

Deep Nonlinear Sufficient Dimension Reduction

Zhou Yu

School of Statistics
East China Normal University

November 6, 2025



Outline

- 1 Introduction
- 2 Method
- 3 Theoretical Analysis
- 4 Simulation Studies
- 5 Real data
- 6 Summary



Outline

- 1 Introduction
- 2 Method
- 3 Theoretical Analysis
- 4 Simulation Studies
- 5 Real data
- 6 Summary



Review of Sufficient Dimension Reduction

Ker-Chau Li (1991) introduced the following model:

$$Y = h(\beta_1^\top X, \beta_2^\top X, \dots, \beta_d^\top X, \varepsilon). \quad (1)$$

- Y is an univariate output variable;
- The dimension of X is p ;
- The random error ε is independent of X ;
- The space \mathcal{B} generated by β_1, \dots, β_d is called the sufficient dimension reduction subspace.
- Model (1) is equivalent to:
 - 1 The conditional distribution of Y given X depends on X only through the d dimensional variable $\beta_1^\top X, \dots, \beta_d^\top X$.
 - 2 Conditioning on $\beta_1^\top X, \dots, \beta_d^\top X$, Y and X are independent.



Review of Sufficient Dimension Reduction

(1) can be characterized via conditional independence

$$Y \perp\!\!\!\perp X | \beta_1^\top X, \beta_2^\top X, \dots, \beta_d^\top X, \quad (2)$$

which is defined as the linear sufficient dimension reduction (SDR).

- Approaching the SDR problem as an inverse regression problem is the most mainstream. But it requires making assumptions on the probability distribution of X , which is difficult to justify.
- Rather than the inverse regression framework, [Fukumizu et al., 2009] proposed the kernel dimension reduction, which involves the use of conditional covariance operators on reproducing kernel Hilbert spaces (RKHSs) and does not impose strong assumptions on the probability distribution of X .



Review of Sufficient Dimension Reduction

Nonlinear sufficient dimension reduction can be modeled as

$$Y \perp\!\!\!\perp X | f_1(X), f_2(X), \dots, f_d(X) \quad (3)$$

- From the perspective of RKHS \mathcal{H}_κ generated by the kernel $\kappa(\cdot, X)$, the reproducing property can reformulate $f(X)$ as $\langle f, \kappa(\cdot, X) \rangle_{\mathcal{H}_\kappa}$, which is essentially a direct generalization of the inner product $\beta^\top X$ in (2).
- (3) is equivalent to:
 - The conditional distribution of Y given X depends on X only through the d nonlinear functions $f_1(X), \dots, f_d(X)$.
 - Conditioning on $f_1(X), \dots, f_d(X)$, Y and X are independent.



Review of Sufficient Dimension Reduction

- Consider the model

$$Y = (X_1 + X_2)^2 \log(X_1^2 + X_2^2 + 0.001) + \varepsilon$$

- For linear SDR, $d = 2$, $\beta_1 = (1, 0, \dots, 0)^\top$ and $\beta_2 = (0, 1, 0, \dots, 0)^\top$. (β_1, β_2) is not identifiable.
- The concern is $\text{Span}(\beta_1, \beta_2)$ rather than (β_1, β_2) itself.
- For nonlinear SDR, $d = 1$, $f_0^{(1)}(X) = (X_1 + X_2)^2 \log(X_1^2 + X_2^2 + 0.001)$. $f_0^{(1)}(X)$ is not identifiable.
- The concern of nonlinear SDR is the σ -field generated by $f_0^{(1)}(X)$.



Review of Sufficient Dimension Reduction

Further, [Lee et al., 2013] developed a general theory for SDR at the level of σ -fields. A sub σ -field \mathcal{G} of $\sigma(X)$ is called the sufficient dimension reduction σ -field if

$$Y \perp\!\!\!\perp X | \mathcal{G}. \quad (4)$$

The minimal sufficient dimension reduction σ -field is called central σ -field, denoted by $\mathcal{G}_{Y|X}$, and the set of all $\mathcal{G}_{Y|X}$ -measurable, square-integrable functions is referred to as the central class $\mathfrak{M}_{Y|X}$.

- The definition in (2) and (3) can be incorporated into this framework by regarding \mathcal{G} as $\sigma(\beta^\top X)$ and $\sigma\{f(X)\}$ respectively.
- GSIR (Generalized Sliced Inverse Regression) uses $\overline{\text{ran}}(E_{X|Y}^* E_{X|Y})$ to recover the central class. The $E_{X|Y}$ is a conditional expectation operator from $L_2(P_X)$ to $L_2(P_Y)$.



- A common strategy to achieve nonlinear SDR is to combine the kernel trick with the existing linear SDR methods. The drawback of these methods is that they require the computation of the eigenvectors or inverse of a $n \times n$ matrix.
- As the field of machine learning flourishes, many scholars have embarked on the application of machine learning to other statistical fields.
- The nonlinear representation of the neural network function class is more intuitive and powerful compared to that of RKHS.
- The dependence measures of variables play a crucial role in variable screening and SDR.



Outline

- 1 Introduction
- 2 Method**
- 3 Theoretical Analysis
- 4 Simulation Studies
- 5 Real data
- 6 Summary



Nonlinear SDR at the Level of σ -fields

Let $X \in \mathbb{R}^p$ and $Y \in \mathbb{R}^q$ be two random vectors. Consider the following nonlinear SDR problem

$$Y \perp\!\!\!\perp X | \mathbf{f}_0(X), \quad (5)$$

where $\mathbf{f}_0(\cdot)$ is a \mathbb{R}^d -valued function ($d \leq p$).

Assumption 1

Assume that the family of probability measures $\{P_{X|Y}(\cdot|y) : y \in \Omega\}$ is dominated by a σ -finite measure. Then there exists a unique minimal sufficient σ -field (or central σ -field) $\mathcal{G}_{Y|X}$ such that

$$Y \perp\!\!\!\perp X | \mathcal{G}_{Y|X}.$$

Moreover, \mathbf{f}_0 in (5) satisfies $E\{\mathbf{f}_0(X)\} = 0$, $\text{Var}\{\mathbf{f}_0(X)\} = I_d$ and

$$\sigma\{\mathbf{f}_0(X)\} = \mathcal{G}_{Y|X}.$$

Another Look of Nonlinear SDR

Let $\mathfrak{M}_{\sigma\{f_0(X)\}}$ be the central class corresponding to the central σ -field $\sigma\{f_0(X)\}$. And let $L_2(P_X)$ be the function class of all square-integrable functions with respect to the measure P_X . It can be expressed as the orthogonal direct sum

$$L_2(P_X) = \mathfrak{M}_{\sigma\{f_0(X)\}} \oplus \mathfrak{M}_{\sigma\{f_0(X)\}}^\perp,$$

where $\mathfrak{M}_{\sigma\{f_0(X)\}}^\perp$ is orthogonal complement of $\mathfrak{M}_{\sigma\{f_0(X)\}}$. Without loss of generality, consider $g \in \mathfrak{M}_{\sigma\{f_0(X)\}}^\perp$ and $E(g(X)) = 0$, it holds that

$$\begin{aligned} E\{g(X)|Y\} &= E[E\{g(X)|f_0(X), Y\}|Y] \\ &= E[E\{g(X)|f_0(X)\}|Y] \\ &= 0 \\ &\stackrel{a.s.}{=} E\{g(X)\}. \end{aligned}$$



Generalized Martingale Difference Divergence

- μ is a Lévy measure on $\mathbb{R}^p \setminus \{0\}$ satisfying $\int_{\mathbb{R}^p \setminus \{0\}} (1 \wedge \|s\|^2) \mu(ds) < \infty$.
- Based on a Lévy measure μ , GMDD of $U \in \mathbb{R}^q$ given $V \in \mathbb{R}^p$ is defined as

$$\text{GMDD}(U|V) = \int_{\mathbb{R}^p} \|g_{U,V}(s) - g_U g_V(s)\|^2 \mu(ds),$$

where $g_{U,V}(s) = EUe^{is^\top V}$, $g_U = EU$, $g_V(s) = Ee^{is^\top V}$.

- $\text{GMDD}(U|V) \geq 0$. Moreover,

$$\text{GMDD}(U|V) = 0 \iff E(U|V) \stackrel{a.s.}{=} E(U)$$

- Let $\Theta(x) = \int \{1 - \cos(s^\top x)\} \mu(ds)$ and $\kappa(x_1, x_2) = \Theta(x_1 - x_2)$.

$$\text{GMDD}(U|V) = -E \left\{ (U - EU)^\top (\tilde{U} - E\tilde{U}) \kappa(V, \tilde{V}) \right\},$$

where (\tilde{U}, \tilde{V}) is an independent copy of (U, V) .



