

Learning Distributions Generated by One-Layer ReLU Networks

Shanshan Wu, Alex Dimakis, Sujay Sanghavi

[NeurIPS 2019]

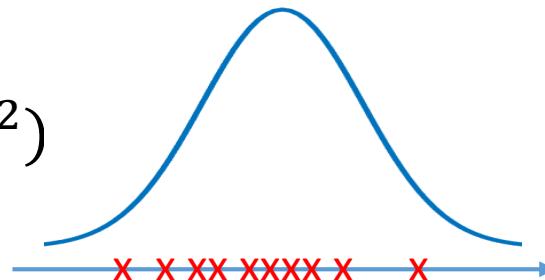
Background: Density Estimation

Input: i.i.d. samples \sim unknown distribution



Goal: estimate the distribution

E.g., samples $\sim N(\mu, \sigma^2)$



Estimate μ, σ^2

E.g., samples \sim natural images



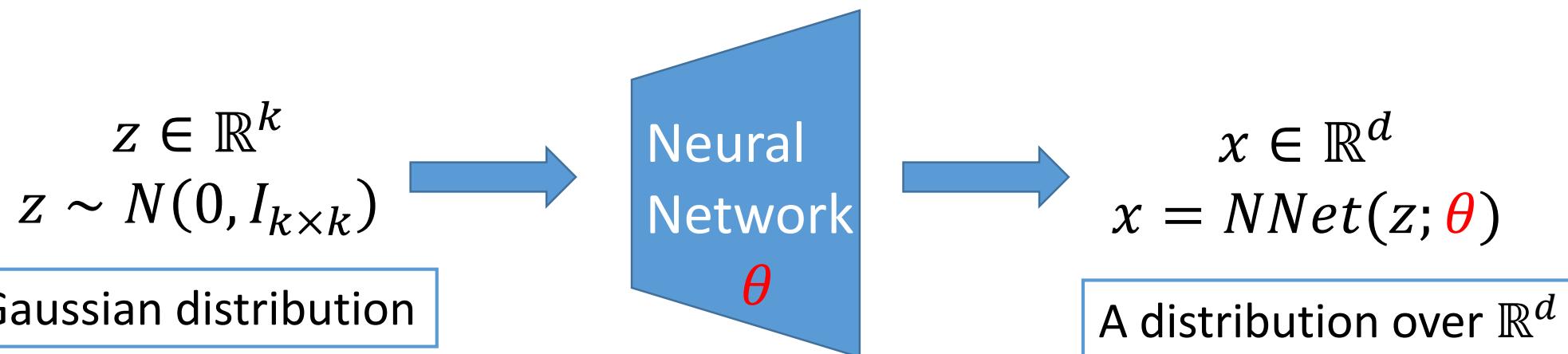
ImageNet



Estimate the parameters
of a **generative model**

Background: Generative Model

- A popular way to model **high-dimensional complex distributions**:



- Given i.i.d. samples, θ can be learned by

- GAN [Goodfellow et al.'2014]
- VAE [Kingma and Welling'2013]
-

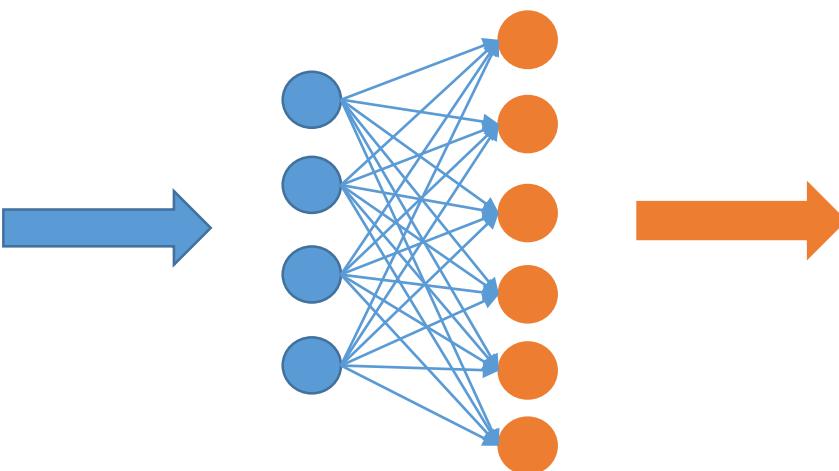
No guarantees & sample complexity



Problem Formulation

- One-layer ReLU generative model with params $W \in \mathbb{R}^{d \times k}$, $b \in \mathbb{R}^d$:

$$\begin{aligned} z &\in \mathbb{R}^k \\ z &\sim N(0, I_{k \times k}) \end{aligned}$$



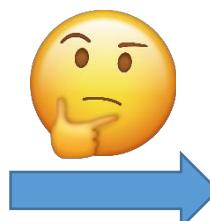
$$\text{ReLU}(Wz + b) \in \mathbb{R}^d$$

A distribution over \mathbb{R}^d

- Our problem:

Given n i.i.d. samples

$$x_1, x_2, \dots, x_n \sim \text{ReLU}(Wz + b)$$

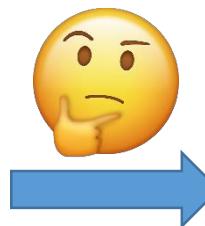


Can we estimate W, b ?

Problem Formulation

- Our problem is **unsupervised learning**:

Given n i.i.d. samples in \mathbb{R}^d
 $x_1, x_2, \dots, x_n \sim \text{ReLU}(Wz + b)$



Can we estimate W, b ?

- Different from **supervised learning** problems [Ge et al., Goel et al., ...]

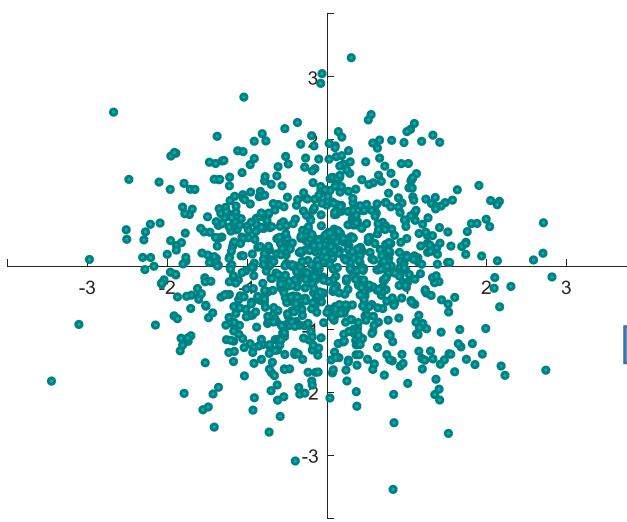
Given n (**input, output**) pairs
 $(z_1, x_1), \dots, (z_n, x_n)$ where
 $x_i = \text{ReLU}(Wz_i + b)$



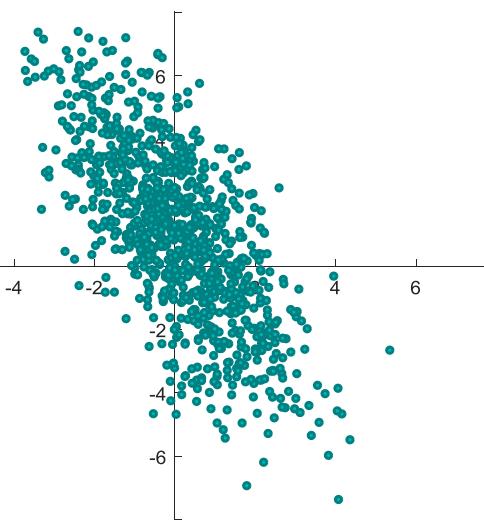
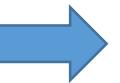
Can we estimate W, b ?

Visualization of Our Problem

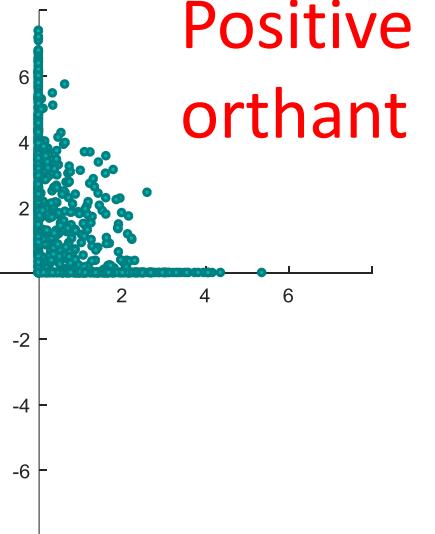
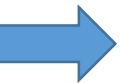
- $W^* = \begin{bmatrix} -0.9 & -1.2 \\ 0.4 & 2.7 \end{bmatrix}$, $b^* = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$, distribution of $\text{ReLU}(W^*z + b^*)$:



$z \sim N(0, I_{2 \times 2})$



$W^*z + b^*$



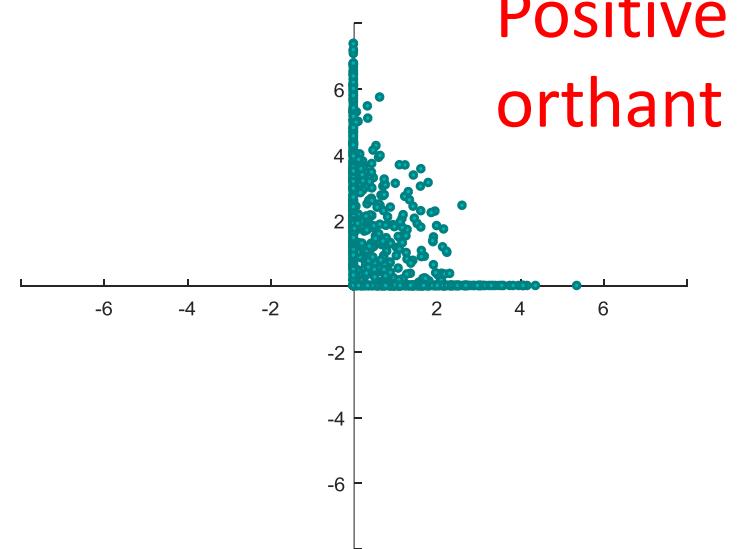
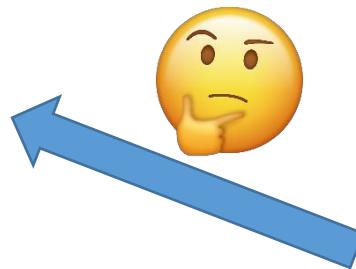
$\text{ReLU}(W^*z + b^*)$

Positive
orthant

Visualization of Our Problem

- $W^* = \begin{bmatrix} -0.9 & -1.2 \\ 0.4 & 2.7 \end{bmatrix}, b^* = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$

Can we recover W^*, b^* ?



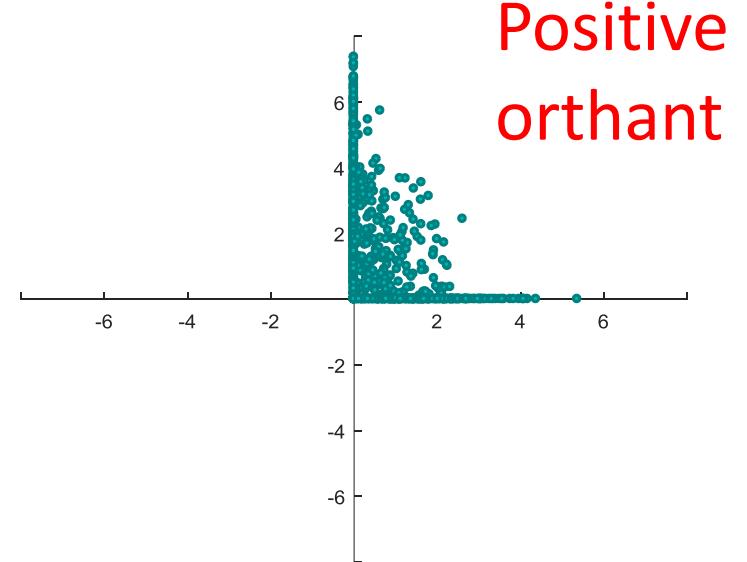
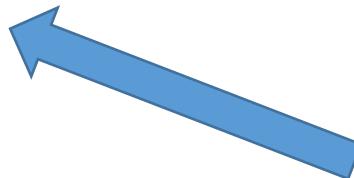
Samples $\sim \text{ReLU}(W^*z + b^*)$

Visualization of Our Problem

- $W^* = \begin{bmatrix} -0.9 & -1.2 \\ 0.4 & 2.7 \end{bmatrix}, b^* = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$

Can we recover W^*, b^* ?

Are they identifiable?



Samples $\sim \text{ReLU}(W^*z + b^*)$

Identifiability

- Is $W \in \mathbb{R}^{k \times d}$ identifiable from the **distribution** $\text{ReLU}(Wz + b)$? No!

Only WW^T can be possibly identified.

- Fact:

$$W_1 W_1^T = W_2 W_2^T$$



$$\text{distr of } \text{ReLU}(W_1 z + b) = \text{distr of } \text{ReLU}(W_2 z + b)$$

Identifiability

- Is $W \in \mathbb{R}^{k \times d}$ identifiable from the **distribution** $\text{ReLU}(Wz + b)$? No!

Only WW^T can be possibly identified.

- Fact:

$$W_1 W_1^T = W_2 W_2^T$$



$$\text{distr of } \text{ReLU}(W_1 z + b) = \text{distr of } \text{ReLU}(W_2 z + b)$$
$$\underbrace{\qquad\qquad\qquad}_{N(b, W_1 W_1^T)} \qquad\qquad\qquad \underbrace{\qquad\qquad\qquad}_{N(b, W_2 W_2^T)}$$

Identifiability

- Is $b \in \mathbb{R}^d$ identifiable from the distribution $\text{ReLU}(Wz + b)$? Yes, but...

If b is negative



Estimate b needs
 $\Omega(\exp(\|b\|_\infty^2))$
samples

Identifiability

- Is $b \in \mathbb{R}^d$ identifiable from the **distribution** $\text{ReLU}(Wz + b)$? **Yes, but...**

If b is **negative**



Estimate b needs
 $\Omega(\exp(\|b\|_\infty^2))$
samples

Claim:

Distinguish $\text{ReLU}(z - |b|)$ & $\text{ReLU}(z - |b| - 0.1)$
needs $\Omega(\exp(b^2))$ samples.

Identifiability

- Is $b \in \mathbb{R}^d$ identifiable from the distribution $\text{ReLU}(Wz + b)$? Yes, but...

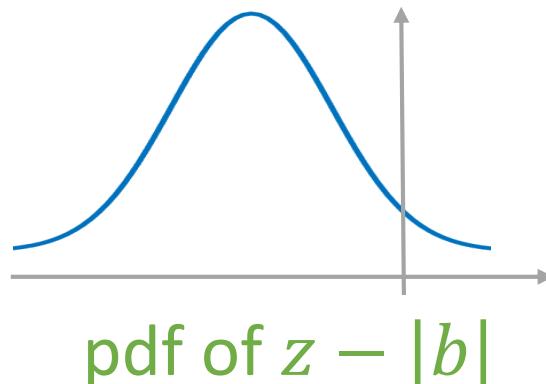
If b is negative



Estimate b needs
 $\Omega(\exp(\|b\|_\infty^2))$
samples

Claim:

Distinguish $\text{ReLU}(z - |b|)$ & $\text{ReLU}(z - |b| - 0.1)$
needs $\Omega(\exp(b^2))$ samples.



Identifiability

- Is $b \in \mathbb{R}^d$ identifiable from the **distribution** $\text{ReLU}(Wz + b)$? Yes, but...

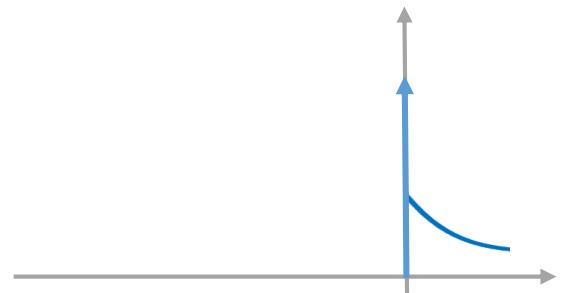
If b is negative



Estimate b needs
 $\Omega(\exp(\|b\|_\infty^2))$
samples

Claim:

Distinguish $\text{ReLU}(z - |b|)$ & $\text{ReLU}(z - |b| - 0.1)$
needs $\Omega(\exp(b^2))$ samples.



pdf of $\text{ReLU}(z - |b|)$

$$\mathbb{P}[\text{ReLU}(z - |b|) > 0] \leq \exp(-b^2)$$

Identifiability

- Is $b \in \mathbb{R}^d$ identifiable from the **distribution** $\text{ReLU}(Wz + b)$? Yes, but...

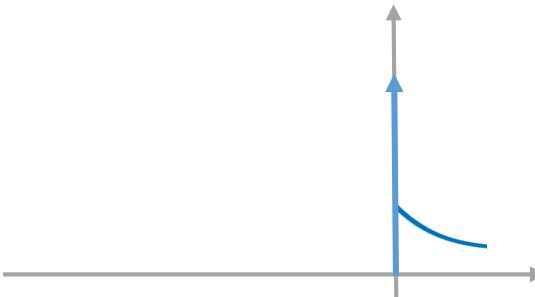
If b is negative



Estimate b needs
 $\Omega(\exp(\|b\|_\infty^2))$
samples

Claim:

Distinguish $\text{ReLU}(z - |b|)$ & $\text{ReLU}(z - |b| - 0.1)$
needs $\Omega(\exp(b^2))$ samples.



$$\mathbb{P}[\text{ReLU}(z - |b| - 0.1) > 0] \leq \exp(-b^2)$$

pdf of $\text{ReLU}(z - |b| - 0.1)$

Identifiability

- Is $b \in \mathbb{R}^d$ identifiable from the **distribution** $\text{ReLU}(Wz + b)$? Yes, but...

