



# Prospectus Presentation

## Minimax Theory of Constrained Hypothesis Testing

Presenter: Yikun Li

Committee members: Prof. Matey Neykov (chair), Prof. Miklos Racz, Prof. Feng Ruan



## Content

### Introduction

- Problem Formulation
- Literature Review

### Main Results

- Lower Bounds
- Upper Bounds
- Discussion

### Experiments

- Numerical Simulation
- Discussion

### Summary

- Summary
- Open Problems
- Plan for Future Works



## Content

### Introduction

- Problem Formulation
- Literature Review

### Main Results

- Lower Bounds
- Upper Bounds
- Discussion

### Experiments

- Numerical Simulation
- Discussion

### Summary

- Summary
- Open Problems
- Plan for Future Works

## Problem Formulation

Problem studied: hypothesis testing

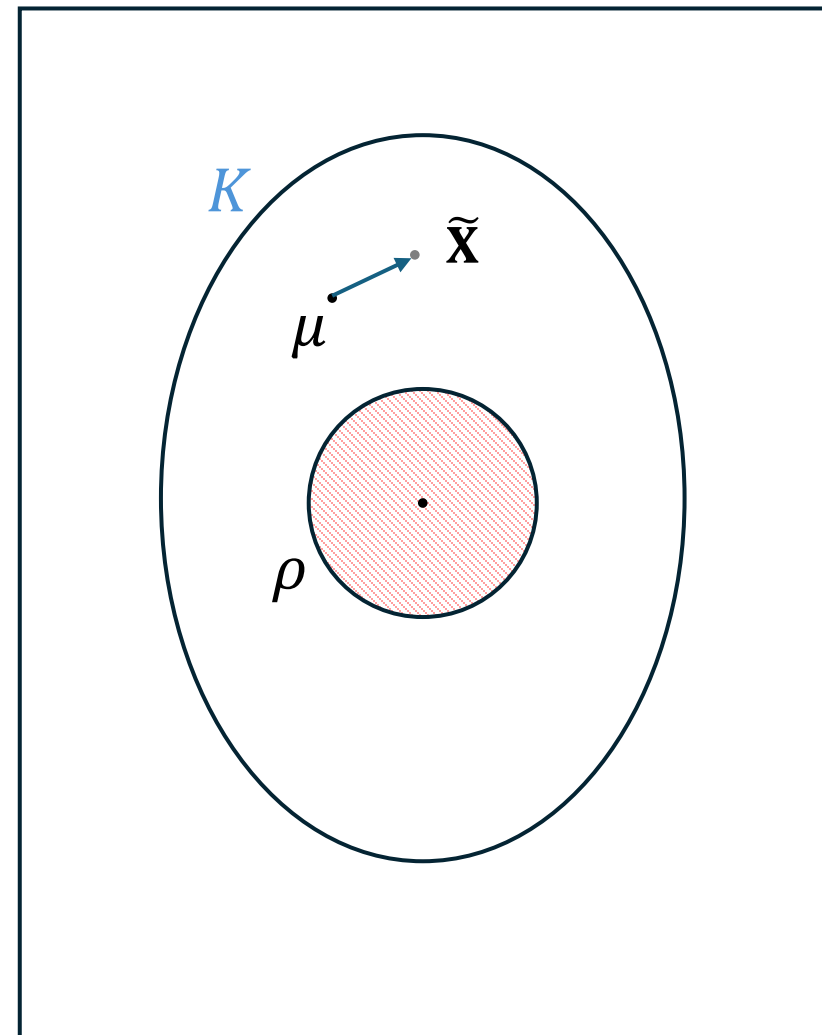
Assumptions:

- Gaussian sequence model (true distribution of the data);
- Known covariance matrix  $\mathbf{I}_d$ , mean vector  $\mu \in K$ ; (prior constraints)

Authentic samples  $\tilde{\mathbf{X}} := \{\tilde{X}_1, \dots, \tilde{X}_N\}$  i.i.d. from Gaussian  $\mathcal{N}(\mu, \mathbf{I}_d)$

$$H_0: \mu = \mathbf{0};$$

$$H_1: \|\mu\|_2 \geq \rho, \mu \in K.$$

 $\mathbb{R}^d$ 

## Problem Formulation

Problem studied: hypothesis testing

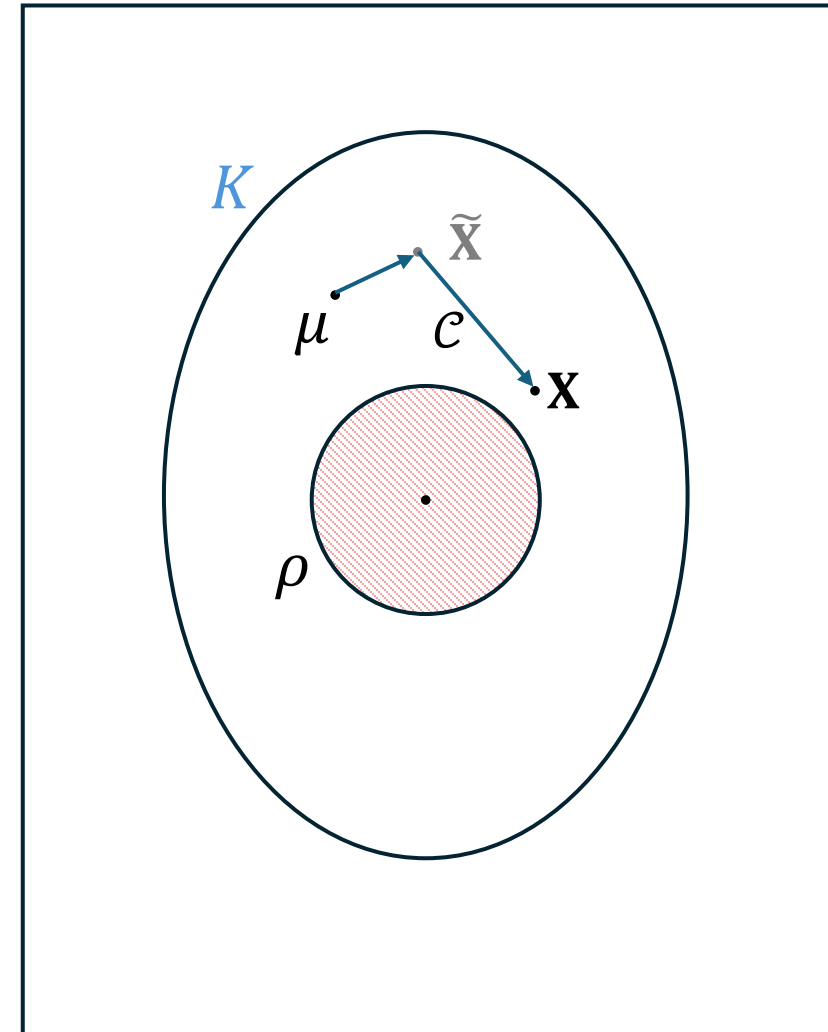
Additional Assumptions:

- Strong contamination adversary  $\mathcal{C}$  (robustness)
- **Powerful** but **limited** contamination fraction  $\epsilon$

$$\tilde{\mathbf{X}} \text{ (unobserved)} \longrightarrow \mathbf{X} := \mathcal{C}(\tilde{\mathbf{X}}) \text{ (observed)}$$

$$H_0: \mu = \mathbf{0};$$

$$H_1: \|\mu\|_2 \geq \rho, \mu \in K.$$

 $\mathbb{R}^d$ 

## Problem Formulation

Problem studied: hypothesis testing

- $\rho$  overly small  $\longrightarrow$  large Type II errors
- $\rho$  overly large  $\longrightarrow$  information not fully utilized

Question: **what is the exact rate of  $\rho$  ?**

Acceptable tests:  $A_s := \left\{ \phi : \sup_{\mathcal{C}} \mathbb{P}_0 \left( \phi(\mathcal{C}(\tilde{\mathbf{X}})) = 1 \right) \leq \alpha \right\},$

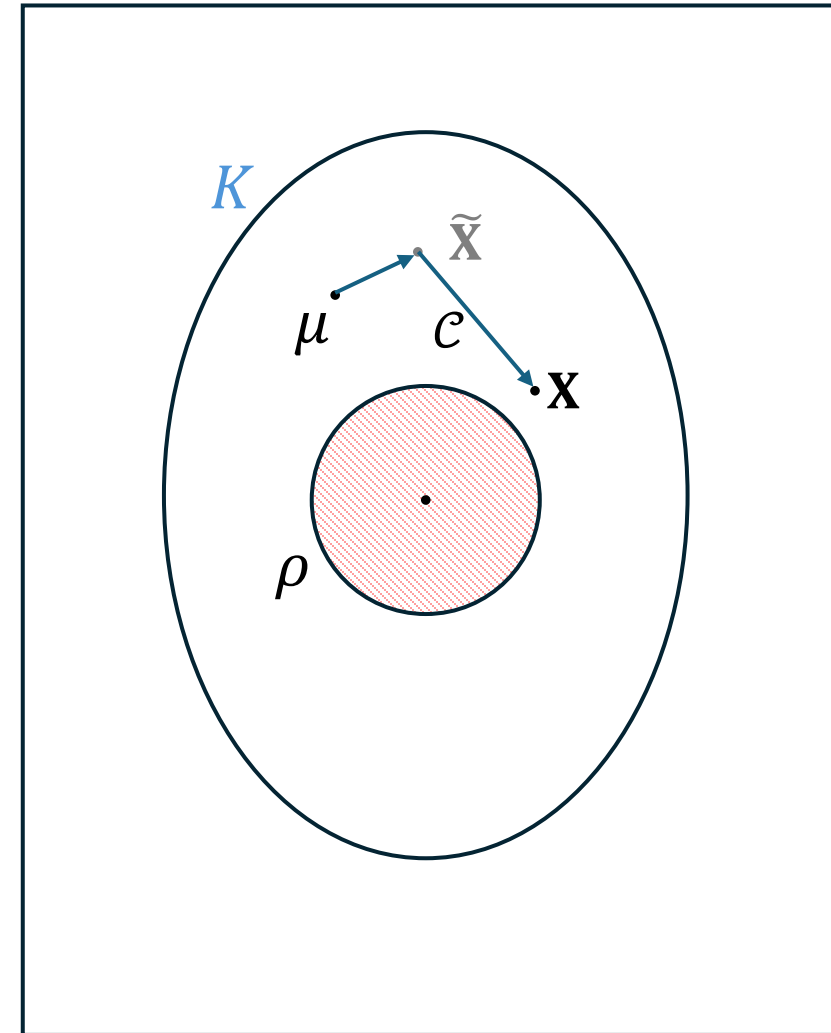
Exact rate:  $\rho_{\text{critical}} := \inf_{\rho} \left\{ \rho : \rho > 0, \sup_{\|\mu\|_2 \geq \rho, \mu \in K} \sup_{\mathcal{C}} \mathbb{P}_{\mu} \left( \phi(\mathcal{C}(\mathbf{X})) = 0 \right) \leq \alpha, \phi \in A_s \right\}$

optimize over the worst case

worst case of parameter

worst case of adversary

$\mathbb{R}^d$





# Literature Review

## Information-theoretic aspect

1973, Le Cam. Le Cam's approach for minimax lower bound.

~1990, Ingster, Suslina, et al. Foundational methodology and rates established for paradigm testing problem.

1992, Donoho, et al. Minimax rate established for constrained estimation problem

2002, Baraud. Establish the minimax rate of testing under ellipsoid constraints

~2020, Canonne, Narayanan, et al. Minimax rate under robust settings

2024 Neykov, et al. Minimax rate for estimation under star-shaped constraints and robust settings

arXiv:2412.03832v2 [math.ST] 12 Jun 2025

### Information theoretic limits of robust sub-Gaussian mean estimation under star-shaped constraints

Akshay Prasad and Matey Neykov

Department of Statistics & Data Science, Carnegie Mellon University  
Department of Statistics and Data Science, Northwestern University

aprasada@andrew.cmu.edu, mneykov@northwestern.edu

#### Abstract

We obtain the minimax rate for a mean location model with a bounded star-shaped set  $K \subseteq \mathbb{R}^n$  constraint on the mean, in an adversarially corrupted data setting with Gaussian noise. We assume an unknown fraction  $\epsilon \leq 1/2 - \kappa$  for some fixed  $\kappa \in (0, 1/2]$  of  $N$  observations are arbitrarily corrupted. We obtain a minimax risk up to proportionality constants under the squared  $\ell_2$  loss of  $\max(\eta^2, \sigma^2 \epsilon^2) \wedge d^2$  with

$$\eta^* = \sup \left\{ \eta > 0 : \frac{N\eta^2}{\sigma^2} \leq \log \mathcal{M}_K^{\text{loc}}(\eta, \epsilon) \right\},$$

where  $\log \mathcal{M}_K^{\text{loc}}(\eta, \epsilon)$  denotes the local entropy of the set  $K$ ,  $d$  is the diameter of  $K$ ,  $\sigma^2$  is the variance, and  $\epsilon$  is some sufficiently large absolute constant. A variant of our algorithm achieves the same rate for settings with known or symmetric sub-Gaussian noise, with a smaller breakdown point, still of constant order. We further study the case of unknown sub-Gaussian noise and show that the rate is slightly slower:  $\max(\eta^2, \sigma^2 \epsilon^2 \log(1/\epsilon)) \wedge d^2$ . We generalize our results to the case when  $K$  is star-shaped but unbounded.

#### Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Related Literature	3
1.2	Notation and Definitions	5
1.3	Organization	7
<b>2</b>	<b>Lower Bounds</b>	<b>9</b>
2.1	Unknown Sub-Gaussian Noise Lower Bound	9
<b>3</b>	<b>Upper Bound with i.i.d. Gaussian Noise</b>	<b>10</b>
3.1	Constructing an Infinite Tree of Points in $K$	10
3.2	Robust Algorithm	12
3.3	Bounding the Error of our Algorithm	13
<b>4</b>	<b>Sub-Gaussian case</b>	<b>15</b>
4.1	Sign-Symmetric/Known Noise	16
4.2	Unknown sub-Gaussian Noise	17
<b>5</b>	<b>Extension to unbounded sets</b>	<b>21</b>
5.1	Lower Bound	21
5.2	Upper Bound	21
5.3	Example: Sparse Robust Mean Estimation	26



# Literature Review

## Computational aspect

2016, Diakonikolas et al. First general computationally efficient framework under robust settings for estimation

2023, Canonne, Narayanan, et al. Efficient framework transferred to robust testing scenario without constraint

### The Full Landscape of Robust Mean Testing: Sharp Separations between Oblivious and Adaptive Contamination

Clément Canonne\*  
University of Sydney

Samuel B. Hopkins†  
MIT

Jerry Li‡  
Microsoft Research

Allen Liu§  
MIT

Shyam Narayanan¶  
MIT

July 21, 2023

#### Abstract

We consider the question of Gaussian mean testing, a fundamental task in high-dimensional distribution testing and signal processing, subject to adversarial corruptions of the samples. We focus on the relative power of different adversaries, and show that, in contrast to the common wisdom in robust statistics, there exists a strict separation between adaptive adversaries (strong contamination) and oblivious ones (weak contamination) for this task. Specifically, we resolve both the information-theoretic and computational landscapes for robust mean testing. In the exponential-time setting, we establish the tight sample complexity of testing  $\mathcal{N}(0, I)$  against  $\mathcal{N}(\alpha v, I)$ , where  $\|v\|_2 = 1$ , with an  $\varepsilon$ -fraction of adversarial corruptions, to be

$$\tilde{\Theta}\left(\max\left(\frac{\sqrt{d}}{\alpha^2}, \frac{d\varepsilon^3}{\alpha^4}, \min\left(\frac{d^{2/3}\varepsilon^{2/3}}{\alpha^{8/3}}, \frac{d\varepsilon}{\alpha^2}\right)\right)\right),$$

while the complexity against adaptive adversaries is

$$\tilde{\Theta}\left(\max\left(\frac{\sqrt{d}}{\alpha^2}, \frac{d\varepsilon^2}{\alpha^4}\right)\right),$$

which is strictly worse for a large range of vanishing  $\varepsilon, \alpha$ . To the best of our knowledge, ours is the first separation in sample complexity between the strong and weak contamination models.

In the polynomial-time setting, we close a gap in the literature by providing a polynomial-time algorithm against adaptive adversaries achieving the above sample complexity  $\tilde{\Theta}(\max(\sqrt{d}/\alpha^2, d\varepsilon^2/\alpha^4))$ , and a low-degree lower bound (which complements an existing reduction from planted clique) suggesting that all efficient algorithms require this many samples, even in the oblivious-adversary setting.

\*clement.canonne@sydney.edu.au. Supported by an ARC DECRA (DE230101329) and an unrestricted gift from Google Research.

†samhop@mit.edu. Supported by NSF Award No. 2238080 and MLA@CSAIL.

‡jerrli@microsoft.com.

§cliu568@mit.edu Supported by an NSF Graduate Research Fellowship and a Fannie and John Hertz Foundation Fellowship.

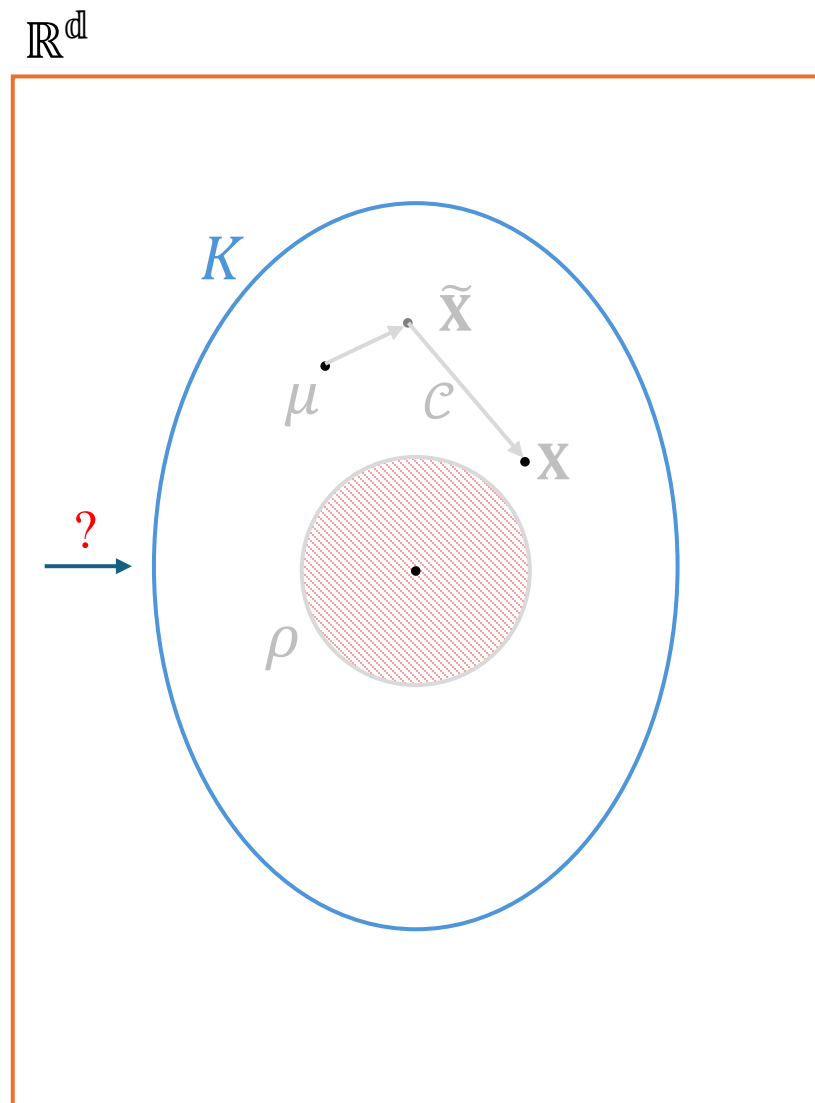
¶shyamsn@mit.edu. Supported by an NSF Graduate Fellowship and a Google Fellowship.

arXiv:2307.10273v1 [cs.DS] 18 Jul 2023

## Literature Review

### Overall Motivation

- Under certain geometric constraints, how does the minimax lower bound change?
- Does the matching upper bound still exist?
- Is the algorithm computationally efficient?





## Content

### Introduction

- Problem Formulation
- Literature Review

### Main Results

- Lower Bounds
- Upper Bounds
- Discussion

### Experiments

- Numerical Simulation
- Discussion

### Summary

- Summary
- Open Problems
- Plan for Future Works



## Main Results

### Preliminaries (Quadratically convex orthosymmetric (QCO) set)

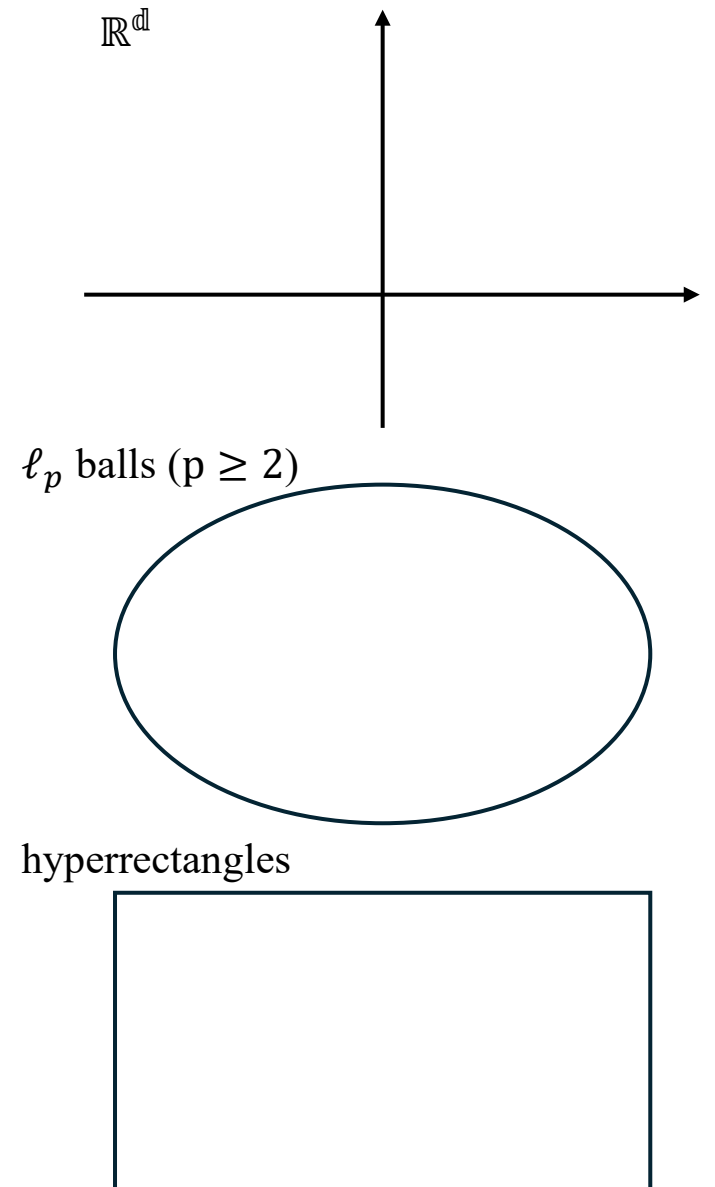
**Definition 2.1** (Quadratically convex orthosymmetric (QCO) set). *Given a set  $K \subset \mathbb{R}^d$ , we say  $K$  is a quadratically convex orthosymmetric (QCO) set if it satisfies the following conditions:*

(1),  $K$  is convex;

(2),  $K$  is quadratically convex, which means that  $K^2$  is also convex, where  $K^2$  is defined as

$$K^2 := \left\{ (\theta_1^2, \dots, \theta_d^2)^\top \mid (\theta_1, \dots, \theta_d)^\top \in K \right\};$$

(3),  $K$  is orthosymmetric, which means that if  $\theta = (\theta_1, \dots, \theta_d)^\top \in K$ , then  $\theta_\eta := (\eta_1\theta_1, \dots, \eta_d\theta_d)^\top \in K$ , where  $\eta_i \in \{-1, 1\}$ ,  $1 \leq i \leq d$ .



# Main Results

## Preliminaries (Kolmogorov widths)

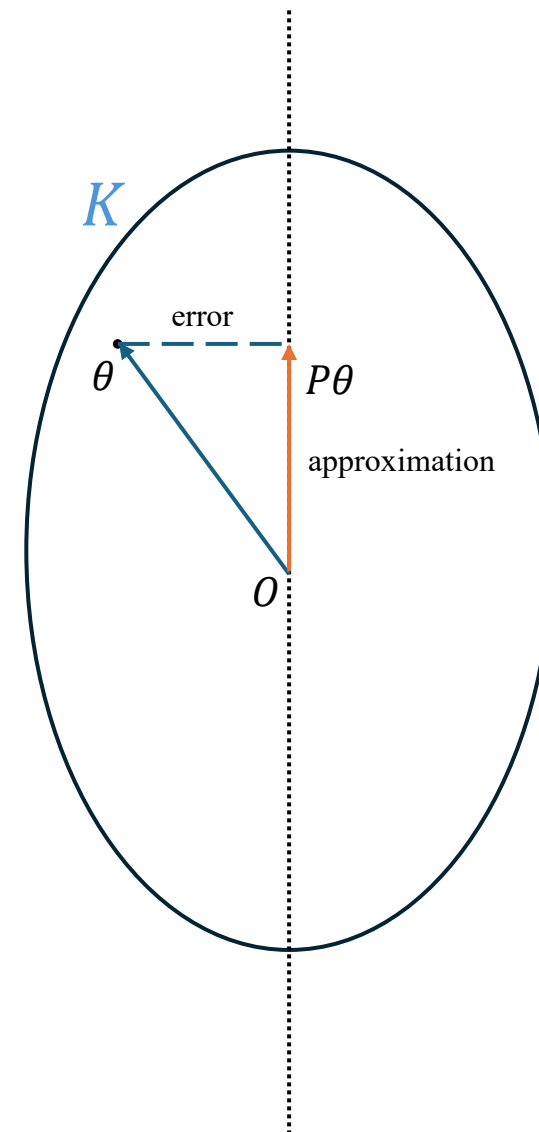
**Definition 2.2** (Kolmogorov  $k$ -width, also known as Kolmogorov  $N$ -width). *Let  $\mathcal{X}$  be a Banach space equipped with the norm  $\|\cdot\|$ , and  $K \subset \mathcal{X}$  is a subset. The Kolmogorov  $k$ -width is defined as*

$$D_k(K) = \inf_{V_k} \sup_{\theta \in K} \min_{\tilde{\theta} \in V_k} \|\theta - \tilde{\theta}\|_2, \tag{2.1}$$

where  $V_k \subset \mathbb{R}^d$  is any  $k$ -dimensional subspace of  $\mathbb{R}^d$ .

