Semi-Supervised Deep Sobolev Regression: Estimation, Variable Selection, and Beyond

Chenguang Duan

Wuhan University Joint work with Zhao Ding, Yuling Jiao, and Jerry Zhijian Yang

July 11, 2025





Introduction and Background

- Ø Sobolev-Penalized Regression
- Semi-Supervised Deep Sobolev Regression
- Applications and Numerical Experiments
 Derivative Estimations
 Nonparametric Variable Selection
 Inverse Problems

6 Conclusions



$$Y_i = f_0(X_i) + \xi_i$$

- $\blacktriangleright X_1,\ldots,X_n\sim^{\mathrm{i.i.d.}}\mu.$
- $\blacktriangleright \xi_1,\ldots,\xi_n\sim^{\text{i.i.d.}} N(0,\sigma^2 I_d).$
- $\xi_{1:n}$ is independent of $X_{1:n}$.

The goal of the regression is to find the unknown function f_0 using labeled data $\{(X_i, Y_i)\}_{i=1}^n$.



$$Y_i = f_0(X_i) + \xi_i$$

- $\blacktriangleright X_1,\ldots,X_n\sim^{\mathrm{i.i.d.}}\mu.$
- $\blacktriangleright \xi_1,\ldots,\xi_n\sim^{\text{i.i.d.}} N(0,\sigma^2 I_d).$
- $\xi_{1:n}$ is independent of $X_{1:n}$.

The goal of the regression is to find the unknown function f_0 using labeled data $\{(X_i, Y_i)\}_{i=1}^n$.

Least-squares regression

$$f_0 = \underset{f \text{ measurable}}{\arg \min} L(f) := \mathbb{E}_{(X,Y)} [(Y - f(X))^2]$$



$$Y_i = f_0(X_i) + \xi_i$$

- $\blacktriangleright X_1,\ldots,X_n\sim^{\mathrm{i.i.d.}}\mu.$
- $\blacktriangleright \ \xi_1,\ldots,\xi_n\sim^{\text{i.i.d.}} N(0,\sigma^2 I_d).$
- $\xi_{1:n}$ is independent of $X_{1:n}$.

The goal of the regression is to find the unknown function f_0 using labeled data $\{(X_i, Y_i)\}_{i=1}^n$.

Least-squares regression

$$f_0 = \underset{f \text{ measurable}}{\arg \min} L(f) := \mathbb{E}_{(X,Y)} [(Y - f(X))^2]$$

- The expectation is intractable.
- Minimization over the measurable function class is intractable.



$$Y_i = f_0(X_i) + \xi_i$$

- $\blacktriangleright X_1,\ldots,X_n\sim^{\mathrm{i.i.d.}}\mu.$
- $\blacktriangleright \ \xi_1,\ldots,\xi_n\sim^{\text{i.i.d.}} N(0,\sigma^2 I_d).$
- $\xi_{1:n}$ is independent of $X_{1:n}$.

The goal of the regression is to find the unknown function f_0 using labeled data $\{(X_i, Y_i)\}_{i=1}^n$.

Least-squares regression

$$f_0 = \underset{f \text{ measurable}}{\arg \min} L(f) := \mathbb{E}_{(X,Y)} [(Y - f(X))^2]$$

- The expectation is intractable.
- Minimization over the measurable function class is intractable.

Monte Carlo Approximation Parameterization



$$\widehat{f}_n \in \operatorname*{arg\,min}_{f \in \mathcal{F}} \widehat{L}_n(f) := \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2$$



$$\widehat{f}_{n} \in \underset{f \in \mathcal{F}}{\arg\min} \widehat{L}_{n}(f) := \frac{1}{n} \sum_{i=1}^{n} (Y_{i} - f(X_{i}))^{2}$$
sion
$$\mathbb{E}_{S} \left[\|\widehat{f}_{n} - f_{0}\|_{2}^{2} \right] \leq \underset{f \in \mathcal{F}}{\inf \|f - f_{0}\|_{L^{2}(\Omega)}}_{Approximation \ error} + \underbrace{\mathbb{E}_{S} \left[\underset{f \in \mathcal{F}}{\sup} (L(f) - \widehat{L}_{n}(f)) \right]}_{Generalization \ error}$$

$$\underbrace{\mathsf{Approximation}_{error \ (Bias)}}_{\substack{\mathsf{Population \ risk}\\ \mininizer}} \xrightarrow{\mathsf{Population \ risk}} \operatorname{Bias-Variance \ trade-off}$$

Error decomposition

C. Duan (WHU)



Choose the hypothesis class as a **deep neural network** class.