If b is negative

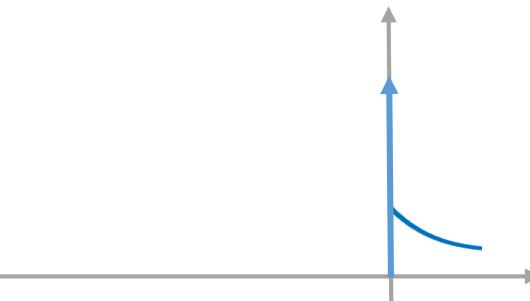


Estimate b needs
 $\Omega(\exp(\|b\|_\infty^2))$
samples

Assumption: b is non-negative

Claim:

Distinguish $\text{ReLU}(z - |b|)$ & $\text{ReLU}(z - |b| - 0.1)$
needs $\Omega(\exp(b^2))$ samples.



$$\mathbb{P}[\text{ReLU}(z - |b| - 0.1) > 0] \leq \exp(-b^2)$$

pdf of $\text{ReLU}(z - |b| - 0.1)$

Our Problem

Observe $x_1, x_2, \dots, x_n \sim \text{ReLU}(W^*z + b^*)$



$z \sim N(0, I_{k \times k})$

Unknown $\mathbb{R}^{d \times k}$

Unknown \mathbb{R}^d & non-negative

Goal: Estimate W^*W^{*T} and b^*

Initial Attempts

- **Idea 1:** Maximum likelihood estimation
- Problem: 1) no closed-form
2) non-convex



- **Idea 2:** Matrix completion

$$x_1, x_2, \dots, x_n \\ d \quad | \quad | \quad | \quad \dots \quad | \quad n \\ = \text{ReLU} \left(\begin{matrix} k & z_1, z_2, \dots, z_n \\ d \begin{bmatrix} W \end{bmatrix} & | \quad | \quad | \quad \dots \quad | \end{matrix} \right)$$

Observe positive entries of rank- k matrix



Experiments
do not work!

Idea 3: Truncated Gaussian

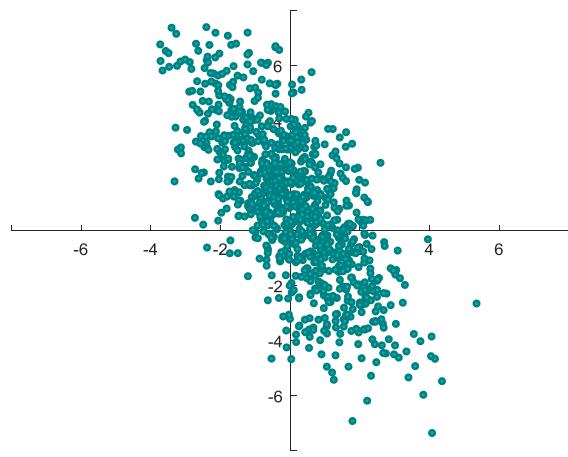
- Given a set $S \subset \mathbb{R}^d$, the pdf of an S -truncated Gaussian distribution

$$q(x; \mu, \Sigma, S) = \begin{cases} \frac{p(x; \mu, \Sigma)}{\int_S p(y; \mu, \Sigma) dy}, & \text{if } x \in S \\ 0, & \text{if } x \notin S \end{cases}$$

↑
pdf of $N(\mu, \Sigma)$

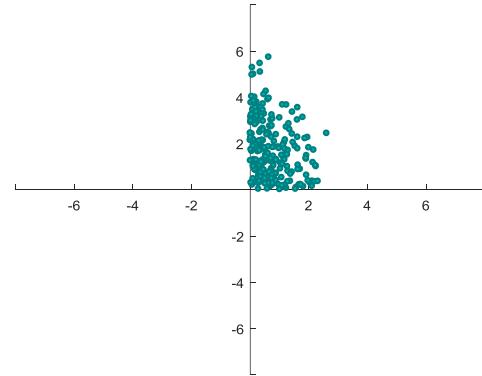
- Given truncated samples, estimate μ, Σ [Daskalakis et al., 2018]

Truncation versus ReLU

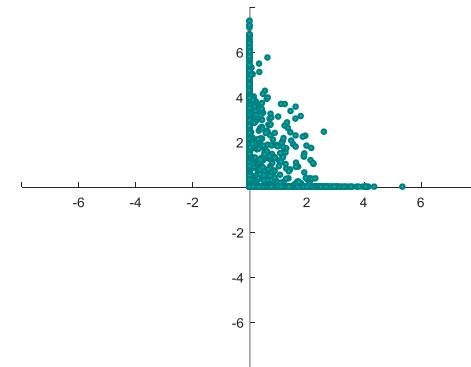


$\text{samples} \sim N(\mu, \Sigma)$

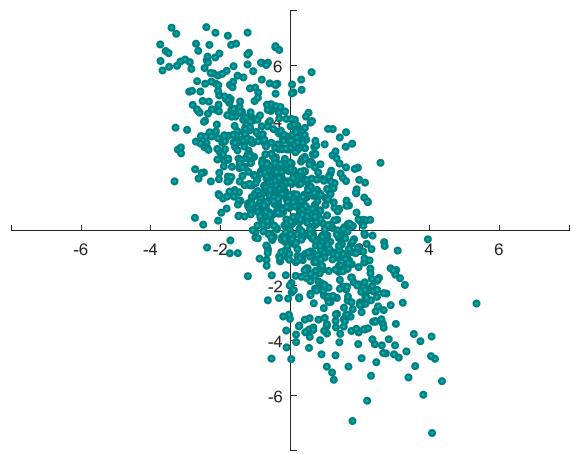
Truncation with
 $S = \{x \in \mathbb{R}^d : x > 0\}$



ReLU(samples)

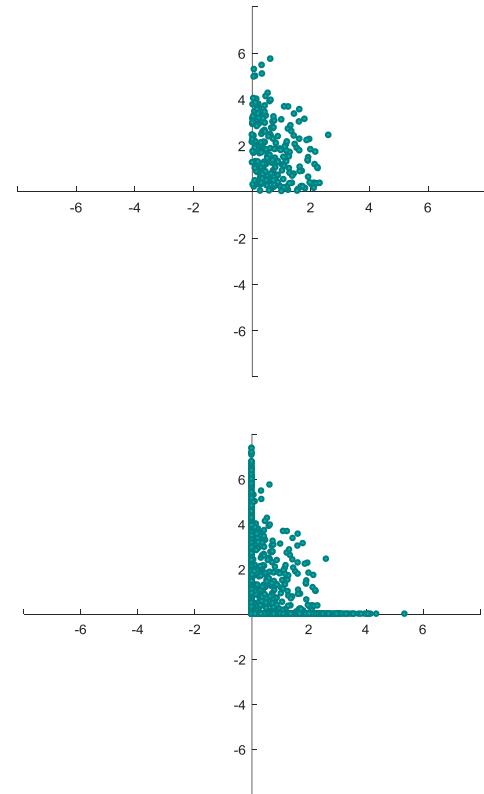


Convert ReLU to Truncated Samples



$\text{samples} \sim N(\mu, \Sigma)$

ReLU(samples)
Truncation with
 $S = \{x \in \mathbb{R}^d : x > 0\}$



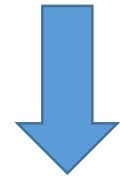
Convert ReLU to Truncated Samples

Samples $\sim \text{ReLU}(\mathbf{W}z + \mathbf{b})$



Remove the samples outside
 $S = \{x \in \mathbb{R}^d : x > 0\}$

Samples $\sim S\text{-truncated } N(\mathbf{b}, \mathbf{W}\mathbf{W}^T)$

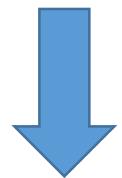


Run [Daskalakis et al., 2018]

Estimate $\mathbf{b}, \mathbf{W}\mathbf{W}^T$

Convert ReLU to Truncated Samples

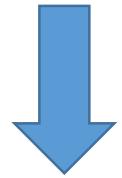
Samples $\sim \text{ReLU}(\mathbf{W}z + \mathbf{b})$



Remove the samples outside
 $S = \{x \in \mathbb{R}^d : x > 0\}$



Samples $\sim S$ -truncated $N(\mathbf{b}, \mathbf{W}\mathbf{W}^T)$



Run [Daskalakis et al., 2018]

Estimate $\mathbf{b}, \mathbf{W}\mathbf{W}^T$

$$x = \begin{pmatrix} \text{ReLU}(\nu^T z) \\ \text{ReLU}(-\nu^T z) \end{pmatrix}$$



$$\mathbb{P}[x \in S] = 0$$

Our Algorithm: Overview

- Our problem:

Given $x_1, \dots, x_n \sim \text{ReLU}(Wz + b)$, recover $WW^T \in \mathbb{R}^{d \times d}$, $b \in \mathbb{R}^d$.