- For  $\ell_2$  norm approximation,  $\tilde{\theta}$  can be expressed as  $\tilde{\theta} = P\theta$ , where  $P$  is the projection operator onto  $V_k$

*Motivation: the minimax risk of  $k$ -dimensional linear approximations*





## Main Results

Preliminaries (Optimal dimensions and projections)

First optimal dimensions:  $k_1^*(K, \sigma, N) := \max \left\{ j \mid 0 \leq j \leq d, D_{j-1}(K) > \frac{j^{\frac{1}{4}}}{\sqrt{N}} \sigma \right\}$

Second optimal dimensions:  $k_2^*(K, \sigma, N, \epsilon) := \max \left\{ j \mid 0 \leq j \leq d, D_{j-1}(K) > \frac{j^{\frac{1}{4}} \sqrt{\epsilon}}{N^{\frac{1}{4}}} \sigma \right\}$

The corresponding projections are first/second optimal projections  $P_1^*, P_2^*$

*Motivation: interaction between the geometry and the data properties*



## Main Results

### Preliminaries (Notations)

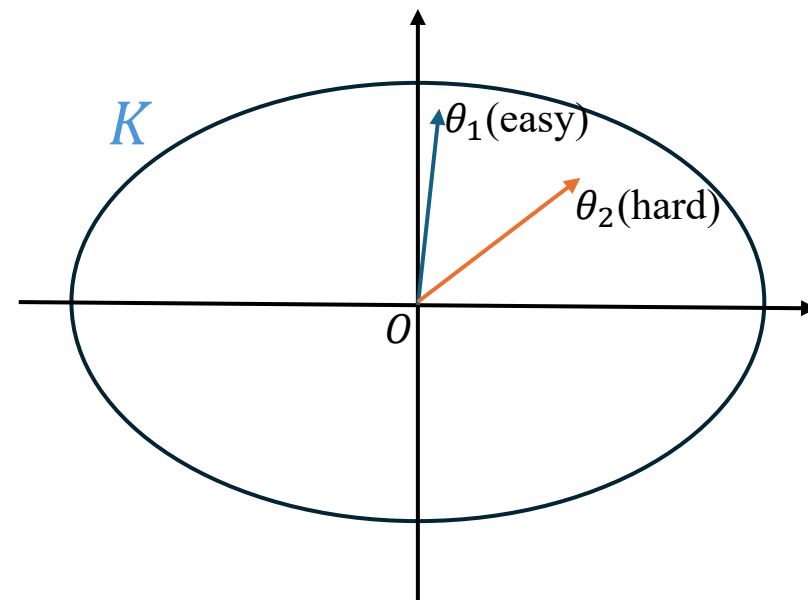
- $d$ : full dimension
- $N$ : sample size
- $\mu$ : mean vector
- $K$ : geometric constraints of  $\mu$
- $\alpha$ : fixed Type I and Type II errors tolerance
  
- $\tilde{\mathbf{X}} := \{\tilde{X}_1, \dots, \tilde{X}_N\}$ : original, unobserved samples from  $\mathcal{N}(\mu, \mathbf{I}_d)$
- $\mathbf{X} := \{X_1, \dots, X_N\}$ : contaminated samples from  $\mathcal{C}$  and  $\tilde{\mathbf{X}}$
- $\mathcal{C}$ : index set of the corrupted samples
  
- $D_k(K)$ : Kolmogorov  $k$ -width of  $K$
- $k_1^*, k_2^*$ : first/second optimal dimensions
- $P_1^*, P_2^*$ : first/second optimal projections

## Main Results (First Lower Bounds)

**Important lemma.** Let  $K \subset \mathcal{X}$  be a QCO set and  $c$  be an arbitrary positive constant. Suppose that  $D_{k-1}(K) > c\sigma$  for some  $k \geq 1$ . Then there exists a vector  $\theta \in K$  such that  $\|\theta\|_2 = c\sigma$  while  $\|\theta\|_\infty \leq \frac{c}{\sqrt{k}}\sigma$ .

### Motivation:

- The output vector  $\theta$  is “smeared” with non-trivial total energy
- Construction of *least favor distribution*





## Main Results (First Lower Bounds)

**First lower bound.** If  $\rho \leq c(\alpha) \frac{(k_1^*)^{1/4}}{\sqrt{N}} \sigma$ , where  $c(\alpha)$  is a constant that only depends on  $\alpha$ , or  $k_1^* = 0$ , then

$$\inf_{\psi: \mathbb{P}_0(\psi=1) \leq \alpha} \sup_{\theta \in K, \|\theta\| \geq \rho} \mathbb{P}_\theta(\psi = 0) \geq \alpha$$

### High-Level Proof:

- Small  $\chi^2$  divergence between  $\mathbb{P}_0^{\otimes N}$  and  $\mathbb{E}_{\mu \sim \gamma} \mathbb{P}_\mu^{\otimes N}$  leads to impossibility of the test via the standard Le Cam's fuzzy hypothesis method
- Control the  $\chi^2$  divergence via the output vector  $\theta$

$\gamma$  is any prior on the mean  $\mu$



## Main Results (First Lower Bounds)

**First lower bound.** If  $\rho \leq c(\alpha) \frac{(k_1^*)^{1/4}}{\sqrt{N}} \sigma$ , where  $c(\alpha)$  is a constant that only depends on  $\alpha$ , or  $k_1^* = 0$ , then

$$\inf_{\psi: \mathbb{P}_0(\psi=1) \leq \alpha} \sup_{\theta \in K, \|\theta\| \geq \rho} \mathbb{P}_\theta(\psi = 0) \geq \alpha$$

Remark:

- No adversary appears here?
- Classic rate  $\asymp \frac{d^{1/4}}{\sqrt{N}} \sigma$  (Ingster)
- One intrinsic dimensional parameter is  $k_1^*$  instead of  $d$  for the constraint  $K$ !
- $k_1^* = 0 \longrightarrow \sup_{\theta \in K} \|\theta\|_2 \leq \frac{1}{\sqrt{N}} \sigma \longrightarrow$  the constraint  $K$  is insufficient rich for the test!



## Main Results (Second Lower Bounds)

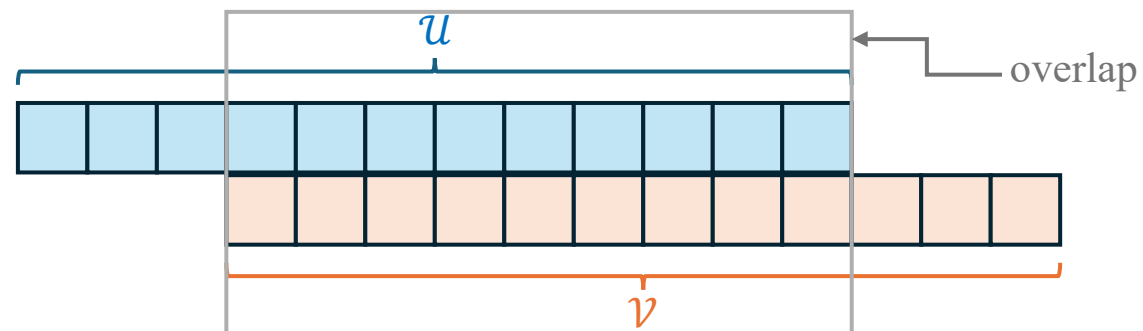
**Important lemma.** Suppose  $\mathbf{X}, \mathbf{X}' \in \mathbb{R}^{N \times d}$ . Assume the laws of  $\mathbf{X}, \mathbf{X}'$  are  $\mathcal{U}$  and  $\mathcal{V}$ , respectively. If there is a coupling between  $\mathcal{U}$  and  $\mathcal{V}$  such that the expectation of the Hamming distance between  $\mathbf{X}$  and  $\mathbf{X}'$  is  $\mathcal{O}(\epsilon N)$ :

$$\mathbb{E}_{(\mathbf{X}, \mathbf{X}') \sim (\mathcal{U}, \mathcal{V})} d_H(\mathbf{X}, \mathbf{X}') \lesssim \epsilon N.$$

Then, there is no robust test to distinguish between  $\mathcal{U}$  and  $\mathcal{V}$  with  $\epsilon$  contamination and small errors.

Why is it useful?

A special coupling between  $\mathbb{P}_0^{\otimes N}$  and  $\mathbb{P}_\mu^{\otimes N} \longrightarrow$  impossibility for a robust test!







## Main Results (Second Lower Bounds)

**Second lower bound.** If  $\rho \lesssim \frac{(k_2^*)^{1/4} \sqrt{\epsilon}}{N^{1/4}} \sigma$ , or  $k_2^* = 0$ , then it is impossible to find a robust testing with uniformly small Type I and Type II errors. Here the omitted constant only depends on  $\alpha$ .

### High-Level Proof

- Interpolate between  $\mathbb{P}_{\emptyset}^{\otimes N}$  and  $\mathbb{E}_{\mu \sim \gamma} \mathbb{P}_{\mu}^{\otimes N}$  via a chain of optimal coupling  

- Expectation of the Hamming distance  calculation of the total variance
- Bound the total variation by KL divergence
- Compute and control the KL divergence using the output vector  $\theta$



## Main Results (Second Lower Bounds)

**Second lower bound.** If  $\rho \lesssim \frac{(k_2^*)^{1/4} \sqrt{\epsilon}}{N^{1/4}} \sigma$ , or  $k_2^* = 0$ , then it is impossible to find a robust testing with uniformly small Type I and Type II errors. Here the omitted constant only depends on  $\alpha$ .

Remark:

- When  $\epsilon = 0$ : the lower bound collapses to zero
- Rate from Canonne, Narayanan, et al (2023):  $\frac{d^{1/4} \sqrt{\epsilon}}{N^{1/4}} \sigma$
- Another intrinsic dimensional parameter is  $k_2^*$ !



## Main Results (Third Lower Bounds)

**Third lower bound.** If  $\epsilon \gtrsim \frac{1}{\sqrt{N}}$  and  $\rho \lesssim \epsilon\sigma$ , then it is impossible to find a robust testing with uniformly small Type I and Type II errors. Here the omitted constants only depend on  $\alpha$ .

### Remark:

- This lower bound comes from the counterpart estimation problem
- Irrelevant with sample size  $N$
- The condition  $\epsilon \gtrsim \frac{1}{\sqrt{N}}$  can be safely removed (See later discussion)



## Main Results (Synthesized Lower Bounds)

**Synthesized lower bound.** Given the hypothesis testing problem previously formulated, the following condition is necessary for the existence of the desired robust testing:

$$\rho_{\text{critical}}^2 \gtrsim \sigma^2 \max \left\{ \underbrace{\frac{\sqrt{k_1^*}}{N}}_{\text{geometric terms}}, \underbrace{\epsilon \sqrt{\frac{k_2^*}{N}}, \epsilon^2}_{\text{corruption terms}} \right\}$$

Remark:

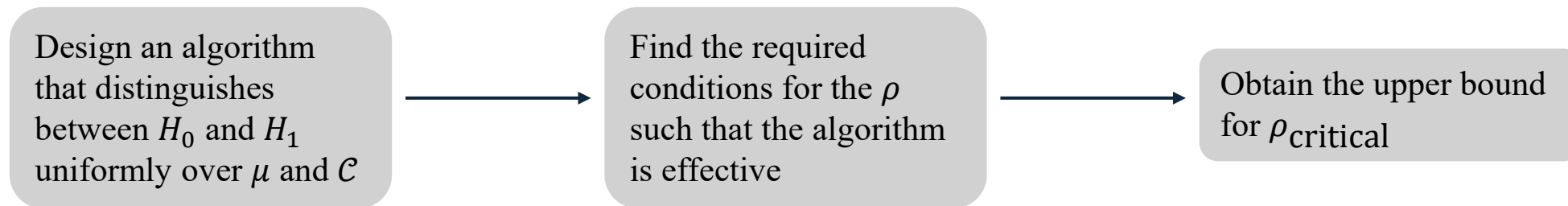
- The first two terms involve the geometry of the constraint
- The last two terms involve the corruption process



