

# Proximal Mean Field Learning in Shallow Neural Networks

Alexis M.H. Teter

Department of Applied Mathematics  
University of California, Santa Cruz

**Joint work with**



Abhishek Halder  
Iowa State University

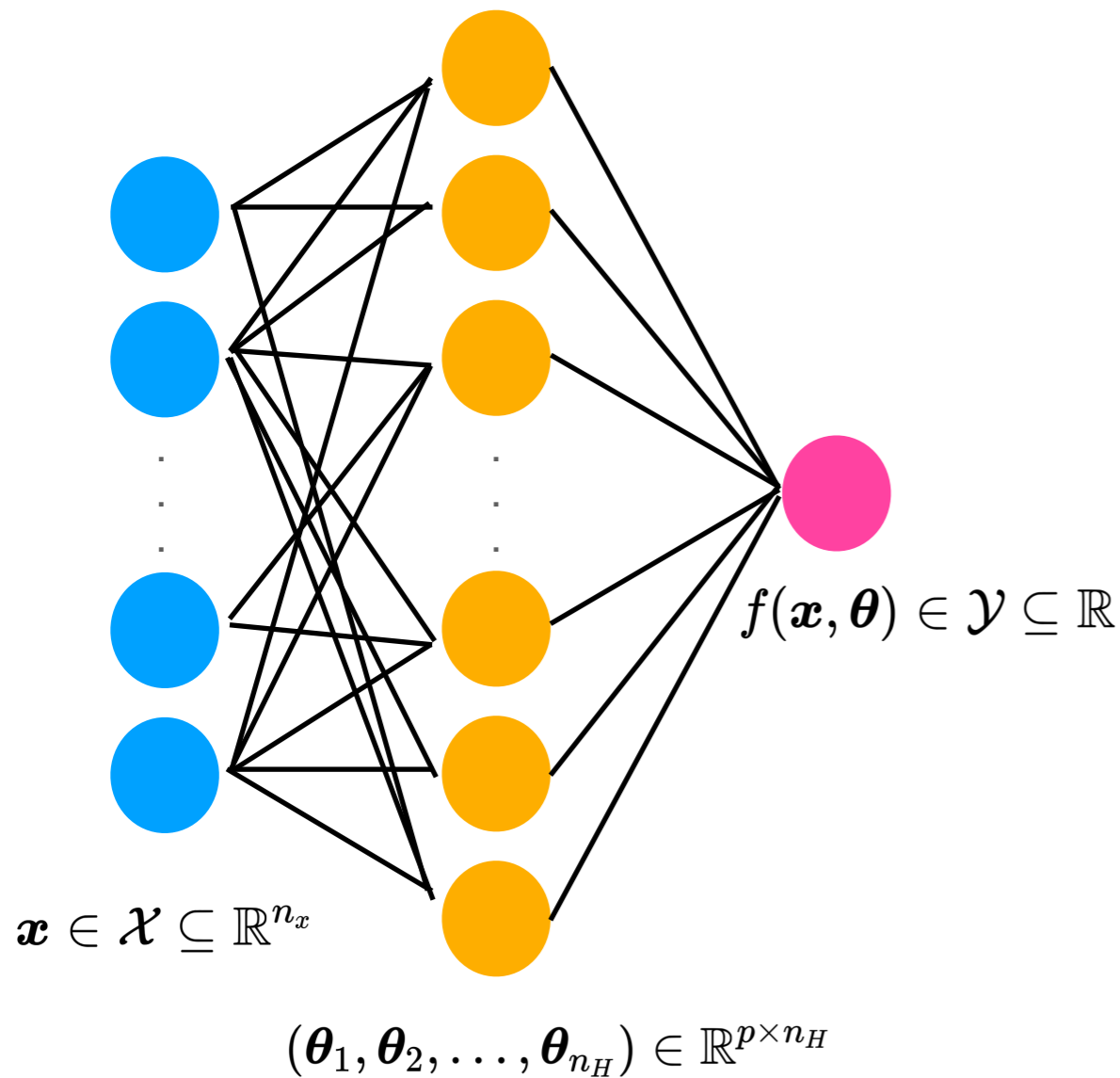


Iman Nodozi  
University of California, Santa Cruz

**SIAM Minisymposium: Recent advances in optimization for DNNs**

**2024 SIAM Annual Meeting, Online, July 18, 2024**

# Structure of shallow neural network



$$f(\mathbf{x}, \boldsymbol{\theta}) := \frac{1}{n_H} \sum_{i=1}^{n_H} \Phi(\mathbf{x}, \boldsymbol{\theta}_i)$$

$$\Phi(\mathbf{x}, \boldsymbol{\theta}_i) := a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle + b_i)$$

# Risk

## Population Risk

$$R(f) := \mathbb{E}_{(\mathbf{x}, y) \sim \gamma} [(y - f(\mathbf{x}, \bar{\boldsymbol{\theta}}))^2]$$

## Empirical Risk

$$R(f) \approx \frac{1}{n_{\text{data}}} \sum_{j=1}^{n_{\text{data}}} (y_j - f(\mathbf{x}_j, \bar{\boldsymbol{\theta}}))^2$$

