

"Optimal" Nonparametric Regression on Low Dimensional Manifolds using Deep ReLU Neural Networks



Minshuo Chen, Haoming Jiang, Wenjing Liao, Tuo Zhao

Georgia Tech

Georgia Tech



Nonparametric Regression

Given $\mathcal{S} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ with $\mathbf{x}_i \in \mathbb{R}^D$ i.i.d. sampled from a distribution \mathcal{D} , we observe

$$y_i = f^*(\mathbf{x}_i) + \epsilon,$$

- $\epsilon_1, \dots, \epsilon_n$ are i.i.d. from $N(0, 1)$;
- $f^* \in \mathcal{F}$ with \mathcal{F} being a class of smooth functions, e.g., Hölder, Sobolev and Besov spaces.

Information-Theoretic Lower Bound:

$$\inf_{\hat{f}} \sup_{f^* \in \mathcal{F}} \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\mathbf{X}|\mathcal{S}} (\hat{f}(\mathbf{X}) - f^*(\mathbf{X}))^2 \asymp n^{-\frac{2(s+\alpha)}{2(s+\alpha)+D}},$$

- $\mathbf{X} \sim \mathcal{D}$ and \mathcal{F} is Hölder class $\mathcal{H}^{s,\alpha}$.

Theory VS Practice

Object Recognition on ImageNet:

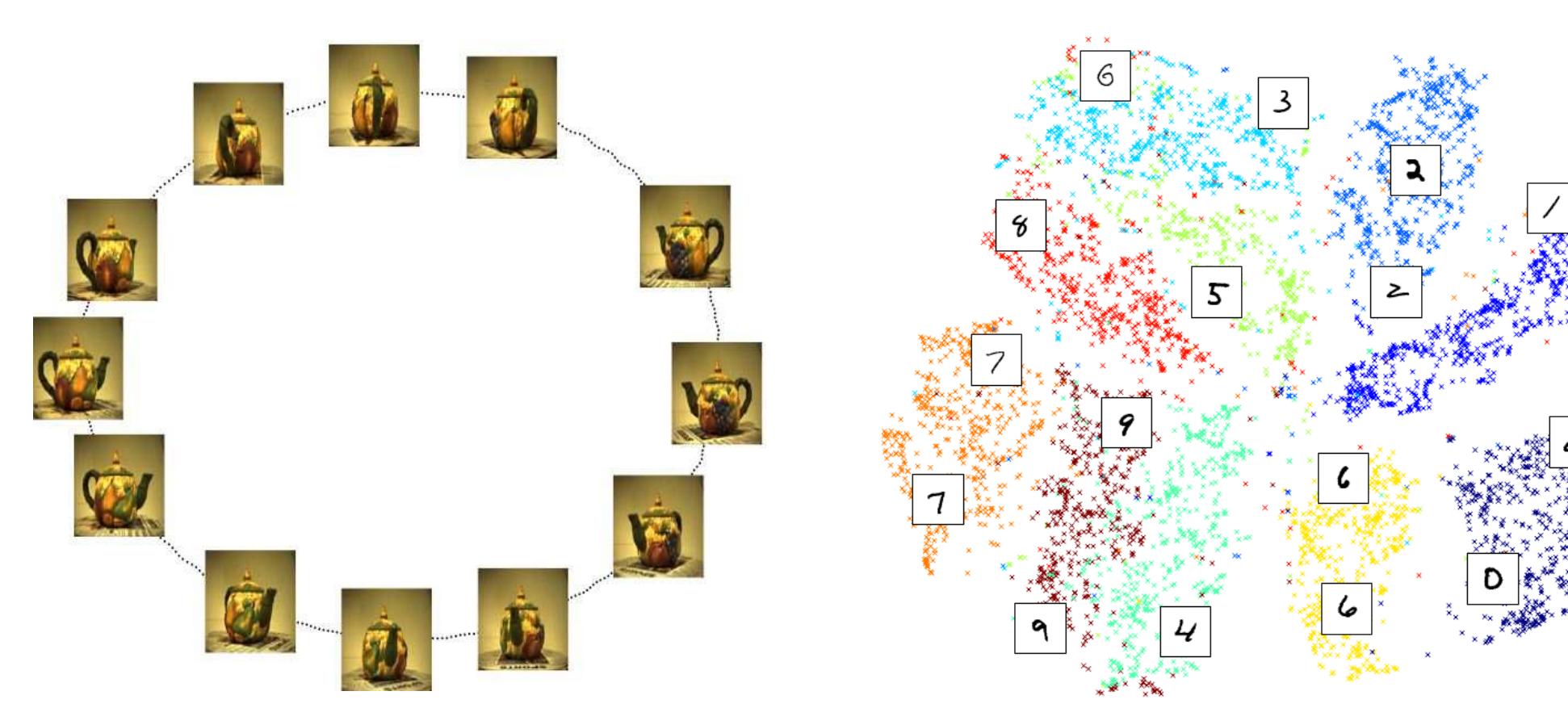
- ImageNet: 1000000 images in 1000 categories;
- Image Resolution: 384×384 ;
- Empirical Performance: Top 1 Accuracy 86.4% and Top 5 Accuracy 98.0%;
- Theoretical Bound: $n_{\text{theory}} \gtrsim \epsilon^{-D/(s+\alpha)}$;
- For moderate s and α , $n_{\text{theory}} \gg 1,000,000$.

