

Gaussian Mean Testing under Truncation

Clément L. Canonne*, Themis Gouleakis, Yuhao Wang, Joy Qiping Yang*

*University of Sydney, Nanyang Technological University, National University of Singapore

Introduction

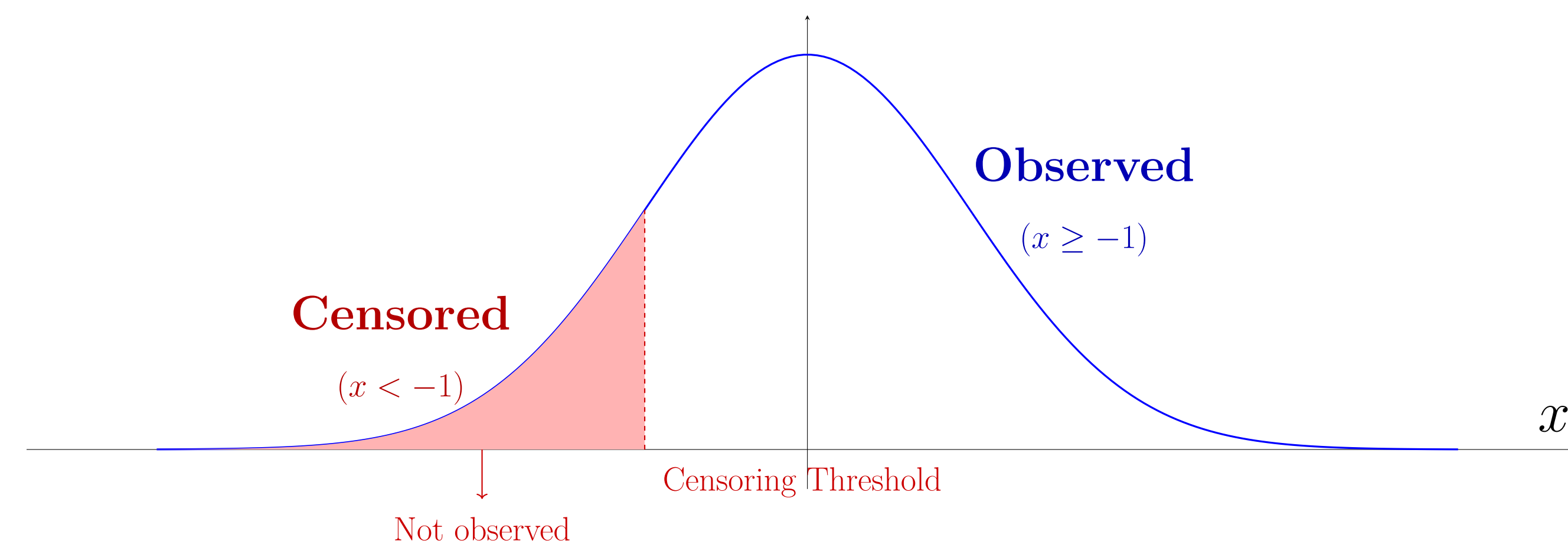
Problem: Learning a high-dimensional Gaussian distribution when data is **missing not at random (MNAR)**.

- **Self-Censoring:** Each coordinate is missing depending on its own.
- **Linear Thresholding:** Missingness depends linearly on the full data vector.

Self-Censoring Model

Model:

- Observations are censored versions of samples from $\mathcal{N}(\mu^*, \Sigma^*)$.
- Coordinate i is observed if $y_i \in S_i$.



Main Theorem (Recovery under Self-Censoring)

- If any two coordinates are observed together with probability at least α .
- Algorithm learns μ^* and Σ^* up to TV distance ε .
- Sample complexity: $O(d^2(\lambda_{\max}/\lambda_{\min})^2/\alpha\varepsilon^2)$.

Algorithm: Truncation_PSGD

Algorithm 1: Recovery of Gaussian parameters under self-censoring via truncated PSGD.