Large dimensional, non convex optimization problem

$$\min_{\bar{\boldsymbol{\theta}} \in \mathbb{R}^{p \times n_H}} R(f)$$

# Mean field limit

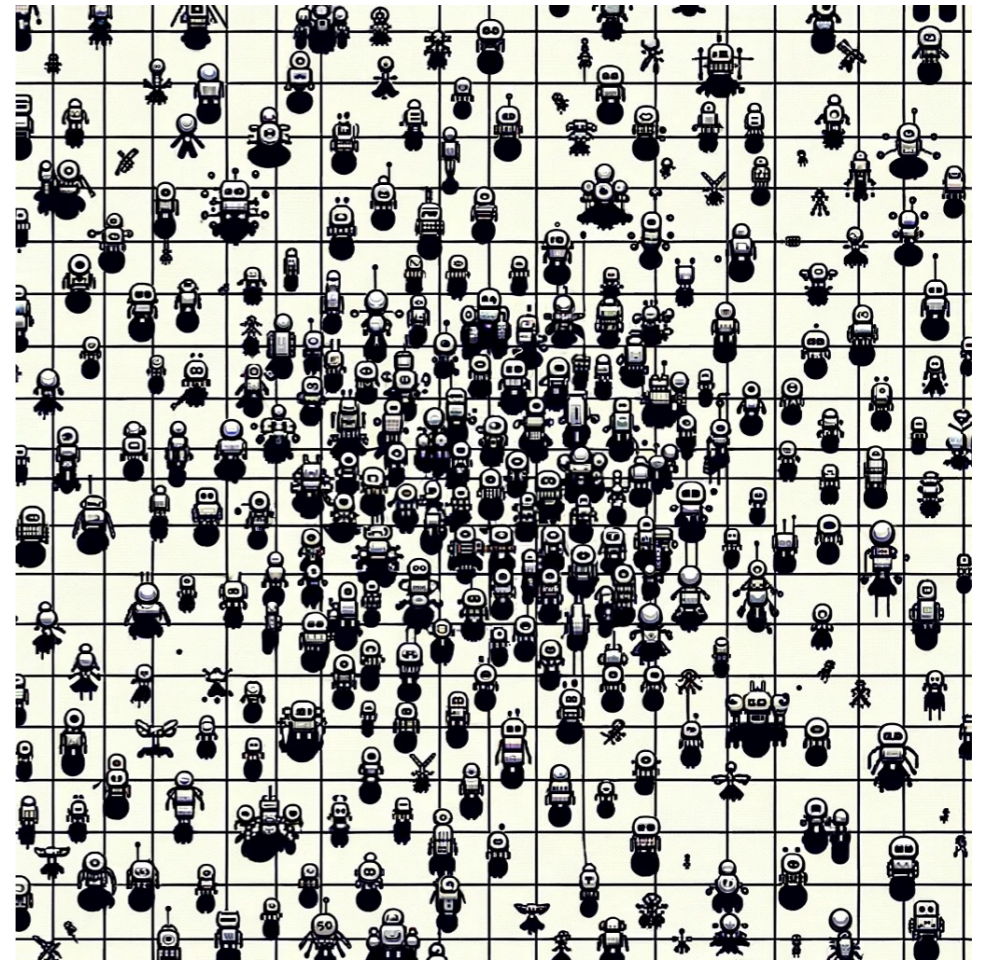
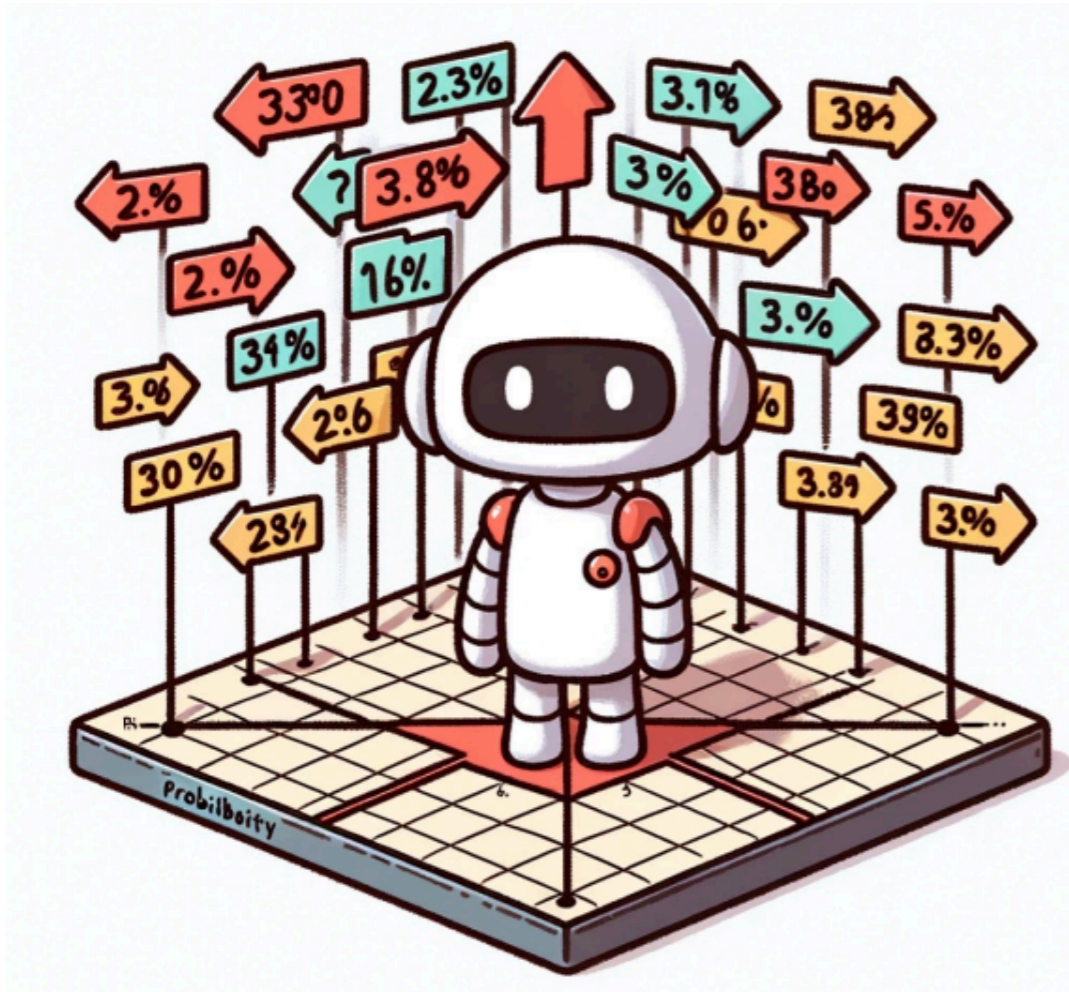
$$f(\mathbf{x}, \boldsymbol{\theta}) := \frac{1}{n_H} \sum_{i=1}^{n_H} \Phi(\mathbf{x}, \boldsymbol{\theta}_i)$$

Let  $n_H \rightarrow \infty$

$$f_{\text{MeanField}} := \int_{\mathbb{R}^p} \Phi(\mathbf{x}, \boldsymbol{\theta}) d\mu(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}}[\Phi(\mathbf{x}, \boldsymbol{\theta})]$$