Question:

- Why does there exist such a **huge gap** between theory and practice?

Low Dimensional Structures

▷ **Practical Motivation:** Images and acoustic signals exhibit low dimensional structures.

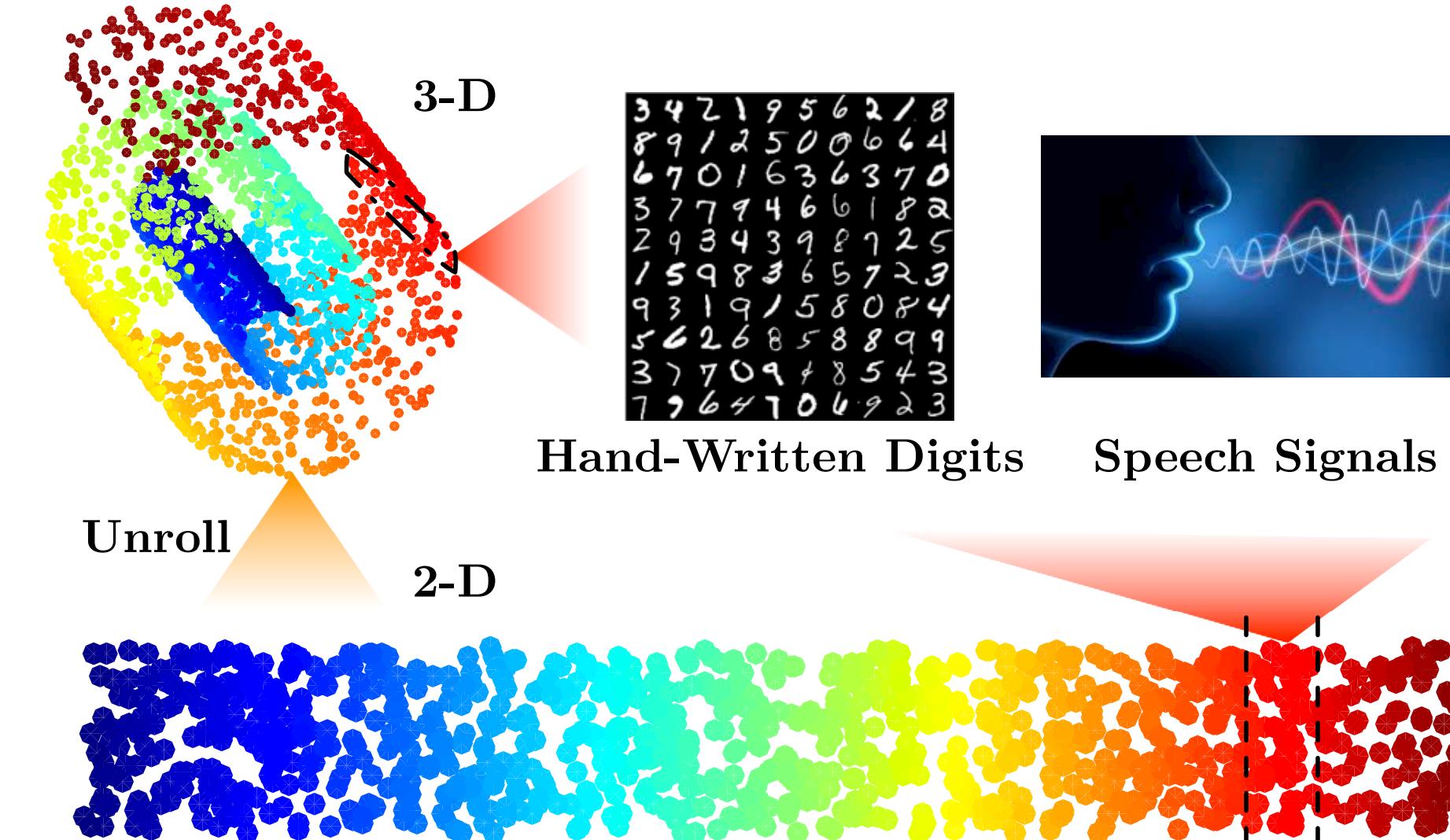


▷ **Key Observation:** Data intrinsic dimension is much smaller than the ambient dimension — making statistical estimation manageable.

Model data using a low dimensional manifold.