Generalized Martingale Difference Divergence

Table: Some real-valued cndfs on \mathbb{R}^p and the corresponding Lévy measures.

cndf $\Theta(x)$	Lévy measure $\mu(ds)$
$\int \left(1 - \cos(s^\top x)\right) \mu(ds)$	μ finite measure
$\frac{ x _p^2}{\lambda^2 + x _p^2}, \lambda > 0, x \in \mathbb{R}^p$	$2\lambda^{-1} e^{-\lambda s } ds$
$\frac{ x _p^2}{\lambda^2 + x _p^2}, \lambda > 0, x \in \mathbb{R}^p$	$\int_0^\infty e^{\lambda^2 \nu - \frac{ s ^2}{4\nu}} \frac{\lambda^2}{(4\pi\nu)^{p/2}} d\nu ds$
$1 - e^{-\frac{1}{2} x _p^2}$	$(2\pi)^{-p/2} e^{-\frac{1}{2} s ^2} ds$
$ x _p$	$\frac{\Gamma(\frac{1+p}{2})}{\pi^{\frac{1+p}{2}}} \frac{ds}{ s ^{p+1}}$
$\frac{1}{2} x _p^2$	no Lévy measure
$ x _p^\alpha, \alpha \in (0, 2)$	$\frac{\alpha 2^{\alpha-1} \Gamma(\frac{\alpha+p}{2})}{\pi^{p/2} \Gamma(1-\frac{\alpha}{2})} \frac{ds}{ s ^{p+\alpha}}$
$\ln\left(1 + \frac{x^2}{2}\right), x \in \mathbb{R}$	$\frac{1}{ s } e^{-\sqrt{\frac{1}{2}} s } ds$
$\ln\left(1 + \frac{ x _p^2}{2}\right), x \in \mathbb{R}^p$	$\int_0^\infty \int_0^1 \frac{1}{\lambda} (2\pi\lambda\rho)^{-p/2} e^{- s ^2/(2\lambda\rho)} e^{-\rho} d\lambda dp ds$
$\ln \cosh(x), x \in \mathbb{R}$	$\frac{ds}{2s \sinh(\pi s/2)}$
$\ln \cosh(x _p), x \in \mathbb{R}^p$	$\frac{\Gamma(\frac{1+p}{2})}{\pi^{\frac{1+p}{2}}} \frac{ds}{ s ^{p+1}} + \iint_0^1 \int_0^\infty \frac{1}{\lambda} e^{\frac{\lambda\rho}{2}} e^{-\rho} \sum_{n=1}^\infty \frac{(-\frac{\lambda\rho}{2})^n}{n!(2\pi)^p} e^{-2n r } e^{ir^\top s} dp d\lambda dr ds$



Lemma 1

Let g be any square-integrable function that lies in the orthogonal complement space of $\mathfrak{M}_{\sigma\{f_0(X)\}}$, then $\text{GMDD}(g(X)|Y) = 0$.

For any $f \in L_2(P_X)$, the property of GMDD tells

$$-E\left([f(X) - E\{f(X)\}][f(\tilde{X}) - E\{f(\tilde{X})\}]\kappa(Y, \tilde{Y})\right) \geq 0.$$

Inspired by this relation and Lemma 1, we speculate that the function from $\mathfrak{M}_{\sigma\{f_0(X)\}}$ may maximize the GMDD index. The following theorem gives an affirmative answer.



Theorem 2

Provided assumption 1 holds, let $m \in \mathbb{N}^+$ be any positive integer and $\mathbf{f}^* = (f_1^*, \dots, f_m^*)$ be an optimal solution of the following objective function

$$\begin{aligned} \max_{\mathbf{f} \in \{L_2(P_X)\}^m} & -E\left([\mathbf{f}(X) - E\{\mathbf{f}(X)\}]^\top [\mathbf{f}(\tilde{X}) - E\{\mathbf{f}(\tilde{X})\}]\kappa(Y, \tilde{Y})\right) \\ \text{subject to} & \quad \text{Var}\{\mathbf{f}(X)\} = I_m. \end{aligned} \quad (6)$$

Then $f_j^* \in \mathfrak{M}_{\sigma\{f_0(X)\}}$ for all $1 \leq j \leq m$.

Note that the dimensionality of $\mathfrak{M}_{\sigma\{f_0(X)\}}$ is infinite. However, $\mathbf{f}_0 : \mathbb{R}^p \mapsto \mathbb{R}^d$ in (5) indicates the most suitable choice $m = d$. Theorem 2 establishes the unbiasedness of nonlinear SDR method (6) since each component of any optimal solution \mathbf{f}^* is $\sigma\{f_0(X)\}$ measurable. In other words, $\sigma\{\mathbf{f}^*(X)\} \subseteq \sigma\{f_0(X)\}$.



The uniform framework of linear and nonlinear SDR

- Consider the linear SDR case of $\mathbf{f}(X) = \beta^\top X$, then (6) becomes

$$\begin{aligned} \max_{\beta \in \mathbb{R}^{p \times d}} & -\text{tr} \left(\beta^\top E \left[\{X - E(X)\} \{\tilde{X} - E(\tilde{X})\}^\top \kappa(Y, \tilde{Y}) \right] \beta \right) \\ \text{subject to} & \quad \beta^\top \Sigma \beta = I_d, \end{aligned}$$

where $\Sigma = \text{Var}(X)$.

- Classical Sliced Inverse Regression can be reformulated similarly as

$$\begin{aligned} \max_{\beta \in \mathbb{R}^{p \times d}} & \text{tr} \left[\beta^\top \cdot \text{Var}(E(X|Y)) \cdot \beta \right] \\ \text{subject to} & \quad \beta^\top \Sigma \beta = I_d. \end{aligned}$$



Two Unconstrained Objectives

- To solve the optimization problem (6) with constraints, we propose the successive direction extraction method.
- Let $\lambda = (\lambda_1, \dots, \lambda_d)$ with $\lambda_j > 0$, for $1 \leq j \leq d$ and first consider

$$L_1(\lambda_1, f_1) = E\left([f_1(X) - E\{f_1(X)\}][f_1(\tilde{X}) - E\{f_1(\tilde{X})\}]\kappa(Y, \tilde{Y})\right) + \lambda_1[\text{Var}\{f_1(X)\} - 1]^2,$$

which gives

$$f_{1,\lambda_1}^* = \operatorname{argmin}_{f_1 \in L_2(P_X)} L_1(\lambda_1, f_1).$$



Two Unconstrained Objectives

- For $j \geq 2$, let

$$\begin{aligned} & L_j(\lambda_j, (\mathbf{f}_{[j-1]}^*, f_j)) \\ &= E\left([f_j(X) - E\{f_j(X)\}][f_j(\tilde{X}) - E\{f_j(\tilde{X})\}]\kappa(Y, \tilde{Y})\right) \\ &+ \lambda_j[\text{Var}\{f_j(X)\} - 1]^2 + \tilde{\lambda}\left\{\sum_{i=1}^{j-1} \text{Cov}(f_{i,\lambda_i}^*(X), f_j(X))\right\}^2, \end{aligned} \quad (7)$$

where $\mathbf{f}_{[j-1]}^* = (f_{1,\lambda_1}^*, \dots, f_{j-1,\lambda_{j-1}}^*)$, $\tilde{\lambda} > 0$ is predefined constant. Then it follows that

$$f_{j,\lambda_j}^* = \underset{f_j \in L_2(P_X)}{\text{argmin}} L_j(\lambda_j, (\mathbf{f}_{[j-1]}^*, f_j)). \quad (8)$$



Two Unconstrained Objectives

A1 There is a strict gap among λ_i^* in Theorem 2 for $m = d + 1$, i.e.,

$$\lambda_i^* = E\bar{f}_i^*(X)\bar{f}_i^*(\tilde{X})\kappa(Y, \tilde{Y}) < E\bar{f}_j^*(X)\bar{f}_j^*(\tilde{X})\kappa(Y, \tilde{Y}) = \lambda_j^* < 0,$$

for $1 \leq i < j \leq d + 1$. Additionally, we assume $\tilde{\lambda} > -\lambda_1^*$, where $\tilde{\lambda}$ is the constant in (7).

For linear SDR, this assumption necessitates the distinctness of the first d eigenvalues of generalized eigendecomposition problem:

$$\text{GEV}(E[\{X - E(X)\}\{\tilde{X} - E(\tilde{X})\}^\top \kappa(Y, \tilde{Y})], \Sigma).$$

Assumption A1 guarantees that f_1^*, \dots, f_d^* in Theorem 2 is unique.