Input: Data $\mathbf{x} \in \mathbb{R}^{n \times d}$, where $n = \frac{1}{\alpha\varepsilon^2}$

- 1 for $i \leftarrow 1$ to d do
- 2 $\hat{\mu}_i, \hat{\Sigma}_{ii} \leftarrow \text{Uni_SGD_trunc}(\mathbf{x}_i, S_i)$
- 3 for $i \leftarrow 1$ to $d-1$ do
- 4 for $j \leftarrow i+1$ to d do
- 5 $\hat{\Sigma}_{ij}^2 \leftarrow \text{Biv_SGD_trunc}(\mathbf{x}_i, \mathbf{x}_j, S_i \times S_j)$
- 6 $\hat{\mu} \leftarrow [\hat{\mu}_1, \dots, \hat{\mu}_d]$
- 7 for $i \leftarrow 1$ to $d-1$ do
- 8 for $j \leftarrow i+1$ to d do
- 9 $\hat{\Sigma}_{ij}, \hat{\Sigma}_{ji} \leftarrow \hat{\Sigma}_{ij}^2$
- 10 return $(\hat{\mu}, \hat{\Sigma})$

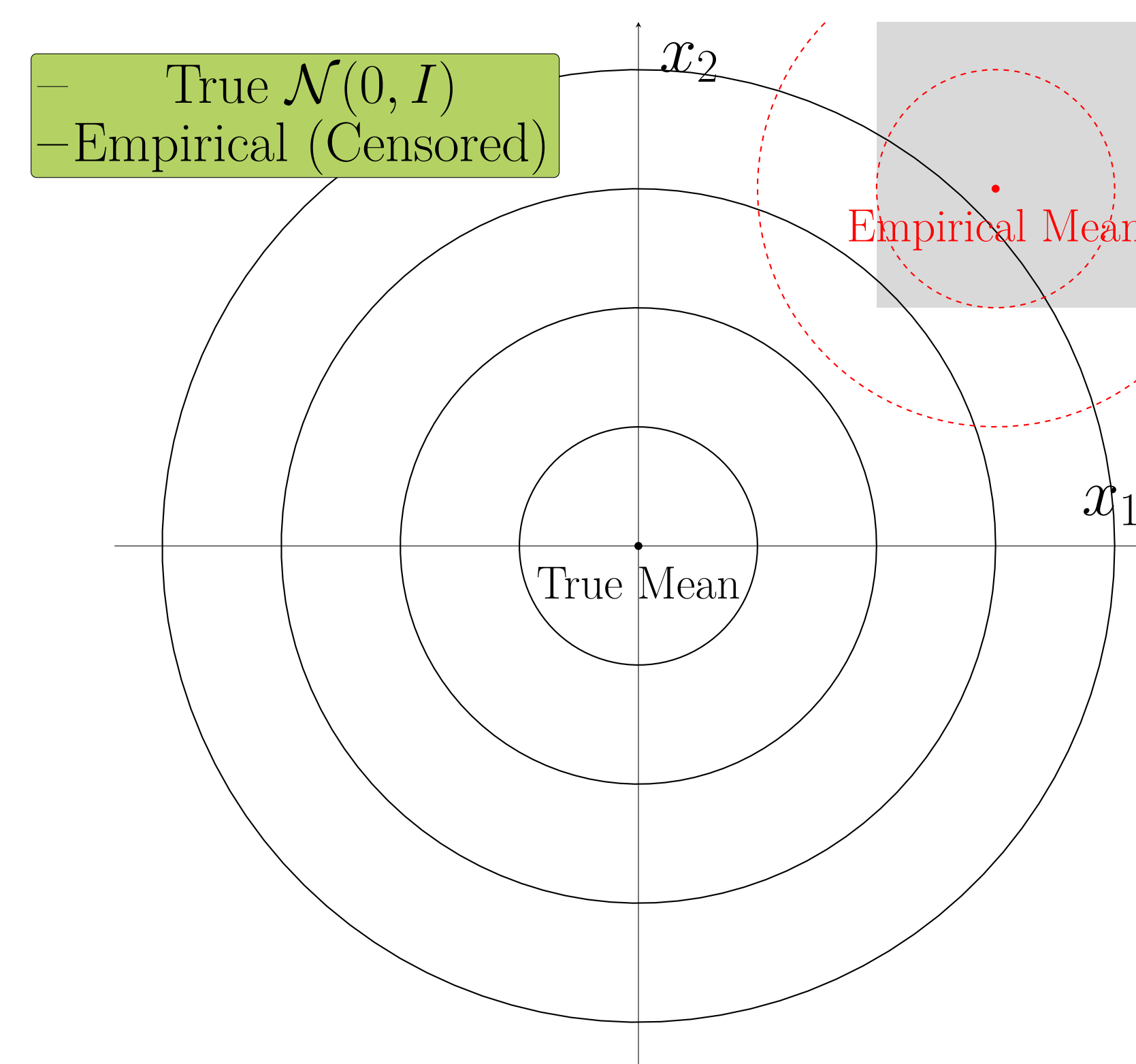
Folklore: Learning a Gaussian Without Missingness