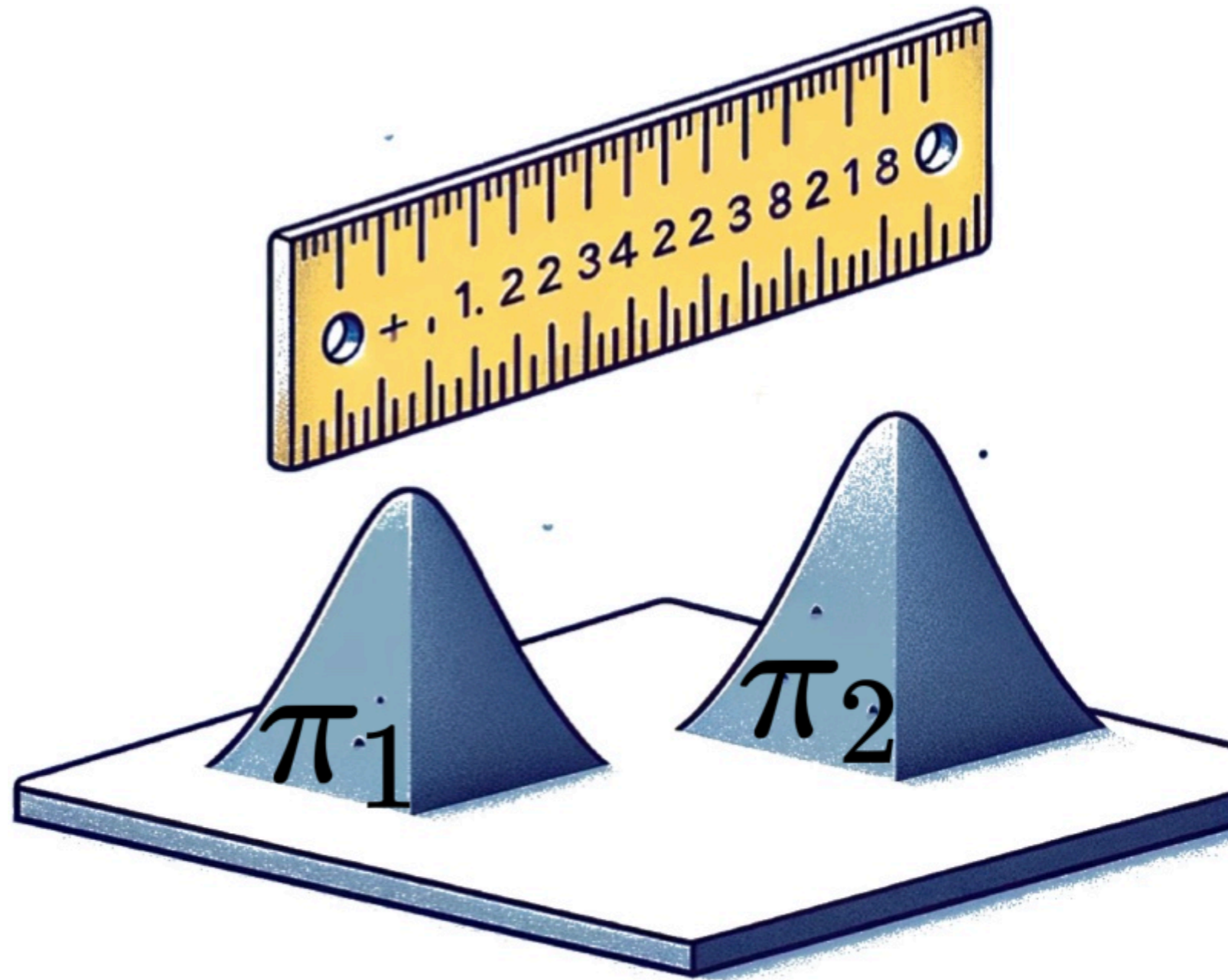
where

$$d\mu(\boldsymbol{\theta}) = \rho(\boldsymbol{\theta}) d\boldsymbol{\theta}$$



$$d\mu = \rho dx$$

# Wasserstein metric

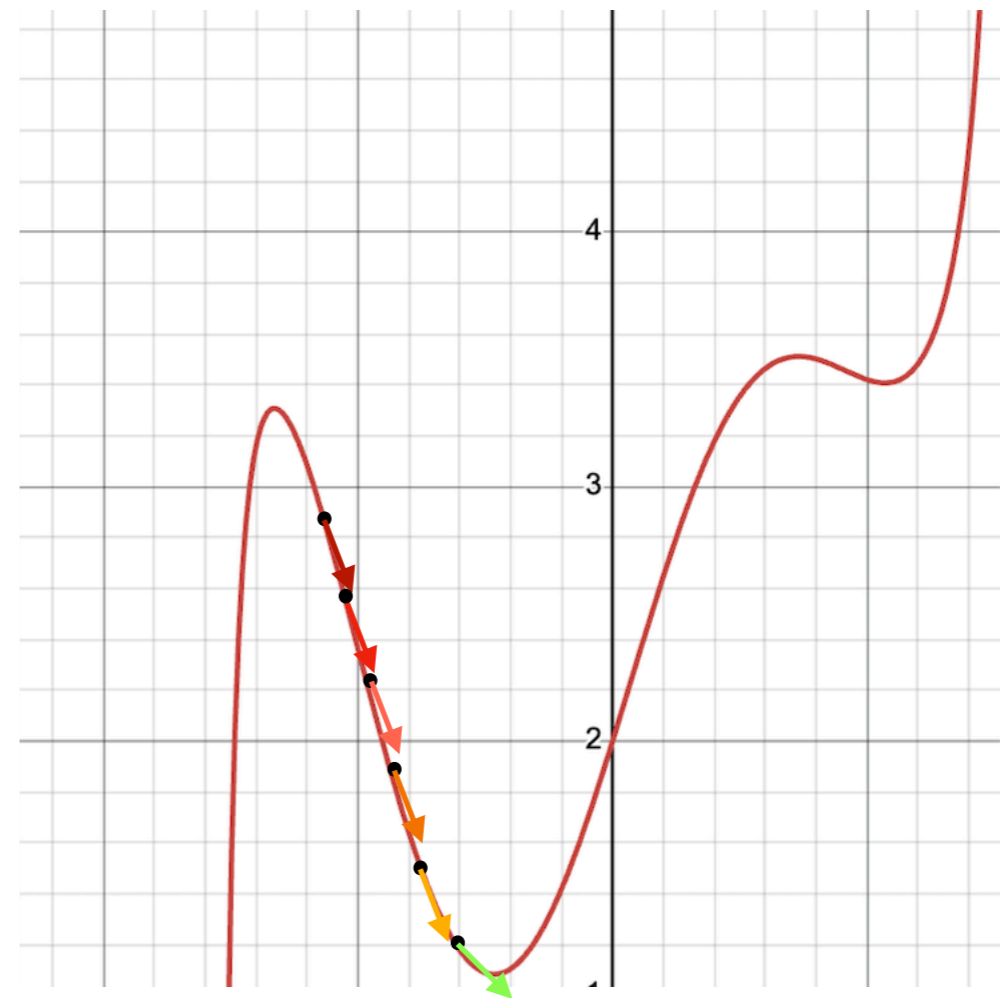


$$W_2^2(\pi_1, \pi_2) := \inf_{\pi \in \Pi(\pi_1, \pi_2)} \int_{\mathcal{Z}_1 \times \mathcal{Z}_2} \|z_1 - z_2\|_2^2 d\pi(z_1, z_2)$$

# Gradient Flows

## Gradient Flow

$$\frac{d\mathbf{x}}{dt} = -\nabla f(\mathbf{x}), \quad \mathbf{x}(0) = \mathbf{x}_0$$



## Recursion

$$\mathbf{x}_k = \mathbf{x}_{k-1} - h\nabla f(\mathbf{x}_k) = \arg \min_{\mathbf{x} \in \mathcal{R}^n} \left( \frac{1}{2} \|\mathbf{x} - \mathbf{x}_{k-1}\|_2^2 + hf(\mathbf{x}) \right) = \text{prox}_{hF}^{|\cdot|_2}(\mathbf{x}_{k-1})$$

proximal update

$$= \arg \inf_{\text{decision variable}} \left\{ \frac{1}{2} \text{dist}^2(\text{decision variable, input}) + \text{time step} \times \text{functional}(\text{decision variable}) \right\}$$

# Gradient Flows

## Gradient Flow

$$\frac{\partial \rho}{\partial t} = -\nabla^{W_2} F(\rho) := -\nabla \cdot \left( \rho \nabla \frac{\delta F}{\delta \rho} \right) \quad \text{where} \quad \rho(\boldsymbol{\theta}, 0) = \rho_0(\boldsymbol{\theta})$$

## Recursion

$$\rho_k = \rho(\cdot, t = kh) = \arg \min_{\rho \in \mathcal{P}_2(\mathbb{R}^p)} \left( \frac{1}{2} W_2^2(\rho, \rho_{k-1}) + hF(\rho) \right) = \text{prox}_{hF}^{W_2}(\rho_{k-1})$$

proximal update

$$= \arg \inf_{\text{decision variable}} \left\{ \frac{1}{2} \text{dist}^2(\text{decision variable, input}) + \text{time step} \times \text{functional}(\text{decision variable}) \right\}$$



# Risk functional

$$R(f_{\text{Mean Field}}(\mathbf{x}, \rho)) = \mathbb{E}_{(\mathbf{x}, y)} \left( y - \int_{\mathbb{R}^p} \Phi(\mathbf{x}, \theta) \rho(\theta) d\theta \right)^2$$

$$R(f_{\text{Mean Field}}(\mathbf{x}, \rho)) = F_0 + \int_{\mathbb{R}^p} V(\theta) \rho(\theta) d\theta + \int_{\mathbb{R}^{2p}} U(\theta, \tilde{\theta}) \rho(\theta) \rho(\tilde{\theta}) d\theta d\tilde{\theta}$$

**Drift potential**

**Interaction potential**

$$F_0 := \mathbb{E}_{(\mathbf{x}, y)} [y^2]$$

$$V(\theta) := \mathbb{E}_{(\mathbf{x}, y)} [-2y \Phi(\mathbf{x}, \theta)]$$

$$U(\theta, \tilde{\theta}) := \mathbb{E}_{(\mathbf{x}, y)} [\Phi(\mathbf{x}, \theta) \Phi(\mathbf{x}, \tilde{\theta})]$$

# Supervised learning in mean field limit

Supervised learning problem

$$\min_{\rho} F(\rho) := \min_{\rho} R(f_{\text{Mean Field}}(\mathbf{x}, \rho))$$

$$R(f_{\text{Mean Field}}(\mathbf{x}, \rho)) = F_0 + \int_{\mathbb{R}^p} V(\boldsymbol{\theta}) \rho(\boldsymbol{\theta}) d\boldsymbol{\theta} + \int_{\mathbb{R}^{2p}} U(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) \rho(\boldsymbol{\theta}) \rho(\tilde{\boldsymbol{\theta}}) d\boldsymbol{\theta} d\tilde{\boldsymbol{\theta}}$$

independent  
of  $\rho$

potential  
energy;  
linear in  $\rho$

bilinear  
interaction  
energy;  
nonlinear in  $\rho$

Regularized risk functional

$$F_{\beta}(\rho) := F(\rho) + \beta^{-1} \int_{\mathbb{R}^p} \rho \log \rho d\boldsymbol{\theta}, \quad \beta > 0$$