## Main Results (Upper Bound)

The upper bound should be derived constructively

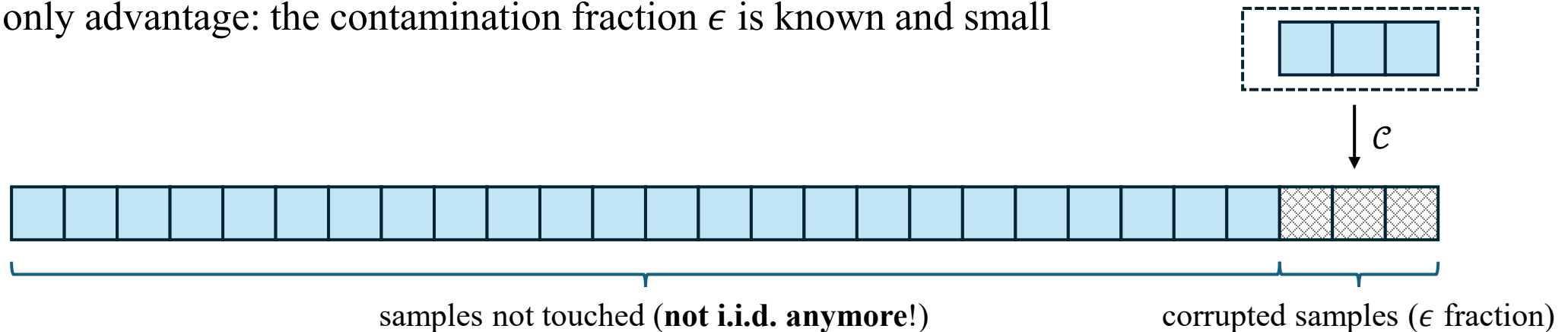
Overall roadmap:



## Main Results (Theoretical Algorithm)

A natural motivation: conceal the influence of the adversary  $\mathcal{C}$

- However,  $\mathcal{C}$  is too powerful  $\longrightarrow$  not easy to work with the contaminated samples
- Our only advantage: the contamination fraction  $\epsilon$  is known and small



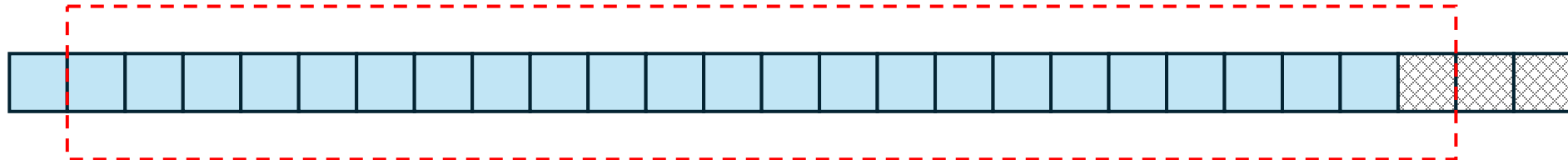


## Main Results (Theoretical Algorithm)

Introduce “consistent subset”

**Definition (consistent subset).** A subset  $S \subset \mathbf{X} = \{X_1, \dots, X_N\}$  is called a consistent subset regarding  $\epsilon$  and a test  $\phi: 2^S \rightarrow \{0,1\}$  if:

- (i),  $|S| \geq (1 - \epsilon)N$ ,
- (ii),  $\phi(S') = \phi(S)$  for any  $S' \subset S$  with  $|S'| \geq (1 - 2\epsilon)N$ .





## Main Results (Theoretical Algorithm)

**Definition (consistent subset).** A subset  $S \subset \mathbf{X} = \{X_1, \dots, X_N\}$  is called a consistent subset regarding  $\epsilon$  and a test  $\phi: 2^S \rightarrow \{0,1\}$  if:

- (i),  $|S| \geq (1 - \epsilon)N$ ,
- (ii),  $\phi(S') = \phi(S)$  for any  $S' \subset S$  with  $|S'| \geq (1 - 2\epsilon)N$ .

**Why consistent subset? It can help recover the testing result on the unobserved original samples  $\tilde{\mathbf{X}}$ !**

Suppose the original samples  $\tilde{\mathbf{X}}$  is consistent w.r.t. a valid test  $\phi_e$ , then:

- It is guaranteed that there exists at least one consistent subset within the observation  $\mathbf{X}$ ! (At least  $[N] \setminus C$ )
- Therefore:

$$\phi_e(\mathbf{X}_S) = \phi_e(\mathbf{X}_{S \cap [N] \setminus C}) = \phi_e(\mathbf{X}_{[N] \setminus C}) = \phi_e(\tilde{\mathbf{X}})$$



## Main Results (Theoretical Algorithm)

Define the testing event  $E_e := \left\{ \left\| P^* \sum_{i \in T} X_i \right\|_2^2 - \overset{\text{any dimension } 1 \leq k^* \leq d}{k^* |T| \sigma^2} \geq c |T|^2 E^2(\epsilon, K, N, \sigma^2) \right\}$

Testing function  $\phi_e = \mathbf{1}_{E_e}$    
  $\uparrow$  optimal projection w.r.t.  $k^*$

- $E^2(\epsilon, K, N, \sigma^2)$  is defined as  $\sigma^2 \max \left\{ \frac{\sqrt{k}}{N}, \epsilon^2 \ln \left( \frac{1}{\epsilon} \right), \sqrt{\frac{\epsilon^2 \ln \left( \frac{1}{\epsilon} \right) k}{N}} \right\}$
- $\phi_e$  is a valid test with small errors and  $\tilde{\mathbf{X}}$  is consistent with  $\phi_e$  w.h.p. when  $\|P^* \mu\|_2^2 \gtrsim E^2(\epsilon, K, N, \sigma^2)$

Requirement on  $\rho$ ?

$$\|P^* \mu\|_2^2 = \rho^2 - \|\mu - P^* \mu\|_2^2 \geq \rho^2 - D_{k^*}^2(K)$$

$$\rho^2 \gtrsim \underbrace{D_k^2(K)}_{\text{approximation gap}} + \underbrace{\sigma^2 \max \left\{ \frac{\sqrt{k}}{N}, \epsilon^2 \ln \left( \frac{1}{\epsilon} \right), \sqrt{\frac{\epsilon^2 \ln \left( \frac{1}{\epsilon} \right) k}{N}} \right\}}_{\text{testing gap}}$$



## Main Results (Theoretical Algorithm)

**Theoretical upper bound.** For the hypothesis testing problem previously formulated, if the following condition is satisfied, then there exists a test  $\phi_e$  to distinguish between  $H_0$  and  $H_1$  with uniformly small errors over  $\mathcal{C}$  and  $\mu$ .

$$\rho^2 \gtrsim D_k^2(K) + \sigma^2 \max \left\{ \frac{\sqrt{k}}{N}, \epsilon^2 \ln \left( \frac{1}{\epsilon} \right), \sqrt{\frac{\epsilon^2 \ln \left( \frac{1}{\epsilon} \right) k}{N}} \right\}$$

**Corollary.** Optimizing over  $k$  in the upper bound above, we obtain the sharpest upper bound as:

$$\rho_{\text{critical}}^2 \lesssim \sigma^2 \max \left\{ \frac{\sqrt{\min \{k_1^*, k_2^*\}}}{N}, \epsilon^2 \ln \left( \frac{1}{\epsilon} \right), \sqrt{\frac{\epsilon^2 \ln \left( \frac{1}{\epsilon} \right) \min \{k_1^*, k_2^*\}}{N}} \right\}$$

where  $k_1^*$  and  $k_2^*$  are the first and second optimal dimensions.

**Limitation: require scan of subset of  $X$   $\longrightarrow$  not efficient!**



## Main Results (Polynomial-time Algorithm)

### From first-order statistic to second-order statistic

- The theoretical algorithm is based on the sum (first-order)
- What if we base on the covariance (second-order)

## Filtering technique (Diakonikolas et al., 2016)

arXiv:1604.06443v2 [cs.DS] 15 Mar 2019

### Robust Estimators in High Dimensions without the Computational Intractability

Ilias Diakonikolas\*   Gautam Kamath<sup>†</sup>   Daniel M. Kane<sup>‡</sup>   Jerry Li<sup>§</sup>   Ankur Moitra<sup>¶</sup>  
Alistair Stewart<sup>||</sup>  
March 18, 2019

#### Abstract

We study high-dimensional distribution learning in an agnostic setting where an adversary is allowed to arbitrarily corrupt an  $\varepsilon$ -fraction of the samples. Such questions have a rich history spanning statistics, machine learning and theoretical computer science. Even in the most basic settings, the only known approaches are either computationally inefficient or lose dimension-dependent factors in their error guarantees. This raises the following question: Is high-dimensional agnostic distribution learning even possible, algorithmically?

In this work, we obtain the first computationally efficient algorithms with dimension-independent error guarantees for agnostically learning several fundamental classes of high-dimensional distributions: (1) a single Gaussian, (2) a product distribution on the hypercube, (3) mixtures of two product distributions (under a natural balancedness condition), and (4) mixtures of spherical Gaussians. Our algorithms achieve error that is independent of the dimension, and in many cases scales nearly-linearly with the fraction of adversarially corrupted samples. Moreover, we develop a general recipe for detecting and correcting corruptions in high-dimensions that may be applicable to many other problems.

\*University of Southern California. Supported by NSF Award CCF-1652862 (CAREER) and a Sloan Research Fellowship. Part of this work was performed while the author was at the University of Edinburgh, supported in part by EPSRC grant EP/L021749/1 and a Marie Curie Career Integration Grant. [diakonik@usc.edu](mailto:diakonik@usc.edu)

<sup>†</sup>Simons Institute for the Theory of Computing. Supported by NSF Award CCF-0953960 (CAREER) and ONR grant N00014-12-1-0999. This work was done in part while the author was an intern at Microsoft Research Cambridge, visiting the Simons Institute for the Theory of Computing, and a graduate student at MIT. [gkamath@csail.mit.edu](mailto:gkamath@csail.mit.edu)

<sup>‡</sup>University of California, San Diego. Part of this work was performed while visiting the University of Edinburgh. [dakab@cs.ucsd.edu](mailto:dakab@cs.ucsd.edu)

<sup>§</sup>Microsoft Research AI. Supported by NSF CAREER Award CCF-1453261, a Google Faculty Research Award, and an NSF Fellowship. This work was done in part while the author was an intern at Microsoft Research Cambridge and a graduate student at MIT. [jerfl@microsoft.com](mailto:jerfl@microsoft.com)

<sup>¶</sup>Massachusetts Institute of Technology. Supported by NSF CAREER Award CCF-1453261, a grant from the MIT NEC Corporation, and a Google Faculty Research Award. [moitra@mit.edu](mailto:moitra@mit.edu)

<sup>||</sup>Web3 Foundation. Part of this work was performed while the author was at the University of Edinburgh and the University of Southern California. Research supported in part by EPSRC grant EP/L021749/1. [stewart.ai@gmail.com](mailto:stewart.ai@gmail.com)

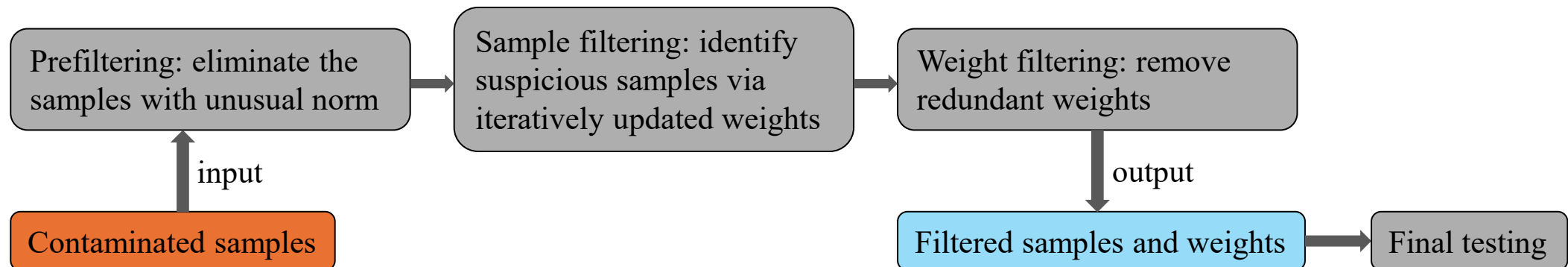


## Main Results (Polynomial-time Algorithm)

Filtering technique (Diakonikolas et al., 2016)

High-level steps:

- Assign each observed sample with an initial equal weight
- Detect the disturbance in the sample covariance matrix to identify suspicious sample and update the weights
- Test based on the filtered samples and weights





## Main Results (Polynomial-time Algorithm)

**Polynomial-time upper bound.** For the hypothesis testing problem previously formulated, if the following condition is satisfied, then there exists a test  $\phi_e$  to distinguish between  $H_0$  and  $H_1$  with uniformly small errors over  $\mathcal{C}$  and  $\mu$ .

$$\rho^2 \gtrsim D_k^2(K) + \sigma^2 \max \left\{ \frac{\epsilon \ln \left( \frac{N}{\alpha} \right)}{\sqrt{N}}, \epsilon^2 \ln \left( \frac{N}{\alpha} \right), \sqrt{\frac{\epsilon^2 k \ln \left( \frac{N}{\alpha} \right)}{N}}, \frac{\sqrt{k} \ln \left( \frac{1}{\alpha} \right)}{N} \right\}$$

**Corollary.** Optimizing over  $k$  in the upper bound above, we obtain the sharpest upper bound as:

$$\rho_{\text{critical}}^2 \lesssim \sigma^2 \max \left\{ \frac{\epsilon \ln \left( \frac{N}{\alpha} \right)}{\sqrt{N}}, \epsilon^2 \ln \left( \frac{N}{\alpha} \right), \sqrt{\frac{\epsilon^2 \min \{k_1^*, k_2^*\} \ln \left( \frac{N}{\alpha} \right)}{N}}, \frac{\sqrt{\min \{k_1^*, k_2^*\} \ln \left( \frac{1}{\alpha} \right)}}{N} \right\}$$

where again  $k_1^*$  and  $k_2^*$  are the first and second optimal dimensions.

