# Representation and regression problems in neural networks

## Mean-field relaxation, generalization, and numerics

Kang Liu

joint work with Enrique Zuazua

August 2024

Friedrich-Alexander-Universität
DYNAMICS, CONTROL,
MACHINE LEARNING
AND NUMERICS

# Table of Contents

# Table of Contents

# A diagram of classification task by NNs

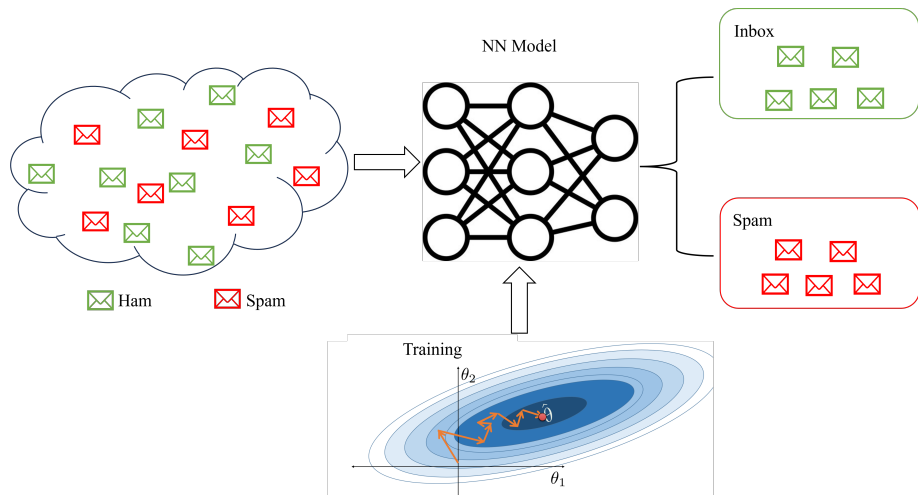# A diagram of classification task by NNs



Key Points: **Data**, **Neural Network Model**, <span style="color:red">**Training**</span>.

# General framework of training problems

- **Data**: $\{(x_i, y_i) \in \mathbb{R}^{d+1}\}_{i=1}^{N}$.

# General framework of training problems

- **Data**: $\{(x_i, y_i) \in \mathbb{R}^{d+1}\}_{i=1}^{N}$.
- **NN architecture**:

$$f: \mathbb{R}^d \times \mathbb{R}^p \to \mathbb{R}, \ (x, \Theta) \mapsto f(x, \Theta), \quad \text{where}$$

$x$ : feature (input), $\quad \Theta$ : parameter (control), $\quad f(x, \Theta)$ : prediction (output).

# General framework of training problems

- **Data**: $\{(x_i, y_i) \in \mathbb{R}^{d+1}\}_{i=1}^{N}$.
- **NN architecture**:

$$f : \mathbb{R}^d \times \mathbb{R}^p \to \mathbb{R}, \ (x, \Theta) \mapsto f(x, \Theta), \quad \text{where}$$

$x :$ feature (input), $\quad \Theta :$ parameter (control), $\quad f(x, \Theta) :$ prediction (output).

- **Three training scenarios**:

# General framework of training problems

- **Data**: $\{(x_i, y_i) \in \mathbb{R}^{d+1}\}_{i=1}^{N}$.

- **NN architecture**:

$$f: \mathbb{R}^d \times \mathbb{R}^p \to \mathbb{R}, \ (x, \Theta) \mapsto f(x, \Theta), \quad \text{where}$$

$x$ : feature (input), $\quad \Theta$ : parameter (control), $\quad f(x, \Theta)$ : prediction (output).

- **Three training scenarios**:

  1. **Exact representation**:

$$f(x_i, \Theta) = y_i, \quad \text{for } i = 1, \ldots, N.$$

# General framework of training problems

- **Data**: $\{(x_i, y_i) \in \mathbb{R}^{d+1}\}_{i=1}^{N}$.
- **NN architecture**:

$$f : \mathbb{R}^d \times \mathbb{R}^p \to \mathbb{R}, \ (x, \Theta) \mapsto f(x, \Theta), \quad \text{where}$$

$x$ : feature (input), $\quad \Theta$ : parameter (control), $\quad f(x, \Theta)$ : prediction (output).

- **Three training scenarios**:

  1. **Exact representation**:

  $$f(x_i, \Theta) = y_i, \quad \text{for } i = 1, \dots, N.$$

  2. **Approximate representation**:

  $$\|f(x_i, \Theta) - y_i\| \le \epsilon, \quad \text{for } i = 1, \dots, N.$$

# General framework of training problems

- **Data**: $\{(x_i, y_i) \in \mathbb{R}^{d+1}\}_{i=1}^{N}$.
- **NN architecture**:

$$f : \mathbb{R}^d \times \mathbb{R}^p \to \mathbb{R}, \ (x, \Theta) \mapsto f(x, \Theta), \quad \text{where}$$

$x$ : feature (input), $\quad \Theta$ : parameter (control), $\quad f(x, \Theta)$ : prediction (output).

- **Three training scenarios**:

  1. **Exact representation**:

  $$f(x_i, \Theta) = y_i, \quad \text{for } i = 1, \ldots, N.$$

  2. **Approximate representation**:

  $$\|f(x_i, \Theta) - y_i\| \leq \epsilon, \quad \text{for } i = 1, \ldots, N.$$

  3. **Regression**:

  $$\inf_{\Theta} \frac{1}{N} \sum_{i=1}^{N} \ell(f(x_i, \Theta) - y_i).$$

# General framework of training problems

- **Data**: $\{(x_i, y_i) \in \mathbb{R}^{d+1}\}_{i=1}^{N}$.

- **NN architecture**:

$$f: \mathbb{R}^d \times \mathbb{R}^p \to \mathbb{R}, \ (x, \Theta) \mapsto f(x, \Theta), \quad \text{where}$$

$x$ : feature (input), $\quad \Theta$ : parameter (control), $\quad f(x, \Theta)$ : prediction (output).

- **Three training scenarios**:

  1. **Exact representation**:

  $$f(x_i, \Theta) = y_i, \quad \text{for } i = 1, \ldots, N.$$

  2. **Approximate representation**:

  $$\|f(x_i, \Theta) - y_i\| \leq \epsilon, \quad \text{for } i = 1, \ldots, N.$$

  3. **Regression**:

  $$\inf_{\Theta} \frac{1}{N} \sum_{i=1}^{N} \ell(f(x_i, \Theta) - y_i).$$

# General framework of training problems

- **Data**: $\{(x_i, y_i) \in \mathbb{R}^{d+1}\}_{i=1}^{N}$.
- **NN architecture**:

$$f : \mathbb{R}^d \times \mathbb{R}^p \to \mathbb{R}, \ (x, \Theta) \mapsto f(x, \Theta), \quad \text{where}$$

$x :$ feature (input), $\quad \Theta :$ parameter (control), $\quad f(x, \Theta) :$ prediction (output).