A neural network $\phi: \mathbb{R}^d o \mathbb{R}$ is a function defined by

$$f(x) = T_L(\rho(T_{L-1}(\cdots \rho(T_0(x))\cdots))),$$
 (DNN)

where the activation function ρ is applied component-wisely and $T_{\ell}(x) = A_{\ell}x + b_{\ell}$ is an affine transformation with $A_{\ell} \in \mathbb{R}^{N_{\ell+1} \times N_{\ell}}$ and $b_{\ell} \in \mathbb{R}^{N_{\ell}}$ for $\ell = 0, ..., L$. Then the ρ -activated deep neural network class $N_{\rho}(L, S, B)$ is defined as

$$N_{\rho}(L,S,B) = \left\{ f(x) \text{ has the form of (DNN)} : \sum_{0 \le \ell \le L+1} (\|A_{\ell}\|_{0} + \|b_{\ell}\|_{2}) \le S, \, \|f\|_{L^{\infty}(\mathbb{R}^{d})} \le B \right\},$$

where L is called the depth of neural networks, S represents the number of non-zero parameters, and B is the uniform bound of neural networks.

Approximation error of DNN

► Approximation in *L^p*-norm

(Dmitry Yarotsky 2017, Zuowei Shen et al. 2019, Johannes Schmidt-Hieber 2020, Jianfeng Lu et al. 2020, Yuling Jiao et al. 2023)

- Approximation in Sobolev norms (Ingo Gühring et al. 2021, Chenguang Duan et al. 2022)
- Approximation with Lipschitz constraint (Jian Huang et al. 2022, Yuling Jiao et al. 2023, Zhao Ding et al. 2024)

Generalization error of DNN

Empirical process + Sample complexity of DNN

- ► (Local) Rademacher complexity (Peter L. Bartlett et al. 2002, 2005)
- Size-independent Rademacher complexity of DNN (Noah Golowich et al. 2018)
- VC-dimention bound for DNN (Martin Anthony et al. 1999, Peter L. Bartlett et al. 1998, 2019)



Suppose that $f_0 \in \mathcal{H}^s(\Omega)$, then the minimax optimal rate of en estimator $\widehat{f_n}$ is given as:

$$\mathbb{E}_{S}\left[\|\widehat{f}_{n}-f_{0}\|_{L^{2}(\Omega)}^{2}\right]=\mathcal{O}\left(n^{-\frac{2s}{d+2s}}\right).$$



Suppose that $f_0 \in \mathcal{H}^s(\Omega)$, then the minimax optimal rate of en estimator $\widehat{f_n}$ is given as:

$$\mathbb{E}_{S}\left[\|\widehat{f}_{n}-f_{0}\|_{L^{2}(\Omega)}^{2}\right]=\mathcal{O}\left(n^{-\frac{2s}{d+2s}}\right).$$

The convergence in L^2 -norm can **NOT** imply the convergence of gradient.

► For example,

$$\widehat{f}_n(x) = rac{\sin(nx)\log n}{n} \stackrel{L^p}{\longrightarrow} 0, \quad ext{for each } 1 \leq p \leq \infty.$$

However, $\widehat{f}'_n(x) = \cos(nx)\log n \to +\infty$.



Suppose that $f_0 \in \mathcal{H}^s(\Omega)$, then the minimax optimal rate of en estimator $\widehat{f_n}$ is given as:

$$\mathbb{E}_{S}\left[\|\widehat{f}_{n}-f_{0}\|_{L^{2}(\Omega)}^{2}\right]=\mathcal{O}\left(n^{-\frac{2s}{d+2s}}\right).$$

The convergence in L^2 -norm can **NOT** imply the convergence of gradient.

For example,

$$\widehat{f}_n(x) = \frac{\sin(nx)\log n}{n} \xrightarrow{L^p} 0$$
, for each $1 \le p \le \infty$.

However, $\widehat{f}'_n(x) = \cos(nx)\log n \to +\infty$.

How can we simultaneously estimate both the regression function f_0 and its gradient ∇f_0 ?

Applications in inverse problems, nonparametric variable selection, generative learning ...