**Advantage: the filtering and testing process are finished in polynomial time!**



## Main Results (Discussion)

**Theorem.** The synthesized lower bound, theoretical upper bound, and polynomial-time upper bound all match, except for a universal constant and logarithmic factors in  $\alpha, N, \frac{1}{\epsilon}$ .

Information-theoretic lower bound:  $\rho_c^2 \gtrsim \sigma^2 \max \left\{ \frac{\sqrt{k_1^*}}{N}, \epsilon \sqrt{\frac{k_2^*}{N}}, \epsilon^2 \right\}$



Theoretical upper bound:  $\rho_c^2 \lesssim \sigma^2 \max \left\{ \frac{\sqrt{\min \{k_1^*, k_2^*\}}}{N}, \epsilon^2 \ln \left( \frac{1}{\epsilon} \right), \sqrt{\frac{\epsilon^2 \ln \left( \frac{1}{\epsilon} \right) \min \{k_1^*, k_2^*\}}{N}} \right\}$



Polynomial upper bound:  $\rho_c^2 \lesssim \sigma^2 \max \left\{ \frac{\epsilon \ln \left( \frac{N}{\alpha} \right)}{\sqrt{N}}, \epsilon^2 \ln \left( \frac{N}{\alpha} \right), \sqrt{\frac{\epsilon^2 \min \{k_1^*, k_2^*\} \ln \left( \frac{N}{\alpha} \right)}{N}}, \frac{\sqrt{\min \{k_1^*, k_2^*\} \ln \left( \frac{1}{\alpha} \right)}}{N} \right\}$



## Main Results (Discussion)

Minimax rate for the constrained hypothesis testing problem with adversary:

$$\rho_c^2 \gtrsim \sigma^2 \max \left\{ \frac{\sqrt{k_1^*}}{N}, \epsilon \sqrt{\frac{k_2^*}{N}}, \epsilon^2 \right\}$$

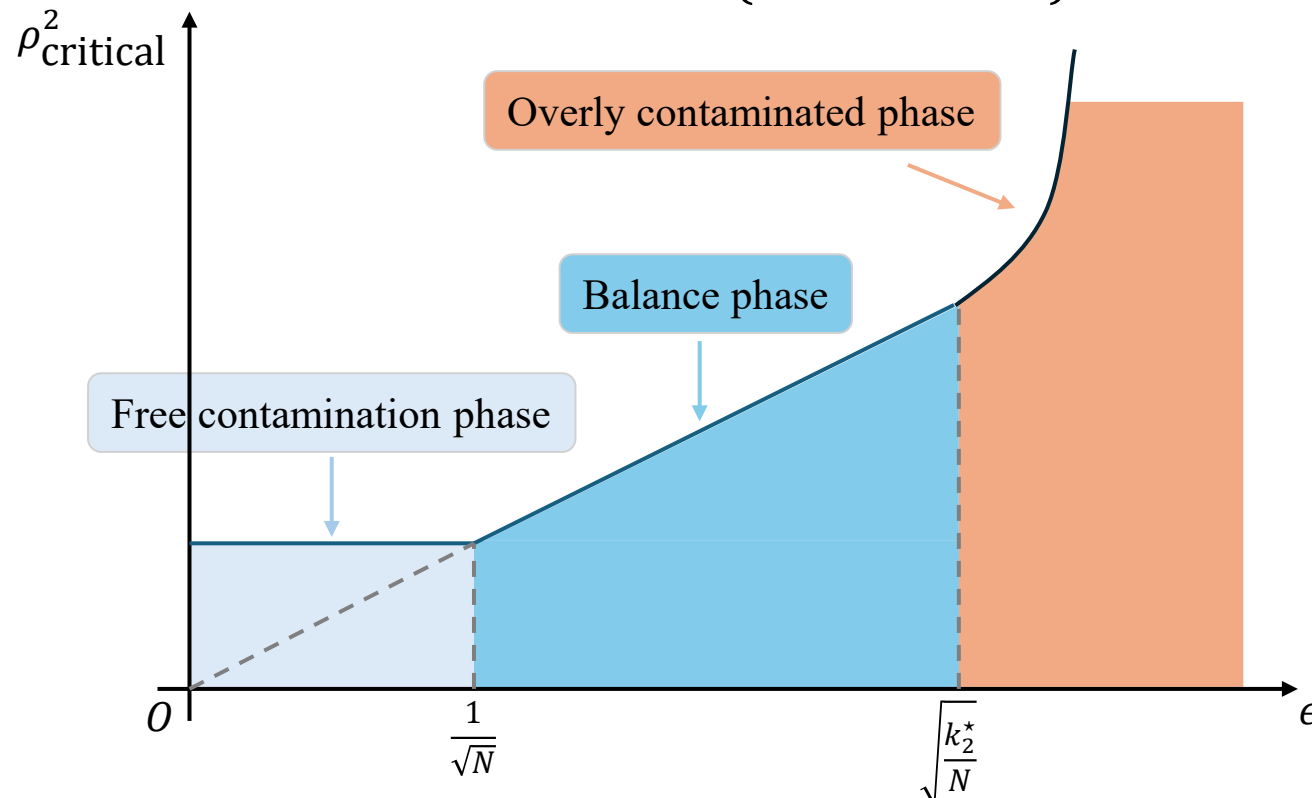
- Full dimension  $d$  does not appear --- completely replaced by  $k_1^*$  and  $k_2^*$ !
- When  $0 \leq \epsilon \lesssim \frac{1}{\sqrt{N}}$ : **the first term** dominates (Free contamination phase)
- When  $\frac{1}{\sqrt{N}} \lesssim \epsilon \lesssim \sqrt{\frac{k_2^*}{N}}$ : **the second term** dominates (Balance phase)
- When  $\epsilon \gtrsim \sqrt{\frac{k_2^*}{N}}$ : **the third term** dominates (Overly contaminated phase)



## Main Results (Discussion)

Minimax rate for the constrained hypothesis testing problem with adversary:

$$\rho_c^2 \gtrsim \sigma^2 \max \left\{ \frac{\sqrt{k_1^*}}{N}, \epsilon \sqrt{\frac{k_2^*}{N}}, \epsilon^2 \right\}$$





## Content

### Introduction

- Problem Formulation
- Literature Review

### Main Results

- Lower Bounds
- Upper Bounds
- Discussion

### Experiments

- Numerical Simulation
- Discussion

### Summary

- Summary
- Open Problems
- Plan for Future Works



## Experiment

### Experiment settings

- i.i.d. samples  $\tilde{\mathbf{X}} := \{\tilde{\mathbf{X}}_1, \dots, \tilde{\mathbf{X}}_N\}$  from  $\mathcal{N}(\mu, \mathbf{I}_d)$
- Constraint  $\mu \in K_e := \left\{v \mid v \in \mathbb{R}^d, \sum_{i=1}^n \frac{v_i^2}{a_i} \leq 1, a_i = \sqrt{d} \cdot i^{-2}\right\}$

### Configuration for the parameters:

- Classical settings:  $(N, d, \epsilon) \in \{200, 2000\} \times \{5, 10, 20, 50\} \times \{0.01, 0.015, 0.02, 0.025\}$
- High-dimensional settings:  $(N, d, \epsilon) \in \{100, 500\} \times \{100, 500\} \times \{0.01, 0.015, 0.02, 0.025\}$

# Experiment

## Experiment settings

Configuration for the adversary  $\mathcal{C}$ :

---

**Algorithm 1:** strategy 1 of the adversary  $\mathcal{C}$

---

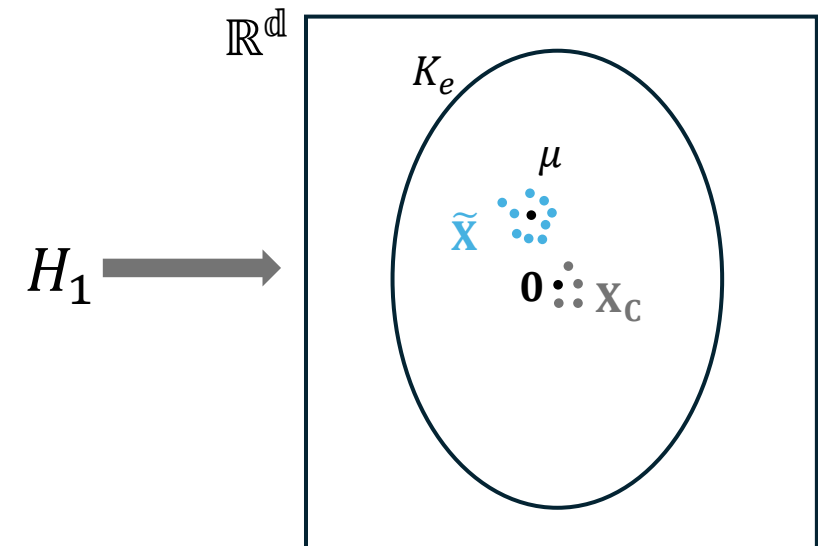
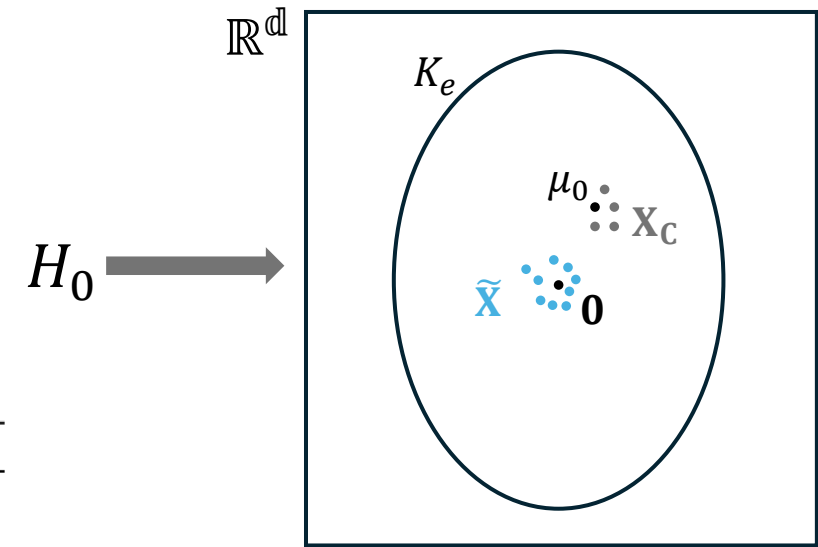
```

1 import  $\tilde{\mathbf{X}}, N, d, \sigma = 1$ 

2 Compute  $\mathbf{arr} \in \mathbb{R}^N$ ,  $\mathbf{arr}_i = \|\tilde{\mathbf{X}}_i\|_2, 1 \leq i \leq N$ 
3 if  $H_1$  is true then
4    $C =$  indices of the largest  $\epsilon N$  entries in  $\mathbf{arr}$ 
5    $\mathbf{X}_C = \epsilon N$  i.i.d. fake data from  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ 
6    $\tilde{\mathbf{X}}_C = \mathbf{X}_C$ 
7 else
8   import  $\mu_0$ 
9    $C =$  indices of the smallest  $\epsilon N$  entries in  $\mathbf{arr}$ 
10   $\mathbf{X}_C = \epsilon N$  i.i.d. fake data from  $\mathcal{N}(\mu_0, \mathbf{I}_d)$ 
11   $\tilde{\mathbf{X}}_C = \mathbf{X}_C$ 
12 return  $\mathbf{X} = \mathcal{C}(\tilde{\mathbf{X}})$ 