Regression on Low Dimensional Manifolds

▷ Low Dimensional Smooth Manifolds

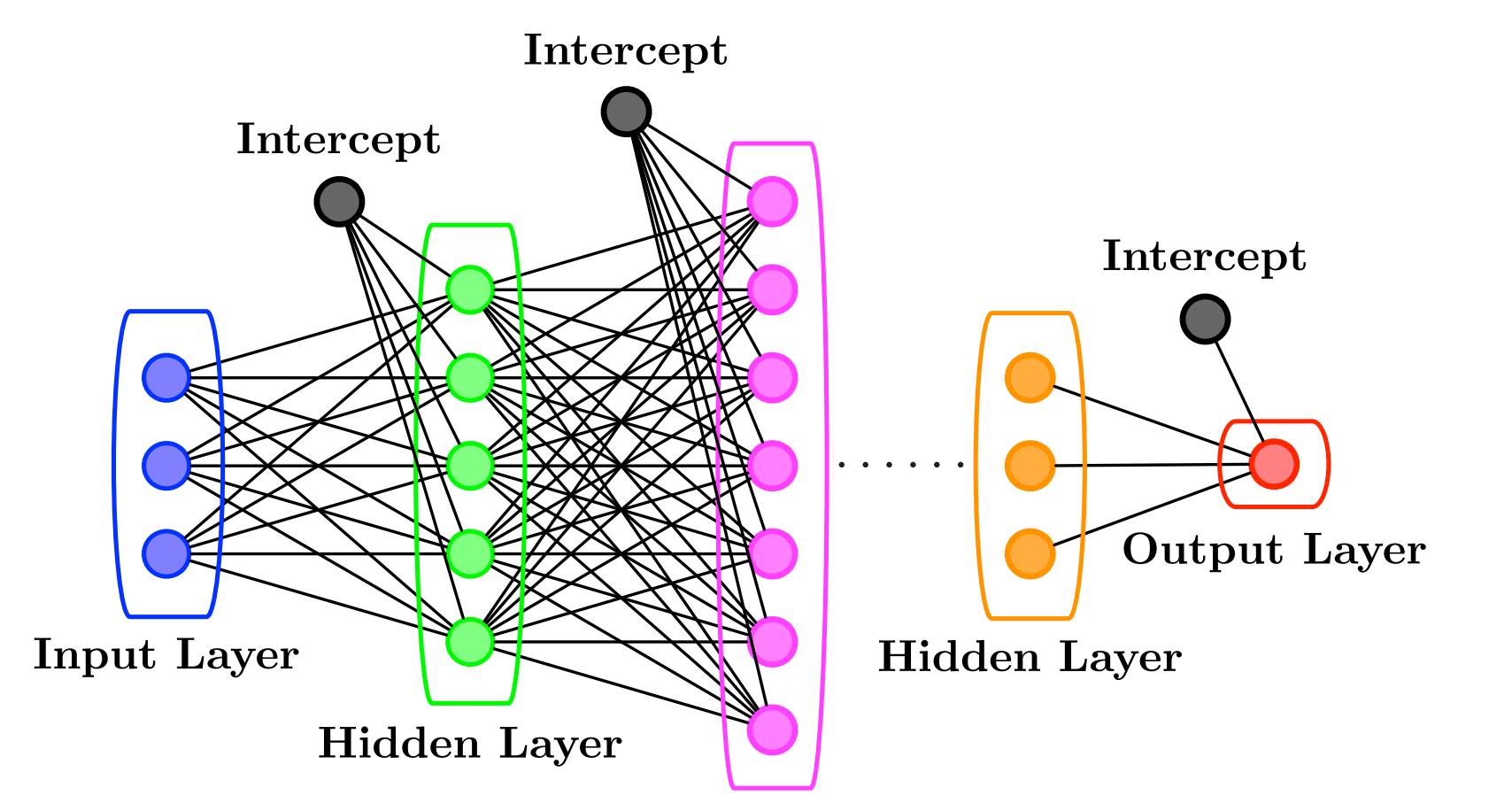


▷ Regression on Low Dimensional Manifolds

- \mathbf{x}_i sampled from a d -dimensional manifold \mathcal{M} ($d \ll D$);