$\rho^{\text{opt}}(t, \boldsymbol{\theta})$

Nonlinear PDE IVP

$$\frac{\partial \rho}{\partial t} = \nabla_{\boldsymbol{\theta}} \cdot \left( \rho \left( V(\boldsymbol{\theta}) + \int_{\mathbb{R}^p} U(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) \rho(\tilde{\boldsymbol{\theta}}) d\tilde{\boldsymbol{\theta}} \right) \right) + \beta^{-1} \Delta_{\boldsymbol{\theta}} \rho$$

where  $\rho(t = 0, \boldsymbol{\theta}) = \rho_0$

# Proximal recursions

## Proximal recursion

$$\varrho_k = \text{prox}_{hF_\beta}^{W_2}(\varrho_{k-1}) := \underset{\varrho \in \mathcal{P}_2(\mathbb{R}^p)}{\text{arginf}} \left\{ \frac{1}{2} (W_2(\varrho, \varrho_{k-1}))^2 + h F_\beta(\varrho) \right\}$$

where  $\varrho_{k-1}(\cdot) := \varrho(\cdot, t_{k-1})$

$$\varrho_0 \equiv \rho_0$$

Approximate bilinear term as...

$$\int_{\mathbb{R}^{2p}} U(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) \varrho(\boldsymbol{\theta}) \varrho(\tilde{\boldsymbol{\theta}}) d\boldsymbol{\theta} d\tilde{\boldsymbol{\theta}} \approx \int_{\mathbb{R}^{2p}} U(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) \varrho(\boldsymbol{\theta}) \varrho_{k-1}(\tilde{\boldsymbol{\theta}}) d\boldsymbol{\theta} d\tilde{\boldsymbol{\theta}}$$

# Proximal recursions

## Proximal recursion

$$\varrho_k = \text{prox}_{hF_\beta}^{W_2}(\varrho_{k-1}) := \operatorname{arginf}_{\varrho \in \mathcal{P}_2(\mathbb{R}^p)} \left\{ \frac{1}{2} (W_2(\varrho, \varrho_{k-1}))^2 + h \hat{F}_\beta(\varrho) \right\}$$

where  $\varrho_{k-1}(\cdot) := \varrho(\cdot, t_{k-1})$

$$\varrho_0 \equiv \rho_0$$

## Approximation of regularized risk functional

$$\hat{F}_\beta(\varrho, \varrho_{k-1}) := \int_{\mathbb{R}^p} \left( F_0 + V(\boldsymbol{\theta}) + \left( \int_{\mathbb{R}^p} U(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) \varrho_{k-1}(\tilde{\boldsymbol{\theta}}) d\tilde{\boldsymbol{\theta}} \right) + \beta^{-1} \log \varrho(\boldsymbol{\theta}) \right) \varrho(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

# Proximal recursions

## Proximal recursion (semi-implicit variant)

$$\varrho_k = \text{prox}_{hF_\beta}^{W_2}(\varrho_{k-1}) := \operatorname{arginf}_{\varrho \in \mathcal{P}_2(\mathbb{R}^p)} \left\{ \frac{1}{2} (W_2(\varrho, \varrho_{k-1}))^2 + h \hat{F}_\beta(\varrho) \right\}$$

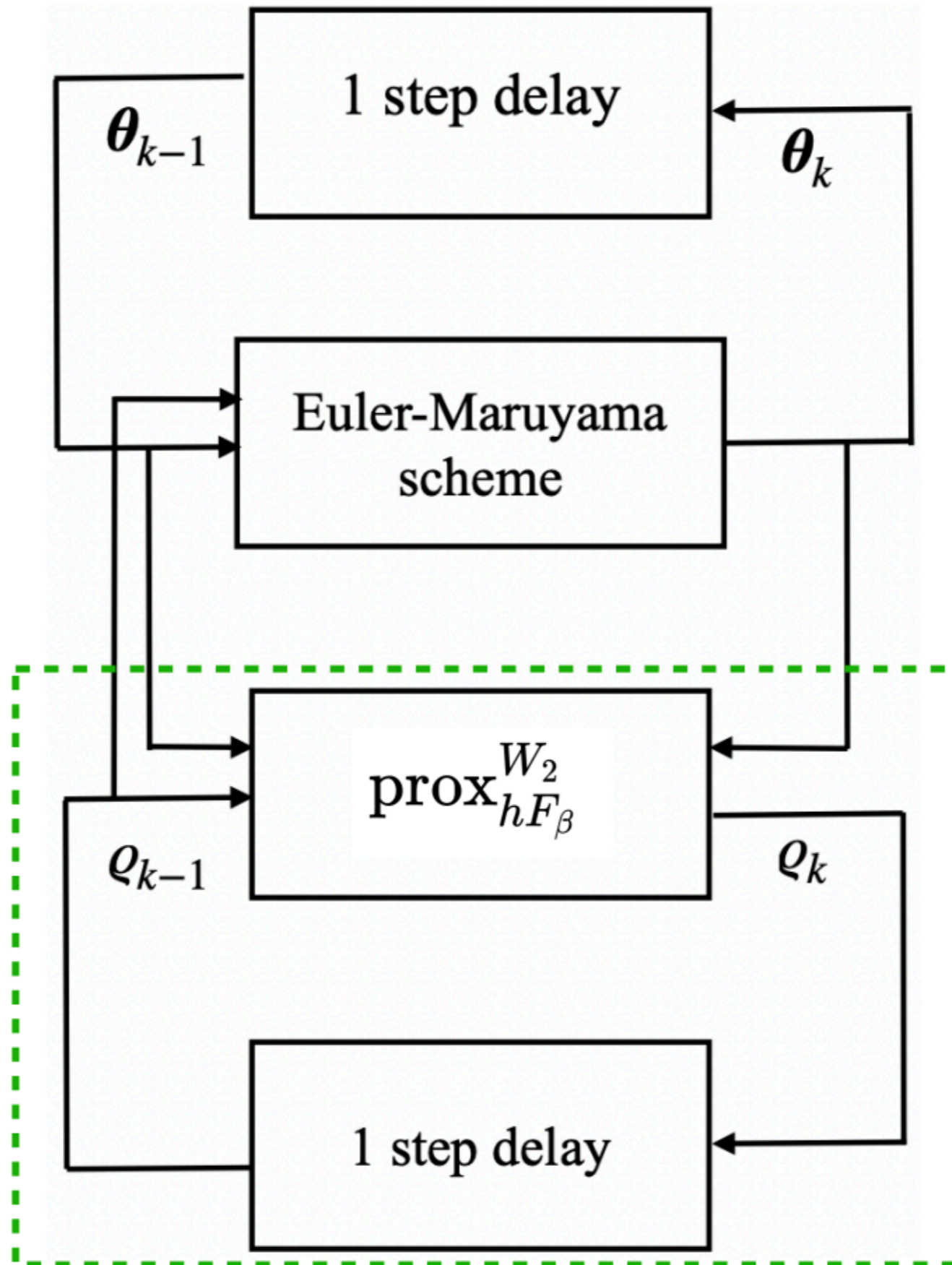
where  $\varrho_{k-1}(\cdot) := \varrho(\cdot, t_{k-1})$

$$\varrho_0 \equiv \rho_0$$

### Thm. 1:

As  $h \rightarrow 0$ , proximal updates converge to solution to PDE IVP.