- (i, j) -th entry of WW^T is

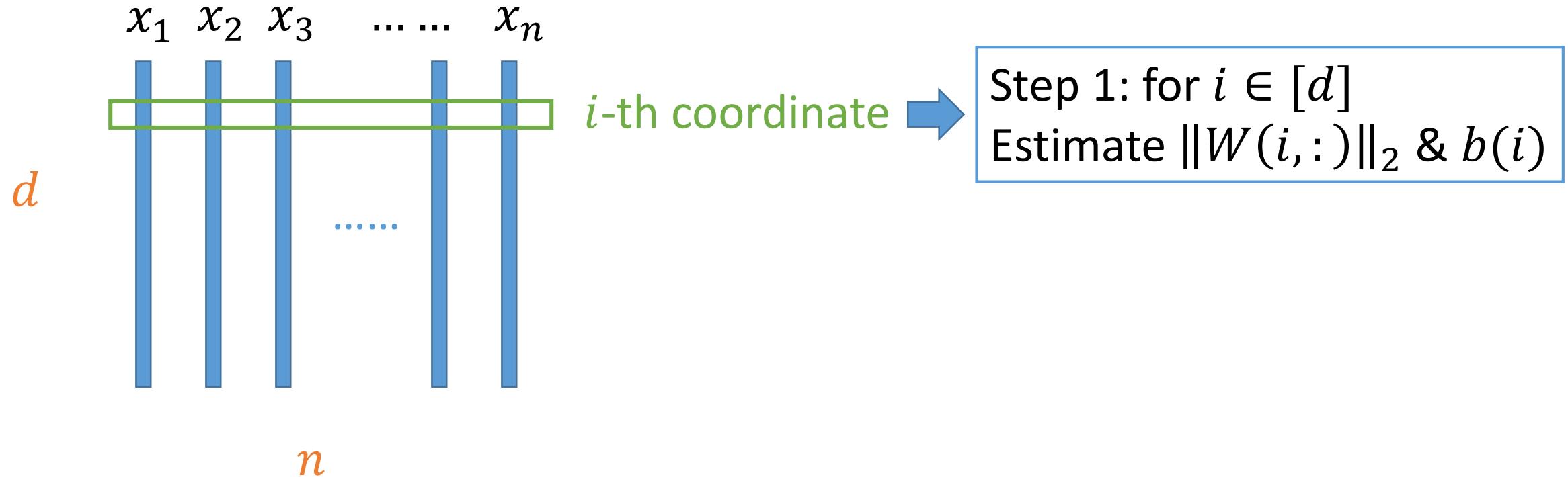
$$\langle W(i, :), W(j, :) \rangle = \|W(i, :)\|_2 \|W(j, :)\|_2 \cos \theta_{ij}$$

- Our algorithm:

- Step 1: estimate $\|W(i, :)\|_2$ & $b(i)$ for $i \in [d]$
- Step 2: estimate θ_{ij} for $i \neq j \in [d]$

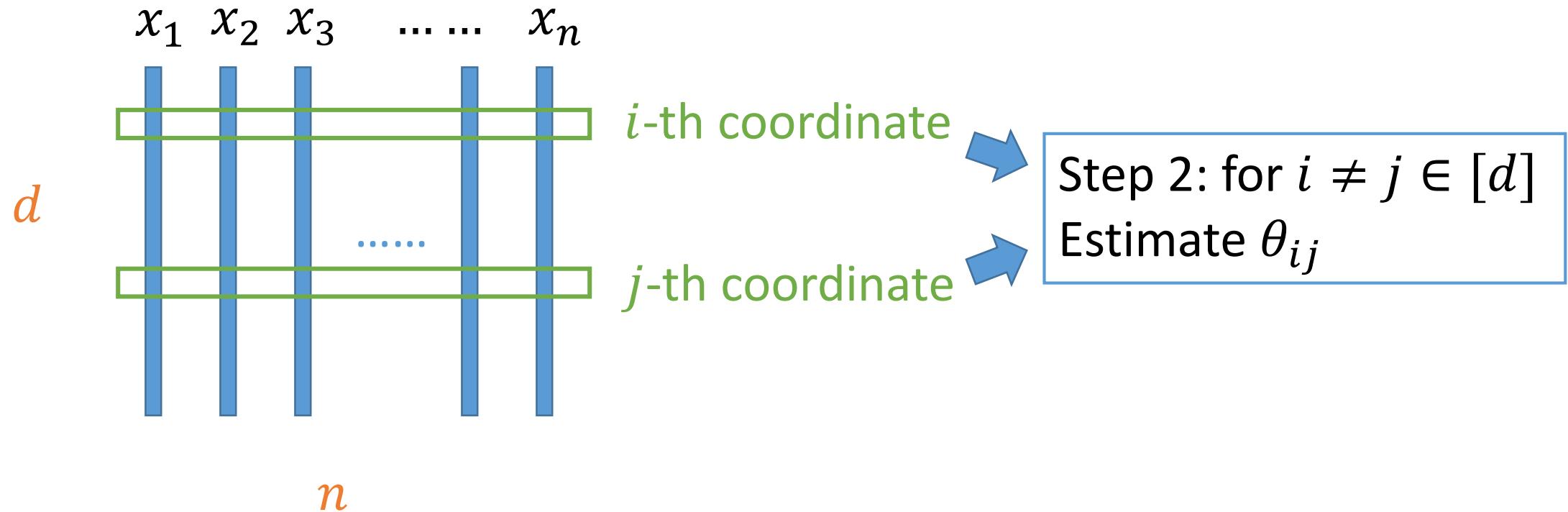
Our Algorithm: Overview

- Our problem: Given $x_1, \dots, x_n \sim \text{ReLU}(Wz + b)$, recover WW^T , b .

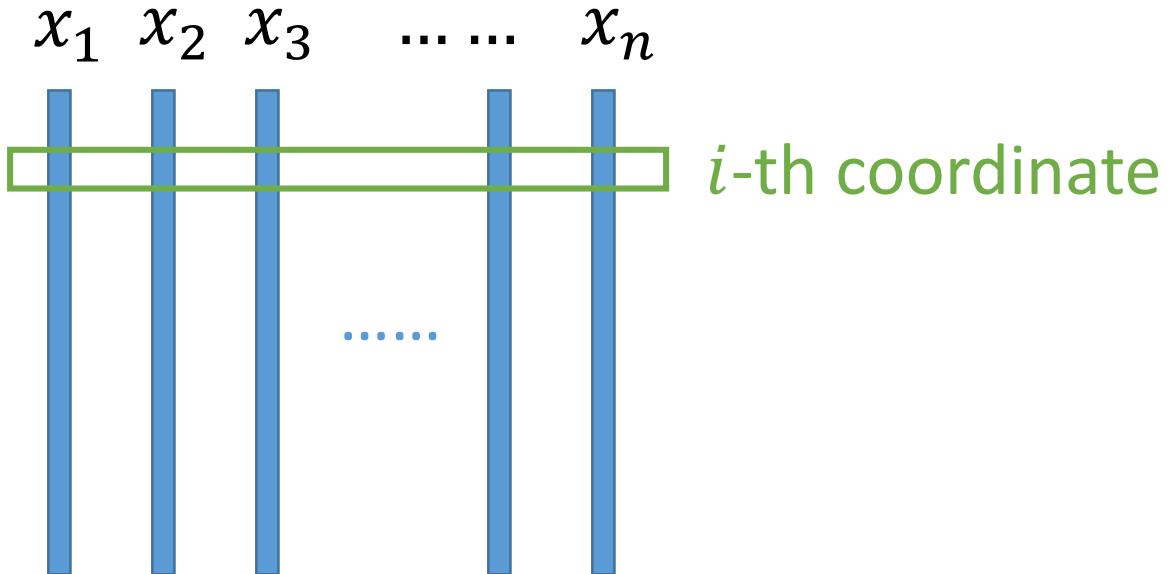


Our Algorithm: Overview

- Our problem: Given $x_1, \dots, x_n \sim \text{ReLU}(Wz + b)$, recover WW^T , b .