Two Unconstrained Objectives

Theorem 3

Under the assumption A1, the solution of the successive direction extraction method is proportional to the optimal direction f_j^ given by (6) in Theorem 2, i.e.,*

$$f_{j,\lambda_j}^* = \pm \sqrt{1 - \frac{\lambda_j^*}{2\lambda_j}} f_j^*, \quad 1 \leq j \leq d.$$



Two Unconstrained Objectives

To alleviate the computational cost, we use a slack variable λ of scalar type called the norm multiplier instead of a matrix, regardless of the dimension d . Now consider another lagrangian function

$$L_F(\lambda, \mathbf{f}) = E\left([\mathbf{f}(X) - E\{\mathbf{f}(X)\}]^\top [\mathbf{f}(\tilde{X}) - E\{\mathbf{f}(\tilde{X})\}] \kappa(Y, \tilde{Y})\right) + \lambda \| \text{Var}\{\mathbf{f}(X)\} - I_d \|_F. \quad (9)$$

Theorem 4

Let $\lambda > 0$ and $\mathbf{f}_{\lambda, F}^*$ be an optimal solution of the optimization target

$$\min_{\mathbf{f} \in \{L_2(P_X)\}^d} L_F(\lambda, \mathbf{f}),$$

then $f_{\lambda, F, j}^* \in \mathfrak{M}_{\sigma\{f_0(X)\}}$ for all $1 \leq j \leq d$. If we further assume $\lambda > |\lambda_1^* + \dots + \lambda_d^*|$, the minimization of $L_F(\lambda, \mathbf{f})$ is equivalent to (6) in Theorem 2.

Define two sets

$$\begin{aligned}\mathcal{I}_{n,k} &= \{(\pi_1, \dots, \pi_k) \in \{1, \dots, i\}^n : \pi_j \neq \pi_l, \text{ for } j \neq l\}; \\ \mathcal{J}_{n,k} &= \{(\pi_1, \dots, \pi_k) \in \{1, \dots, i\}^n : \pi_j < \pi_l, \text{ for } j < l\}.\end{aligned}$$

Let $Z_i = (X_i, Y_i)$. Given the data

$$\mathcal{D}_n = \{Z_1, Z_2, \dots, Z_n\},$$

the most intuitive empirical form of the successive direction extraction method follows as



$$L_{1,n}(\lambda_1, f_1) = \binom{n}{4}^{-1} \sum_{(\pi_1, \pi_2, \pi_3, \pi_4) \in \mathcal{J}_{n,4}} h_{\lambda_1, f_1}(Z_{\pi_1}, Z_{\pi_2}, Z_{\pi_3}, Z_{\pi_4}),$$

$$\hat{f}_{1, \lambda_1} = \operatorname{argmin}_{f_1 \in \mathcal{F}_n} L_{1,n}(\lambda_1, f_1),$$

$$\begin{aligned} L_{j,n}(\lambda_j, (\hat{\mathbf{f}}_{[j-1]}, f_j)) &= \binom{n}{4}^{-1} \sum_{(\pi_1, \pi_2, \pi_3, \pi_4) \in \mathcal{J}_{n,4}} \{h_{\lambda_j, f_j}(Z_{\pi_1}, Z_{\pi_2}, Z_{\pi_3}, Z_{\pi_4}) \\ &\quad + \sum_{s=1}^{j-1} h_{\tilde{\lambda}, \hat{f}_s, \lambda_s, f_j}(Z_{\pi_1}, Z_{\pi_2}, Z_{\pi_3}, Z_{\pi_4})\}, \end{aligned}$$

$$\hat{f}_{j, \lambda_j} = \operatorname{argmin}_{f_j \in \mathcal{F}_n} L_{j,n}(\lambda_j, (\hat{\mathbf{f}}_{[j-1]}, f_j)),$$



where $\hat{\mathbf{f}}_{[j-1]} = (\hat{f}_{1,\lambda_1}, \dots, \hat{f}_{j-1,\lambda_{j-1}})$ is the collection of previous estimated directions, and the two symmetric kernel functions are

$$\begin{aligned} & h_{\lambda,f}(Z_1, Z_2, Z_3, Z_4) \\ &= \frac{1}{12} \sum_{(i,j) \in \mathcal{I}_{4,2}} f(X_i)f(X_j)\kappa(Y_i, Y_j) - \frac{1}{12} \sum_{(i,j,k) \in \mathcal{I}_{4,3}} f(X_i)f(X_j)\kappa(Y_j, Y_k) \\ &+ \frac{1}{24} \sum_{(i,j,k,l) \in \mathcal{I}_{4,4}} f(X_i)f(X_j)\kappa(Y_k, Y_l) + \lambda h_{1,f,f} \\ &- 2\lambda \left\{ \frac{1}{4} \sum_{(i) \in \mathcal{I}_{4,1}} f(X_i)f(X_i) + \frac{1}{12} \sum_{(i,j) \in \mathcal{I}_{4,2}} f(X_i)f(X_j) \right\} + \lambda \end{aligned}$$



$$\begin{aligned} & h_{\tilde{\lambda}, f_s, f}(Z_1, Z_2, Z_3, Z_4) \\ &= \tilde{\lambda} \left\{ \frac{1}{12} \sum_{(i,j) \in \mathcal{I}_{4,2}} f(X_i) f_s(X_i) f(X_j) f_s(X_j) \right. \\ &\quad - \frac{1}{12} \sum_{(i,j,k) \in \mathcal{I}_{4,3}} f(X_i) f_s(X_i) f(X_j) f_s(X_k) \\ &\quad \left. + \frac{1}{24} \sum_{(i,j,k,l) \in \mathcal{I}_{4,4}} f(X_i) f_s(X_j) f(X_k) f_s(X_l) \right\}. \end{aligned}$$



Estimation

- Choose a suitable class \mathcal{F}_n to approximate $L_2(P_X)$.
- we optimize $L_{j,n}(\lambda_j, (\hat{\mathbf{f}}_{[j-1]}, f_j))$ sequentially from 1 to d (when $j = 1$, $L_{j,n}(\lambda_j, (\hat{\mathbf{f}}_{[j-1]}, f_j))$ degenerates to $L_{1,n}(\lambda_1, f_1)$), and all estimations from previous steps should fill into $\hat{\mathbf{f}}_{[j-1]}$ in preparation for the next direction of estimation.
- Under the guidance of Theorem 3, we divide the sample \mathcal{D}_n into a learning sample \mathcal{D}_{n_l} of size n_{n_l} and a testing sample \mathcal{D}_{n_t} of n_t , where $n = n_l + n_t$. For any $\lambda \in \Lambda$, a pre-given set of the optional hyperparameter, define

$$\hat{f}_{j,\lambda_j,n_t} \in \operatorname{argmin}_{f_j \in \mathcal{F}_{n_t}} L_{j,n_t}(\lambda_j, (\hat{\mathbf{f}}_{[j-1]}, f_j))$$

based on the learning sample \mathcal{D}_{n_l} . Then we choose $\lambda^* \in \Lambda$ maximizing the empirical distance correlation ([Böttcher et al., 2018]) of \hat{f}_{λ,n_t} and Y on the testing sample \mathcal{D}_{n_t} .



The determination of dimension d

1. Calculate the estimator

$$\hat{\mathbf{f}}_{\lambda} \in \underset{\mathbf{f} \in \mathcal{F}_n^r}{\operatorname{argmin}} L_F(\lambda, \mathbf{f})$$

with $\lambda = 1$ and a relatively large r such as $r = \min(p, 32)$.