```

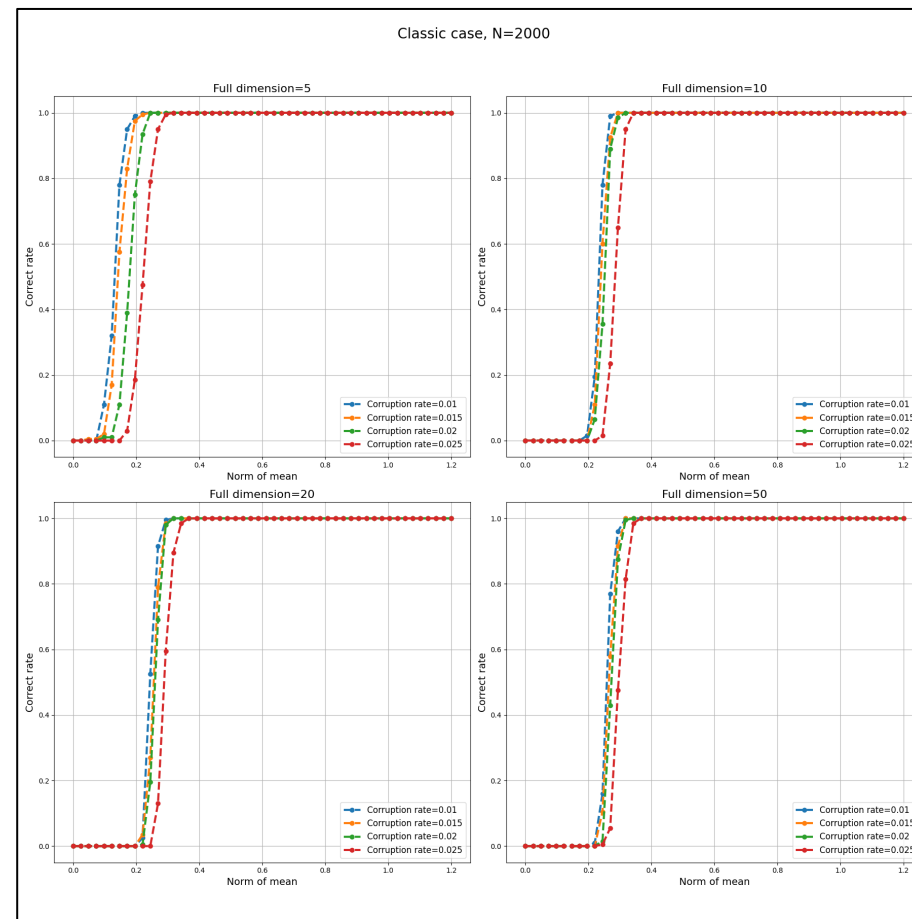
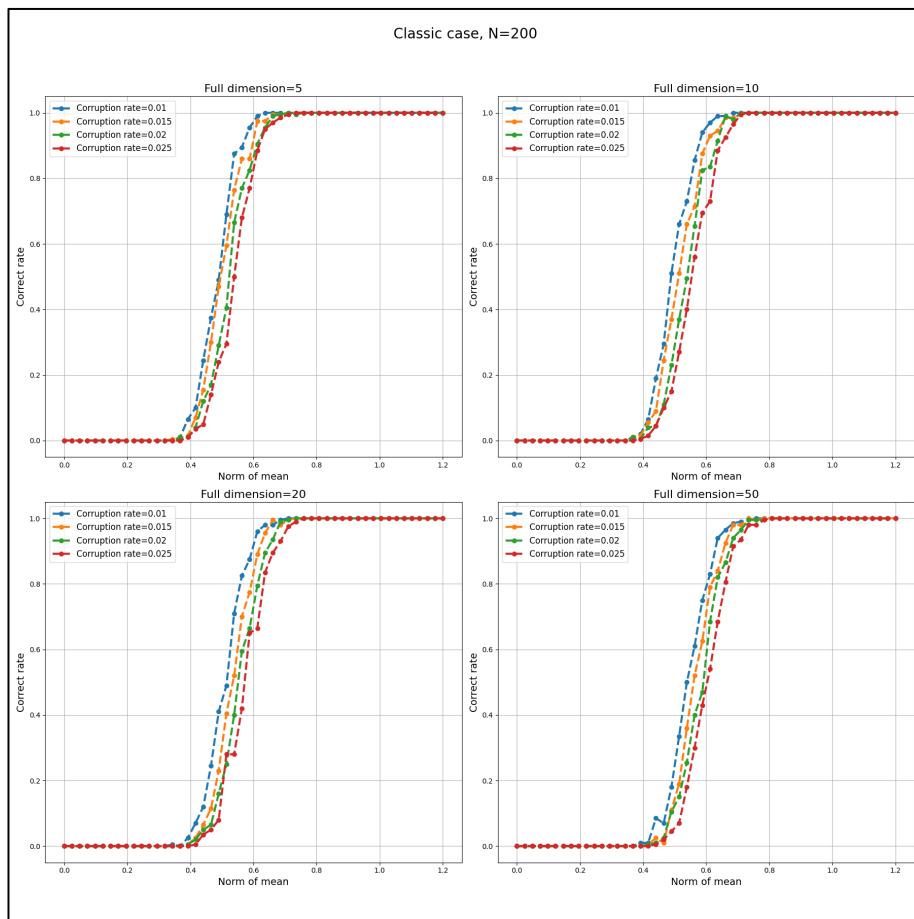
---





# Experiment

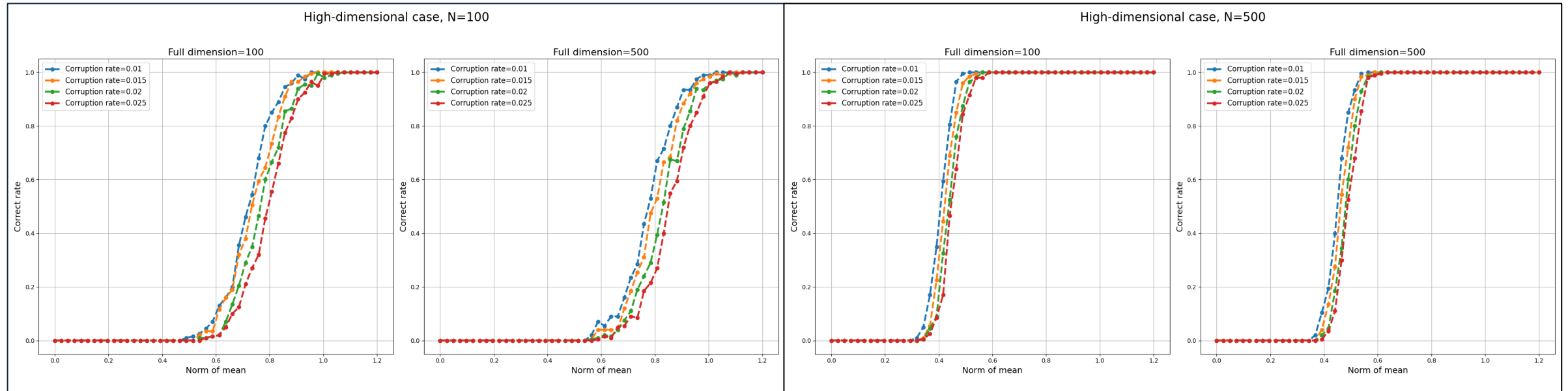
## Experiment results (classic scenario)





# Experiment

## Experiment results (high-dimensional scenario)

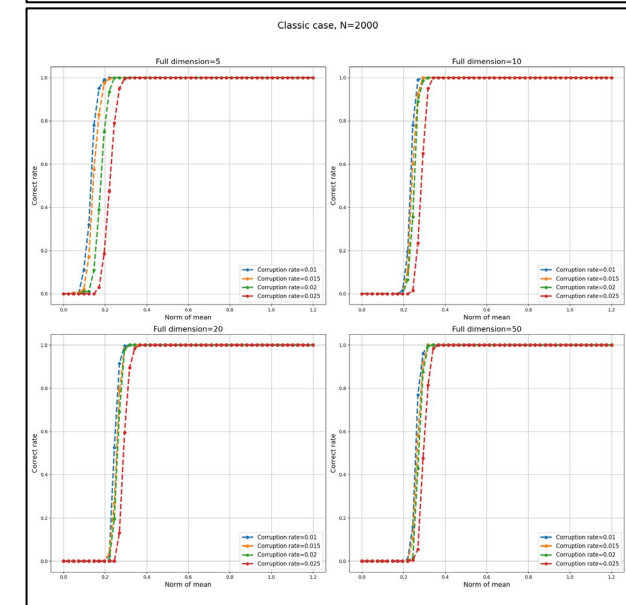
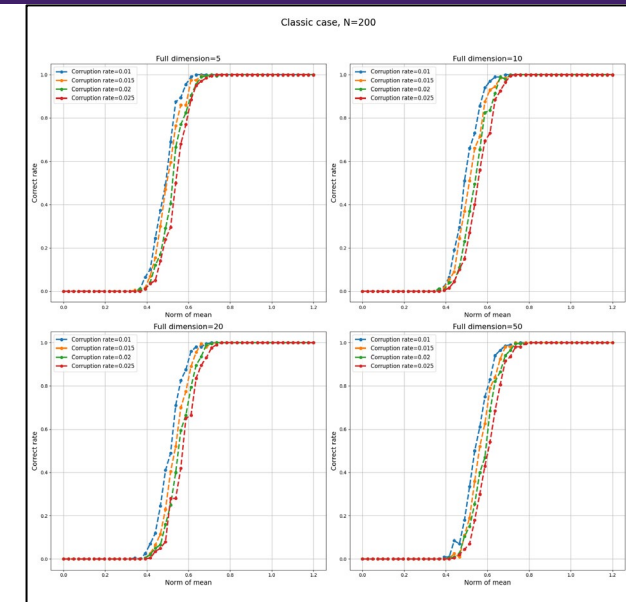




# Experiment

## Discussion

- The algorithm performs reasonably well across different parameter configurations
- The strategy of the adversary is effective
- **The algorithm is less sensitive to  $d$  as  $d$  increases**  
$$[5,10,20,50] \xrightarrow{\min\{k_1^*, k_2^*\}} [5,9,10,12]$$
- The algorithm is less sensitive to  $\epsilon$  as  $N$  increases





## Content

### Introduction

- Problem Formulation
- Literature Review

### Main Results

- Lower Bounds
- Upper Bounds
- Discussion

### Experiments

- Numerical Simulation
- Discussion

### Summary

- Summary
- Open Problems
- Plan for Future Works

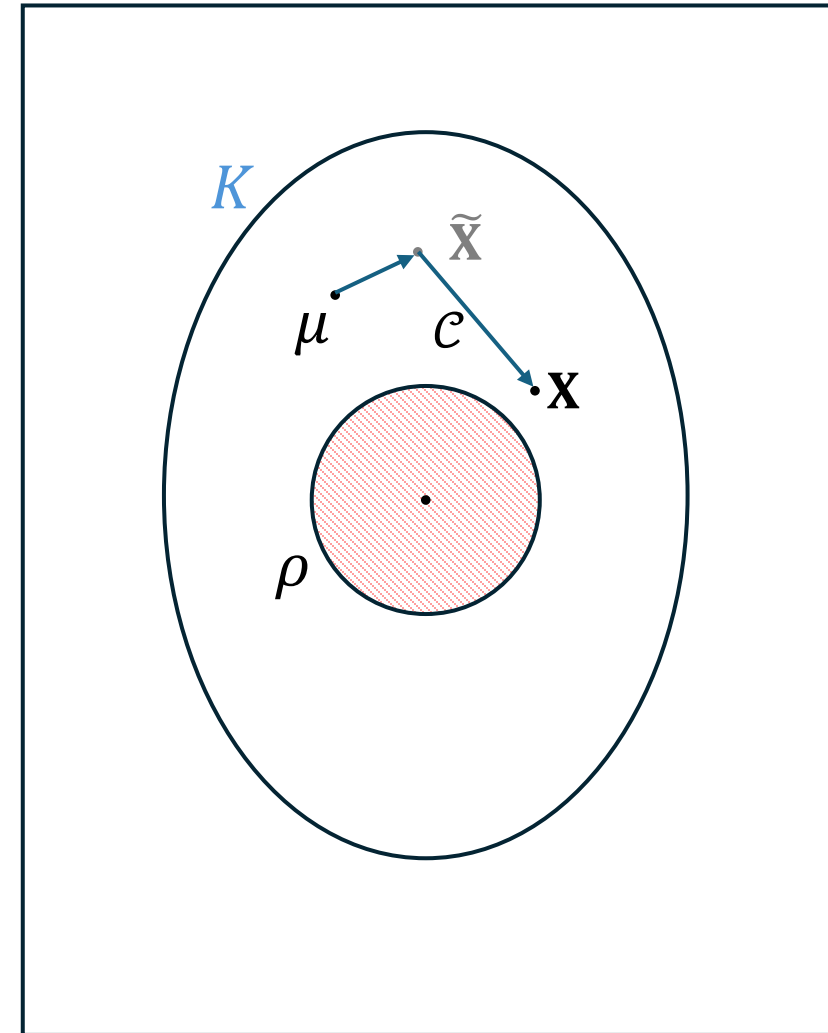


## Summary

### Conclusions & Takeaways

- Hypothesis testing problem with adversary and constraints
- Information-theoretic lower bounds
- Theoretical and polynomial-time algorithms
- Computational efficiency