Step 1: estimate $\|W(i, :) \|_2$ & $b(i)$

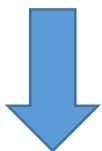


Goal:
Estimate $\|W(i, :) \|_2$ & $b(i)$

Step 1: estimate $\|W(i, :) \|_2$ & $b(i)$

$$x_1 \ x_2 \ x_3 \ \dots \dots \ x_n$$

 *i*-th coordinate

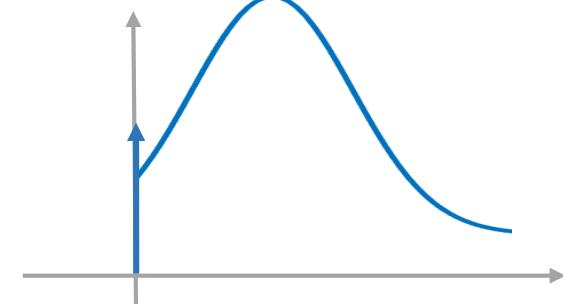


$$x(i) \sim \text{ReLU}(W(i, :) z + b(i))$$

$$\underbrace{\qquad\qquad\qquad}_{N(b(i), \|W(i, :) \|_2^2)}$$

$$N(b(i), \|W(i, :) \|_2^2)$$

Goal:
Estimate $\|W(i, :) \|_2$ & $b(i)$



Pdf of $x(i)$

Step 1: estimate $\|W(i, :) \|_2$ & $b(i)$

$x_1 \ x_2 \ x_3 \ \dots \dots \ x_n$



i-th coordinate



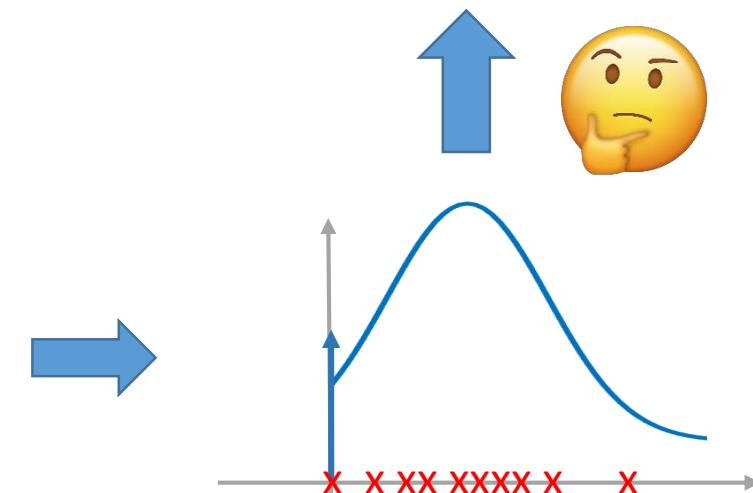
$$x(i) \sim \text{ReLU}(W(i, :)z + b(i))$$

$$\underbrace{\qquad\qquad\qquad}_{N(b(i), \|W(i, :) \|_2^2)}$$

$$N(b(i), \|W(i, :) \|_2^2)$$

Goal:

Estimate $\|W(i, :) \|_2$ & $b(i)$



Samples $\sim x(i)$

Step 1: estimate $\|W(i, :) \|_2$ & $b(i)$

$x_1 \ x_2 \ x_3 \ \dots \dots \ x_n$



i-th coordinate



$$x(i) \sim \text{ReLU}(W(i, :) z + b(i))$$

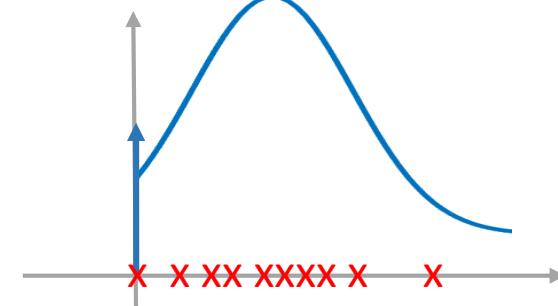
$$\underbrace{\quad}_{\quad}$$

$$N(b(i), \|W(i, :) \|_2^2)$$

Goal:

Estimate $\|W(i, :) \|_2$ & $b(i)$

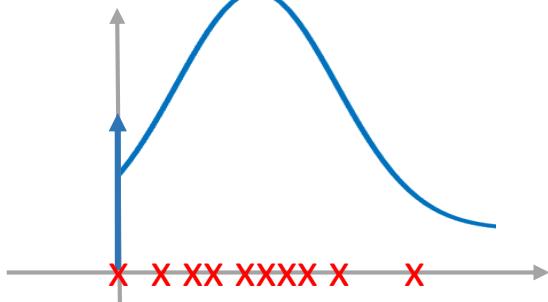
Truncated
Gaussian



Samples $\sim x(i)$

Step 1: estimate $\|W(i, :) \|_2$ & $b(i)$

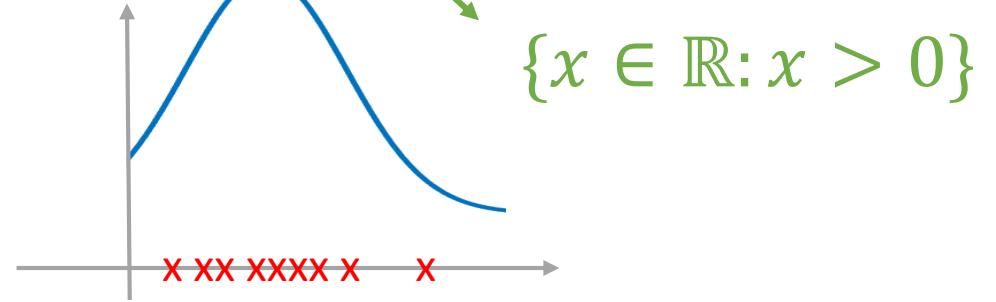
Samples $\sim x(i)$



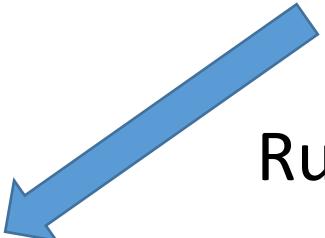
Remove zeros



Samples \sim truncated $N(b(i), \|W(i, :) \|_2^2)$

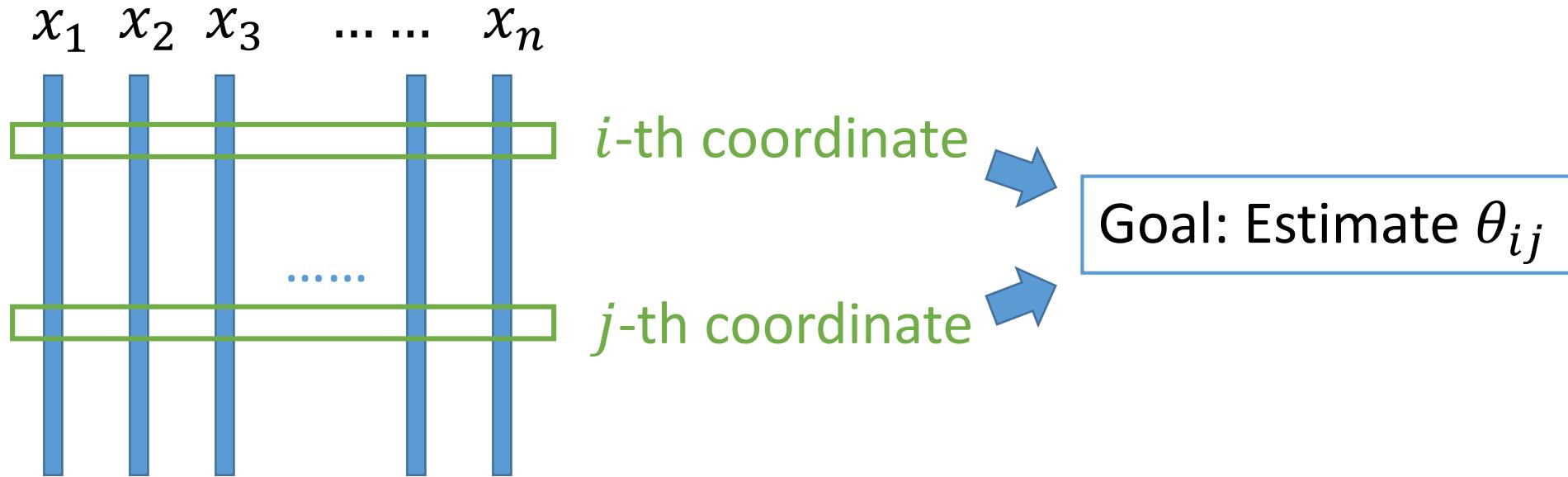


Estimate $b(i), \|W(i, :) \|_2$



Run algo [Daskalakis et al., 2018]

Step 2: estimate θ_{ij}



Step 2: estimate θ_{ij}

- $x(i) \sim \text{ReLU}(W(i,:)z + b(i))$
- $x(j) \sim \text{ReLU}(W(j,:)z + b(j))$



Goal: Estimate
 $\theta_{ij} := \text{angle b/t } W(i,:) \text{ & } W(j,:)$

Step 2: estimate θ_{ij}

- $x(i) \sim \text{ReLU}(W(i,:)z + b(i))$
- $x(j) \sim \text{ReLU}(W(j,:)z + b(j))$

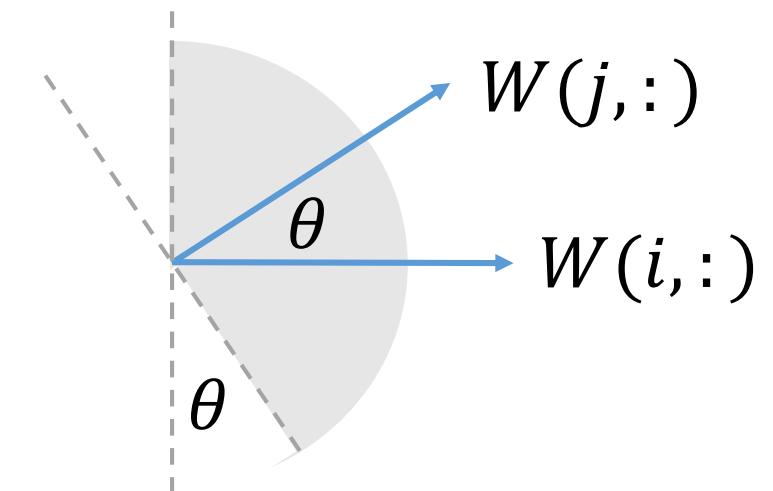


Goal: Estimate
 $\theta_{ij} := \text{angle b/t } W(i,:) \text{ & } W(j,:)$

- Fact [Williamson and Shmoys'11]:

$$\mathbb{P}_z[W(i,:)z > 0 \text{ & } W(j,:)z > 0]$$

$$= \frac{\pi - \theta_{ij}}{2\pi}$$



Step 2: estimate θ_{ij}

- $x(i) \sim \text{ReLU}(W(i,:)z + b(i))$
- $x(j) \sim \text{ReLU}(W(j,:)z + b(j))$