Introduction and Background

2 Sobolev-Penalized Regression

Semi-Supervised Deep Sobolev Regression

Applications and Numerical Experiments
 Derivative Estimations
 Nonparametric Variable Selection
 Inverse Problems

6 Conclusions



Least-squares regression with gradient penalty

$$L^{\lambda}(f) := \underbrace{\mathbb{E}_{(X,Y)}\left[(Y - f(X))^2\right]}_{\text{least-squares}} + \underbrace{\lambda |f|^2_{H^1(\Omega)}}_{\text{gradient penalty}} = \|f - f_0\|^2_{L^2(\Omega)} + \sigma^2 + \lambda |f|_{H^1(\Omega)}.$$



Least-squares regression with gradient penalty

$$L^{\lambda}(f) := \underbrace{\mathbb{E}_{(X,Y)}\left[(Y - f(X))^2\right]}_{\text{least-squares}} + \underbrace{\lambda |f|_{H^1(\Omega)}^2}_{\text{gradient penalty}} = \|f - f_0\|_{L^2(\Omega)}^2 + \sigma^2 + \lambda |f|_{H^1(\Omega)}$$

Variational problem

Find $f^{\lambda} \in H^1(\Omega)$, such that $\delta L^{\lambda}(f^{\lambda},v) = 0$ for each $v \in H^1(\Omega)$, that is,

$$(f^{\lambda} - f_0, v)_{L^2(\Omega)} + \lambda (\nabla (f^{\lambda} - f_0), \nabla v)_{L^2(\Omega)} = -\lambda (\nabla f_0, \nabla v)_{L^2(\Omega)}.$$



Least-squares regression with gradient penalty

$$L^{\lambda}(f) := \underbrace{\mathbb{E}_{(X,Y)}\left[(Y - f(X))^2\right]}_{\text{least-squares}} + \underbrace{\lambda |f|_{H^1(\Omega)}^2}_{\text{gradient penalty}} = \|f - f_0\|_{L^2(\Omega)}^2 + \sigma^2 + \lambda |f|_{H^1(\Omega)}$$

Variational problem

Find $f^\lambda\in H^1(\Omega)$, such that $\delta L^\lambda(f^\lambda,v)=0$ for each $v\in H^1(\Omega)$, that is,

$$(f^{\lambda} - f_0, v)_{L^2(\Omega)} + \lambda (\nabla (f^{\lambda} - f_0), \nabla v)_{L^2(\Omega)} = -\lambda (\nabla f_0, \nabla v)_{L^2(\Omega)}.$$

Notice that f^{λ} is the solution to the elliptic equation

$$\begin{cases} -\lambda \Delta f + f = f_0, & \text{in } \Omega, \\ \nabla f \cdot \mathbf{n} = 0, & \text{on } \partial \Omega. \end{cases}$$

Substituting $v = f^{\lambda} - f_0$ into $\delta L^{\lambda}(f^{\lambda}, v) = 0$ implies $\overbrace{\|f^{\lambda} - f_0\|_{L^2(\Omega)}^2}^{\text{interior } L^2 \text{ error}} + \lambda \overbrace{\|f^{\lambda} - f_0\|_{H^1(\Omega)}^2}^{\text{interior gradient error}} = \lambda (\Delta f_0, f_0 - f^{\lambda})_{L^2(\Omega)} \le \lambda \|\Delta f_0\|_{L^2(\Omega)} \|f^{\lambda} - f_0\|_{L^2(\Omega)}.$

Substituting
$$v = f^{\lambda} - f_0$$
 into $\delta L^{\lambda}(f^{\lambda}, v) = 0$ implies
interior L^2 error
 $\|f^{\lambda} - f_0\|_{L^2(\Omega)}^2 + \lambda \quad f_0^{\lambda} - f_0|_{H^1(\Omega)}^2$
 $= \lambda (\Delta f_0, f_0 - f^{\lambda})_{L^2(\Omega)} \le \lambda \|\Delta f_0\|_{L^2(\Omega)} \|f^{\lambda} - f_0\|_{L^2(\Omega)}.$

• Interior L^2 error:

$$\|f^{\lambda} - f_0\|_{L^2(\Omega)}^2 \le \lambda^2 \|\Delta f_0\|_{L^2(\Omega)}^2.$$

Substituting
$$v = f^{\lambda} - f_0$$
 into $\delta L^{\lambda}(f^{\lambda}, v) = 0$ implies
interior L^2 error
 $\|f^{\lambda} - f_0\|_{L^2(\Omega)}^2 + \lambda$ interior gradient error
 $\|f^{\lambda} - f_0\|_{H^1(\Omega)}^2$
 $= \lambda (\Delta f_0, f_0 - f^{\lambda})_{L^2(\Omega)} \le \lambda \|\Delta f_0\|_{L^2(\Omega)} \|f^{\lambda} - f_0\|_{L^2(\Omega)}.$

• Interior L^2 error:

$$||f^{\lambda} - f_0||^2_{L^2(\Omega)} \le \lambda^2 ||\Delta f_0||^2_{L^2(\Omega)}.$$

Interior gradient error:

$$|f^{\lambda} - f_0|^2_{H^1(\Omega)} \le \lambda \|\Delta f_0\|^2_{L^2(\Omega)}.$$