$\mathbb{R}^d$



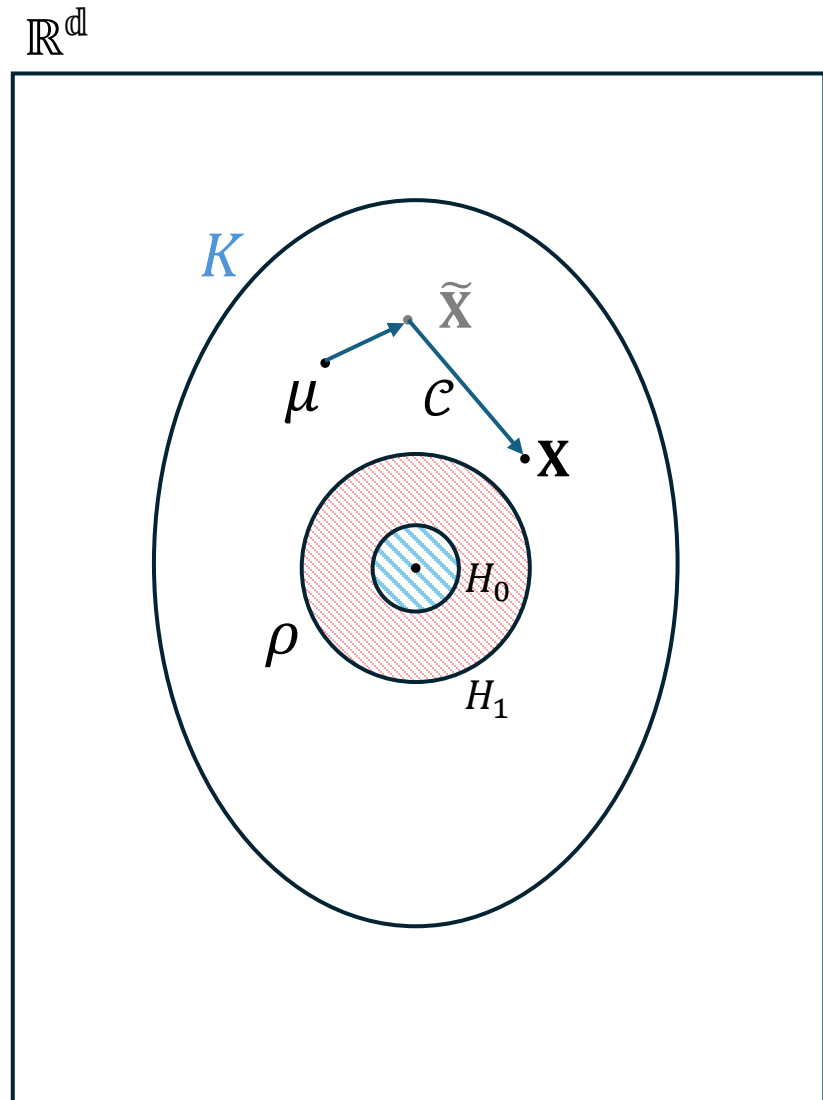
## Summary

### Open problems and possible future directions

- Kolmogorov width approximations (currently working)
- Imprecise  $H_0$  (Kania, et al., 2025)

$$H_0: \|\mu\|_2 \leq \rho_0$$
$$H_1: \|\mu\|_2 \geq \rho_1$$

- $\ell_p$  norm separation ( $p \geq 2$ ) ( $\chi^p$  test)
- More general constraints (Convex? Star-shaped? ...?)





Thank you!  
Q & A



## Appendix ( $\omega$ -regularity)

**Definition.** Given a weight vector  $\omega = (\omega_1, \dots, \omega_N) \in \mathbb{R}^N$  and any  $1 \leq k^* \leq d$ , samples  $\mathbf{Y} \in \mathbb{R}^{N \times d}$  is said to be  $(\epsilon, \beta_1, \beta_2)$ -regular if for all subsets  $S \in [N]$  with  $|S| \leq \epsilon N$ , we have:

$$(i), \left| \sum_{i \in S} \|Y_i\|_2^2 - |S|k^* \right| \leq c\beta_1;$$

$$(ii), \left| \left\| \sum_{i \in S} \sqrt{\omega_i} Y_i \right\|_2^2 - \|\omega_S\|_1 k^* \right| \leq c\beta_2;$$

$$(iii), \left| \left\langle \sum_{i \in S} \sqrt{\omega_i} Y_i, \sum_{j \in [N]} \sqrt{\omega_j} Y_j \right\rangle - \|\omega_S\|_1 k^* \right| \leq c\sqrt{N}\beta_1$$

where  $\omega_S \in \mathbb{R}^N$  means the truncated  $\omega$  on the subset  $S$  (fill truncated entries with zero).

## Appendix (prefiltering, sample filtering, and weight filtering)

---

Algorithm 4: Prefiltering.

---

```
1 Set  $\gamma_1 = c \left[ \sqrt{k^* \ln \left( \frac{N}{\alpha} \right) + \ln \left( \frac{N}{\alpha} \right)} \right]$ ,  $\text{count} = 0, i = 0$ 
2 while  $i < N$  do
3   if  $\left| \|Y_i\|_2^2 - k^* \right| > \gamma_1$  then
4      $\text{count} = \text{count} + 1$ 
5     if  $\text{count} > \epsilon N$  then
6       return None
7     Delete  $Y_i$  from  $\mathbf{Y}$ 
8    $i = i + 1$ 
9 return  $\mathbf{Y}$ 
```

---

## Appendix (prefiltering, sample filtering, and weight filtering)

---

**Algorithm 5:** Sample filtering when  $N > k^*$ .

---

- 1 Set  $\gamma_2$  as (C.13), and  $\lambda = \|\mathbf{Y}^\top D(\omega)\mathbf{Y} - N\mathbf{I}_{k^*}\|_2$ . ( $\omega$  is initialized as  $\mathbf{1}$ .)
- 2 **while**  $\lambda \geq \gamma_2$  **do**
- 3     Set  $v$  to be the unit singular vector associated with  $\lambda$
- 4     Compute  $\tau_i = \langle v, Y_i \rangle^2 \mathbf{1}_{\{\omega_i > 0\}}$  for  $1 \leq i \leq N$
- 5     Set  $I$  be the smallest index such that  $\sum_{i=1}^I \omega_i \geq 2\epsilon N$
- 6     Update  $w$  according to (C.14)
- 7     **if**  $\|\omega\|_1 < N(1 - 2\epsilon)$  **then**
- 8         **return** None
- 9     Set  $\lambda = \|\mathbf{Y}^\top D(\omega)\mathbf{Y} - N\mathbf{I}_{k^*}\|_2$
- 10 **return**  $\omega$

---

$$\omega_i^{(t+1)} = \begin{cases} \left(1 - \frac{\tau_i^{(t)}}{\tau_1^{(t)}}\right) \omega_i^{(t)} & \text{if } i \leq I, \\ \omega_i^{(t)} & \text{if } i > I. \end{cases}$$

## Appendix (prefiltering, sample filtering, and weight filtering)

---

**Algorithm 6:** Sample filtering when  $N \leq k^*$ .

---

- 1 Set  $\gamma_3$  as (C.21), and  $\lambda = \left\| \sqrt{D(\omega)} \mathbf{Y} \mathbf{Y}^\top \sqrt{D(\omega)} - k^* D(\omega) \right\|_2$ . ( $\omega$  is initialized as  $\mathbf{1}$ .)
  - 2 while  $\lambda \geq \gamma_3$  do
    - 3 Set  $v$  to be the unit singular vector associated with  $\lambda$
    - 4 Compute  $\tau_i = \frac{v_i^2}{\omega_i} \mathbf{1}_{\{w_i > 0\}}$
    - 5 Update  $w$  according to (C.22)
    - 6 if  $\|\omega\|_1 < N(1 - 6\epsilon)$  then
      - 7 return None
    - 8 Set  $\lambda = \left\| \sqrt{D(\omega)} \mathbf{Y} \mathbf{Y}^\top \sqrt{D(\omega)} - k^* D(\omega) \right\|_2$
  - 9 return  $\omega$
- 

$$\omega_i = \left( 1 - \frac{\tau_i}{\max_i \tau_i} \right) \omega_i, 1 \leq i \leq N$$



## Appendix (prefiltering, sample filtering, and weight filtering)

---

Algorithm 7: Weight filtering.

---

- 1 Compute  $\tau_i = \left| \left\langle \sqrt{\omega_i} Y_i, \sum_{j=1}^N \sqrt{\omega_j} Y_j \right\rangle - \omega_i k^* \right| \mathbf{1}_{\{\omega_i > 0\}}, 1 \leq i \leq N$
  - 2 Sort  $\tau_i$  by decreasing order and find the indices  $\{i_1, i_2, \dots, i_{\epsilon N}\}$  corresponding to the first  $\epsilon N$  maximal  $\tau_i$
  - 3 Set  $\omega_{i_1}, \omega_{i_2}, \dots, \omega_{i_{\epsilon N}}$  to zero
  - 4 return  $\omega$
-

## Appendix (alternative strategy for $\mathcal{C}$ )

---

**Algorithm 2:** strategy 2 of the adversary  $\mathcal{C}$

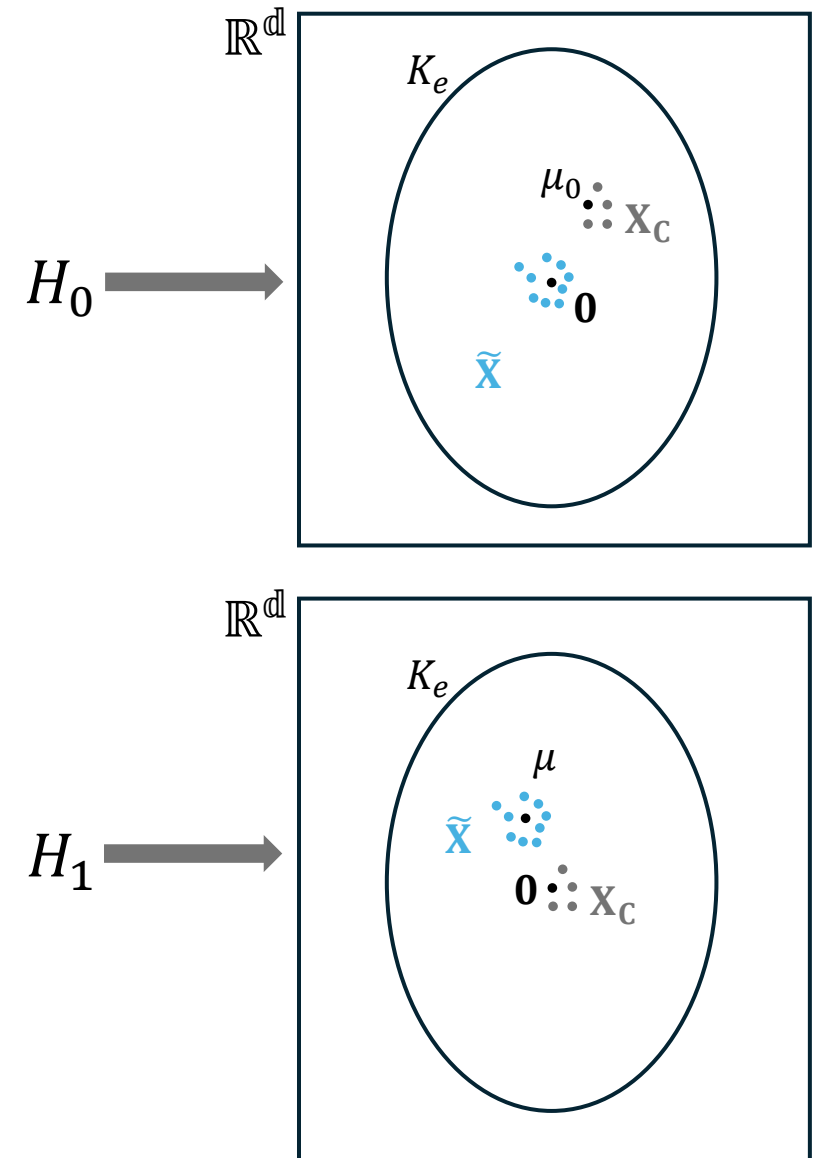
---

```

1 import  $\tilde{\mathbf{X}}, N, d, \sigma = 1$ 
2 if  $H_1$  is true then
3   Compute  $\mathbf{dist} \in \mathbb{R}^N$ ,  $\mathbf{dist}_i = \|\tilde{\mathbf{X}}_i\|_2, 1 \leq i \leq N$ 
4    $C =$  indices of the largest  $\epsilon N$  entries in  $\mathbf{dist}$ 
5    $\mathbf{X}_C = \epsilon N$  i.i.d. fake data from  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ 
6    $\tilde{\mathbf{X}}_C = \mathbf{X}_C$ 
7 else
8   Compute  $\mu_0$  according to the constraint type
9   Compute  $\mathbf{dist} \in \mathbb{R}^N$ ,  $\mathbf{dist}_i = \|\tilde{\mathbf{X}}_i - \mu_0\|_2$ 
10   $C =$  indices of the largest  $\epsilon N$  entries in  $\mathbf{dist}$ 
11   $\mathbf{X}_C = \epsilon N$  i.i.d. fake data from  $\mathcal{N}(\mu_0, \mathbf{I}_d)$ 
12   $\tilde{\mathbf{X}}_C = \mathbf{X}_C$ 
13 return  $\mathbf{X} = \mathcal{C}(\tilde{\mathbf{X}})$ 