Goal: Estimate
 $\theta_{ij} := \text{angle b/t } W(i,:) \text{ & } W(j,:)$

- Suppose we know $b(i), b(j)$:

$$\mathbb{P}_x[x(i) > b(i) \& x(j) > b(j)] = \mathbb{P}_z [W(i,:)z > 0 \& W(j,:)z > 0]$$

Assumption:
 $b(i), b(j) \geq 0$

$$= \frac{\pi - \theta_{ij}}{2\pi}$$

Step 2: estimate θ_{ij}

- $x(i) \sim \text{ReLU}(W(i,:)z + b(i))$
- $x(j) \sim \text{ReLU}(W(j,:)z + b(j))$



Goal: Estimate
 $\theta_{ij} := \text{angle b/t } W(i,:) \text{ & } W(j,:)$

- Suppose we know $b(i), b(j)$:

$$\mathbb{P}_x[x(i) > b(i) \& x(j) > b(j)] = \mathbb{P}_z [W(i,:)z > 0 \& W(j,:)z > 0]$$

Assumption:

$b(i), b(j) \geq 0$

$$= \frac{\pi - \theta_{ij}}{2\pi}$$

If $b(i) < 0$, since $x(i) \geq 0$, $\mathbb{P}_x[x(i) > b(i)] = 1$

Step 2: estimate θ_{ij}

- $x(i) \sim \text{ReLU}(W(i,:)z + b(i))$
- $x(j) \sim \text{ReLU}(W(j,:)z + b(j))$



Goal: Estimate
 $\theta_{ij} := \text{angle b/t } W(i,:) \text{ & } W(j,:)$

- Given $\hat{b}(i)$ and $\hat{b}(j)$ from Step 1:

$$\hat{b}(i) \approx b(i), \hat{b}(j) \approx b(j)$$



$$\mathbb{P}_x[x(i) > \hat{b}(i) \& x(j) > \hat{b}(j)] \approx \frac{\pi - \theta_{ij}}{2\pi}$$

Step 2: estimate θ_{ij}

- $x(i) \sim \text{ReLU}(W(i,:)z + b(i))$
- $x(j) \sim \text{ReLU}(W(j,:)z + b(j))$



Goal: Estimate
 $\theta_{ij} := \text{angle b/t } W(i,:) \text{ & } W(j,:)$

- Given $\hat{b}(i)$ and $\hat{b}(j)$ from Step 1:

$$\hat{b}(i) \approx b(i), \hat{b}(j) \approx b(j)$$



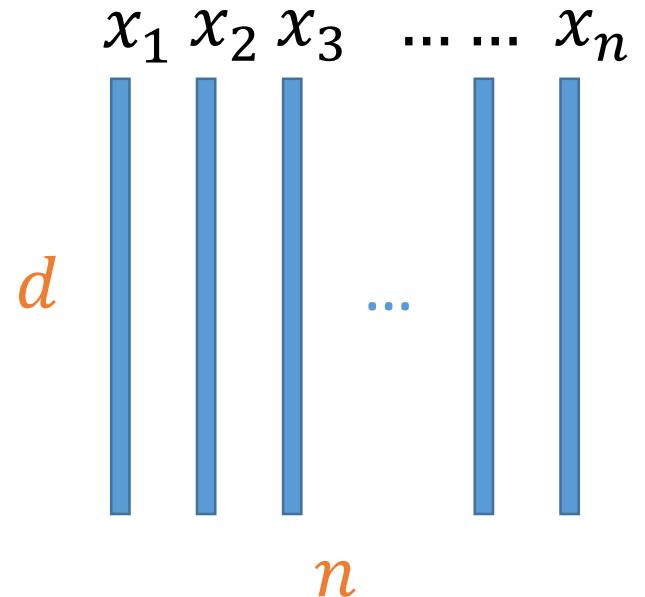
$$\mathbb{P}_x[x(i) > \hat{b}(i) \& x(j) > \hat{b}(j)] \approx \frac{\pi - \theta_{ij}}{2\pi}$$



$$\hat{\theta}_{ij} = \pi - 2\pi \frac{\sum_{k=1}^n \mathbf{1}[x_k(i) > \hat{b}(i) \& x_k(j) > \hat{b}(j)]}{n}$$

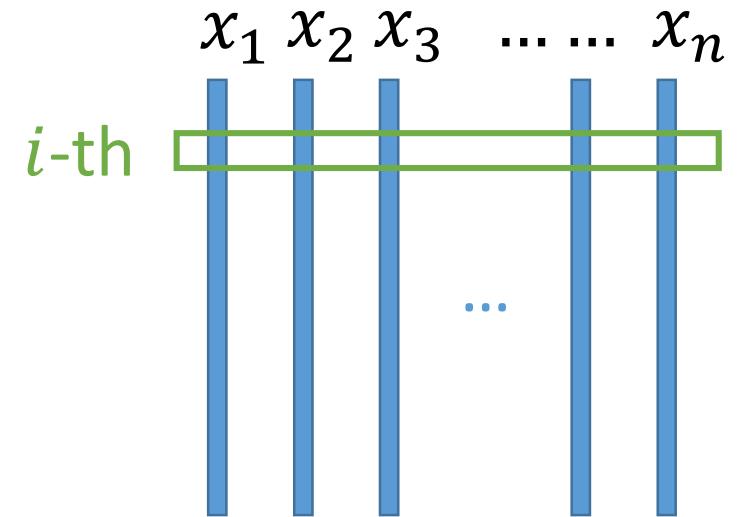
Our Algorithm: Overall

- **Input:** $x_1, \dots, x_n \sim \text{ReLU}(Wz + b)$, b non-negative
- **Output:** $\widehat{WW^T} \in \mathbb{R}^{d \times d}$, $\widehat{b} \in \mathbb{R}^d$
- **For** $i = 1, 2, \dots, d$
 - Remove zero samples from i -th coordinates
 - $\widehat{b}(i), \|\widehat{W(i,:)}\|_2^2 \leftarrow \text{Run [Daskalakis et al.'2018]}$
- **For** $i \neq j \in [d]$
 - $\widehat{\theta}_{ij} \leftarrow \pi - 2\pi \frac{\sum_{k=1}^n \mathbf{1}[x_k(i) > \widehat{b}(i) \& x_k(j) > \widehat{b}(j)]}{n}$
 - $\widehat{WW^T}(i, j) = \|\widehat{W(i,:)}\|_2 \|\widehat{W(j,:)}\|_2 \cos \widehat{\theta}_{ij}$



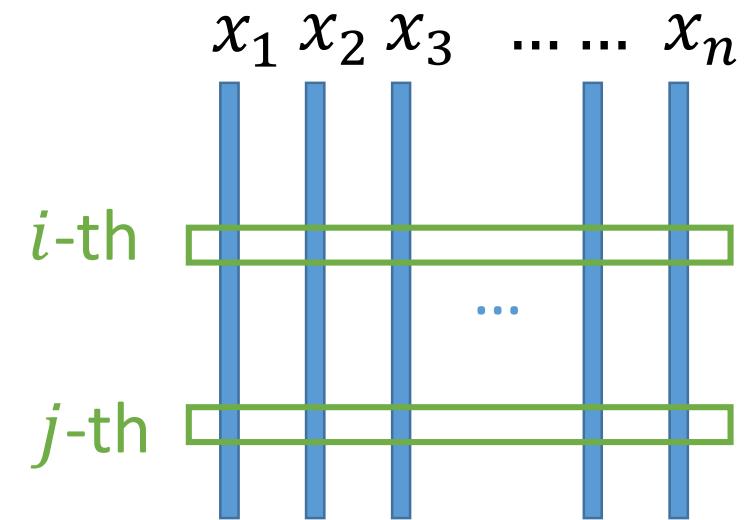
Our Algorithm: Overall

- **Input:** $x_1, \dots, x_n \sim \text{ReLU}(Wz + b)$, b non-negative
- **Output:** $\widehat{WW^T} \in \mathbb{R}^{d \times d}$, $\widehat{b} \in \mathbb{R}^d$
- **For** $i = 1, 2, \dots, d$
 - Remove zero samples from i -th coordinates
 - $\widehat{b}(i), \|\widehat{W(i,:)}\|_2^2 \leftarrow \text{Run [Daskalakis et al.'2018]}$
- **For** $i \neq j \in [d]$
 - $\widehat{\theta}_{ij} \leftarrow \pi - 2\pi \frac{\sum_{k=1}^n \mathbf{1}[x_k(i) > \widehat{b}(i) \& x_k(j) > \widehat{b}(j)]}{n}$
 - $\widehat{WW^T}(i, j) = \|\widehat{W(i,:)}\|_2 \|\widehat{W(j,:)}\|_2 \cos \widehat{\theta}_{ij}$