Substituting
$$v = f^{\lambda} - f_0$$
 into $\delta L^{\lambda}(f^{\lambda}, v) = 0$ implies
interior L^2 error
 $\|f^{\lambda} - f_0\|_{L^2(\Omega)}^2 + \lambda \quad |f^{\lambda} - f_0|_{H^1(\Omega)}^2$
 $= \lambda (\Delta f_0, f_0 - f^{\lambda})_{L^2(\Omega)} \le \lambda \|\Delta f_0\|_{L^2(\Omega)} \|f^{\lambda} - f_0\|_{L^2(\Omega)}.$

► Interior *L*² error:

$$||f^{\lambda} - f_0||^2_{L^2(\Omega)} \le \lambda^2 ||\Delta f_0||^2_{L^2(\Omega)}.$$

Interior gradient error:

$$|f^{\lambda} - f_0|^2_{H^1(\Omega)} \le \lambda ||\Delta f_0||^2_{L^2(\Omega)}.$$

The Sobolev-penalized regressor f^{λ} is an estimator of f_0 in both L^2 -norm and H^1 -semi-norm.



- Introduction and Background
- Ø Sobolev-Penalized Regression

Semi-Supervised Deep Sobolev Regression

Applications and Numerical Experiments
 Derivative Estimations
 Nonparametric Variable Selection
 Inverse Problems

6 Conclusions



Sobolev-Penalized Risk

$$L^{\lambda}(f) := \mathbb{E}_{(X,Y)} \left[(Y - f(X))^2 \right] + \lambda |f|^2_{H^1(\Omega)}$$

- ► The expectation is intractable.
- Minimization over the measurable function class is intractable.

Monte Carlo Approximation Parameterization



Sobolev-Penalized Risk

$$L^{\lambda}(f) := \mathbb{E}_{(X,Y)}\left[(Y - f(X))^2 \right] + \lambda |f|^2_{H^1(\Omega)}$$

- ► The expectation is intractable.
- Minimization over the measurable function class is intractable.

Monte Carlo Approximation

Parameterization

Empirical Sobolev-Penalized Risk Minimization

$$\widehat{f}_{n,m}^{\lambda} \in \operatorname*{arg\,min}_{f \in \operatorname{conv}(\mathcal{F})} \widehat{L}_{n,m}^{\lambda}(f) = \frac{1}{n} \sum_{i=1}^{n} (Y_i - f(X_i))^2 + \frac{\lambda}{m} \sum_{j=1}^{m} \|\nabla f(Z_j)\|_2^2.$$

- ► Labeled data: $(X_i, Y_i) \sim^{i.i.d.} P$
- Unlabeled data: $Z_j \sim^{i.i.d.} unif(\Omega)$



Sobolev-Penalized Risk

$$L^{\lambda}(f) := \mathbb{E}_{(X,Y)}\left[(Y - f(X))^2 \right] + \lambda |f|^2_{H^1(\Omega)}$$

- ► The expectation is intractable.
- Minimization over the measurable function class is intractable.

Monte Carlo Approximation

Parameterization

Empirical Sobolev-Penalized Risk Minimization

$$\widehat{f}_{n,m}^{\lambda} \in \operatorname*{arg\,min}_{f \in \operatorname{conv}(\mathcal{F})} \widehat{L}_{n,m}^{\lambda}(f) = \frac{1}{n} \sum_{i=1}^{n} (Y_i - f(X_i))^2 + \frac{\lambda}{m} \sum_{j=1}^{m} \|\nabla f(Z_j)\|_2^2.$$

• Labeled data:
$$(X_i, Y_i) \sim^{\text{i.i.d.}} P$$

• Unlabeled data: $Z_j \sim^{i.i.d.} unif(\Omega)$

Semi-supervised framework



A1. Sub-Gaussian noise. The noise ξ is sub-Gaussian with mean 0 and finite variance proxy σ^2 .



- ▶ **A1. Sub-Gaussian noise.** The noise ξ is sub-Gaussian with mean 0 and finite variance proxy σ^2 .
- ► **A2. Bounded hypothesis.** There exists an absolute positive constant B_0 , such that $\sup_{x \in \Omega} |f_0(x)| \le B_0$. Further, functions in hypothesis class \mathcal{F} are also bounded, that is, $\sup_{x \in \Omega} |f(x)| \le B_0$.