2. Give an estimator of GMDD matrix by

$$\hat{\Delta} = \frac{2}{n(n-1)} \sum_{j < k} \left\{ \hat{\mathbf{f}}_{\lambda}(X_j) - \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{f}}_{\lambda}(X_i) \right\} \left\{ \hat{\mathbf{f}}_{\lambda}(X_k) - \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{f}}_{\lambda}(X_i) \right\}$$

3. Calculate eigenvalues $\{\hat{\lambda}_i\}_{i=1}^r$ of $\hat{\Delta}$ and sort them in descending order.
4. Define the eigenvalue ratio as $\frac{\hat{\lambda}_i}{\sum_{i=1}^r \hat{\lambda}_i}$. And we determine the ultimate dimension d for SDR to be the minimum integer k such that the cumulative eigenvalue ratio

$$\sum_{j=1}^k \hat{\lambda}_j / \sum_{i=1}^r \hat{\lambda}_i \geq 0.9.$$



Deep Neural Networks

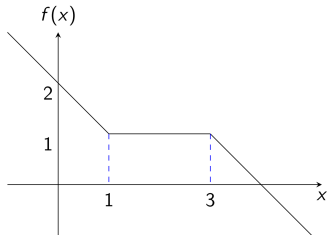
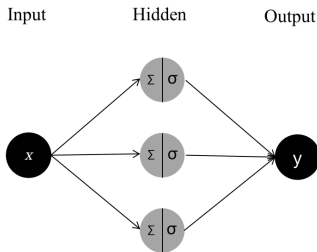
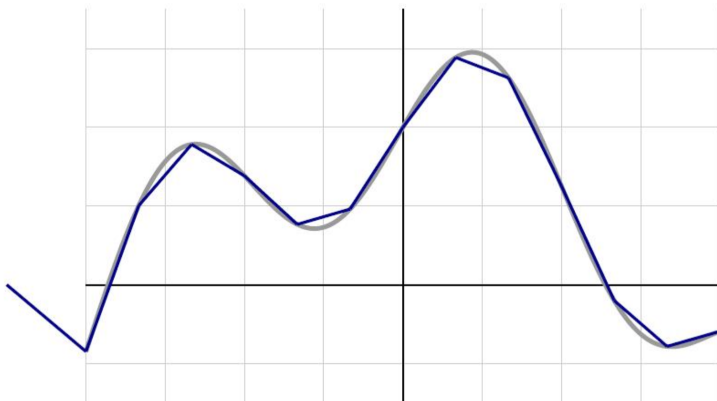


Figure: Single layer ReLU neural network.

Figure: Piecewise linear approximation, $f(x) = \sigma(1) + \sigma(-x+1) - \sigma(x-3)$.



Deep Neural Networks



Deep Neural Networks

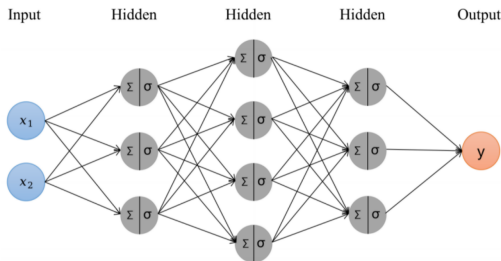


Figure: A feedforward neural network with width $N = 4$, depth $L = 3$, and size $S = 44$.

- Depth \mathcal{D} : the number of all hidden layers;
- Width \mathcal{W} : the maximum width of all hidden layers;
- Size \mathcal{S} : the total number of parameters in the network;
- Number of neurons \mathcal{U} : the number of nodes in hidden layers.



There is a lot of literature investigating the convergence properties of nonparametric regression using feedforward neural networks.

- [[Györfi et al., 2002](#)] gave a discussion about consistency and prediction error bound of deep neural networks based on the least squares objective function.
- [[Farrell et al., 2021](#)] established non-asymptotic bounds for deep feed-forward architecture net for general Lipschitz loss functions with the demand that the regression function lies in a Sobolev ball.
- [[Schmidt-Hieber, 2020](#)] proved sparsely connected and well-designed neural networks with ReLU activation function can achieve the nearly minimax rates of convergence under a general composition assumption on the regression function.



- [[Bauer and Kohler, 2019](#)] derived the prediction error bound when considering to use specially structured sigmoid activated multilayer neural networks to fit a similar structured regression function.
- [[Kohler et al., 2022](#)] gave similar results but took into account a different form of regression function and a simpler structured network.
- [[Jiao et al., 2021](#)] further improved the prefactor of optimal bounds on the prediction error to make it more meaningful in the high-dimensional setting.



Outline

- 1 Introduction
- 2 Method
- 3 Theoretical Analysis**
- 4 Simulation Studies
- 5 Real data
- 6 Summary



We define the excess risk for the j th step estimation

$$R_j(\lambda_j, \hat{f}_{j,\lambda_j}) = L_j(\lambda_j, (\mathbf{f}_{[j-1]}^*, \hat{f}_{j,\lambda_j})) - L_j(\lambda_j, (\mathbf{f}_{[j-1]}^*, f_{j,\lambda_j}^*))$$

and the total excess risk

$$R(\boldsymbol{\lambda}, \hat{\mathbf{f}}) = \sum_{j=1}^d R_j(\lambda_j, \hat{f}_{j,\lambda_j}),$$

where $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_d)$, $\hat{\mathbf{f}} = \hat{\mathbf{f}}_{[d]} = (\hat{f}_{1,\lambda_1}, \dots, \hat{f}_{d,\lambda_d})$ is obtained from (10) and $\mathbf{f}_{[j-1]}^* = (f_{1,\lambda_1}^*, \dots, f_{j-1,\lambda_{j-1}}^*)$.



- Consider the simplest case where $d = 1$.

$$\begin{aligned} R(\lambda, \hat{f}) &= R_1(\lambda_1, \hat{f}_{1,\lambda_1}) = L_1(\lambda_1, \hat{f}_{1,\lambda_1}) - L_1(\lambda_1, f_{1,\lambda_1}^*) \\ &= \underbrace{\left\{ L_1(\lambda_1, \hat{f}_{1,\lambda_1}) - \inf_{f_1 \in \mathcal{F}_n} L_1(\lambda_1, f_1) \right\}}_{\text{the statistical error}} \\ &\quad + \underbrace{\left\{ \inf_{f_1 \in \mathcal{F}_n} L_1(\lambda_1, f_1) - L_1(\lambda_1, f_{1,\lambda_1}^*) \right\}}_{\text{the approximation error}}. \end{aligned}$$

- A common strategy for tackling the first term is to reduce the U - process related problem to an empirical process problem by introducing Rademacher random variables.
- The approximation error can be solved by the functional approximation theory.



Assumptions

We make the following assumptions.

- A2 There exists an absolute constant $B > 1$ such that $\|f_{j,\lambda_j}^*\|_\infty \leq B$, $\|f\|_\infty \leq B$, and $\kappa(y, \tilde{y}) \leq B$ for any $1 \leq j \leq d$, $f \in \mathcal{F}_n$, and $y, \tilde{y} \in \mathbb{R}^q$.
- A3 f_{j,λ_j}^* is β -Hölder continuous for all $1 \leq j \leq d$, i.e., $f_{j,\lambda_j}^* \in C^\beta([0, 1]^p)$.
- A4 There exist an universal constant C and a parameter V depending on \mathcal{F}_n such that

$$\log N(\epsilon, \mathcal{F}_n, \|\cdot\|_{L_2(Q)}) \leq CV\{1 + \log(1/\epsilon)\},$$

where Q is the probability measure P_X or P_n , the empirical probability measure of X .



Assumptions

- By Theorem 2.6.7 in [Vaart and Wellner, 1996] and Theorem 7 in [Bartlett] the log covering number of the scalar-valued ReLU neural network class \mathcal{F}_n with respect to L_2 -norm for any probability measure Q can be bounded by

$$\begin{aligned}\log N(\epsilon, \mathcal{F}_n, \|\cdot\|_{L_2(Q)}) &\leq K_1 \cdot \text{VC}(\mathcal{F}_n) \{1 + \log(1/\epsilon)\} \\ &\leq K_2 \cdot \mathcal{SL} \log(\mathcal{S}) \{1 + \log(1/\epsilon)\}.\end{aligned}$$

Assumption (A4) is satisfied with $V = \text{VC}(\mathcal{F}_n)$ for ReLU neural networks.

- Regarding the approximation error, the β -Hölder continuous function class $C^\beta([0, 1]^p)$ can be approximated by the ReLU networks with error (see Theorem 1.1 in [Shen, 2020])

$$\sup_{f \in \mathcal{F}_n; h \in C^\beta([0, 1]^p)} \|f - h\|_\infty = \mathcal{O}\left((\mathcal{NL})^{-\frac{2\beta}{p}}\right).$$



Define the Rademacher average

$$R_n = \sup_{f_1 \in \mathcal{F}_n} \frac{1}{\lfloor n/4 \rfloor} \left| \sum_{i=1}^{\lfloor n/4 \rfloor} \epsilon_i h_{\lambda_1, f_1}(Z_i, Z_{\lfloor n/4 \rfloor + i}, Z_{2\lfloor n/4 \rfloor + i}, Z_{3\lfloor n/4 \rfloor + i}) \right|,$$

where $\epsilon_1, \dots, \epsilon_n$ are i.i.d. Rademacher random variables (i.e., random symmetric sign variables) and $Z = (X, Y)$.

Lemma 5

Assume \hat{f}_{1, λ_1} is a minimizer of $L_{1, n}(\lambda_1, f_1)$, then for any convex nondecreasing function ψ ,

$$E\psi\left(L_1(\lambda_1, \hat{f}_{1, \lambda_1}) - \inf_{f_1 \in \mathcal{F}_n} L_1(\lambda_1, f_1)\right) \leq E\psi(4R_n).$$



By setting $\psi(x) = e^x$, we can obtain the moment generating function of $L_1(\lambda_1, \hat{f}_{1,\lambda_1}) - \inf_{f_1 \in \mathcal{F}_n} L_1(\lambda_1, f_1)$, which is the essential component to deriving the non-asymptotic bound.