- \mathcal{F} is Hölder on \mathcal{M} (distinguish from Hölder on \mathbb{R}^D).
- We learn \hat{f} by minimizing the empirical ℓ_2 risk, i.e.,

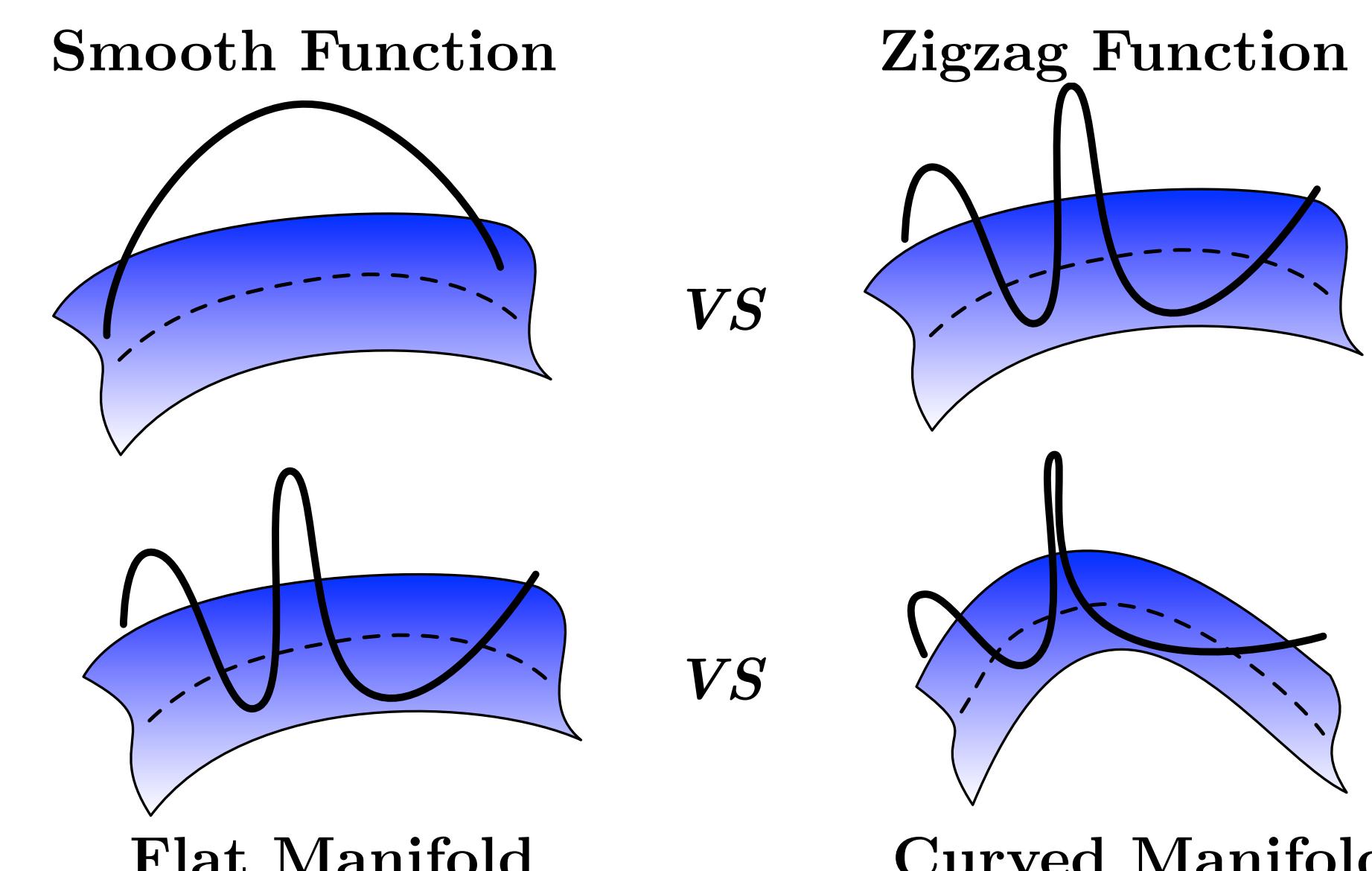
$$\hat{f} = \underset{f \in \mathcal{F}_{\text{NN}}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2, \quad (1)$$