- **Three training scenarios**:

  1. **Exact representation**:

  $$f(x_i, \Theta) = y_i, \quad \text{for } i = 1, \dots, N.$$

  2. **Approximate representation**:

  $$\|f(x_i, \Theta) - y_i\| \le \epsilon, \quad \text{for } i = 1, \dots, N.$$

  3. **Regression**:

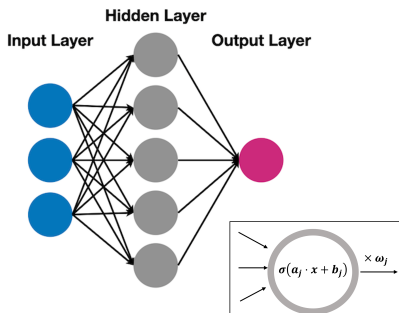  $$\inf_{\Theta} \frac{1}{N} \sum_{i=1}^{N} \ell(f(x_i, \Theta) - y_i).$$

## Problems

Existence, design of loss function, generalization property, numerical algorithms...

# Shallow Neural Network

**Shallow NNs with $P$ neurons**

$$f_{\text{shallow}}(x, \Theta) = \sum_{j=1}^{P} \omega_j \sigma(\langle a_j, x \rangle + b_j),$$

where $\Theta = (\omega_j, a_j, b_j)_{j=1}^{P}$, with $\omega_j \in \mathbb{R}$ and $(a_j, b_j) \in \mathbb{R}^{d+1}$.
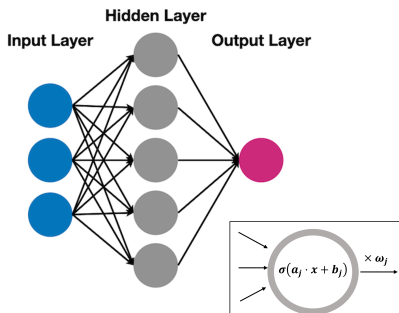
# Shallow Neural Network

**Shallow NNs with $P$ neurons**

$$f_{\text{shallow}}(x, \Theta) = \sum_{j=1}^{P} \omega_j \sigma(\langle a_j, x \rangle + b_j),$$

where $\Theta = (\omega_j, a_j, b_j)_{j=1}^{P}$, with
$\omega_j \in \mathbb{R}$ and $(a_j, b_j) \in \mathbb{R}^{d+1}$.



**Input Layer**   **Hidden Layer**   **Output Layer**

$\sigma(a_j \cdot x + b_j)$   $\times \omega_j$

## Why shallow NNs

- Simple structure;

# Shallow Neural Network

**Shallow NNs with $P$ neurons**

$$f_{\text{shallow}}(x, \Theta) = \sum_{j=1}^{P} \omega_j \sigma(\langle a_j, x \rangle + b_j),$$

where $\Theta = (\omega_j, a_j, b_j)_{j=1}^{P}$, with
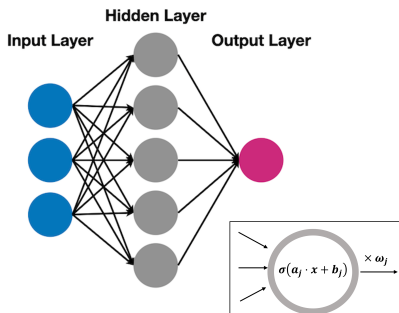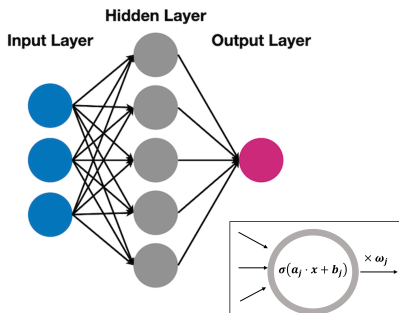$\omega_j \in \mathbb{R}$ and $(a_j, b_j) \in \mathbb{R}^{d+1}$.



**Why shallow NNs**

- Simple structure;
- Universal approximation property [Cybenko, 1989];

# Shallow Neural Network

**Shallow NNs with $P$ neurons**

$$f_{\text{shallow}}(x, \Theta) = \sum_{j=1}^{P} \omega_j \sigma(\langle a_j, x \rangle + b_j),$$

where $\Theta = (\omega_j, a_j, b_j)_{j=1}^{P}$, with $\omega_j \in \mathbb{R}$ and $(a_j, b_j) \in \mathbb{R}^{d+1}$.
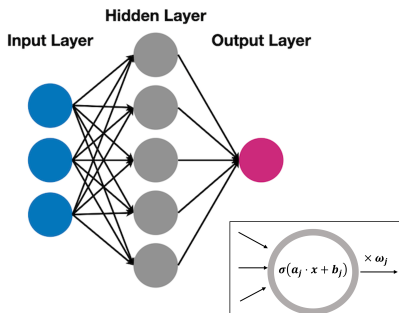


**Why shallow NNs**

- **Simple** structure;
- **Universal approximation** property [Cybenko, 1989];
- **Finite-sample representation** property [Pinkus, 1999];

# Shallow Neural Network

**Shallow NNs with $P$ neurons**

$$f_{\text{shallow}}(x, \Theta) = \sum_{j=1}^{P} \omega_j \sigma(\langle a_j, x \rangle + b_j),$$

where $\Theta = (\omega_j, a_j, b_j)_{j=1}^{P}$, with $\omega_j \in \mathbb{R}$ and $(a_j, b_j) \in \mathbb{R}^{d+1}$.



**Why shallow NNs**

- Simple structure;
- Universal approximation property [Cybenko, 1989];
- Finite-sample representation property [Pinkus, 1999];
- "Convergence" of the SGD algorithm [Chizat-Bach, 2018].

# Finite-sample representation property

Recall that

$$f_{\text{shallow}}(x, \Theta) = \sum_{j=1}^{P} \omega_j \sigma(\langle a_j, x \rangle + b_j).$$

### Finite-sample representation property [Pinkus 1999]

Assume that $P \geq N$ and $m = 1$. If $\sigma$ is non-polynomial, then for any distinct dataset $\{x_i, y_i\}_{i=1}^{N}$, there exists $\Theta$ such that

$$f_{\text{shallow}}(x_i, \Theta) = y_i, \quad \text{for } i = 1, \ldots, N.$$

# Finite-sample representation property

Recall that

$$f_{\text{shallow}}(x, \Theta) = \sum_{j=1}^{P} \omega_j \sigma(\langle a_j, x \rangle + b_j).$$

> **Finite-sample representation property [Pinkus 1999]**
>
> Assume that $P \geq N$ and $m = 1$. If $\sigma$ is non-polynomial, then for any distinct dataset $\{x_i, y_i\}_{i=1}^{N}$, there exists $\Theta$ such that
>
> $$f_{\text{shallow}}(x_i, \Theta) = y_i, \quad \text{for } i = 1, \ldots, N.$$

We extend in [L.-Zuazua, 2024] the previous result to the case where $y_i$ is in high dimension and $(a_j, b_j)$ are within a compact set. The proof is by induction and the application of the Hahn-Banach Theorem.