Given n i.i.d. samples $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n)} \sim \mathcal{N}(\mu^*, \Sigma^*)$:

- Estimate the mean: $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}^{(i)}$
- Estimate the covariance: $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\mathbf{y}^{(i)} - \hat{\mu})(\mathbf{y}^{(i)} - \hat{\mu})^\top$
- With $n = \tilde{O}(d^2/\varepsilon^2)$ samples, we recover:

$$\|\mu^* - \hat{\mu}\| = \mathcal{O}(\varepsilon), \quad \|\Sigma^* - \hat{\Sigma}\|_F = \mathcal{O}(\varepsilon)$$

Censoring



Linear Thresholding Model

Model: Each coordinate is observed if a random linear inequality is satisfied. For coordinate $i \in [d]$, we observe y_i only if:

$$v_i^\top \mathbf{y} \leq b_i$$

where $v_i \in \mathbb{R}^d$ and $b_i \in \mathbb{R}$ are (possibly random) threshold parameters.

Goal: Recover the true mean μ^* given covariance matrix Σ^* of a Gaussian $\mathcal{N}(\mu^*, \Sigma^*)$ from censored samples.

Assumptions:

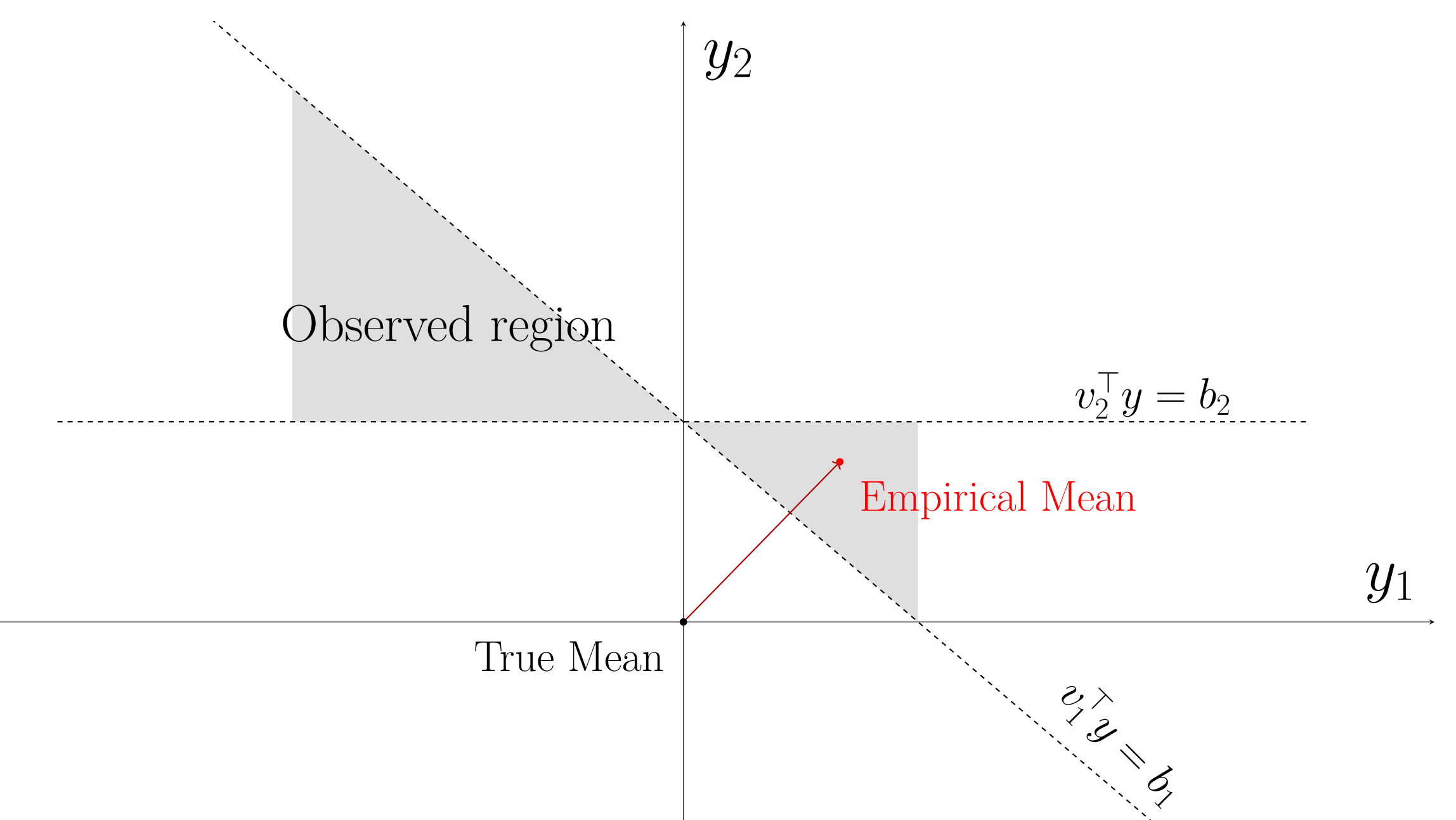
- Small subsets (up to βd coordinates) are observed together with probability at least α .
- A stable anchoring subset C is always observed (parameter γ).