where \mathcal{F}_{NN} denotes a class of ReLU neural networks.



$$f(\mathbf{x}) = \mathbf{W}_L \sigma(\cdots \sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2) \cdots) + \mathbf{b}_L, \quad \sigma(\cdot) = \max\{0, \cdot\}$$

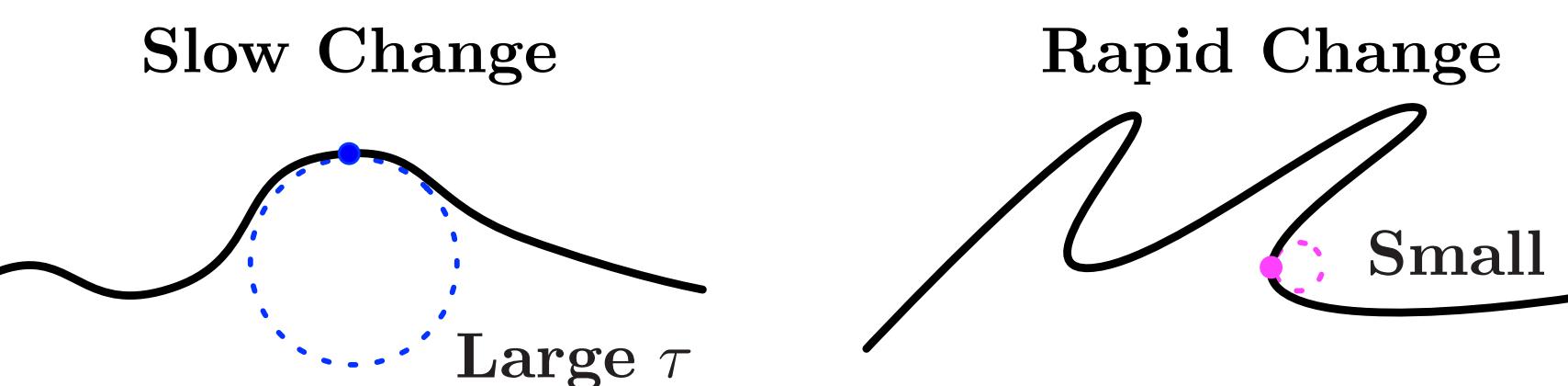
▷ Size of ReLU Neural Networks in \mathcal{F}_{NN}



Highly depends on **smoothness** of f^* , **curvature** of \mathcal{M} .