# Design of loss function/regularization

- A well-known principle [1] in machine learning is the following:
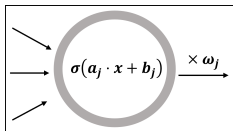
  "sparsity" mitigates "overfitting".

----

[1]Srivastava et al. "Dropout: A simple way to prevent Neural Networks from overfitting". In JMLR, 2014.

[2]Candes and Romberg. "Quantitative robust uncertainty principles and optimally sparse decompositions". In FOCM, 2006.

# Design of loss function/regularization

- A well-known principle [1] in machine learning is the following:

    "sparsity" mitigates "overfitting".

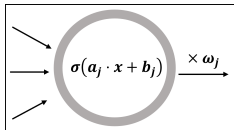- In shallow NNs, the number of activated neurons is $\|\omega\|_{\ell^0}$.



---

[1] Srivastava et al. "Dropout: A simple way to prevent Neural Networks from overfitting". In JMLR, 2014.

[2] Candes and Romberg. "Quantitative robust uncertainty principles and optimally sparse decompositions". In FOCM, 2006.

# Design of loss function/regularization

- A well-known principle [1] in machine learning is the following:

  "sparsity" mitigates "overfitting".

- In shallow NNs, the number of activated neurons is $\|\omega\|_{\ell^0}$.



- The function $\|\omega\|_{\ell^0}$ is non-convex. A practical replacement from compressed sensing [2]:

$$\|\omega\|_{\ell^0} \mapsto \|\omega\|_{\ell^1}.$$

---

[1] Srivastava et al. "Dropout: A simple way to prevent Neural Networks from overfitting". In JMLR, 2014.

[2] Candes and Romberg. "Quantitative robust uncertainty principles and optimally sparse decompositions". In FOCM, 2006.

# Primal problems

Let $\Omega$ be a compact subset of $\mathbb{R}^{d+1}$. Note $\Theta = (\omega_j, a_j, b_j)_{j=1}^{P}$.

- The sparse **exact representation** problem:

$$\inf_{\Theta \in (\mathbb{R} \times \Omega)^P} \|\omega\|_{\ell^1}, \quad \text{s.t.} \sum_{j=1}^{P} \omega_j \sigma(\langle a_j, x_i \rangle + b_j) = y_i, \quad \text{for } i = 1, \ldots, N. \quad (\text{P}_0)$$

# Primal problems

Let $\Omega$ be a compact subset of $\mathbb{R}^{d+1}$. Note $\Theta = (\omega_j, a_j, b_j)_{j=1}^P$.

- The sparse **exact representation** problem:

$$\inf_{\Theta \in (\mathbb{R} \times \Omega)^P} \|\omega\|_{\ell^1}, \quad \text{s.t.} \sum_{j=1}^P \omega_j \sigma(\langle a_j, x_i \rangle + b_j) = y_i, \quad \text{for } i = 1, \ldots, N. \quad (\text{P}_0)$$

- The sparse **approximate representation** problem:

$$\inf_{\Theta \in (\mathbb{R} \times \Omega)^P} \|\omega\|_{\ell^1}, \quad \text{s.t.} \left| \sum_{j=1}^P \omega_j \sigma(\langle a_j, x_i \rangle + b_j) - y_i \right| \leq \epsilon, \quad \text{for } i = 1, \ldots, N,$$

$$(\text{P}_\epsilon)$$

where $\epsilon > 0$ is a hyperparameter.

# Primal problems

Let $\Omega$ be a compact subset of $\mathbb{R}^{d+1}$. Note $\Theta = (\omega_j, a_j, b_j)_{j=1}^P$.

- The sparse **exact representation** problem:

$$\inf_{\Theta \in (\mathbb{R} \times \Omega)^P} \|\omega\|_{\ell^1}, \quad \text{s.t.} \sum_{j=1}^P \omega_j \sigma(\langle a_j, x_i \rangle + b_j) = y_i, \quad \text{for } i = 1, \ldots, N. \quad (\text{P}_0)$$

- The sparse **approximate representation** problem:

$$\inf_{\Theta \in (\mathbb{R} \times \Omega)^P} \|\omega\|_{\ell^1}, \quad \text{s.t.} \left| \sum_{j=1}^P \omega_j \sigma(\langle a_j, x_i \rangle + b_j) - y_i \right| \leq \epsilon, \quad \text{for } i = 1, \ldots, N,$$
$$(\text{P}_\epsilon)$$

where $\epsilon > 0$ is a hyperparameter.

- The sparse **regression** problem:

$$\inf_{\Theta \in (\mathbb{R} \times \Omega)^P} \|\omega\|_{\ell^1} + \frac{\lambda}{N} \sum_{i=1}^N \ell \left( \sum_{j=1}^P \omega_j \sigma(\langle a_j, x_i \rangle + b_j) - y_i \right), \quad (\text{P}_\lambda^{\text{reg}})$$

where $\lambda > 0$ is a hyperparameter.

# Problem

- How can we address these high-dimensional and non-convex optimization problems?

# Table of Contents

# Mean-field relaxation

Primal problems $(P_0)$, $(P_\epsilon)$, and $(P_\lambda^{\text{reg}})$ are non-convex optimization problems, where the non-convexity is from the non-linearity of shallow NNs, e.g.,

$$\left\{ \Theta \,\middle|\, \sum_{j=1}^{P} \omega_j \sigma(\langle a_j, x_i \rangle + b_j) = y_i, \,\forall i = 1, \ldots, N \right\} \text{ is a non-convex set.}$$

# Mean-field relaxation

Primal problems (P$_0$), (P$_\epsilon$), and (P$_\lambda^{\text{reg}}$) are non-convex optimization problems, where the non-convexity is from the non-linearity of shallow NNs, e.g.,

$$\left\{ \Theta \,\Big|\, \sum_{j=1}^{P} \omega_j \sigma(\langle a_j, x_i \rangle + b_j) = y_i, \, \forall i = 1, \ldots, N \right\} \text{ is a non-convex set.}$$

The mean-field relaxation technique is commonly employed in shallow NNs, see [Mei-Montanari-Nguyen, 2018] and [Chizat-Bach, 2018].