Theorem 6

Under the assumption A1, A2 and A4, if $n \geq V$, for any $\delta > 0$ with probability at least $1 - \delta$, it holds that

$$\begin{aligned} & L_1(\lambda_1, \hat{f}_{1,\lambda_1}) - L_1(\lambda_1, f_{1,\lambda_1}^*) \\ &= \mathcal{O}\left(\sqrt{\frac{V}{n}} + \sqrt{\frac{\log(1/\delta)}{n-3}} + \inf_{f_1 \in \mathcal{F}_n} \|f_1 - f_{1,\lambda_1}^*\|_\infty^2\right). \end{aligned}$$



Slow Rate For Neural Networks

Corollary 7

For any arbitrary $N \in \mathbb{N}^+$, consider \mathcal{F}_n as scalar-valued ReLU neural network classes with **the width $\mathcal{N} = 3^{p+1} \max(p \lfloor N^{1/p} \rfloor, N + 1)$ and depth $\mathcal{L} = 12n^{\frac{p}{2(p+4\beta)}} + 14 + 2p$** . Under the assumption **A1**, **A2** and **A3**, if $n \geq \mathcal{S} \mathcal{L} \log(\mathcal{S})$, the following result holds for any $\delta > 0$ with probability at least $1 - \delta$,

$$L_1(\lambda_1, \hat{f}_{1,\lambda_1}) - L_1(\lambda_1, f_{1,\lambda_1}^*) = \mathcal{O}\left(n^{-\frac{2\beta}{p+4\beta}} \log^{\frac{1}{2}} n + \sqrt{\frac{\log(1/\delta)}{n}}\right).$$

- [Shen, 2020] provided a quantitative and non-asymptotic approximation rate of deep ReLU neural networks in terms of width and depth.
- [Bartlett et al., 2019] established the nearly tight VC dimension bound of the function class constructed by feedforward neural networks the piecewise-linear activation function.



Theorem 8

Under the assumption A1, A2 and A4, the following holds for any $\delta > 0$ with probability at least $1 - \delta$,

$$L_1(\lambda_1, \hat{f}_{1,\lambda_1}) - L_1(\lambda_1, f_{1,\lambda_1}^*) = \mathcal{O}\left(\frac{V}{n} + \frac{V \log(1/\delta)}{n} + \inf_{f_1 \in \mathcal{F}_n} \|f_1 - f_{1,\lambda_1}^*\|_\infty^2\right)$$
$$\rho^2(\hat{f}_{1,\lambda_1}, f_{1,\lambda_1}^*) = \mathcal{O}\left(\frac{V}{n} + \frac{V \log(1/\delta)}{n} + \inf_{f_1 \in \mathcal{F}_n} \|f_1 - f_{1,\lambda_1}^*\|_\infty^2\right).$$



Corollary 9

For any arbitrary $N \in \mathbb{N}^+$, consider \mathcal{F}_n^d as the ReLu neural network classes with **the width** $\mathcal{N} = 3^{p+1} \max(p \lfloor N^{1/p} \rfloor, N+1)$ and **depth** $\mathcal{L} = 12n^{\frac{p}{2(p+2\beta)}} + 14 + 2p$. Under the assumption **A1**, **A2** and **A3**, if $n \geq \mathcal{SL} \log(\mathcal{S})$, the following holds for any $\delta > 0$ with probability at least $1 - \delta$,

$$L_1(\lambda_1, \hat{f}_{1,\lambda_1}) - L_1(\lambda_1, f_{1,\lambda_1}^*) = \mathcal{O}\left(n^{-\frac{2\beta}{p+2\beta}} \{1 + \log(1/\delta)\} \log n\right),$$

$$\rho^2(\hat{f}_{1,\lambda_1}, f_{1,\lambda_1}^*) = \mathcal{O}\left(n^{-\frac{2\beta}{p+2\beta}} \{1 + \log(1/\delta)\} \log n\right).$$



- The implication of Corollary 9 is

$$E\rho^2(\hat{f}_{1,\lambda_1}, f_{1,\lambda_1}^*) = \mathcal{O}\left(n^{-\frac{2\beta}{p+2\beta}} \log n\right).$$

- The implication of Corollary 9 is

$$E\rho^2(\hat{f}_{1,\lambda_1}, f_{1,\lambda_1}^*) = \mathcal{O}\left(n^{-\frac{2\beta}{p+2\beta}} \log n\right).$$

- Remarkably, our nonlinear SDR estimator, leveraging deep neural networks, achieves this minimax optimal rate up to $\log n$.
- Minimax optimal rate up to $\log n$ also is obtained in [Schmidt-Hieber, for nonparametric regression problem with ReLU networks.



Theorem 10

Under the same assumptions of Theorem 8, the following holds for any $\delta > 0$ with probability at least $1 - (4j - 3)\delta$,

$$\begin{aligned} L_j(\lambda_j, (\mathbf{f}_{[j-1]}^*, \hat{f}_{j,\lambda_j})) - L_j(\lambda_j, (\mathbf{f}_{[j-1]}^*, f_{j,\lambda_j}^*)) &= \mathcal{O}(A(n, V, \delta)), \\ \rho^2(\hat{f}_{j,\lambda_j}, f_{j,\lambda_j}^*) &= \mathcal{O}(A(n, V, \delta)), \end{aligned}$$

where

$$A(n, V, \delta) = \frac{V}{n} + \frac{V \log(1/\delta)}{n} + \sup_{1 \leq j \leq d} \inf_{f_j \in \mathcal{F}_n} \|f_j - f_{j,\lambda_j}^*\|_\infty^2.$$

Moreover, if the assumptions in Corollary 9 hold true, then for $1 \leq j \leq d$,

$$E\rho^2(\hat{f}_{j,\lambda_j}, f_{j,\lambda_j}^*) = \mathcal{O}\left(n^{-\frac{2\beta}{p+2\beta}} \log n\right).$$

Fast Rate: Hoeffding Decomposition

Set

$$q_{\lambda_1, f_1}(Z_1, Z_2, Z_3, Z_4) = h_{\lambda_1, f_1}(Z_1, Z_2, Z_3, Z_4) - h_{\lambda_1, f_{1, \lambda_1}^*}(Z_1, Z_2, Z_3, Z_4),$$

then $R_1(\lambda_1, f_1) = L_1(\lambda_1, f_1) - L_1(\lambda_1, f_{1, \lambda_1}^*) = E\{q_{\lambda_1, f_1}(Z_1, Z_2, Z_3, Z_4)\}$
and the corresponding sample estimate of $R_1(\lambda_1, f_1)$ is

$$R_{1,n}(\lambda_1, f_1) = \frac{(n-4)!}{n!} \sum_{(\pi_1, \pi_2, \pi_3, \pi_4) \in \mathcal{I}_{n,4}} q_{\lambda_1, f_1}(Z_{\pi_1}, Z_{\pi_2}, Z_{\pi_3}, Z_{\pi_4}).$$

Define the conditional kernel as

$$\begin{aligned} q_{\lambda_1, f_1}^k(z_1, \dots, z_k) &= E\{q_{\lambda_1, f_1}(Z_1, \dots, Z_4) | Z_1 = z_1, \dots, Z_k = z_k\} \\ &= \int q_{\lambda_1, f_1}(z_1, \dots, z_4) dP(z_{k+1}) \dots dP(z_4), \quad k = 0, \dots, \end{aligned}$$

and an operator $U_{n,i}$ which maps a i th order kernel q to a U -statistic as

$$U_{n,i}(q) = \frac{(n-i)!}{n!} \sum_{(\pi_1, \dots, \pi_i) \in \mathcal{I}_{n,i}} q(Z_{\pi_1}, \dots, Z_{\pi_i}).$$



Fast Rate: Hoeffding Decomposition

Then the Hoeffding decomposition of $R_{1,n}(\lambda_1, f_1)$ is

$$\begin{aligned} R_{1,n}(\lambda_1, f_1) &= A_{1,1,n}(\lambda_1, f_1) + A_{2,1,n}(\lambda_1, f_1) \\ &= \sum_{i=0}^1 \binom{4}{i} U_{n,i}(q_{\lambda_1, f_1}^{(i)}) + \sum_{i=2}^4 \binom{4}{i} U_{n,i}(q_{\lambda_1, f_1}^{(i)}), \end{aligned}$$

where $\binom{4}{i} = \frac{4!}{i!(4-i)!}$ and

$$q_{\lambda_1, f_1}^{(i)}(x_1, \dots, x_i) = \sum_{k=0}^i (-1)^{i-k} \sum_{\mathcal{J}_{i,k}} q_{\lambda_1, f_1}^k(x_{\pi_1}, \dots, x_{\pi_k}).$$