Assumptions

- \mathcal{M} is compact with reach bounded by τ .



- f^* belongs to the Hölder ball $\mathcal{H}_{\mathcal{M}}^{s,\alpha}$ of radius 1 on \mathcal{M} .

$$\left| D^s(f^* \circ \phi^{-1})|_{\phi(\mathbf{x}_1)} - D^s(f^* \circ \phi^{-1})|_{\phi(\mathbf{x}_2)} \right| \leq \|\phi(\mathbf{x}_1) - \phi(\mathbf{x}_2)\|_2^\alpha.$$

Statistical Recovery

▷ Architecture: Sparse ReLU Network

Theorem 1 (Bias Characterization). For any given $\epsilon \in (0, 1)$, there exists a ReLU network architecture with

- no more than $\mathcal{O}(\log \frac{1}{\epsilon} + \log D)$ layers;
- at most $\mathcal{O}(\epsilon^{-\frac{d}{s+\alpha}} \log \frac{1}{\epsilon} + D \log \frac{1}{\epsilon} + D \log D)$ neurons and weight parameters,

such that for any $f^* \in \mathcal{H}_{\mathcal{M}}^{s,\alpha}$, the network with properly chosen parameters yields \tilde{f} satisfying

$$\|\tilde{f} - f^*\|_\infty \leq \epsilon.$$

↓ Bias-Variance Trade-off

▷ Statistical Recovery Guarantee

Theorem 2 (Estimation Error). Suppose that \hat{f}_n minimizes the empirical risk (1), where

$\mathcal{F}_{\text{NN}} = \{f \mid f \text{ is an } L\text{-layer ReLU network with width bounded by } p, \|f\|_\infty \leq R, \|\mathbf{W}_j\|_\infty \leq \kappa, \|\mathbf{b}_j\|_\infty \leq \kappa \text{ and } \sum_{j=1}^L \|\mathbf{W}_j\|_0 + \|\mathbf{b}_j\|_0 \leq K\}.$

Let $L = \mathcal{O}(\frac{(s+\alpha) \log n}{d+2s+2\alpha})$, $p = \mathcal{O}(n^{\frac{d}{d+2s+2\alpha}})$ and $K = \mathcal{O}(Lp)$,

$$\mathbb{E}_{\mathcal{S}} \mathbb{E}_{\mathbf{X}|\mathcal{S}} (\hat{f}_n(\mathbf{X}) - f^*(\mathbf{X}))^2 = \mathcal{O}\left(n^{-\frac{2(s+\alpha)}{2s+2\alpha+d}} \log^3 n\right).$$

↓ Optimal?