- ▶ **A1. Sub-Gaussian noise.** The noise ξ is sub-Gaussian with mean 0 and finite variance proxy σ^2 .
- ► **A2. Bounded hypothesis.** There exists an absolute positive constant B_0 , such that $\sup_{x \in \Omega} |f_0(x)| \le B_0$. Further, functions in hypothesis class \mathcal{F} are also bounded, that is, $\sup_{x \in \Omega} |f(x)| \le B_0$.
- ► A3. Bounded derivatives of hypothesis. There exists positive constants $\{B_{1,k}\}_{k=1}^d$, such that $\sup_{x \in \Omega} |D_k f_0(x)| \le B_{1,k}$ for $1 \le k \le d$. Further, the first-order partial derivatives of functions in hypothesis class \mathcal{F} are also bounded, i.e., $\sup_{x \in \Omega} |D_k f(x)| \le B_{1,k}$ for each $1 \le k \le d$ and $f \in \mathcal{F}$. Denote by $B_1^2 := \sum_{k=1}^d B_{1,k}^2$.



- ▶ **A1. Sub-Gaussian noise.** The noise ξ is sub-Gaussian with mean 0 and finite variance proxy σ^2 .
- ► **A2. Bounded hypothesis.** There exists an absolute positive constant B_0 , such that $\sup_{x \in \Omega} |f_0(x)| \le B_0$. Further, functions in hypothesis class \mathcal{F} are also bounded, that is, $\sup_{x \in \Omega} |f(x)| \le B_0$.
- ► A3. Bounded derivatives of hypothesis. There exists positive constants $\{B_{1,k}\}_{k=1}^d$, such that $\sup_{x \in \Omega} |D_k f_0(x)| \le B_{1,k}$ for $1 \le k \le d$. Further, the first-order partial derivatives of functions in hypothesis class \mathcal{F} are also bounded, i.e., $\sup_{x \in \Omega} |D_k f(x)| \le B_{1,k}$ for each $1 \le k \le d$ and $f \in \mathcal{F}$. Denote by $B_1^2 := \sum_{k=1}^d B_{1,k}^2$.
- ► **A4. Regularity of regression function.** The regression function satisfies $\Delta f_0 \in L^2(\Omega)$ and $\nabla f_0 \cdot \mathbf{n} = 0$ a.e. on $\partial \Omega$, where n is the unit normal to the boundary.



Oracle inequality

Under A1 to A4. Suppose $n \ge \log N(B_0 \delta, \mathcal{F}, L^2(\mathcal{D}))$ and $m \ge \max_k \log N(B_{1,k} \delta, D_k \mathcal{F}, L^2(\mathcal{S}))$. Then

$$\begin{split} & \mathbb{E}\left[\|\widehat{f}_{n,m}^{\lambda} - f_{0}\|_{L^{2}(\Omega)}^{2}\right] \lesssim \beta\lambda^{2} + \varepsilon_{\mathrm{app}}(\mathcal{F},\lambda) + \varepsilon_{\mathrm{gen}}(\mathcal{F},n) + \varepsilon_{\mathrm{gen}}^{\mathrm{reg}}(\nabla\mathcal{F},m), \\ & \mathbb{E}\left[\|\nabla(\widehat{f}_{n,m}^{\lambda} - f_{0})\|_{L^{2}(\Omega)}^{2}\right] \lesssim \beta\lambda + \lambda^{-1}\varepsilon_{\mathrm{app}}(\mathcal{F},\lambda) + \lambda^{-1}\varepsilon_{\mathrm{gen}}(\mathcal{F},n) + \lambda^{-1}\varepsilon_{\mathrm{gen}}^{\mathrm{reg}}(\nabla\mathcal{F},m), \end{split}$$

where $\beta = \|\Delta f_0\|_{L^2(\Omega)}^2 + B_1^2$. The approximation error $\varepsilon_{app}(\mathcal{F}, \lambda)$, the generalization errors $\varepsilon_{gen}(\mathcal{F}, n)$ and $\varepsilon_{gen}^{reg}(\nabla \mathcal{F}, m)$ are defined, respectively, as

$$\begin{split} \varepsilon_{\mathrm{app}}(\mathcal{F},\lambda) &= \inf_{f \in \mathcal{F}} \Big\{ \|f - f_0\|_{L^2(\Omega)}^2 + \lambda \|\nabla(f - f_0)\|_{L^2(\Omega)}^2 \Big\}, \\ \varepsilon_{\mathrm{gen}}(\mathcal{F},n) &= (B_0^2 + \sigma^2)(\log n) \inf_{\delta > 0} \Big\{ \Big(\frac{2\log N(B_0\delta,\mathcal{F},L^2(\mathcal{D}))}{n}\Big)^{\frac{1}{2}} + \delta \Big\} \\ \varepsilon_{\mathrm{gen}}^{\mathrm{reg}}(\nabla\mathcal{F},m) &= B_1^2 \inf_{\delta > 0} \Big\{ \max_{1 \le k \le d} \frac{\log N(B_{1,k}\delta, D_k\mathcal{F}, L^2(\mathcal{S}))}{m} + \delta \Big\}. \end{split}$$