# Mean-field relaxation

Primal problems $(P_0)$, $(P_\epsilon)$, and $(P_\lambda^{reg})$ are non-convex optimization problems, where the non-convexity is from the non-linearity of shallow NNs, e.g.,

$$\left\{ \Theta \;\Big|\; \sum_{j=1}^{P} \omega_j \sigma(\langle a_j, x_i \rangle + b_j) = y_i, \; \forall i = 1, \ldots, N \right\} \text{ is a non-convex set.}$$

The mean-field relaxation technique is commonly employed in shallow NNs, see [Mei-Montanari-Nguyen, 2018] and [Chizat-Bach, 2018].

## Shallow NN

The original shallow NN writes:

$$\sum_{j=1}^{P} \omega_j \sigma(\langle a_j, x \rangle + b_j),$$

where $(\omega_j, a_j, b_j) \in \mathbb{R} \times \Omega$ for all $j$.

**Cost function**: $\|\omega\|_{\ell^1}$.

## Mean-field shallow NN

The mean-field shallow NN writes:

$$\int_\Omega \sigma(\langle a, x \rangle + b) d\mu(a, b),$$

where $\mu \in \mathcal{M}(\Omega)$. The outcome is linear with respect to $\mu$.

**Cost function**: $\|\mu\|_{TV}$.

# Relaxed problems

Let $Y = (y_1, \ldots, y_N)$. Define the following linear mapping:

$$\phi \, \mu := (\phi_i \, \mu)_{i=1}^N = \left( \int_\Omega \sigma(\langle a, x_i \rangle + b) d\mu(a, b) \right)_{i=1}^N$$

# Relaxed problems

Let $Y = (y_1, \ldots, y_N)$. Define the following linear mapping:

$$\phi \mu := (\phi_i \mu)_{i=1}^N = \left( \int_\Omega \sigma(\langle a, x_i \rangle + b) d\mu(a, b) \right)_{i=1}^N$$

Convex relaxations:

- The relaxation of $(\mathrm{P}_0)$:

$$\inf_{\mu \in \mathcal{M}(\Omega)} \|\mu\|_{\mathsf{TV}}, \quad \text{s.t. } \phi \mu = Y. \qquad (\mathrm{PR}_0)$$

# Relaxed problems

Let $Y = (y_1, \ldots, y_N)$. Define the following linear mapping:

$$\phi \, \mu := (\phi_i \, \mu)_{i=1}^N = \left( \int_\Omega \sigma(\langle a, x_i \rangle + b) d\mu(a, b) \right)_{i=1}^N$$

Convex relaxations:

- The relaxation of $(P_0)$:

$$\inf_{\mu \in \mathcal{M}(\Omega)} \|\mu\|_{\mathsf{TV}}, \quad \text{s.t. } \phi \, \mu = Y. \qquad (\mathsf{PR}_0)$$

- The relaxation of $(P_\epsilon)$:

$$\inf_{\mu \in \mathcal{M}(\Omega)} \|\mu\|_{\mathsf{TV}}, \quad \text{s.t. } \|\phi \, \mu - Y\|_{\ell^\infty} \leq \epsilon. \qquad (\mathsf{PR}_\epsilon)$$

# Relaxed problems

Let $Y = (y_1, \ldots, y_N)$. Define the following linear mapping:

$$\phi\,\mu := (\phi_i\,\mu)_{i=1}^N = \left( \int_\Omega \sigma(\langle a, x_i \rangle + b)\,d\mu(a, b) \right)_{i=1}^N$$

Convex relaxations:

- The relaxation of $(P_0)$:

$$\inf_{\mu \in \mathcal{M}(\Omega)} \|\mu\|_{\mathsf{TV}}, \quad \text{s.t. } \phi\,\mu = Y. \qquad (PR_0)$$

- The relaxation of $(P_\epsilon)$:

$$\inf_{\mu \in \mathcal{M}(\Omega)} \|\mu\|_{\mathsf{TV}}, \quad \text{s.t. } \|\phi\,\mu - Y\|_{\ell^\infty} \le \epsilon. \qquad (PR_\epsilon)$$

- The relaxation of $(P_\lambda^{\mathbf{reg}})$:

$$\inf_{\mu \in \mathcal{M}(\Omega)} \|\mu\|_{\mathsf{TV}} + \frac{\lambda}{N} \sum_{i=1}^N \ell\,(\phi_i\,\mu - y_i). \qquad (PR_\lambda^{\mathbf{reg}})$$

# Free of relaxation gap

## Theorem (L.-Zuazua,2024)

*Under mild assumptions [1] on $\sigma$ and $\Omega$, if $P \geq N$, then*

$$\text{val}(P_0) = \text{val}(PR_0); \quad \text{val}(P_\epsilon) = \text{val}(PR_\epsilon); \quad \text{val}(P_\lambda^{\text{reg}}) = \text{val}(PR_\lambda^{\text{reg}}).$$

*Moreover, the extreme points of the solution sets of relaxed problems have the following form:*

$$\mu^* = \sum_{j=1}^{N} \omega_j^* \delta_{(a_j^*, b_j^*)}.$$

---

[1] An example of $(\sigma, \Omega)$: $\sigma$ is the ReLU function and $\Omega$ is the unit ball.

[2] Similar results for particular scenarios of exact representation and regression in ML obtained by representer theorems are studied in [Unser, 2019] and [Dios-Bruna, 2020].

# Free of relaxation gap

## Theorem (L.-Zuazua,2024)

*Under mild assumptions [1] on $\sigma$ and $\Omega$, if $P \geq N$, then*

$$\text{val}(P_0) = \text{val}(PR_0); \quad \text{val}(P_\epsilon) = \text{val}(PR_\epsilon); \quad \text{val}(P_\lambda^{\text{reg}}) = \text{val}(PR_\lambda^{\text{reg}}).$$

*Moreover, the extreme points of the solution sets of relaxed problems have the following form:*

$$\mu^* = \sum_{j=1}^{N} \omega_j^* \delta_{(a_j^*, b_j^*)}.$$

Main techniques in the proof:

- Existence of solutions: finite-sample representation property.
- "Representer Theorem" [2] from [Fisher-Jerome, 1975].

---

[1] An example of $(\sigma, \Omega)$: $\sigma$ is the ReLU function and $\Omega$ is the unit ball.
[2] Similar results for particular scenarios of exact representation and regression in ML obtained by representer theorems are studied in [Unser, 2019] and [Dios-Bruna, 2020].

# Problems

- How should the hyperparameters $\epsilon$ and $\lambda$ be chosen in these problems? (Generalization)

- How can the relaxed problems be solved, and how can solutions of the primal problems be found? (Numerical algorithms)

# Table of Contents

# A generalization bound

- **Training/Testing** dataset: $\{(x_i, y_i)\}_{i=1}^{N}$ / $\{(x_i', y_i')\}_{i=1}^{N'}$.