▷ Minimax Lower Bound

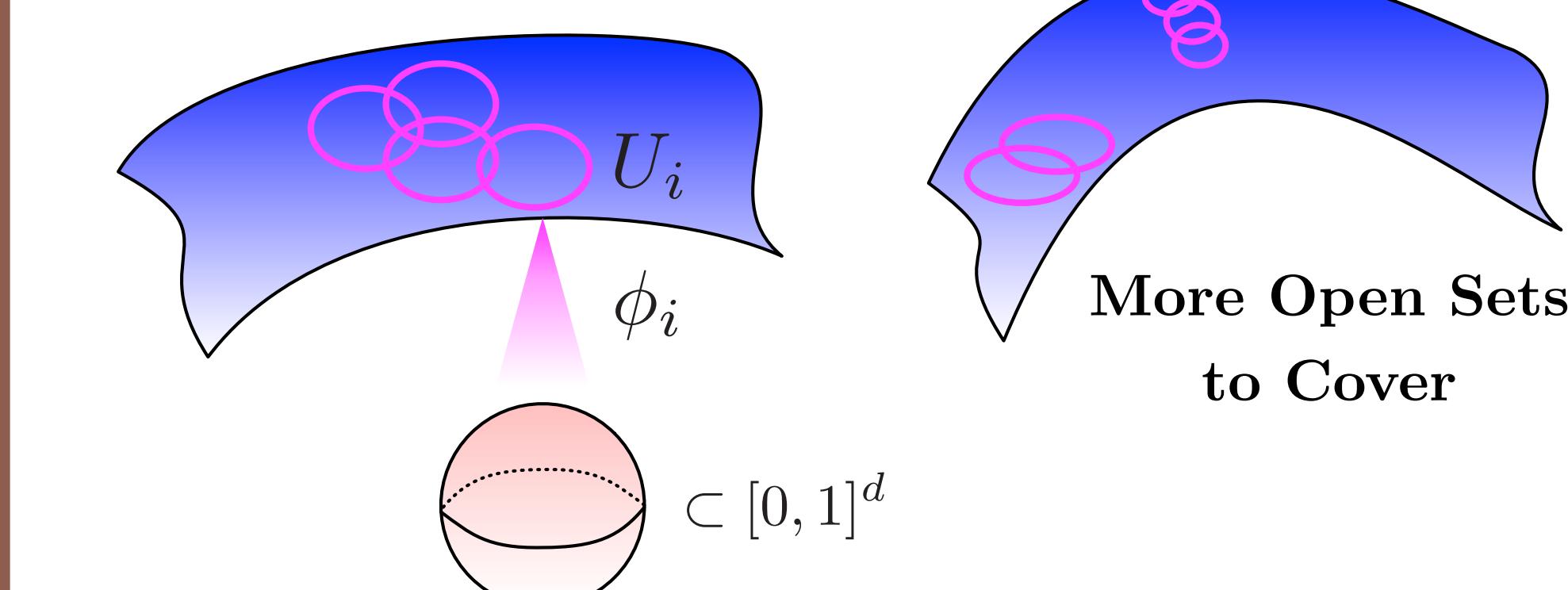
Conjecture 3 (Common Belief). Given $\mathbf{X} \sim \mathcal{D}$ supported on \mathcal{M} , we have

$$\inf_{\hat{f}} \sup_{f^* \in \mathcal{H}_{\mathcal{M}}^{s,\alpha}} \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\mathbf{X}|\mathcal{S}} (\hat{f}(\mathbf{X}) - f^*(\mathbf{X}))^2 = \Omega\left(n^{-\frac{2(s+\alpha)}{2(s+\alpha)+d}}\right).$$

Long lasting **belief** for nonparametric regression on low dimensional manifolds — no rigorous proof.