Our Algorithm: Overall

- **Input:** $x_1, \dots, x_n \sim \text{ReLU}(Wz + b)$, b non-negative
- **Output:** $\widehat{WW^T} \in \mathbb{R}^{d \times d}$, $\widehat{b} \in \mathbb{R}^d$
- **For** $i = 1, 2, \dots, d$
 - Remove zero samples from i -th coordinates
 - $\widehat{b}(i), \|\widehat{W(i,:)}\|_2^2 \leftarrow \text{Run [Daskalakis et al.'2018]}$
- **For** $i \neq j \in [d]$
 - $\widehat{\theta}_{ij} \leftarrow \pi - 2\pi \frac{\sum_{k=1}^n \mathbf{1}[x_k(i) > \widehat{b}(i) \& x_k(j) > \widehat{b}(j)]}{n}$
 - $\widehat{WW^T}(i,j) = \|\widehat{W(i,:)}\|_2 \|\widehat{W(j,:)}\|_2 \cos \widehat{\theta}_{ij}$



Analysis

- **Theorem:** Assuming that $b^* \in \mathbb{R}^d$ is non-negative, then our algorithm takes $\tilde{\mathcal{O}}\left(\frac{1}{\epsilon^2} \ln\left(\frac{d}{\delta}\right)\right)$ samples $\sim \text{ReLU}(W^*z + b^*)$ and outputs

$$\left\| \widehat{WW^T} - W^*W^{*T} \right\|_F \leq \epsilon \|W^*\|_F^2, \quad \left\| \widehat{b} - b^* \right\|_2 \leq \epsilon \|W^*\|_F$$

with probability at least $1 - \delta$.

Analysis: non-degenerate case

- **Corollary:** Let $\kappa := \text{condition number of } W^*W^{*T}$. Our algorithm takes $\tilde{O}\left(\frac{\kappa^2 d^2}{\epsilon^2} \ln\left(\frac{d}{\delta}\right)\right)$ samples $\sim \text{ReLU}(W^*z + b^*)$ and outputs

$$\mathbf{TV}(\text{ReLU}(\hat{W}z + \hat{b}), \text{ReLU}(W^*z + b^*)) \leq \epsilon$$

with probability at least $1 - \delta$.

Sample Complexity Lower Bounds

- Parameter Estimation: $\|\widehat{WW^T} - W^*W^{*T}\|_F \leq \epsilon\|W^*\|_F^2, \quad \|\widehat{b} - b^*\|_2 \leq \epsilon\|W^*\|_F$

Our algorithm	Lower bound
$\tilde{O}\left(\frac{1}{\epsilon^2} \ln\left(\frac{d}{\delta}\right)\right)$	$\Omega\left(\frac{1}{\epsilon^2}\right)$

Our algorithm is optimal

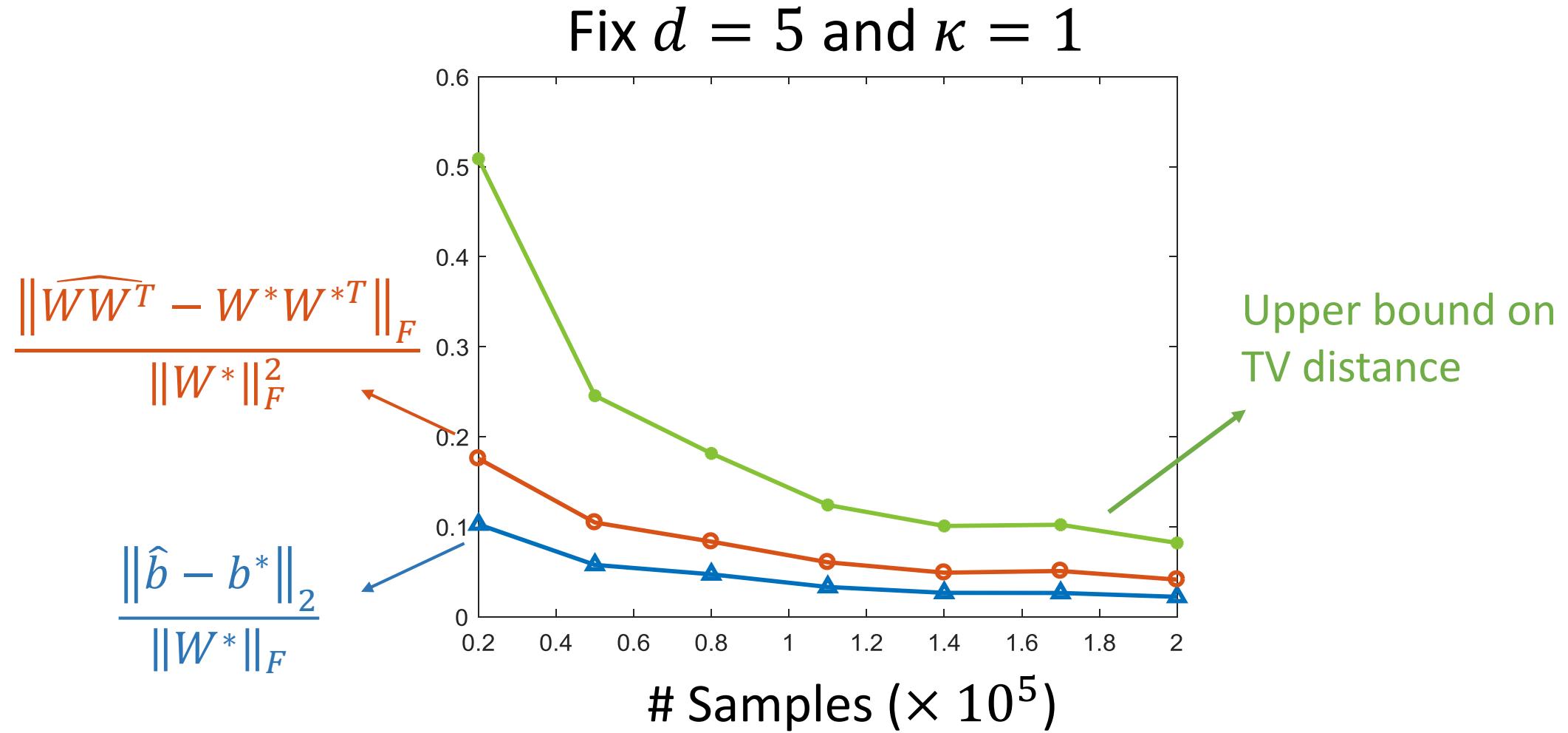
- Total Variation Distance: $\text{TV}(D, \text{ReLU}(W^*z + b^*)) \leq \epsilon$

Our algorithm	Lower bound
$\tilde{O}\left(\frac{\kappa^2 d^2}{\epsilon^2} \ln\left(\frac{d}{\delta}\right)\right)$	$\Omega\left(\frac{d}{\epsilon^2}\right)$

This gap comes from:

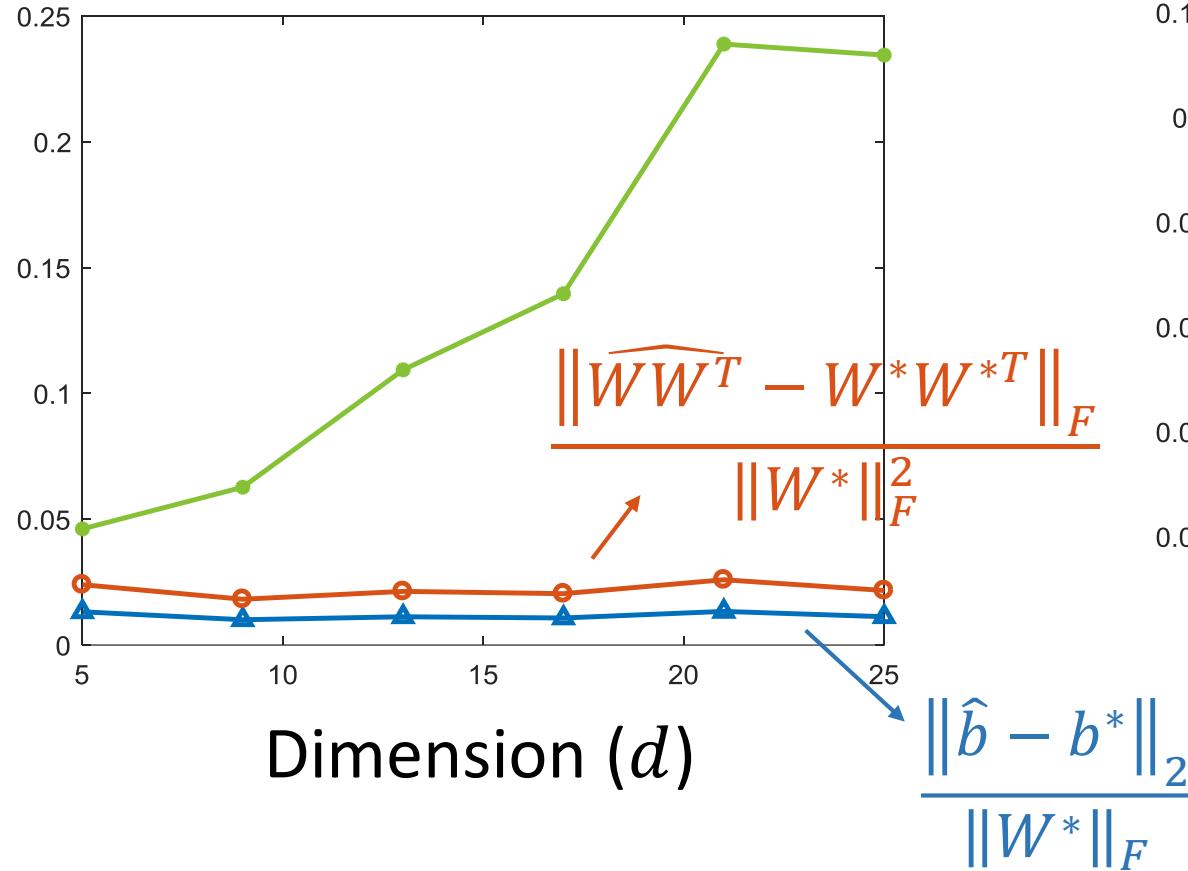
- 1) Param est \rightarrow TV distance
- 2) Lower bound is not tight