# ProxLearn Algorithm



# Euler-Maruyama

$$\boldsymbol{\theta}_k^i = \boldsymbol{\theta}_{k-1}^i - h \nabla (V(\boldsymbol{\theta}_{k-1}^i) + \omega(\boldsymbol{\theta}_{k-1}^i)) + \sqrt{2\beta^{-1}} (\boldsymbol{\eta}_k^i - \boldsymbol{\eta}_{k-1}^i)$$

where  $\omega(\cdot) := \int U(\cdot, \tilde{\boldsymbol{\theta}}) \varrho(\tilde{\boldsymbol{\theta}}) d\tilde{\boldsymbol{\theta}}$

and  $\boldsymbol{\eta}_{k-1}^i := \boldsymbol{\eta}^i(t = (k-1)h)$





# Derivation of ProxLearn

Proximal recursion (semi-implicit variant)

$$\varrho_k = \text{prox}_{h\hat{F}_\beta}^{W_2}(\varrho_{k-1}) := \underset{\varrho \in \mathcal{P}_2(\mathbb{R}^p)}{\text{arginf}} \left\{ \frac{1}{2} (W_2(\varrho, \varrho_{k-1}))^2 + h \hat{F}_\beta(\varrho) \right\}$$

# Derivation of ProxLearn

Discrete version of proximal recursion

$$\mathbf{q}_k = \arg \min_{\mathbf{q}} \left\{ \min_{\mathbf{M} \in \Pi(\mathbf{q}_{k-1}, \mathbf{q})} \frac{1}{2} \langle \mathbf{C}_k, \mathbf{M} \rangle + h \langle \mathbf{v}_{k-1} + \mathbf{U}_{k-1} \mathbf{q}_{k-1} + \beta^{-1} \log \mathbf{q}, \mathbf{q} \rangle \right\}$$

where  $\Pi(\mathbf{q}_{k-1}, \mathbf{q}) := \{\mathbf{M} \in \mathbb{R}^{N \times N} \mid \mathbf{M} \geq \mathbf{0} \text{ (elementwise)}, \mathbf{M}\mathbf{1} = \mathbf{q}_{k-1}, \mathbf{M}^\top \mathbf{1} = \mathbf{q}\}$

# Derivation of ProxLearn

Regularized discrete version of proximal recursion

$$\mathbf{q}_k = \arg \min_{\mathbf{q}} \left\{ \min_{\mathbf{M} \in \Pi(\mathbf{q}_{k-1}, \mathbf{q})} \frac{1}{2} \langle \mathbf{C}_k, \mathbf{M} \rangle + \epsilon \langle \mathbf{M}, \log \mathbf{M} \rangle + h \langle \mathbf{v}_{k-1} + \mathbf{U}_{k-1} \mathbf{q}_{k-1} + \beta^{-1} \log \mathbf{q}, \mathbf{q} \rangle \right\}$$

where  $\Pi(\mathbf{q}_{k-1}, \mathbf{q}) := \{\mathbf{M} \in \mathbb{R}^{N \times N} \mid \mathbf{M} \geq \mathbf{0} \text{ (elementwise)}, \mathbf{M}\mathbf{1} = \mathbf{q}_{k-1}, \mathbf{M}^\top \mathbf{1} = \mathbf{q}\}$

# Derivation of ProxLearn

Regularized discrete version of proximal recursion

$$\mathbf{q}_k = \arg \min_{\mathbf{q}} \left\{ \min_{\mathbf{M} \in \Pi(\mathbf{q}_{k-1}, \mathbf{q})} \frac{1}{2} \langle \mathbf{C}_k, \mathbf{M} \rangle + \epsilon \langle \mathbf{M}, \log \mathbf{M} \rangle + h \langle \mathbf{v}_{k-1} + \mathbf{U}_{k-1} \mathbf{q}_{k-1} + \beta^{-1} \log \mathbf{q}, \mathbf{q} \rangle \right\}$$

where  $\Pi(\mathbf{q}_{k-1}, \mathbf{q}) := \{ \mathbf{M} \in \mathbb{R}^{N \times N} \mid \mathbf{M} \geq \mathbf{0} \text{ (elementwise)}, \mathbf{M} \mathbf{1} = \mathbf{q}_{k-1}, \mathbf{M}^\top \mathbf{1} = \mathbf{q} \}$

Use Lagrange dual problem  
with Lagrange multipliers  $\lambda_0$   
and  $\lambda_1$

Let:

$$\mathbf{z} := \exp(\lambda_1 h / \epsilon)$$

$$\mathbf{q} := \exp(\lambda_0 h / \epsilon)$$

$$\mathbf{\Gamma}_k := \exp(-\mathbf{C}_k / 2\epsilon)$$

$$\boldsymbol{\xi}_{k-1} := \exp(-\beta \mathbf{v}_{k-1} - \beta \mathbf{U}_{k-1} \mathbf{q}_{k-1} - \mathbf{1})$$

Proximal update

$$\mathbf{q}_k = \mathbf{z} \odot \mathbf{\Gamma}_k^\top \mathbf{q}$$

# Proximal algorithm

---

## Algorithm 1 Proximal Algorithm

---

```
1: procedure PROXLEARN( $\boldsymbol{\rho}_{k-1}, \boldsymbol{\Theta}_{k-1}, \beta, h, \varepsilon, N, \mathbf{X}, \mathbf{y}, \delta, L$ )
2:    $\mathbf{v}_{k-1}, \mathbf{U}_{k-1}, \boldsymbol{\Theta}_k \leftarrow \text{EULERMARUYAMA}(h, \beta, \boldsymbol{\Theta}_{k-1}, \mathbf{X}, \mathbf{y}, \boldsymbol{\rho}_{k-1})$   $\triangleright$  Update the location of the samples
3:    $\mathbf{C}_k(i, j) \leftarrow \|\boldsymbol{\theta}_k^i - \boldsymbol{\theta}_{k-1}^j\|_2^2$ 
4:    $\boldsymbol{\Gamma}_k \leftarrow \exp(-\mathbf{C}_k/2\varepsilon)$ 
5:    $\boldsymbol{\xi}_{k-1} \leftarrow \exp(-\beta\mathbf{v}_{k-1} - \beta\mathbf{U}_{k-1}\boldsymbol{\rho}_{k-1} - \mathbf{1})$ 
6:    $\mathbf{z}_0 \leftarrow \text{rand}_{N \times 1}$ 
7:    $\mathbf{z} \leftarrow [\mathbf{z}_0, \mathbf{0}_{N \times (L-1)}]$ 
8:    $\mathbf{q} \leftarrow [\boldsymbol{\rho}_{k-1} \odot (\boldsymbol{\Gamma}_k \mathbf{z}_0), \mathbf{0}_{N \times (L-1)}]$ 
9:    $\ell = 1$ 
10:  while  $\ell \leq L$  do
11:     $\mathbf{z}(:, \ell + 1) \leftarrow (\boldsymbol{\xi}_{k-1} \odot (\boldsymbol{\Gamma}_k^\top \mathbf{q}(:, \ell)))^{\frac{1}{1+\beta\varepsilon/h}}$ 
12:     $\mathbf{q}(:, \ell + 1) \leftarrow \boldsymbol{\rho}_{k-1} \odot (\boldsymbol{\Gamma}_k \mathbf{z}(:, \ell + 1))$ 
13:    if  $\|\mathbf{q}(:, \ell + 1) - \mathbf{q}(:, \ell)\| < \delta$  and  $\|\mathbf{z}(:, \ell + 1) - \mathbf{z}(:, \ell)\| < \delta$  then
14:      Break
15:    else
16:       $\ell \leftarrow \ell + 1$ 
17:    end if
18:  end while
19:  return  $\boldsymbol{\rho}_k \leftarrow \mathbf{z}(:, \ell) \odot (\boldsymbol{\Gamma}_k^\top \mathbf{q}(:, \ell))$ 
20: end procedure
```