```

---

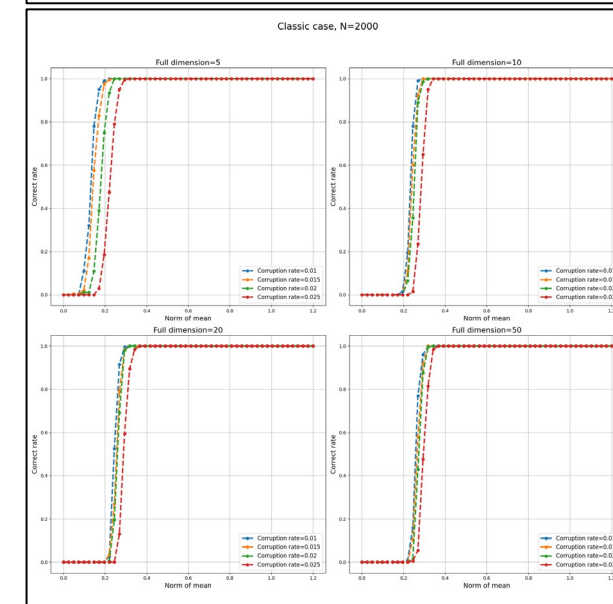
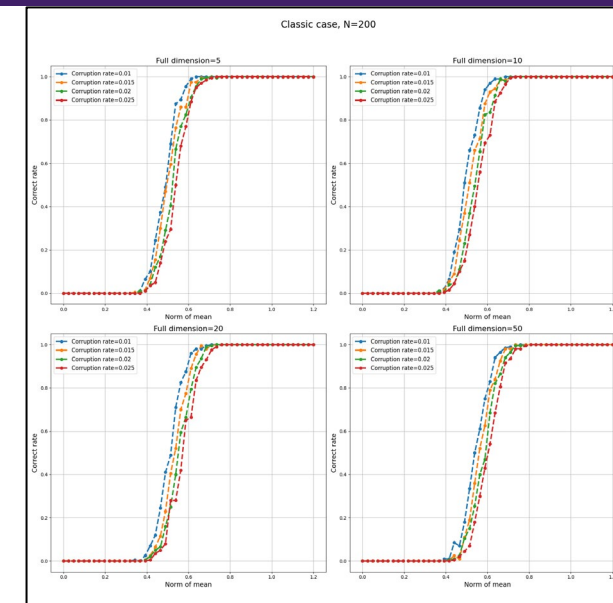




# Experiment

## Limitations

- Unable to verify the match between the lower and upper bounds
- Optimality of the adversary?

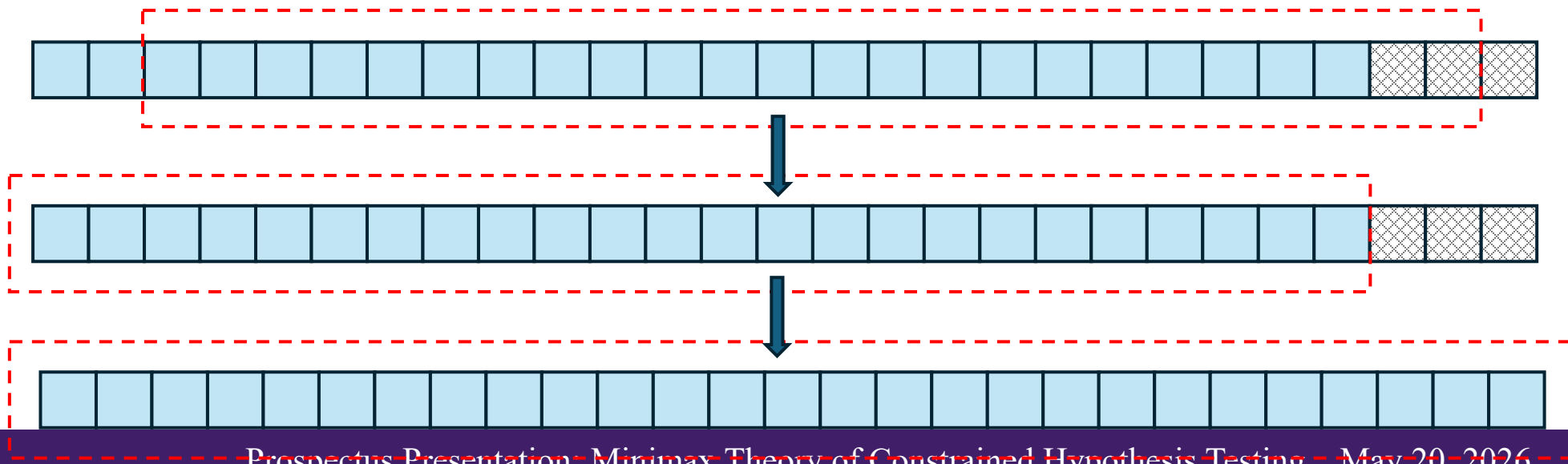


## Main Results (Theoretical Algorithm)

**Definition (consistent subset).** A subset  $S \subset \mathbf{X} = \{X_1, \dots, X_N\}$  is called a consistent subset regarding  $\epsilon$  and a test  $\phi: 2^S \rightarrow \{0,1\}$  if:

- (i),  $|S| \geq (1 - \epsilon)N$ ,
- (ii),  $\phi(S') = \phi(S)$  for any  $S' \subset S$  with  $|S'| \geq (1 - 2\epsilon)N$ .

Motivation: consistent subset can help recover the testing result on the unobserved original samples!



## What is Minimax Theory and Risk?

We interact with the general environment.

- Our decision or action yields rewards, penalties, etc

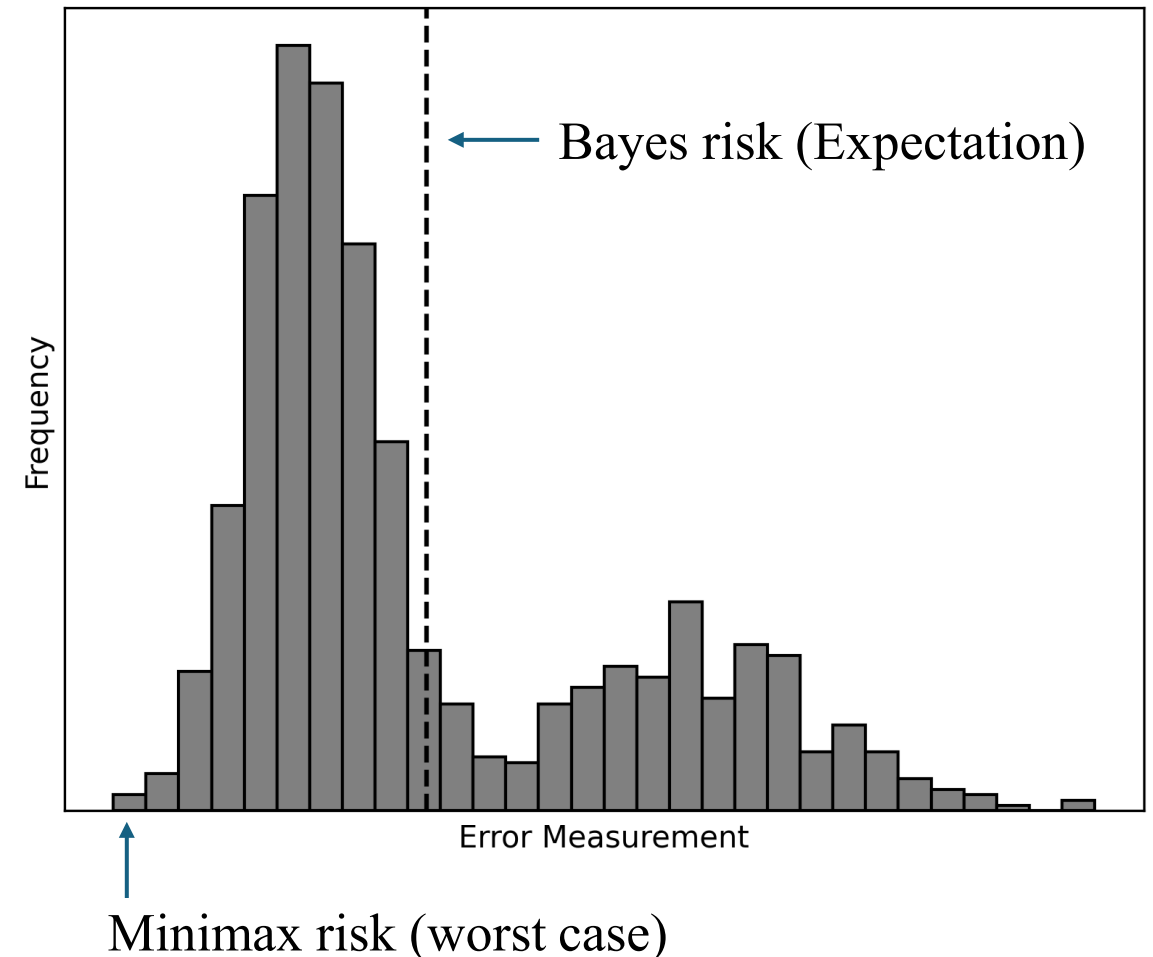
The randomness exists extensively.

- The response is not fixed

How to measure our performance?

- Model the response with a distribution, and calculate the average (Bayes risk)
- Only evaluate from the worst case (**minimax risk**)

*In the worst case, we will...*





## Minimax risk is ubiquitous in real-world applications.

### Robust Control

- LQR controller: noise is well-behaved and predictable
- Robust control ( $H_\infty$  control): zero-sum game with nature
- Nature's goal: find worst-case disturbance (wind, friction loss, sensor noise) that will **maximize** system's error.
- Our goal: design a control policy that **minimizes** the system's error, **even when** Nature plays its best move.
- Example: Autonomous driving, aerospace and flight control, etc.



Source: <https://www.tesla.com/fsd/safety>

# Minimax risk is ubiquitous in real-world applications.

## Robust Finance

- Should we assume that the market has some “desired” properties, or its behaviors can be very “wild”?
- The famous Black-Scholes model assumes that the volatility  $\sigma$  of a stock is constant and known. If volatility spikes unexpectedly, this can be disastrous
- In contrast, modern derivative dealers assume volatility is unknown but bounded within a range  $[\sigma_{min}, \sigma_{max}]$ . We want to find a dynamic hedging strategy that **minimizes the maximum possible hedging**



The front page of the New York Times on September 16, 2008. (Flickr/kbaird)

Source: <https://www.americanprogress.org/article/think-again-the-crisis-from-nowhere/>

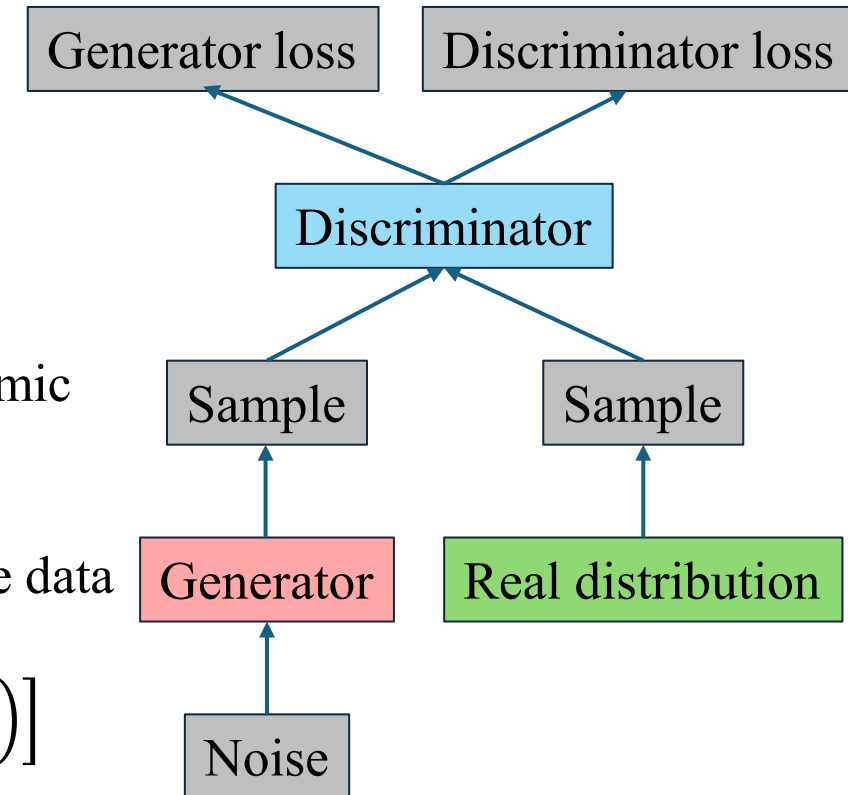


## Minimax risk is ubiquitous in real-world applications.

### Generative Adversarial Network (GAN)

- Introduced by Ian Goodfellow et al. in 2014, GAN is a landmark in generative artificial intelligence
- The **generator** learns to generate plausible data from pure noise to mimic the real data
- The **discriminator** learns to distinguish between the real data and fake data

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim p_z} [\log (1 - D(G(z)))]$$



Source: [https://developers.google.com/machine-learning/gan/gan\\_structure](https://developers.google.com/machine-learning/gan/gan_structure)