# A generalization bound

- **Training/Testing** dataset: $\{(x_i, y_i)\}_{i=1}^{N}$ / $\{(x_i', y_i')\}_{i=1}^{N'}$.
- **Predictions** on testing set by the shallow NN with parameter $\Theta$:

$$\{(x_i', f_{\text{shallow}}(x_i', \Theta))\}_{i=1}^{N'}$$

# A generalization bound

- **Training/Testing** dataset: $\{(x_i, y_i)\}_{i=1}^{N}$ / $\{(x_i', y_i')\}_{i=1}^{N'}$.
- **Predictions** on testing set by the shallow NN with parameter $\Theta$:

$$\{(x_i', f_{\text{shallow}}(x_i', \Theta))\}_{i=1}^{N'}$$

- **Empirical measures**:

$$m_{\text{train}} = \frac{1}{N} \sum_{i=1}^{N} \delta_{(x_i, y_i)}, \quad m_{\text{test}} = \frac{1}{N'} \sum_{i=1}^{N'} \delta_{(x_i', y_i')}, \quad m_{\text{pred}}(\Theta) = \frac{1}{N'} \sum_{i=1}^{N'} \delta_{(x_i', f_{\text{shallow}}(x_i', \Theta))}.$$

# A generalization bound

- **Training/Testing** dataset: $\{(x_i, y_i)\}_{i=1}^{N}$ / $\{(x_i', y_i')\}_{i=1}^{N'}$.
- **Predictions** on testing set by the shallow NN with parameter $\Theta$:

$$\{(x_i', f_{\text{shallow}}(x_i', \Theta))\}_{i=1}^{N'}$$

- **Empirical measures**:

$$m_{\text{train}} = \frac{1}{N} \sum_{i=1}^{N} \delta_{(x_i, y_i)}, \quad m_{\text{test}} = \frac{1}{N'} \sum_{i=1}^{N'} \delta_{(x_i', y_i')}, \quad m_{\text{pred}}(\Theta) = \frac{1}{N'} \sum_{i=1}^{N'} \delta_{(x_i', f_{\text{shallow}}(x_i', \Theta))}.$$

# A generalization bound

- **Training/Testing** dataset: $\{(x_i, y_i)\}_{i=1}^{N}$ / $\{(x_i', y_i')\}_{i=1}^{N'}$.

- **Predictions** on testing set by the shallow NN with parameter $\Theta$:

$$\{(x_i', f_{\text{shallow}}(x_i', \Theta))\}_{i=1}^{N'}$$

- **Empirical measures**:

$$m_{\text{train}} = \frac{1}{N} \sum_{i=1}^{N} \delta_{(x_i, y_i)}, \quad m_{\text{test}} = \frac{1}{N'} \sum_{i=1}^{N'} \delta_{(x_i', y_i')}, \quad m_{\text{pred}}(\Theta) = \frac{1}{N'} \sum_{i=1}^{N'} \delta_{(x_i', f_{\text{shallow}}(x_i', \Theta))}.$$

## Theorem (L.-Zuazua,2024)

*Let $W_1(\cdot, \cdot)$ denote the Wasserstein-1 distance. If $\sigma$ is 1-Lipschitz, then for any $\Theta$,*

$$W_1(m_{\text{test}}, m_{\text{pred}}(\Theta)) \leq \underbrace{2 W_1(m_{\text{train}}, m_{\text{test}})}_{\text{Bias from datasets}} + r(\Theta), \quad \text{where}$$

$$r(\Theta) = \underbrace{\frac{1}{N} \sum_{i=1}^{N} |f_{\text{shallow}}(x_i, \Theta) - y_i|}_{\text{Bias from training}} + \underbrace{W_1(m_{\text{train}}, m_{\text{test}}) \sum_{j=1}^{P} |\omega_j| \|a_j\|}_{\text{``Variance''}}.$$

# Generalization bounds by optimal solutions

Fix the following:

- $\sigma$: ReLU;
- $\Omega$: $B^{d+1}(0,1)$;
- $\ell(\cdot) = |\cdot|$.

Recall that

$$W_1(m_{\text{test}}, m_{\text{pred}}(\Theta)) \leq \underbrace{2W_1(m_{\text{train}}, m_{\text{test}})}_{\text{Bias from datasets}} + r(\Theta).$$

# Generalization bounds by optimal solutions

Fix the following:

- $\sigma$: ReLU;
- $\Omega$: $B^{d+1}(0, 1)$;
- $\ell(\cdot) = |\cdot|$.

Recall that

$$W_1(m_{\text{test}}, m_{\text{pred}}(\Theta)) \leq \underbrace{2W_1(m_{\text{train}}, m_{\text{test}})}_{\text{Bias from datasets}} + r(\Theta).$$

---

### Proposition

*Let $P \geq N$. For any $\epsilon \geq 0$ and $\lambda > 0$, let $\Theta_\epsilon$ and $\Theta_\lambda^{\text{reg}}$ be the solutions of $(P_\epsilon)$ and $(P_\lambda^{\text{reg}})$, respectively. Then,*

$$r(\Theta_\epsilon) \leq \mathcal{U}(\epsilon) := \epsilon + C \, \text{val}(P_\epsilon);$$

$$r(\Theta_\lambda^{\text{reg}}) \leq \mathcal{L}(\lambda) := \max\{\lambda^{-1}, C\} \, \text{val}(PR_\lambda^{\text{reg}}),$$

*where $C = W_1(m_{\text{train}}, m_{\text{test}})$.*

# Optimal hyperparameters

Recall that $C = W_1(m_{\text{train}}, m_{\text{test}})$.

- Optimal value of $\lambda$: $\lambda^* = C^{-1}$.

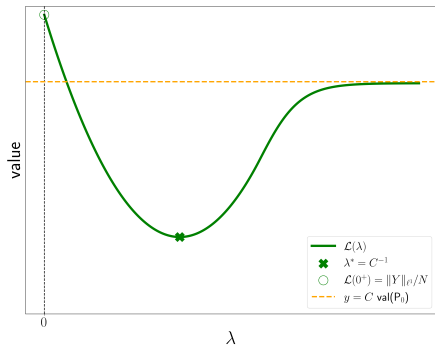# Optimal hyperparameters

Recall that $C = W_1(m_{\text{train}}, m_{\text{test}})$.

- Optimal value of $\lambda$: $\lambda^* = C^{-1}$.
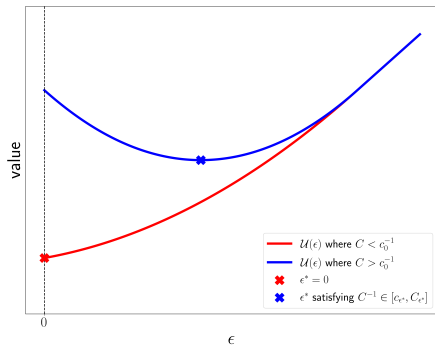- Optimal value of $\epsilon$:

# Optimal hyperparameters

Recall that $C = W_1(m_{\text{train}}, m_{\text{test}})$.