---

# Case Study: Binary Classification on WDBC Data

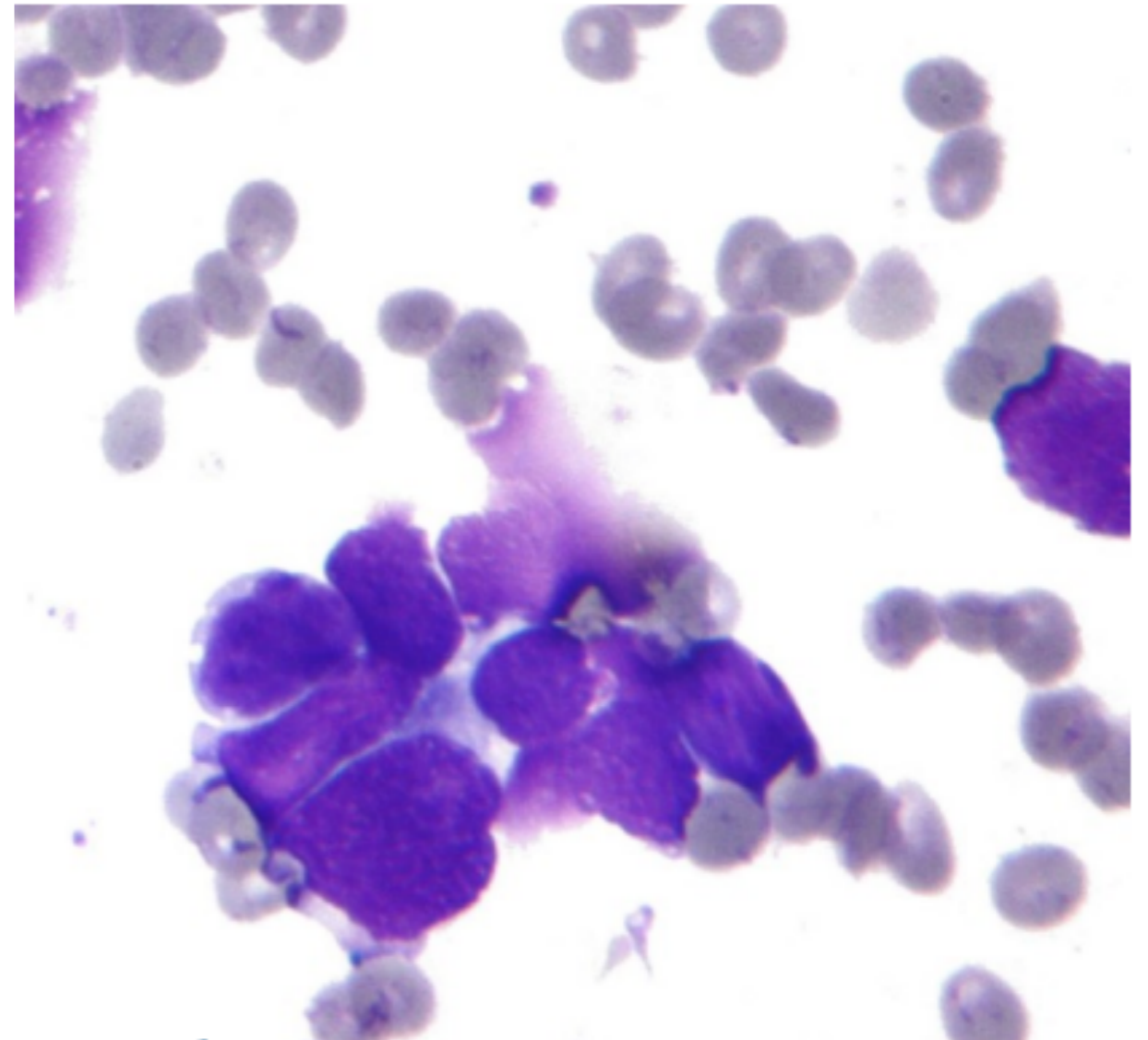
**WDBC:**

**# of features:**

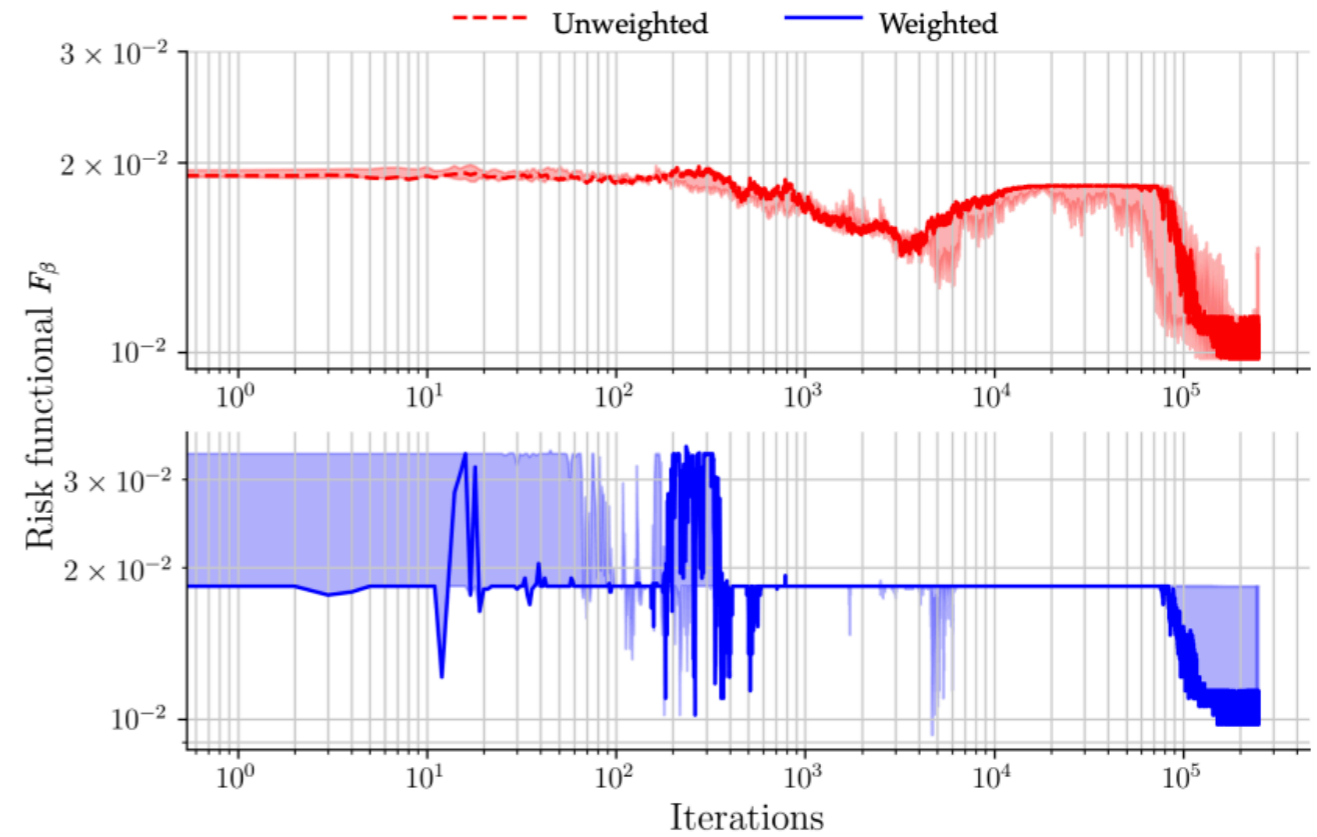
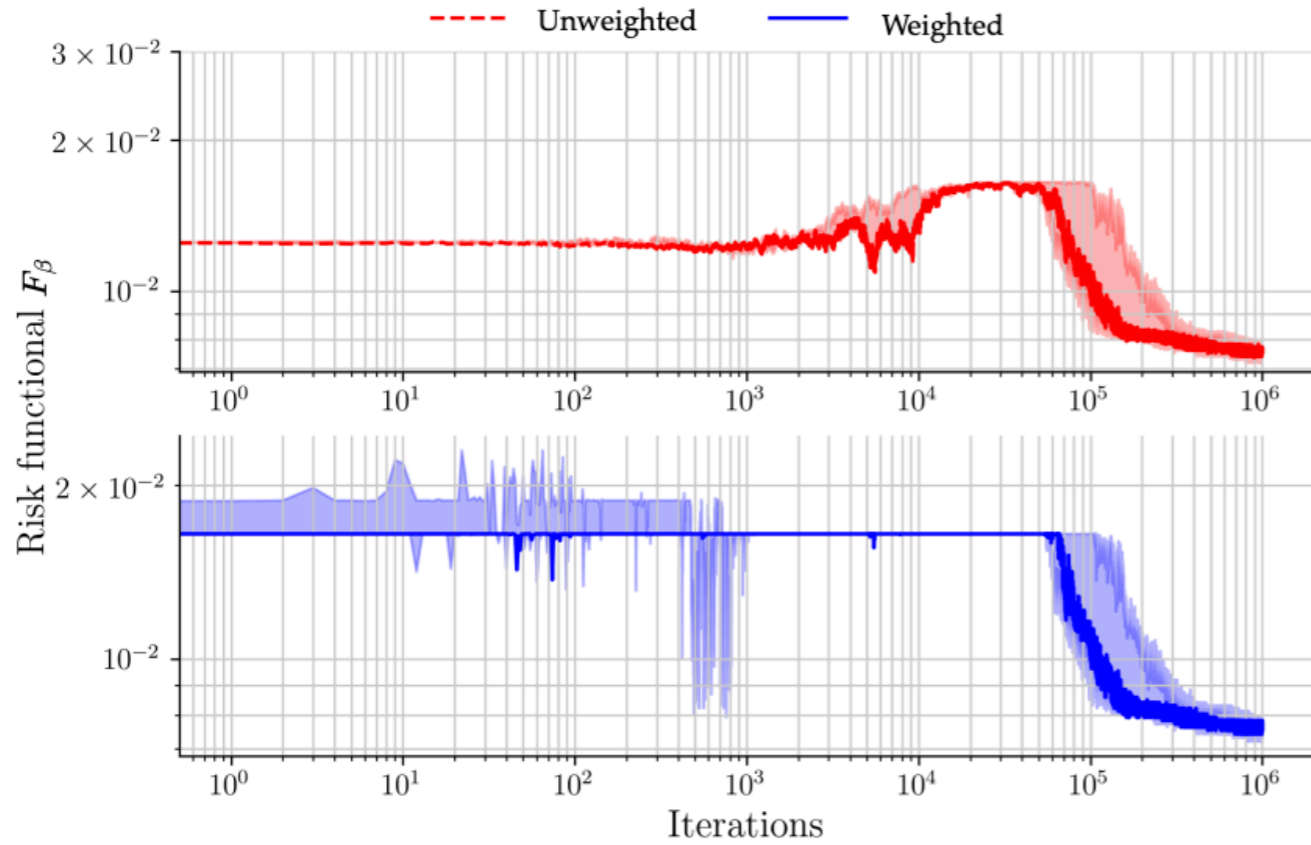
$$n_x = 30$$