Experimental Results

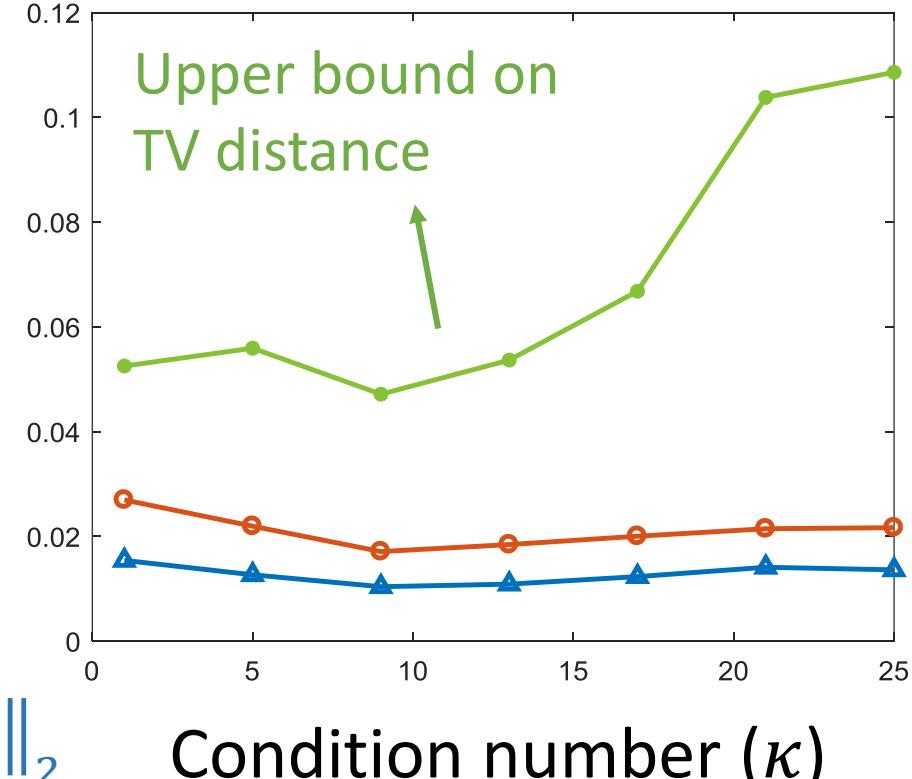


Experimental Results

Fix # samples = 2e5



Fix # samples = 2e5



What if b has negative values?

$$\exists \eta \geq 0, b(i) \geq -\eta \|W(i, :) \|_2, \forall i \in [d]$$



Given $\tilde{O}(\frac{1}{\epsilon^2} \ln(\frac{d}{\delta}))$ samples, our algorithm outputs:

$$\|\widehat{WW^T} - W^*{W^*}^T\|_F \leq \max(\epsilon, \eta) \|W^*\|_F^2$$

$$\|\hat{b} - b^*\|_2 \leq \max(\epsilon, \eta) \|W^*\|_F$$

Related Work

- Learning ReLU neural networks [GKLW'19, GKM'18, LY'17,.....]

Supervised v.s. Unsupervised

- Provably learning a generative model

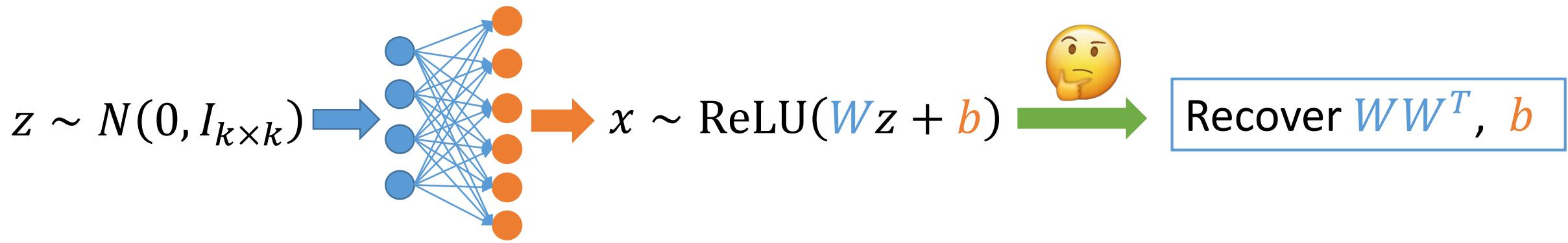
- [Nguyen, Wong, and Hedge'18]

Linear generative model

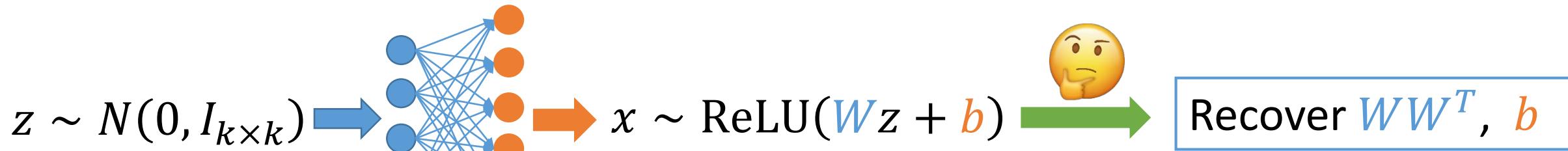
- [Mazumdar and Rawat'19]

Random b , recover W 's column space

Summary



Summary



Our algorithm:

Step 1: est $\|W(i, :) \|_2, b(i)$

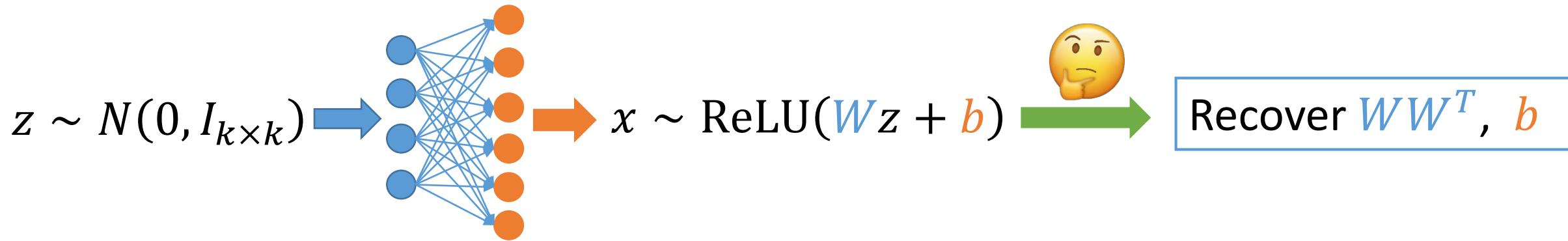
Step 2: est θ_{ij}

	Our algorithm	Lower bound
Param Est	$\tilde{O}\left(\frac{1}{\epsilon^2} \ln\left(\frac{d}{\delta}\right)\right)$	$\Omega\left(\frac{1}{\epsilon^2}\right)$
TV dist	$\tilde{O}\left(\frac{\kappa^2 d^2}{\epsilon^2} \ln\left(\frac{d}{\delta}\right)\right)$	$\Omega\left(\frac{d}{\epsilon^2}\right)$

Assumption:
 b non-negative
Otherwise:
 $\Omega(\exp(\|b\|_\infty^2))$

Experimental results are consistent with analysis.

Summary



- Open problems:
 - Two-layer Generative model
 - Negative bias vector
 - Agnostic setting
 - GAN/VAE.....

When we use neural network to model distributions, what is the structure of those distributions?