Approximation error

Let $\Omega \subseteq K \subseteq \mathbb{R}^d$ be two bounded domain. Set the hypothesis class as a deep ReQU neural network $\mathcal{F} = \mathcal{N}(L, W, S)$ with $L = \mathcal{O}(\log N)$ and $S = \mathcal{O}(N^d)$. Then for each $\phi \in C^s(K)$ with $s \in \mathbb{N}_{\geq 2}$, there exists $f \in \mathcal{F}$ such that

$$\|f - \phi\|_{L^{2}(\Omega)} \leq CN^{-s} \|\phi\|_{C^{s}(K)},$$

$$\|\nabla(f - \phi)\|_{L^{2}(\Omega)} \leq CN^{-(s-1)} \|\phi\|_{C^{s}(K)},$$

where C is a constant independent of N.

Approximation error

Let $\Omega \subseteq K \subseteq \mathbb{R}^d$ be two bounded domain. Set the hypothesis class as a deep ReQU neural network $\mathcal{F} = \mathcal{N}(L, W, S)$ with $L = \mathcal{O}(\log N)$ and $S = \mathcal{O}(N^d)$. Then for each $\phi \in C^s(K)$ with $s \in \mathbb{N}_{\geq 2}$, there exists $f \in \mathcal{F}$ such that

$$\|f - \phi\|_{L^{2}(\Omega)} \leq CN^{-s} \|\phi\|_{C^{s}(K)},$$

$$\|\nabla(f - \phi)\|_{L^{2}(\Omega)} \leq CN^{-(s-1)} \|\phi\|_{C^{s}(K)},$$

where C is a constant independent of N.

Generalization error

Suppose the activation function is piecewise-polynomial. Let $\mathcal{D} = \{X_i\}_{i=1}^n$ and $\mathcal{S} = \{Z_j\}_{j=1}^m$. Then

$$\log N(\delta, N(L, S, B), L^{2}(\mathcal{D})) \lesssim LS \log(S) \log\left(\frac{nB}{\delta}\right),$$
$$\log N(\delta, D_{k}N(L, S, B), L^{2}(\mathcal{S})) \lesssim L^{2}S \log(S) \log\left(\frac{mB}{\delta}\right)$$



Under **A1** to **A4**. Let $\Omega \subseteq K \subseteq \mathbb{R}^d$ be two bounded domain. Assume that $f_0 \in C^s(K)$ with $s \in \mathbb{N}_{\geq 2}$. Set the hypothesis class as a deep ReQU neural network class $\mathcal{F} = N(L, W, S)$ with $L = \mathcal{O}(\log n)$ and $S = \mathcal{O}(n^{\frac{d}{d+4s}})$. Let $\lambda = \mathcal{O}(n^{-\frac{s}{d+4s}}\log^2 n)$. Then

$$\mathbb{E}\left[\|\widehat{f}_{n,m}^{\lambda} - f_{0}\|_{L^{2}(\Omega)}^{2}\right] \leq \mathcal{O}\left(n^{-\frac{2s}{d+4s}}\log^{4}n\right) + \mathcal{O}\left(n^{\frac{d}{d+4s}}\log^{4}nm^{-1}\right), \\
\mathbb{E}\left[\|\nabla(\widehat{f}_{n,m}^{\lambda} - f_{0})\|_{L^{2}(\Omega)}^{2}\right] \leq \mathcal{O}\left(n^{-\frac{s}{d+4s}}\log^{2}n\right) + \mathcal{O}\left(n^{\frac{d+s}{d+4s}}\log^{2}nm^{-1}\right).$$



Under **A1** to **A4**. Let $\Omega \subseteq K \subseteq \mathbb{R}^d$ be two bounded domain. Assume that $f_0 \in C^s(K)$ with $s \in \mathbb{N}_{\geq 2}$. Set the hypothesis class as a deep ReQU neural network class $\mathcal{F} = N(L, W, S)$ with $L = \mathcal{O}(\log n)$ and $S = \mathcal{O}(n^{\frac{d}{d+4s}})$. Let $\lambda = \mathcal{O}(n^{-\frac{s}{d+4s}}\log^2 n)$. Then

$$\mathbb{E}\left[\|\widehat{f}_{n,m}^{\lambda} - f_0\|_{L^2(\Omega)}^2\right] \leq \mathcal{O}\left(n^{-\frac{2s}{d+4s}}\log^4 n\right) + \mathcal{O}\left(n^{\frac{d}{d+4s}}\log^4 nm^{-1}\right),$$
$$\mathbb{E}\left[\|\nabla(\widehat{f}_{n,m}^{\lambda} - f_0)\|_{L^2(\Omega)}^2\right] \leq \mathcal{O}\left(n^{-\frac{s}{d+4s}}\log^2 n\right) + \mathcal{O}\left(n^{\frac{d+4s}{d+4s}}\log^2 nm^{-1}\right).$$

Simultaneously estimate both the regression function f_0 and its gradient ∇f_0 .