Main Result:

- If Σ is known, there is a polynomial-time algorithm to estimate μ^* .
- Sample complexity:

$$\text{poly}(d, 1/\alpha, 1/\beta, 1/\gamma, \lambda_{\max}/\lambda_{\min}, 1/\varepsilon, \log(1/\delta))$$

Linear Thresholding Missingness Model



Algorithm: MissingDescent

Algorithm 2: Mean recovery algorithm given oracle access to incomplete data samples.

Input: Access to data generator \mathcal{O} ; parameters

$\beta, \lambda_{\text{sgd}}, \eta_{\text{lmc}}, R_{\text{lmc}}, r_{\text{proj}}, M_{\text{init}}, M_{\text{sgd}}, M_{\text{grad}}$

- 1 $\mu^{(0)} \leftarrow \text{Initialize}(\mathcal{O}, \beta, M_{\text{init}})$
- 2 for $i \leftarrow 1$ to M_{sgd} do
- 3 Sample $(A^{(i)}, \mathbf{x}^{(i)})$ from \mathcal{O} ;
- 4 $\eta_i \leftarrow \frac{1}{\lambda_{\text{sgd}} \cdot i}$;
- 5 $\mathbf{g}^{(i)} \leftarrow \text{SampleGradient}((A^{(i)}, \mathbf{x}^{(i)}), \mu^{(i-1)}, \eta_{\text{lmc}}, R_{\text{lmc}}, M_{\text{grad}})$;
- 6 $\mathbf{v}^{(i)} \leftarrow \mu^{(i-1)} - \eta_i \mathbf{g}^{(i)}$;
- 7 $\mu^{(i)} \leftarrow \text{ProjectToDomain}(\mu^{(0)}, \mathbf{v}^{(i)}, r_{\text{proj}})$;
- 8 $\bar{\mu} \leftarrow \frac{1}{M_{\text{sgd}}} \sum_{i=1}^{M_{\text{sgd}}} \mu^{(i)}$;
- 9 return $\bar{\mu}$

Future Direction

- Extend recovery guarantees to unknown covariance settings.
- Study more general non-linear missingness patterns.
- Explore robust estimation under adversarial censoring.



Paper Link