Construct Sparse ReLU Network

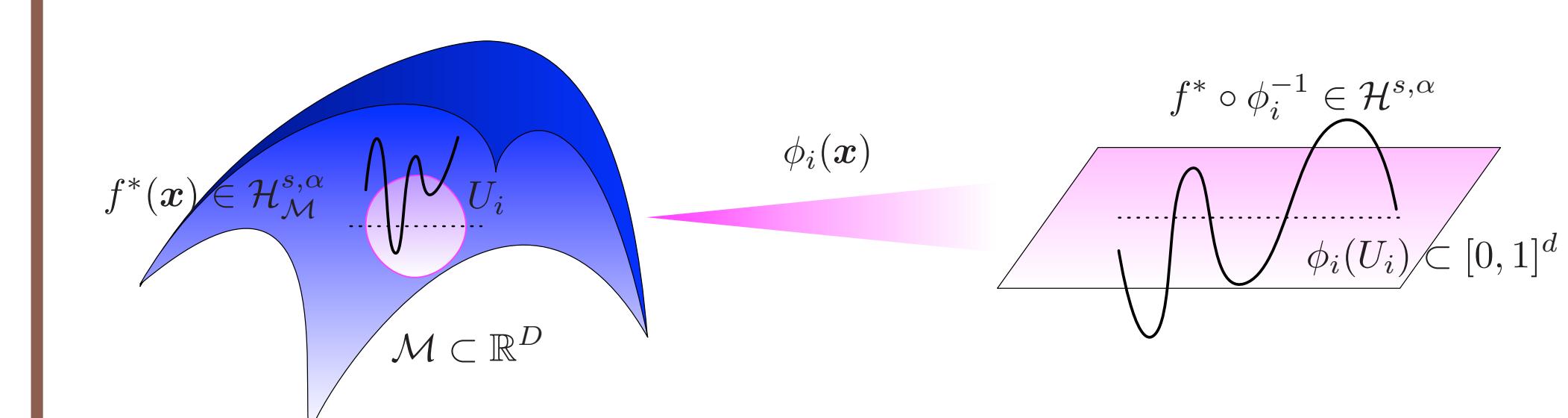
▷ Step I: Chart Construction



- $\mathcal{M} = \bigcup_{i=1}^{C_{\mathcal{M}}} U_i$ with $C_{\mathcal{M}} = \lceil \frac{\text{Vol}(\mathcal{M}) \cdot d \log d}{\tau^d} \rceil$;
- $\phi_i : \mathbb{R}^D \mapsto \mathbb{R}^d$ data transformation.

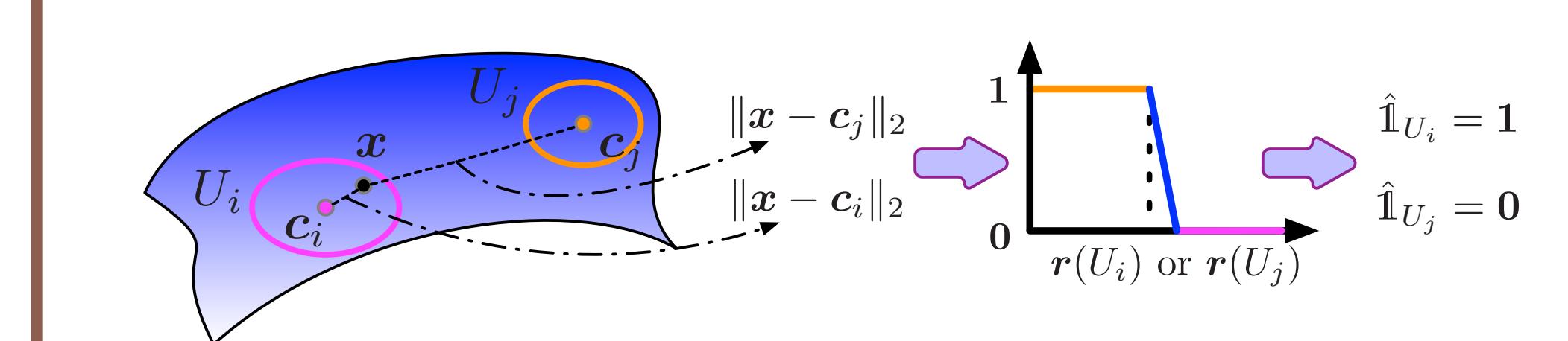
Overlapping Charts ⇒ Balance Contribution (Step III)

▷ Step II: Local Taylor Approximation



Approximation in \mathbb{R}^d ⇒ Ambient Dimension D Free

▷ Step III: Chart Determination and Pairing



Weighted sum to approximate f^* :

$$\tilde{f} = \sum_{i=1}^{C_{\mathcal{M}}} \eta_i \cdot \hat{1}_{U_i} \cdot (P_i \circ \phi_i) \approx f^*,$$

where $\sum_{i=1}^{C_{\mathcal{M}}} \eta_i(\mathbf{x}) = 1$ to balance the contribution of local Taylor approximations P_i 's.

Discussions

▷ ReLU activation vs Smooth activation

- Popular in modern applications;
- Mitigate the gradient vanishing issue.

▷ Feedforward Network vs Convolutional Network

- Conv. filters to extract the low dimensional structures.

▷ Computational Concerns

- Global optimum?
- Overparameterization?

