In Advances in Neural Information Processing Systems. [15] Hajek, B., Wu, Y. and Xu, J. (2014). Computational lower bounds for community detection on random graphs. arXiv preprint arXiv:1406.6625. [16] Johnstone, I. M. (1994). On minimax estimation of a sparse normal mean vector. The Annals of Statistics, 22 271–289. [17] Joulin, A. and Bach, F. R. (2012). A convex relaxation for weakly supervised classifiers. In International Conference on Machine Learning. [18] Kearns, M. (1993). Efficient noise-tolerant learning from statistical queries. In ACM Symposium on Theory of Computing. [19] Ma, Z. and Wu, Y. (2014). Computational barriers in minimax submatrix detection. The Annals of Statistics, 43 1089–1116. [20] Nettleton, D. F., Orriols-Puig, A. and Fornells, A. (2010). A study of the effect of different types of noise on the precision of supervised learning techniques. Artificial Intelligence Review, 33 275–306. [21] Patrini, G., Nielsen, F., Nock, R. and Carioni, M. (2016). Loss factorization, weakly supervised learning and label noise robustness. arXiv preprint arXiv:1602.02450. [22] Ramdas, A., Singh, A. and Wasserman, L. (2016). Classification accuracy as a proxy for two sample testing. arXiv preprint arXiv:1602.02210. [23] Tony Cai, T., Liu, W. and Xia, Y. (2014). Two-sample test of high dimensional means under dependence. Journal of the Royal Statistical Society: Series B, 76 349–372. [24] Tsybakov, A. B. (2008). Introduction to nonparametric estimation. Springer. [25] Vershynin, R. (2010). Introduction to the non-asymptotic analysis of random matrices. arXiv preprint arXiv:1011.3027. [26] Wang, T., Berthet, Q. and Samworth, R. J. (2014). Statistical and computational trade-offs in estimation of sparse principal components. arXiv preprint arXiv:1408.5369. [27] Wang, Z., Gu, Q. and Liu, H. (2015). Sharp computational-statistical phase transitions via oracle computational model. arXiv preprint arXiv:1512.08861. [28] Zhang, Y., Wainwright, M. J. and Jordan, M. I. (2014). Lower bounds on the performance of polynomialtime algorithms for sparse linear regression. In Conference on Learning Theory.