**# of data points:**

$$n = 569$$



# Case Study: Binary Classification on WDBC Data



$\beta$	Unweighted	Weighted
0.03	91.17%	92.35%
0.05	92.94%	92.94%
0.07	78.23%	92.94%

$\beta$	Unweighted	Weighted	Runtime (hr)
0.03	91.18%	91.18%	1.415
0.05	91.18%	92.94%	1.533
0.07	90.59%	91.76%	1.704



# Case Study: Binary Classification

Comparison to Mokrov et al (2021) & Bonet et al (2022)

**Banana:**

# of features:

$$n_x = 2$$

# of data points:

$$n = 5300$$

**Diabetes:**

# of features:

$$n_x = 8$$

# of data points:

$$n = 768$$

**Twonorm:**

# of features:

$$n_x = 20$$

# of data points:

$$n = 7400$$

Dataset	JKO-ICNN	SWGf + RealNVP	ProxLearn, Weighted	ProxLearn, Unweighted
Banana	$0.550 \pm 10^{-2}$	$0.559 \pm 10^{-2}$	$0.551 \pm 10^{-2}$	$0.535 \pm 5 \cdot 10^{-2}$
Diabetes	$0.777 \pm 7 \cdot 10^{-3}$	$0.778 \pm 2 \cdot 10^{-3}$	$0.736 \pm 2 \cdot 10^{-2}$	$0.731 \pm 10^{-2}$
Twonorm	$0.981 \pm 2 \cdot 10^{-4}$	$0.981 \pm 6 \cdot 10^{-4}$	$0.972 \pm 2 \cdot 10^{-3}$	$0.972 \pm 2 \cdot 10^{-3}$

# Case Study: Multi-Class Classification



## Semeion Handwritten Digit Data Set

# of features:

$$n_x = 16 \times 16 = 256$$

# of data points:

$$n = 1593$$

# Case Study: Multi-Class Classification

$$P_{k-1}(j, i) := \Phi(\boldsymbol{\theta}_{k-1}^j, \mathbf{X}(i, :), \mathbf{Y}(i, :))$$

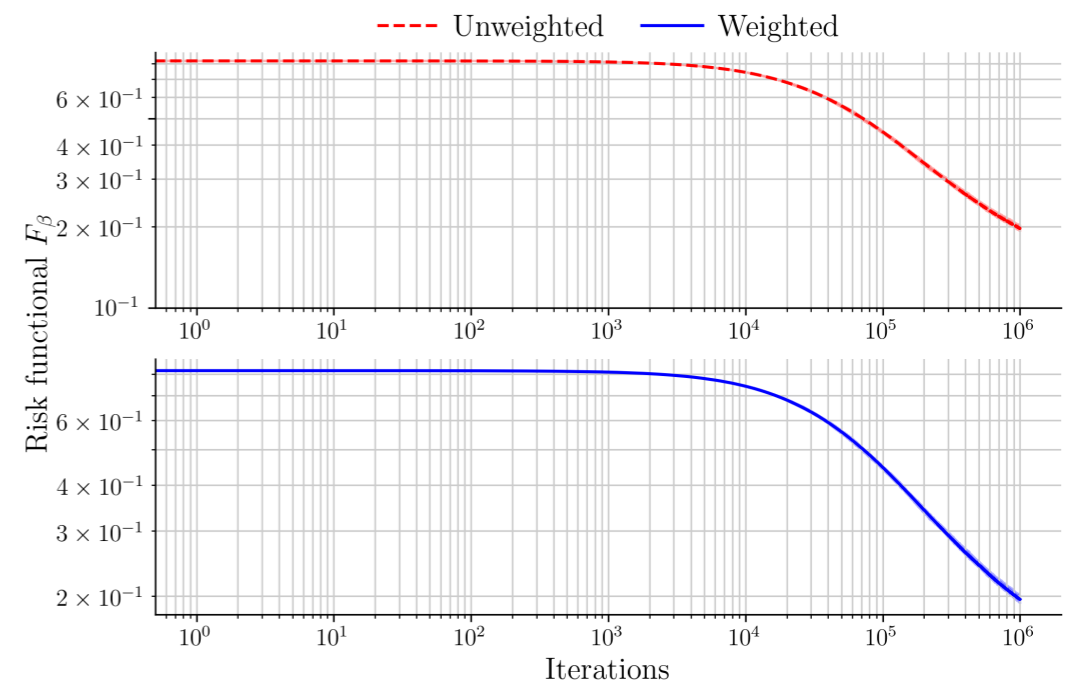
$$:= \left\langle \text{softmax}(\mathbf{X}(i, :)(\boldsymbol{\theta}_{k-1}^j)^\top), (\mathbf{Y}(i, :))^\top \right\rangle$$

## Weighted

$$F_\beta \approx \frac{1}{n_{\text{test}}} \left\| \mathbf{1} - (\mathbf{P}_{k-1}^{\text{test}})^\top \boldsymbol{\rho} \right\|_2^2$$

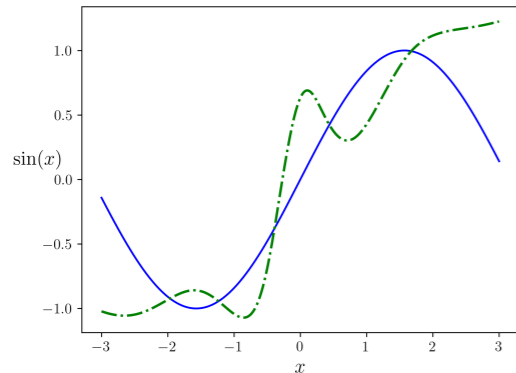
## Unweighted

$$F_\beta \approx \frac{1}{n_{\text{test}}} \left\| \mathbf{1} - \frac{1}{N} (\mathbf{P}_{k-1}^{\text{test}})^\top \mathbf{1} \right\|_2^2$$