- Introduction and Background
- Sobolev-Penalized Regression
- Semi-Supervised Deep Sobolev Regression
- Applications and Numerical Experiments
 Derivative Estimations
 Nonparametric Variable Selection
 Inverse Problems

6 Conclusions





$$f_0(x) = 1 + 36x^2 - 59x^3 + 21x^5 + 0.5\cos(\pi x)$$

July 11, 2025



► A5. Sparsity structure

There exists $f_0^*: \mathbb{R}^{d^*} \to \mathbb{R}$ $(1 \le d^* < d)$ such that for each $x := (x_1, \dots, x_d) \in \mathbb{R}^d$,

$$f_0(x_1,\ldots,x_d) = f_0^*(x_{j_1},\ldots,x_{j_{d^*}}), \quad \{j_1,\ldots,j_{d^*}\} \subseteq [d].$$



A5. Sparsity structure

There exists $f_0^* : \mathbb{R}^{d^*} \to \mathbb{R}$ ($1 \le d^* < d$) such that for each $x := (x_1, \dots, x_d) \in \mathbb{R}^d$,

$$f_0(x_1,\ldots,x_d) = f_0^*(x_{j_1},\ldots,x_{j_{d^*}}), \quad \{j_1,\ldots,j_{d^*}\} \subseteq [d].$$

The derivatives can indicate whether a variable is relevant to the output.



► A5. Sparsity structure

There exists $f_0^* : \mathbb{R}^{d^*} \to \mathbb{R}$ $(1 \le d^* < d)$ such that for each $x := (x_1, \dots, x_d) \in \mathbb{R}^d$,

$$f_0(x_1,\ldots,x_d) = f_0^*(x_{j_1},\ldots,x_{j_{d^*}}), \quad \{j_1,\ldots,j_{d^*}\} \subseteq [d].$$

The derivatives can indicate whether a variable is relevant to the output.

Relevant variable

• A variable $k \in [d]$ is irrelevant for the function f with respect to Lebesgue measure on Ω , if

 $D_k f(X) = 0$ almost surely,

and relevant otherwise. The set of relevant variables is defined as

 $\mathcal{I}(f) = \{k \in [d] : \|D_k f\|_{L^2(\Omega)} > 0\}.$

	Is f_0 sparse?	Convergence rate
$\mathbb{E}[\ \widehat{f}_{n,m}^{\lambda} - f_0\ _{L^2(\Omega)}^2]$	X	$\widetilde{\mathcal{O}}(n^{-rac{2s}{d+4s}})$
$\mathbb{E}[\ \widehat{f}_{n,m}^{\lambda} - f_0\ _{L^2(\Omega)}^2]$	1	$\widetilde{\mathcal{O}}(n^{-rac{2s}{d^*+4s}})$
$\mathbb{E}[\ \nabla(\widehat{f}_{n,m}^{\lambda}-f_0)\ _{L^2(\Omega)}^2]$	×	$\widetilde{\mathcal{O}}(n^{-rac{s}{d+4s}})$
$\mathbb{E}[\ \nabla(\widehat{f}_{n,m}^{\lambda}-f_0)\ _{L^2(\Omega)}^2]$	\checkmark	$\widetilde{\mathcal{O}}(n^{-rac{2s}{d^*+4s}})$

	Is f_0 sparse?	Convergence rate
$\mathbb{E}[\ \widehat{f}_{n,m}^{\lambda} - f_0\ _{L^2(\Omega)}^2]$	X	$\widetilde{\mathcal{O}}(n^{-\frac{2s}{d+4s}})$
$\mathbb{E}[\ \widehat{f}_{n,m}^{\lambda} - f_0\ _{L^2(\Omega)}^2]$	1	$\widetilde{\mathcal{O}}(n^{-rac{2s}{d^*+4s}})$
$\mathbb{E}[\ \nabla(\widehat{f}_{n,m}^{\lambda}-f_0)\ _{L^2(\Omega)}^2]$	×	$\widetilde{\mathcal{O}}(n^{-rac{s}{d+4s}})$
$\mathbb{E}[\ \nabla(\widehat{f}_{n,m}^{\lambda}-f_0)\ _{L^2(\Omega)}^2]$	1	$\widetilde{\mathcal{O}}(n^{-rac{2s}{d^*+4s}})$