Fast Rate: Main Steps

Let

$$\ell_1(\lambda_1, f_1, z) = 4Eh_{\lambda_1, f_1}(z, Z_2, Z_3, Z_4) - 3L_1(\lambda_1, f_1).$$

Consider a new target function and its centered version as

$$v_{1,n}(\lambda_1, f_1) = \frac{1}{n} \sum_{i=1}^n \ell_1(\lambda_1, f_1, Z_i),$$

$$\bar{v}_{1,n}(\lambda_1, f_1) = \frac{1}{n} \sum_{i=1}^n \ell_1(\lambda_1, f_1, Z_i) - L_1(\lambda_1, f_1).$$

Then $A_{1,1,n}(\lambda_1, f_1)$ can be expressed as

$$A_{1,1,n}(\lambda_1, f_1) = \sum_{i=1}^n \ell_1(\lambda_1, f_1, z_i) - \sum_{i=1}^n \ell_1(\lambda_1, f_{1,\lambda_1}^*, z_i),$$



Fast Rate: Main Steps

- Define a distance on $L_2(P_X) \times L_2(P_X)$ as

$$\rho^2(f, g) = C_1 \min \left\{ \|f - g\|_{L_2(P_X)}^2, \|f + g\|_{L_2(P_X)}^2 \right\},$$

where C_1 is appropriate constant.

- By the empirical process theory, we should find a function $\phi : [0, \infty) \mapsto [0, \infty)$ with property:

$$W(\sigma) = E \left[\sup_{f_1 \in \mathcal{F}_n, d^2(f_{1, \mathcal{F}_n}^*, f_1) \leq \sigma^2} \{ \bar{v}_n(\lambda, f_{1, \mathcal{F}_n}^*) - \bar{v}_n(\lambda, f_1) \} \right] \leq \phi(\sigma) / \sqrt{n}$$

where f_{1, \mathcal{F}_n}^* is the minimizer of $L_1(\lambda_1, f_1)$ over \mathcal{F}_n .



Fast Rate: Generalization to $d > 1$

- Define

$$A(n, V, \delta) = \frac{V}{n} + \frac{V \log(1/\delta)}{n} + \sup_{1 \leq j \leq d} \inf_{f_j \in \mathcal{F}_n} \|f_j - f_{j, \lambda_j}^*\|_\infty^2.$$

- Assume that with probability $1 - (4j - 7)\delta$ for any $0 < i \leq j - 1 < d$, it holds that

$$\begin{aligned} L_i(\lambda_i, (f_{[i-1]}^*, \hat{f}_{i, \lambda_i})) - L_i(\lambda_i, (f_{[i-1]}^*, f_{i, \lambda_i}^*)) &= \mathcal{O}(A(n, V, \delta)), \\ \rho^2(\hat{f}_{i, \lambda_i}, f_{i, \lambda_i}^*) &= \mathcal{O}(A(n, V, \delta)). \end{aligned} \quad (11)$$

- Introduce the intermediate function \tilde{f}_{j, λ_j} as the minimizer of

$$L_{j, n}(\lambda_j, (f_{[j-1]}^*, f_j))$$

over $f_j \in \mathcal{F}_n$.



Fast Rate: Generalization to $d > 1$

- All the relationship in $d = 1$ can transfer to $d > 1$ with different universal constant, except the inequality

$$\left| A_{1,1,n}(\lambda_1, \hat{f}_{1,\lambda_1}) - \inf_{f_1} A_{1,1,n}(\lambda_1, f_1) \right| = \mathcal{O} \left(\sup_{f_1} |A_{2,1,n}(\lambda_1, f_1)| \right).$$

Such gap also holds for $A_{1,j,n}$ with \tilde{f}_{j,λ_j} rather than \hat{f}_{j,λ_j} , that is,

$$\begin{aligned} & \left| A_{1,j,n}(\lambda_j, (\mathbf{f}_{[j-1]}^*, \tilde{f}_{j,\lambda_j})) - \inf_{f_j} A_{1,j,n}(\lambda_j, (\mathbf{f}_{[j-1]}^*, f_j)) \right| \\ &= \mathcal{O} \left(\sup_{f_j} |A_{2,j,n}(\lambda_j, (\mathbf{f}_{[j-1]}^*, f_j))| \right), \end{aligned}$$

which implies $\rho^2(\tilde{f}_{j,\lambda_j}, f_{j,\lambda_j}^*) = \mathcal{O}(A(n, V, \delta))$.



Fast Rate: Generalization to $d > 1$

- the bound of $R_j(\lambda_j, \hat{f}_{j,\lambda_j})$ can be derived once we know

$$|A_{1,j,n}(\lambda_j, (\mathbf{f}_{[j-1]}^*, \hat{f}_{j,\lambda_j})) - A_{1,j,n}(\lambda_j, (\mathbf{f}_{[j-1]}^*, \tilde{f}_{j,\lambda_j}))|.$$

- By the same relationship $R_{j,n} = A_{1,j,n} + A_{2,j,n}$ as the case of $d = 1$, we can turn to bound

$$\begin{aligned} & R_{j,n}(\lambda_j, \hat{f}_{j,\lambda_j}) - R_{j,n}(\lambda_j, \tilde{f}_{j,\lambda_j}) \\ &= L_{j,n}(\lambda_j, (\mathbf{f}_{[j-1]}^*, \hat{f}_{j,\lambda_j})) - L_{j,n}(\lambda_j, (\mathbf{f}_{[j-1]}^*, \tilde{f}_{j,\lambda_j})) \\ &= \{L_{j,n}(\lambda_j, (\mathbf{f}_{[j-1]}^*, \hat{f}_{j,\lambda_j})) - L_{j,n}(\lambda_j, (\hat{\mathbf{f}}_{[j-1]}, \hat{f}_{j,\lambda_j}))\} \\ &+ \{L_{j,n}(\lambda_j, (\hat{\mathbf{f}}_{[j-1]}, \hat{f}_{j,\lambda_j})) - L_{j,n}(\lambda_j, (\hat{\mathbf{f}}_{[j-1]}, \tilde{f}_{j,\lambda_j}))\} \\ &+ \{L_{j,n}(\lambda_j, (\hat{\mathbf{f}}_{[j-1]}, \tilde{f}_{j,\lambda_j})) - L_{j,n}(\lambda_j, (\mathbf{f}_{[j-1]}^*, \tilde{f}_{j,\lambda_j}))\} \\ &\triangleq s_1 + s_2 + s_3, \end{aligned}$$



Fast Rate: Generalization to $d > 1$

- by the inequality $(\sum_{i=1}^k x_i)^2 \leq k \sum_{i=1}^k x_i^2$, the orthogonality of $f_{s,\lambda_s}^*, f_{k,\lambda_k}^*$ with $s \neq k$ and (11), we have

$$\begin{aligned} Es_1 &= \mathcal{O} \left(\sum_{i=1}^{j-1} \{ \text{Cov}^2(f_{i,\lambda_i}^*, \hat{f}_{j,\lambda_j}) - \text{Cov}^2(\hat{f}_{i,\lambda_i}, \hat{f}_{j,\lambda_j}) \} \right) \\ &= \mathcal{O} \left(\sum_{i=1}^{j-1} \{ \rho^2(\hat{f}_{i,\lambda_i}, f_{i,\lambda_i}^*) + \rho(\hat{f}_{i,\lambda_i}, f_{i,\lambda_i}^*) \rho(\hat{f}_{j,\lambda_j}, f_{j,\lambda_j}^*) \} \right) \\ &= \mathcal{O} \left(A(n, V, \delta) + \sqrt{A(n, V, \delta)} \rho(\hat{f}_{j,\lambda_j}, f_{j,\lambda_j}^*) \right), \\ Es_3 &= \mathcal{O} \left(\sum_{i=1}^{j-1} \{ \rho^2(\hat{f}_{i,\lambda_i}, f_{i,\lambda_i}^*) + \rho(\hat{f}_{i,\lambda_i}, f_{i,\lambda_i}^*) \rho(\tilde{f}_{j,\lambda_j}, f_{j,\lambda_j}^*) \} \right) \\ &= \mathcal{O} \left(A(n, V, \delta) + \sqrt{A(n, V, \delta)} \rho(\tilde{f}_{j,\lambda_j}, f_{j,\lambda_j}^*) \right) \\ &= \mathcal{O} \left(A(n, V, \delta) \right). \end{aligned}$$