- Optimal value of $\lambda$: $\lambda^* = C^{-1}$.
- Optimal value of $\epsilon$:
  1. if $C < c_0^{-1}$, then $\epsilon^* = 0$;

# Optimal hyperparameters

Recall that $C = W_1(m_{\text{train}}, m_{\text{test}})$.

- Optimal value of $\lambda$: $\lambda^* = C^{-1}$.
- Optimal value of $\epsilon$:
  1. if $C < c_0^{-1}$, then $\epsilon^* = 0$;
  2. if $C \geq c_0^{-1}$, then $\epsilon^*$ satisfies the first-order optimality condition $C^{-1} \in [c_{\epsilon^*}, C_{\epsilon^*}]$.

# Optimal hyperparameters

Recall that $C = W_1(m_{\text{train}}, m_{\text{test}})$.

- Optimal value of $\lambda$: $\lambda^* = C^{-1}$.
- Optimal value of $\epsilon$:
  1. if $C < c_0^{-1}$, then $\epsilon^* = 0$;
  2. if $C \geq c_0^{-1}$, then $\epsilon^*$ satisfies the first-order optimality condition $C^{-1} \in [c_{\epsilon^*}, C_{\epsilon^*}]$.

# Optimal hyperparameters

Recall that $C = W_1(m_{\text{train}}, m_{\text{test}})$.

- Optimal value of $\lambda$: $\lambda^* = C^{-1}$.
- Optimal value of $\epsilon$:
  1. if $C < c_0^{-1}$, then $\epsilon^* = 0$;
  2. if $C \geq c_0^{-1}$, then $\epsilon^*$ satisfies the first-order optimality condition $C^{-1} \in [c_{\epsilon^*}, C_{\epsilon^*}]$.

  Here, $(c_\epsilon, C_\epsilon)$ is related to the solutions of the dual problem of $(\text{PR}_\epsilon)$.

# Optimal hyperparameters

Recall that $C = W_1(m_{\text{train}}, m_{\text{test}})$.

- Optimal value of $\lambda$: $\lambda^* = C^{-1}$.
- Optimal value of $\epsilon$:
  1. if $C < c_0^{-1}$, then $\epsilon^* = 0$;
  2. if $C \geq c_0^{-1}$, then $\epsilon^*$ satisfies the first-order optimality condition $C^{-1} \in [c_{\epsilon^*}, C_{\epsilon^*}]$.

  Here, $(c_\epsilon, C_\epsilon)$ is related to the solutions of the dual problem of $(\text{PR}_\epsilon)$.



(a) Qualitative curve of $\mathcal{L}(\lambda)$.  (b) Two scenarios of $\mathcal{U}(\epsilon)$.

# Table of Contents

# Guideline for numerical algorithms

Relaxed problems are convex, but in an infinite-dimensional space.

$$\inf_{\mu \in \mathcal{M}(\Omega)} \|\mu\|_{\mathsf{TV}}, \quad \text{s.t.} \ \|\phi\,\mu - Y\|_{\ell^\infty} \leq \epsilon. \tag{PR$_\epsilon$}$$

$$\inf_{\mu \in \mathcal{M}(\Omega)} \|\mu\|_{\mathsf{TV}} + \frac{\lambda}{N} \sum_{i=1}^{N} |\phi_i\,\mu - y_i|. \tag{PR$_\lambda^{\mathsf{reg}}$}$$

# Guideline for numerical algorithms

Relaxed problems are **convex**, but in an **infinite-dimensional** space.

$$\inf_{\mu \in \mathcal{M}(\Omega)} \|\mu\|_{\mathsf{TV}}, \quad \text{s.t. } \|\phi\,\mu - Y\|_{\ell^\infty} \leq \epsilon. \tag{$PR_\epsilon$}$$

$$\inf_{\mu \in \mathcal{M}(\Omega)} \|\mu\|_{\mathsf{TV}} + \frac{\lambda}{N} \sum_{i=1}^{N} |\phi_i\,\mu - y_i|. \tag{$PR_\lambda^{\mathsf{reg}}$}$$

A general approach: **Discretization**, then **Optimization**.

# Guideline for numerical algorithms

Relaxed problems are convex, but in an infinite-dimensional space.

$$\inf_{\mu \in \mathcal{M}(\Omega)} \|\mu\|_{\mathsf{TV}}, \quad \text{s.t. } \|\phi\,\mu - Y\|_{\ell^\infty} \leq \epsilon. \qquad (\mathsf{PR}_\epsilon)$$

$$\inf_{\mu \in \mathcal{M}(\Omega)} \|\mu\|_{\mathsf{TV}} + \frac{\lambda}{N} \sum_{i=1}^{N} |\phi_i\,\mu - y_i|. \qquad (\mathsf{PR}_\lambda^{\mathsf{reg}})$$

A general approach: **Discretization**, then **Optimization**.

**Two numerical scenarios**

1. When $\dim(\Omega) = d + 1$ is small, discretize $\Omega$ by a mesh, then optimize by the simplex method.

# Guideline for numerical algorithms

Relaxed problems are convex, but in an infinite-dimensional space.

$$\inf_{\mu \in \mathcal{M}(\Omega)} \|\mu\|_{\mathsf{TV}}, \quad \text{s.t.} \ \|\phi\,\mu - Y\|_{\ell^\infty} \leq \epsilon. \qquad (\mathrm{PR}_\epsilon)$$

$$\inf_{\mu \in \mathcal{M}(\Omega)} \|\mu\|_{\mathsf{TV}} + \frac{\lambda}{N} \sum_{i=1}^{N} |\phi_i\,\mu - y_i|. \qquad (\mathrm{PR}_\lambda^{\mathsf{reg}})$$

A general approach: **Discretization**, then **Optimization**.

**Two numerical scenarios**

1. When $\dim(\Omega) = d + 1$ is small, discretize $\Omega$ by a mesh, then optimize by the simplex method.

2. When $\dim(\Omega) = d + 1$ is great, discretize $(\mathrm{PR}_\lambda^{\mathsf{reg}})$ by an overparameterized version (problem $(\mathrm{P}_\lambda^{\mathsf{reg}})$ with a large $P$), then optimize by the SGD algorithm.