Selection consistency

Under A1 to A5. It follows that

$$\lim_{n\to\infty} \Pr\left\{\mathcal{I}(f_0) = \mathcal{I}(\widehat{f}_{n,m}^{\lambda})\right\} = 1,$$

where $\lambda = O(n^{-\frac{s}{d^*+4s}} \log^2 n)$, and *m* is sufficiently large.



sparse structure





Elliptic equation with unknown source

$$\begin{cases} -\nabla \cdot (a(x)\nabla u(x)) + c(x)u = f(x), & \text{in } \Omega, \\ \nabla u \cdot \mathbf{n} = 0, & \text{on } \partial \Omega. \end{cases}$$

Measurement model

$$Y = \mathcal{S}(f^{\dagger})(X) + \xi.$$

• $u^{\dagger} = S(f^{\dagger})$ is the solution to the elliptic equation.

• $X \sim unif(\Omega)$, and $\xi \sim subG(\sigma^2)$ is the random noise independent of X.

Recovering Procedure

Sobolev-penalized regression using interior measurements

$$\widehat{u}_{n,m}^{\lambda} \in \operatorname*{arg\,min}_{u \in \operatorname{conv}(\mathcal{U})} \widehat{L}_{n,m}^{\lambda}(u) = \frac{1}{n} \sum_{i=1}^{n} (Y_i - u(X_i))^2 + \frac{\lambda}{m} \sum_{j=1}^{m} \|\nabla u(Z_j)\|_2^2.$$

► Interior position-measurement pairs: $\{(X_i, Y_i)\}_{i=1}^n$. expensive ► Positions variables $Z_1, \ldots, Z_m \sim^{i.i.d.} unif(\Omega)$ very cheap

C. Duan (WHU)

Recovering Procedure

Sobolev-penalized regression using interior measurements

$$\widehat{u}_{n,m}^{\lambda} \in \operatorname*{arg\,min}_{u \in \operatorname{conv}(\mathcal{U})} \widehat{L}_{n,m}^{\lambda}(u) = \frac{1}{n} \sum_{i=1}^{n} (Y_i - u(X_i))^2 + \frac{\lambda}{m} \sum_{j=1}^{m} \|\nabla u(Z_j)\|_2^2.$$

- Interior position-measurement pairs: {(X_i, Y_i)}ⁿ_{i=1}.
 Positions variables Z₁,..., Z_m ~^{i.i.d.} unif(Ω)
 very cheap
- Recovering unknown source using gradient estimator Find $\widehat{f}_{n,m}^{\lambda}$, such that for each $v \in H^1(\Omega)$,

$$(a(x)\nabla \widehat{u}_{n,m}^{\lambda}, \nabla v)_{L^{2}(\Omega)} + (c(x)\widehat{u}_{n,m}^{\lambda}, v)_{L^{2}(\Omega)} = (\widehat{f}_{n,m}^{\lambda}, v)_{L^{2}(\Omega)}$$

Since $\widehat{u}_{n,m}^{\lambda} \in H^2(\Omega)$, we find

$$\widehat{f}_{n,m}^{\lambda} = -\nabla \cdot (a(x)\nabla \widehat{u}_{n,m}^{\lambda}) + c(x)\nabla \widehat{u}_{n,m}^{\lambda} \in L^{2}(\Omega).$$

Convergence rate in weak norm

$$\mathbb{E}\Big[\|\widehat{f}_{n,m}^{\lambda} - f^{\dagger}\|_{(H^{1}(\Omega))^{*}}\Big] \lesssim \mathcal{O}\Big(n^{-\frac{s}{2(d+4s)}}\Big).$$





- Introduction and Background
- Ø Sobolev-Penalized Regression
- Semi-Supervised Deep Sobolev Regression
- Applications and Numerical Experiments
 Derivative Estimations
 Nonparametric Variable Selection
 Inverse Problems

6 Conclusions



- Simultaneously estimations of both the regression function and its gradient.
- Nonasymptotic convergence rate.
- ► Applications in nonparametric variable selection and inverse problems.

<u>Reference:</u> Zhao Ding, Chenguang Duan, Yuling Jiao, and Jerry Zhijian Yang. Semi-Supervised Deep Sobolev Regression: Estimation and Variable Selection by ReQU Neural Network. *IEEE Transactions on Information Theory*, 2025.

Thanks for your attention!

```
Homepage:
https://chenguangduan.github.io/
Google Scholar:
https://scholar.google.com/citations?user=RpmGgyMAAAAJ
Email:cgduan.math@gmail.com
```