- Bounded difference inequality entails

$$s_1 = \mathcal{O}\left(A(n, V, \delta) + \sqrt{A(n, V, \delta)}\rho(\hat{f}_{j,\lambda_j}, f_{j,\lambda_j}^*)\right),$$
$$s_3 = \mathcal{O}\left(A(n, V, \delta)\right)$$

with high probability, which gives

$$\sum_{i=1}^3 s_i = \mathcal{O}\left(A(n, V, \delta) + \sqrt{A(n, V, \delta)}\rho(\hat{f}_{j,\lambda_j}, f_{j,\lambda_j}^*)\right).$$



Fast Rate: Generalization to $d > 1$

- Finally, we can bound

$$\begin{aligned} & |A_{1,j,n}(\lambda_j, (\mathbf{f}_{[j-1]}^*, \hat{f}_{j,\lambda_j})) - \inf_{f_j} A_{1,j,n}(\lambda_j, (\mathbf{f}_{[j-1]}^*, f_j))| \\ &= \mathcal{O}\left(A(n, V, \delta) + \sqrt{A(n, V, \delta)}\rho(\hat{f}_{j,\lambda_j}, f_{j,\lambda_j}^*)\right) \end{aligned}$$

- Therefore, we obtain

$$\begin{aligned} \rho^2(\hat{f}_{j,\lambda_j}, f_{j,\lambda_j}^*) &= \mathcal{O}(R_j(\lambda_j, \hat{\mathbf{f}}_{j,\lambda_j})) \\ &= \mathcal{O}\left(A(n, V, \delta) + \sqrt{A(n, V, \delta)}\rho(\hat{f}_{j,\lambda_j}, f_{j,\lambda_j}^*)\right). \end{aligned}$$

It concludes

$$\rho^2(\hat{f}_{j,\lambda_j}, f_{j,\lambda_j}^*) = \mathcal{O}\left(A(n, V, \delta)\right).$$



Outline

- 1 Introduction
- 2 Method
- 3 Theoretical Analysis
- 4 Simulation Studies**
- 5 Real data
- 6 Summary



Simulation of $d = 1$

We compare our method against GSIR. Six different variants of our method with six different distance metric κ for Y is denoted by the numbers 1 – 6.

Setting:

$$\left\{ \begin{array}{l} A: \quad Y = \frac{X_1}{1+(1+X_2)^2} + (1 + X_2)^2 + \epsilon_1, \\ B: \quad Y = \sin((X_1 + X_2)\pi/10) + X_1^2 + \epsilon_2, \\ C: \quad Y = (X_1 + X_2)^2 \log(X_1^2 + X_2^2 + 0.001) + \epsilon_3, \end{array} \right.$$

with $\epsilon \sim N(0, 0.25)$ independent of X , as well as three different distributional scenarios for the predictor vector $X = (X_1, \dots, X_p)^\top$:

$$\left\{ \begin{array}{l} \text{I:} \quad X_i \sim N(0, 1/2) \text{ where } 1 \leq i \leq p \text{ and } X_i, X_j \text{ are mutually independent} \\ \text{II:} \quad X_i \sim -1 + \text{Pois}(1) \text{ where } 1 \leq i \leq p \text{ and } X_i, X_j \text{ are mutually independent} \\ \text{III:} \quad X \sim t_4(\mathbf{0}_p, 0.6 \cdot \mathbf{I}_p + 0.4 \cdot \mathbf{1}_{p \times p}), \end{array} \right.$$



Simulation of $d = 1$

The quality of estimated sufficient predictors is assessed by their distance correlation with the true sufficient predictors.

Table: Distance correlation with true predictor when $p = 50, d = 1$ for GMDDNet-S and GSIR.

(X, Y)	sample size	1	2	3	4	5	6	GSIR
A-I	1k	0.88(0.01)	0.89(0.01)	0.76(0.03)	0.74(0.03)	0.86(0.01)	0.89(0.01)	0.92(0.00)
	2k	0.89(0.01)	0.92(0.01)	0.79(0.01)	0.78(0.02)	0.88(0.01)	0.91(0.01)	0.92(0.01)
	5k	0.93(0.00)	0.95(0.00)	0.81(0.01)	0.81(0.01)	0.91(0.01)	0.94(0.00)	0.93(0.00)
A-II	1k	0.93(0.01)	0.94(0.01)	0.74(0.02)	0.74(0.02)	0.86(0.01)	0.95(0.00)	0.92(0.00)
	2k	0.95(0.01)	0.97(0.00)	0.78(0.02)	0.76(0.01)	0.88(0.01)	0.97(0.00)	0.92(0.00)
	5k	0.97(0.00)	0.98(0.00)	0.82(0.01)	0.80(0.01)	0.92(0.01)	0.98(0.00)	0.92(0.00)
A-III	1k	0.86(0.06)	0.75(0.11)	0.61(0.06)	0.60(0.06)	0.82(0.04)	0.80(0.11)	0.61(0.05)
	2k	0.93(0.03)	0.80(0.13)	0.66(0.03)	0.65(0.03)	0.88(0.02)	0.88(0.08)	0.59(0.12)
	5k	0.96(0.01)	0.89(0.06)	0.69(0.02)	0.67(0.02)	0.88(0.02)	0.95(0.03)	0.61(0.05)



Simulation of $d = 1$

Table: Distance correlation with true predictor when $p = 50, d = 1$ for GMDDNet-S and GSIR.

(X, Y)	sample size	1	2	3	4	5	6	GSIR
B-I	1k	0.50(0.05)	0.50(0.05)	0.35(0.04)	0.37(0.04)	0.50(0.04)	0.50(0.07)	0.38(0.04)
	2k	0.66(0.03)	0.67(0.03)	0.56(0.04)	0.58(0.04)	0.67(0.02)	0.67(0.03)	0.42(0.03)
	5k	0.74(0.01)	0.77(0.02)	0.64(0.02)	0.67(0.02)	0.72(0.01)	0.76(0.02)	0.43(0.02)
B-II	1k	0.81(0.03)	0.80(0.04)	0.38(0.07)	0.31(0.09)	0.74(0.04)	0.83(0.03)	0.70(0.02)
	2k	0.87(0.02)	0.90(0.01)	0.53(0.02)	0.50(0.07)	0.82(0.02)	0.91(0.01)	0.72(0.01)
	5k	0.90(0.02)	0.94(0.01)	0.58(0.02)	0.59(0.02)	0.86(0.01)	0.94(0.01)	0.72(0.02)
B-III	1k	0.88(0.09)	0.72(0.13)	0.66(0.04)	0.66(0.04)	0.86(0.02)	0.80(0.13)	0.57(0.03)
	2k	0.92(0.08)	0.80(0.14)	0.66(0.04)	0.67(0.04)	0.88(0.02)	0.86(0.15)	0.55(0.02)
	5k	0.97(0.01)	0.84(0.14)	0.70(0.02)	0.71(0.02)	0.90(0.01)	0.91(0.06)	0.57(0.03)



Simulation of $d = 1$

Table: Distance correlation with true predictor when $p = 50$, $d = 1$ for GMDDNet-S and GSIR.

(X, Y)	sample size	1	2	3	4	5	6	GSIR
C-I	1k	0.34(0.08)	0.31(0.08)	0.29(0.05)	0.29(0.06)	0.38(0.05)	0.34(0.08)	0.11(0.04)
	2k	0.72(0.05)	0.67(0.10)	0.50(0.05)	0.54(0.04)	0.70(0.06)	0.70(0.10)	0.14(0.02)
	5k	0.88(0.04)	0.90(0.03)	0.61(0.03)	0.67(0.05)	0.81(0.04)	0.88(0.02)	0.19(0.02)
C-II	1k	0.70(0.04)	0.60(0.06)	0.22(0.06)	0.23(0.08)	0.65(0.04)	0.65(0.05)	0.61(0.02)
	2k	0.81(0.03)	0.69(0.03)	0.41(0.08)	0.42(0.09)	0.74(0.02)	0.74(0.04)	0.61(0.02)
	5k	0.95(0.00)	0.82(0.12)	0.58(0.02)	0.56(0.02)	0.91(0.01)	0.92(0.04)	0.62(0.01)
C-III	1k	0.73(0.08)	0.60(0.10)	0.54(0.07)	0.53(0.08)	0.77(0.07)	0.64(0.13)	0.62(0.08)
	2k	0.79(0.06)	0.56(0.11)	0.55(0.03)	0.54(0.03)	0.82(0.02)	0.66(0.11)	0.64(0.03)
	5k	0.85(0.12)	0.66(0.10)	0.54(0.08)	0.53(0.09)	0.73(0.12)	0.73(0.13)	0.59(0.09)