# Low-dimensional scenario

- Discretization of the domain:

$$\Omega \to \Omega_h = \{(a_j, b_j)\}_{j=1}^{M}.$$

# Low-dimensional scenario

- Discretization of the domain:

$$\Omega \to \Omega_h = \{(a_j, b_j)\}_{j=1}^{M}.$$

- Discretized problems:

$$\inf_{\omega \in \mathbb{R}^M} \|\omega\|_{\ell^1}, \quad \text{s.t. } \|A\omega - Y\|_{\ell^\infty} \leq \epsilon, \tag{$PD_\epsilon$}$$

$$\inf_{\omega \in \mathbb{R}^M} \|\omega\|_{\ell^1} + \frac{\lambda}{N}\|A\omega - Y\|_{\ell^1}, \tag{$PD_\lambda^{\text{reg}}$}$$

where $A \in \mathbb{R}^{N \times M}$ with $A_{ij} = \sigma(\langle a_j, x_i \rangle + b_j)$.

# Low-dimensional scenario

- Discretization of the domain:

$$\Omega \to \Omega_h = \{(a_j, b_j)\}_{j=1}^{M}.$$

- Discretized problems:

$$\inf_{\omega \in \mathbb{R}^M} \|\omega\|_{\ell^1}, \quad \text{s.t. } \|A\omega - Y\|_{\ell^\infty} \leq \epsilon, \tag{PD$_\epsilon$}$$

$$\inf_{\omega \in \mathbb{R}^M} \|\omega\|_{\ell^1} + \frac{\lambda}{N}\|A\omega - Y\|_{\ell^1}, \tag{PD$_\lambda^{\text{reg}}$}$$

where $A \in \mathbb{R}^{N \times M}$ with $A_{ij} = \sigma(\langle a_j, x_i\rangle + b_j)$.

- Error estimates:

$$|\text{val}(\text{PD}_\epsilon) - \text{val}(\text{PR}_\epsilon)|, \, |\text{val}(\text{PD}_\lambda^{\text{reg}}) - \text{val}(\text{PR}_\lambda^{\text{reg}})| = \mathcal{O}(d_{\text{Hausdorff}}(\Omega, \Omega_h)).$$

# Low-dimensional scenario

- Discretization of the domain:

$$\Omega \to \Omega_h = \{(a_j, b_j)\}_{j=1}^{M}.$$

- Discretized problems:

$$\inf_{\omega \in \mathbb{R}^M} \|\omega\|_{\ell^1}, \quad \text{s.t.} \ \|A\omega - Y\|_{\ell^\infty} \leq \epsilon, \tag{PD$_\epsilon$}$$

$$\inf_{\omega \in \mathbb{R}^M} \|\omega\|_{\ell^1} + \frac{\lambda}{N}\|A\omega - Y\|_{\ell^1}, \tag{PD$_\lambda^{\text{reg}}$}$$

where $A \in \mathbb{R}^{N \times M}$ with $A_{ij} = \sigma(\langle a_j, x_i \rangle + b_j)$.

- Error estimates:

$$|\text{val}(\text{PD}_\epsilon) - \text{val}(\text{PR}_\epsilon)|, \ |\text{val}(\text{PD}_\lambda^{\text{reg}}) - \text{val}(\text{PR}_\lambda^{\text{reg}})| = \mathcal{O}(d_{\text{Hausdorff}}(\Omega, \Omega_h)).$$

- Equivalent to linear programming problems, solvable using the simplex method.

# Low-dimensional scenario

- Discretization of the domain:

$$\Omega \to \Omega_h = \{(a_j, b_j)\}_{j=1}^{M}.$$

- Discretized problems:

$$\inf_{\omega \in \mathbb{R}^M} \|\omega\|_{\ell^1}, \quad \text{s.t. } \|A\omega - Y\|_{\ell^\infty} \leq \epsilon, \tag{PD$_\epsilon$}$$

$$\inf_{\omega \in \mathbb{R}^M} \|\omega\|_{\ell^1} + \frac{\lambda}{N} \|A\omega - Y\|_{\ell^1}, \tag{PD$_\lambda^{\text{reg}}$}$$

where $A \in \mathbb{R}^{N \times M}$ with $A_{ij} = \sigma(\langle a_j, x_i \rangle + b_j)$.

- Error estimates:

$$|\text{val}(\text{PD}_\epsilon) - \text{val}(\text{PR}_\epsilon)|, \, |\text{val}(\text{PD}_\lambda^{\text{reg}}) - \text{val}(\text{PR}_\lambda^{\text{reg}})| = \mathcal{O}(d_{\text{Hausdorff}}(\Omega, \Omega_h)).$$

- Equivalent to linear programming problems, solvable using the simplex method.
  - **Advantage**: Terminates at an extreme point of the solution set, which corresponds to a solution of the primal problems.

# Low-dimensional scenario

- Discretization of the domain:

$$\Omega \to \Omega_h = \{(a_j, b_j)\}_{j=1}^{M}.$$

- Discretized problems:

$$\inf_{\omega \in \mathbb{R}^M} \|\omega\|_{\ell^1}, \quad \text{s.t.} \ \|A\,\omega - Y\|_{\ell^\infty} \leq \epsilon, \tag{$PD_\epsilon$}$$

$$\inf_{\omega \in \mathbb{R}^M} \|\omega\|_{\ell^1} + \frac{\lambda}{N}\|A\,\omega - Y\|_{\ell^1}, \tag{$PD_\lambda^{\text{reg}}$}$$

where $A \in \mathbb{R}^{N \times M}$ with $A_{ij} = \sigma(\langle a_j, x_i \rangle + b_j)$.

- Error estimates:

$$|\text{val}(PD_\epsilon) - \text{val}(PR_\epsilon)|, \ |\text{val}(PD_\lambda^{\text{reg}}) - \text{val}(PR_\lambda^{\text{reg}})| = \mathcal{O}(d_{\text{Hausdorff}}(\Omega, \Omega_h)).$$

- Equivalent to linear programming problems, solvable using the simplex method.
  - **Advantage**: Terminates at an extreme point of the solution set, which corresponds to a solution of the primal problems.
  - **Limitation**: Suffer from the curse of dimensionality.

# High-dimensional scenario

- Apply the SGD algorithm to the following overparameterized problem:

$$\inf_{\Theta \in (\mathbb{R} \times \Omega)^{\bar{P}}} \|\omega\|_{\ell^1} + \frac{\lambda}{N} \sum_{i=1}^{N} \ell \left( \sum_{j=1}^{\bar{P}} \omega_j \sigma(\langle a_j, x_i \rangle + b_j) - y_i \right),$$
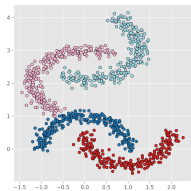
where $\bar{P}$ is large [1].