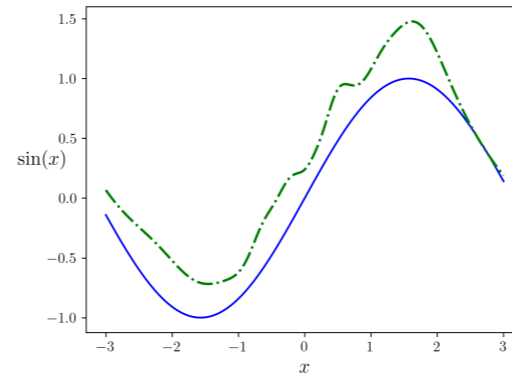


# Learning a sinusoid

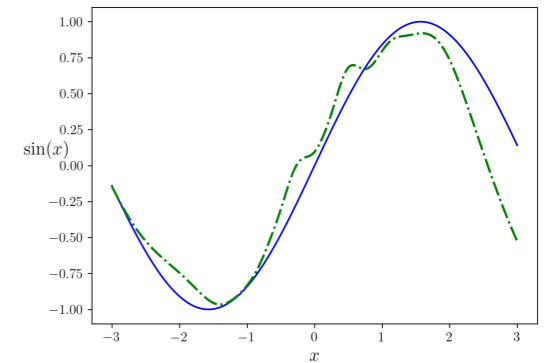
$N$	Final $\hat{F}_\beta$
500	0.01241931
700	0.01075817
1000	0.00806645
2000	0.00762518



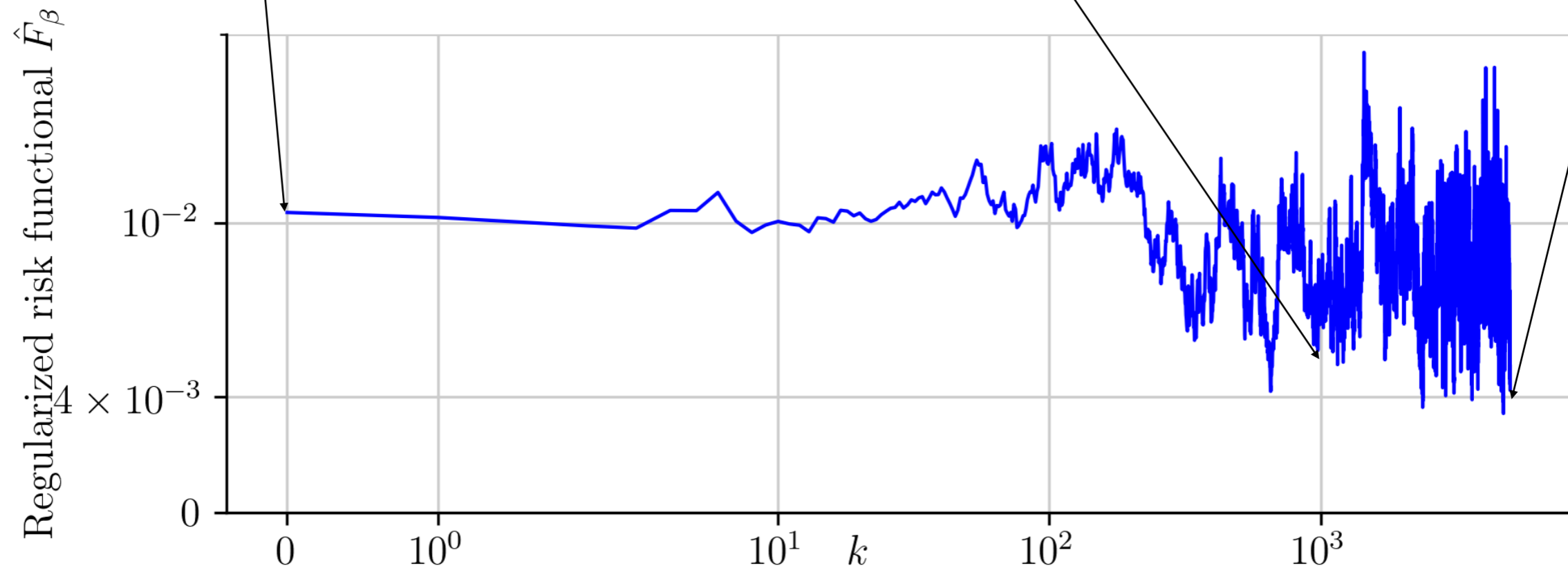
Iteration#1



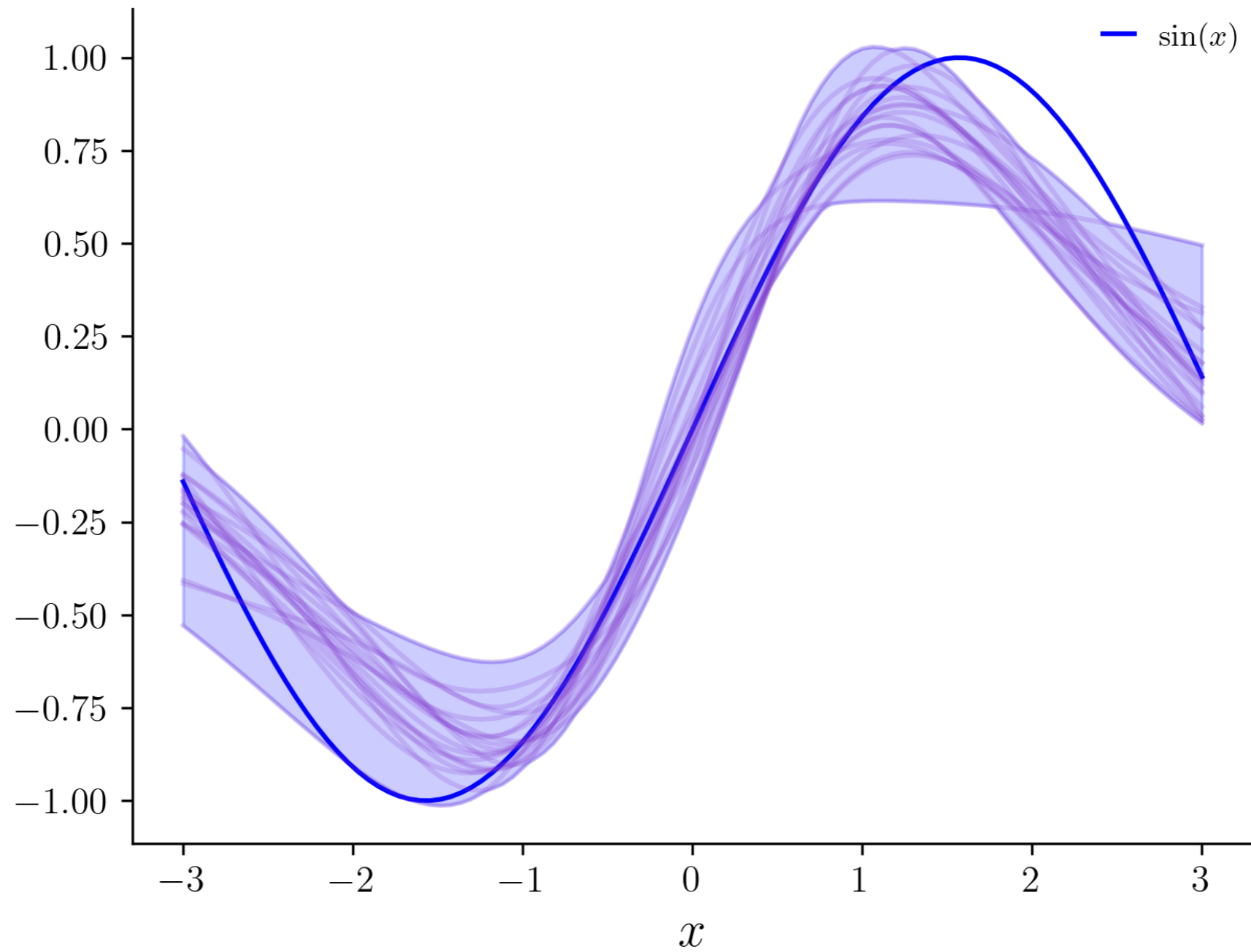
Iteration#1000



Iteration#5000



# Learning a sinusoid



# Additional Avenues of Research

## Multiple hidden layer setting

**(\*) Infinite width limit on one hidden layer; width of other hidden layers held constant**

**(\*) Widths of all hidden layers go to infinity**

# Thank You

**Acknowledgement:**



**2112755**