---

[1]The convergence properties of SGD for the training of overparameterized NNs have been extensively studied recently, including [Chitzat-Bach, 2018], [Zhu-Li-Song, 2019], [Bach, 2024, Chp.12], etc.
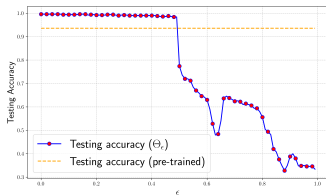
# High-dimensional scenario

- Apply the SGD algorithm to the following overparameterized problem:

$$\inf_{\Theta \in (\mathbb{R} \times \Omega)^{\bar{P}}} \|\omega\|_{\ell^1} + \frac{\lambda}{N} \sum_{i=1}^{N} \ell \left( \sum_{j=1}^{\bar{P}} \omega_j \sigma(\langle a_j, x_i \rangle + b_j) - y_i \right),$$

where $\bar{P}$ is large [1].

- Use the sparsification method developed in [L.-Zuazua, 2024] to filter the previous solution, obtaining one with fewer than $N$ activated neurons.

---

[1] The convergence properties of SGD for the training of overparameterized NNs have been extensively studied recently, including [Chitzat-Bach, 2018], [Zhu-Li-Song, 2019], [Bach, 2024, Chp.12], etc.

# High-dimensional scenario

- Apply the SGD algorithm to the following overparameterized problem:

$$\inf_{\Theta \in (\mathbb{R} \times \Omega)^{\bar{P}}} \|\omega\|_{\ell^1} + \frac{\lambda}{N} \sum_{i=1}^{N} \ell \left( \sum_{j=1}^{\bar{P}} \omega_j \sigma(\langle a_j, x_i \rangle + b_j) - y_i \right),$$

where $\bar{P}$ is large [1].

- Use the sparsification method developed in [L.-Zuazua, 2024] to filter the previous solution, obtaining one with fewer than $N$ activated neurons.

---

[1] The convergence properties of SGD for the training of overparameterized NNs have been extensively studied recently, including [Chitzat-Bach, 2018], [Zhu-Li-Song, 2019], [Bach, 2024, Chp.12], etc.

# High-dimensional scenario

- Apply the SGD algorithm to the following overparameterized problem:

$$\inf_{\Theta \in (\mathbb{R} \times \Omega)^{\bar{P}}} \|\omega\|_{\ell^1} + \frac{\lambda}{N} \sum_{i=1}^{N} \ell \left( \sum_{j=1}^{\bar{P}} \omega_j \sigma(\langle a_j, x_i \rangle + b_j) - y_i \right),$$

where $\bar{P}$ is large [1].

- Use the sparsification method developed in [L.-Zuazua, 2024] to filter the previous solution, obtaining one with fewer than $N$ activated neurons.

This approach is free from the curse of dimensionality but lacks rigorous convergence analysis.

---

[1] The convergence properties of SGD for the training of overparameterized NNs have been extensively studied recently, including [Chitzat-Bach, 2018], [Zhu-Li-Song, 2019], [Bach, 2024, Chp.12], etc.

# Table of Contents

# Classification in 2-D



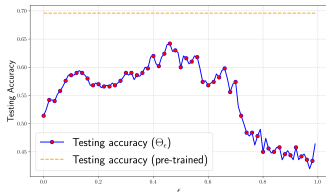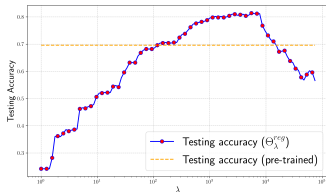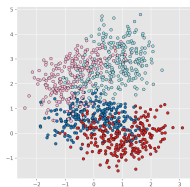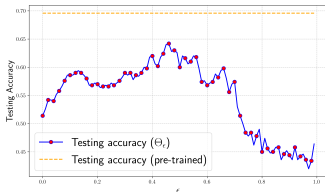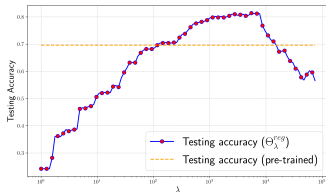(a) Datasets.

(b) Testing accuracy w.r.t. $\epsilon$.

(c) Testing accuracy w.r.t. $\lambda$.

# Classification in 2-D



(a) Datasets.

(b) Testing accuracy w.r.t. $\epsilon$.
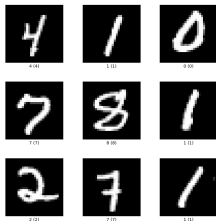
(c) Testing accuracy w.r.t. $\lambda$.

(a) Datasets.

(b) Testing accuracy w.r.t. $\epsilon$.

(c) Testing accuracy w.r.t. $\lambda$.

# Classification in 2-D



(a) Datasets.



(b) Testing accuracy w.r.t. $\epsilon$.



(c) Testing accuracy w.r.t. $\lambda$.

# Classification in 2-D



(a) Datasets.

(b) Testing accuracy w.r.t. $\epsilon$.
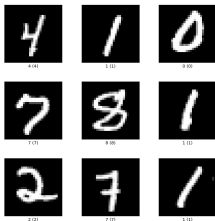
(c) Testing accuracy w.r.t. $\lambda$.

Conclusion:

- If the datasets have clear separable boundaries, consider $(P_0)$, $(P_\epsilon)$ with $\epsilon \to 0$, or $(P_\lambda^{\text{reg}})$ with $\lambda \to \infty$;
- If the datasets have heavily overlapping areas, consider the regression problem $(P_\lambda^{\text{reg}})$ with a particular range of $\lambda \sim W_1^{-1}(m_{\text{train}}, m_{\text{test}})$.
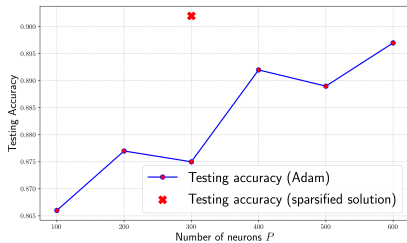
# Classification in a high-dimensional space



- The Mnist dataset, vectors in $\mathbb{R}^{28 \times 28}$.
- Training data: 300 samples of numbers 0, 1, and 2.
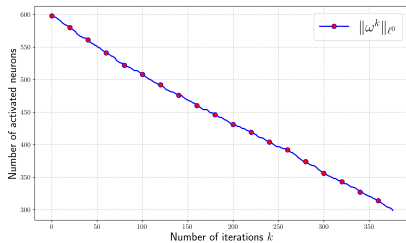- Testing data: 1000 samples of numbers 0, 1, and 2.

# Classification in a high-dimensional space



- The Mnist dataset, vectors in $\mathbb{R}^{28\times 28}$.
- Training data: 300 samples of numbers 0, 1, and 2.
- Testing data: 1000 samples of numbers 0, 1, and 2.



(a) Testing accuracy w.r.t. $P$.



(b) $\|\omega\|_{\ell^0}$ w.r.t. the iteration number.

*Thank you!*