Simulation of $d = 1$

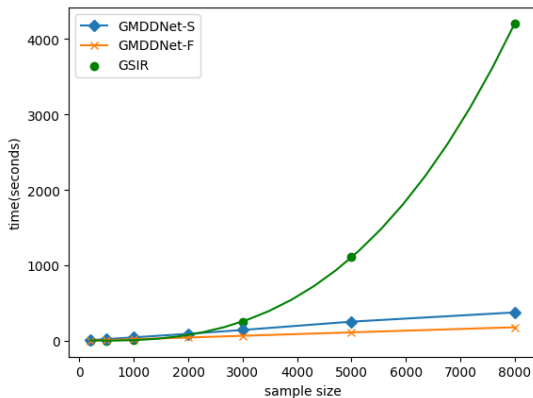


Figure: Time for different methods



Simulation of $d = 2$

Now consider the case of $d = 2$:

$$\begin{cases} D: & Y = \left((1 + X_1)^2 + \epsilon_1, \sin((X_4 + X_5)\pi/10) + X_4^2 + \epsilon_2 \right)^\top, \\ E: & Y = \mathbb{I}\{f_1(X) + \epsilon_3 > 0\} + 2 \cdot \mathbb{I}\{f_2(X) + \epsilon_4 > 0\}, \\ F: & Y = (f_3(X) + \epsilon_5, f_4(X) + \epsilon_6, \epsilon_7, \dots, \epsilon_{14})^\top \in \mathbb{R}^{10}, \end{cases}$$

where

$$f_1(X) = \log \frac{(2X_1 + 1)^2 + (2X_2 + 1)^2}{(2X_6 + 1)^2 + (2X_7 + 1)^2},$$

$$f_2(X) = \log \frac{(2X_4 + 1)^2 + (2X_5 + 1)^2}{(2X_6 + 1)^2 + (2X_7 + 1)^2},$$

$$f_3(X) = \sqrt{(X_1 + 1)^2 + (X_2 + 1)^2 + 1} + \sqrt{X_1^2 + X_2^2 + 1},$$

$$f_4(X) = (1 + X_4)^2 + (1 + X_5)^2,$$

the distributions for the predictor vector X still adhere to the I, II, III aforementioned and the mutually independent $\epsilon_i \sim N(0, 0.25)$.



Simulation of $d = 2$

Table: Distance correlation with true predictor at $p = 50, d = 2$ for GMDDNet-S and GSIR.

(X, Y)	sample size	1	2	3	4	5	6	GSIR
D-I	1k	0.80(0.02)	0.79(0.03)	0.69(0.02)	0.68(0.02)	0.79(0.02)	0.80(0.02)	0.79(0.01)
	2k	0.83(0.01)	0.83(0.01)	0.73(0.02)	0.72(0.01)	0.82(0.01)	0.83(0.01)	0.80(0.01)
	5k	0.85(0.01)	0.87(0.01)	0.78(0.00)	0.77(0.01)	0.84(0.00)	0.88(0.01)	0.81(0.00)
D-II	1k	0.91(0.01)	0.92(0.02)	0.63(0.02)	0.62(0.02)	0.80(0.01)	0.93(0.01)	0.84(0.01)
	2k	0.92(0.01)	0.94(0.01)	0.70(0.02)	0.70(0.02)	0.84(0.01)	0.94(0.01)	0.84(0.01)
	5k	0.94(0.00)	0.97(0.01)	0.74(0.02)	0.74(0.01)	0.89(0.01)	0.97(0.01)	0.85(0.00)
D-III	1k	0.84(0.04)	0.72(0.08)	0.58(0.03)	0.54(0.03)	0.79(0.04)	0.79(0.05)	0.58(0.02)
	2k	0.88(0.05)	0.73(0.11)	0.65(0.01)	0.61(0.02)	0.83(0.02)	0.82(0.07)	0.59(0.02)
	5k	0.93(0.01)	0.83(0.04)	0.68(0.03)	0.65(0.02)	0.85(0.02)	0.91(0.02)	0.59(0.02)



Simulation of $d = 2$

Table: Distance correlation with true predictor at $p = 50, d = 2$ for GMDDNet-S and GSIR.

(X, Y)	sample size	1	2	3	4	5	6	GSIR
E-I	1k	—	—	—	—	—	0.46(0.03)	0.61(0.02)
	2k	—	—	—	—	—	0.55(0.02)	0.62(0.01)
	5k	—	—	—	—	—	0.64(0.02)	0.65(0.02)
E-II	1k	—	—	—	—	—	0.59(0.02)	0.79(0.01)
	2k	—	—	—	—	—	0.61(0.02)	0.81(0.01)
	5k	—	—	—	—	—	0.68(0.01)	0.82(0.01)
E-III	1k	—	—	—	—	—	0.42(0.02)	0.24(0.03)
	2k	—	—	—	—	—	0.55(0.04)	0.28(0.03)
	5k	—	—	—	—	—	0.65(0.01)	0.34(0.04)



Simulation of $d = 2$

Table: Distance correlation with true predictor at $p = 50, d = 2$ for GMDDNet-S and GSIR.

(X, Y)	sample size	1	2	3	4	5	6	GSIR
F-I	1k	0.76(0.02)	0.75(0.02)	0.70(0.02)	0.60(0.03)	0.75(0.02)	0.76(0.02)	0.79(0.01)
	2k	0.81(0.01)	0.81(0.02)	0.75(0.02)	0.65(0.02)	0.80(0.01)	0.81(0.02)	0.80(0.01)
	5k	0.85(0.01)	0.84(0.01)	0.82(0.00)	0.74(0.00)	0.86(0.01)	0.85(0.00)	0.80(0.01)
F-II	1k	0.80(0.02)	0.80(0.02)	0.70(0.02)	0.55(0.02)	0.79(0.02)	0.81(0.02)	0.83(0.01)
	2k	0.87(0.01)	0.84(0.02)	0.75(0.01)	0.59(0.02)	0.84(0.01)	0.86(0.01)	0.84(0.01)
	5k	0.91(0.00)	0.87(0.01)	0.83(0.01)	0.67(0.01)	0.90(0.00)	0.89(0.01)	0.84(0.01)
F-III	1k	0.85(0.03)	0.77(0.05)	0.64(0.03)	0.52(0.03)	0.77(0.02)	0.84(0.03)	0.64(0.03)
	2k	0.88(0.02)	0.80(0.09)	0.69(0.03)	0.56(0.02)	0.81(0.01)	0.88(0.04)	0.65(0.02)
	5k	0.91(0.01)	0.89(0.01)	0.74(0.02)	0.62(0.01)	0.86(0.00)	0.92(0.02)	0.63(0.01)



Simulation of $d = 2$

Table: Distance correlation with true predictor when $p = 50$, $d = 2$ for GMDDNet-F and GSIR.

(X, Y)	sample size	1	2	3	4	5	6	GSIR
D-I	1k	0.74(0.03)	0.73(0.03)	0.70(0.04)	0.70(0.03)	0.75(0.03)	0.74(0.03)	0.79(0.01)
	2k	0.81(0.03)	0.78(0.04)	0.78(0.02)	0.78(0.01)	0.83(0.02)	0.80(0.03)	0.80(0.01)
	5k	0.88(0.02)	0.90(0.01)	0.82(0.01)	0.81(0.01)	0.86(0.01)	0.89(0.02)	0.81(0.00)
D-II	1k	0.84(0.02)	0.83(0.05)	0.70(0.03)	0.68(0.03)	0.82(0.01)	0.84(0.04)	0.84(0.01)
	2k	0.89(0.02)	0.87(0.05)	0.74(0.02)	0.72(0.02)	0.87(0.01)	0.91(0.02)	0.84(0.01)
	5k	0.94(0.01)	0.92(0.04)	0.77(0.01)	0.75(0.02)	0.90(0.00)	0.94(0.02)	0.85(0.00)
D-III	1k	0.78(0.05)	0.62(0.08)	0.65(0.02)	0.62(0.02)	0.76(0.03)	0.68(0.07)	0.58(0.02)
	2k	0.75(0.18)	0.51(0.09)	0.68(0.02)	0.66(0.02)	0.82(0.02)	0.63(0.12)	0.59(0.02)
	5k	0.86(0.03)	0.53(0.09)	0.70(0.03)	0.68(0.02)	0.84(0.02)	0.69(0.05)	0.59(0.02)



Simulation of $d = 2$

Table: Distance correlation with true predictor when $p = 50, d = 2$ for GMDDNet-F and GSIR.

(X, Y)	sample size	1	2	3	4	5	6	GSIR
E-I	1k	—	—	—	—	—	0.48(0.03)	0.61(0.02)
	2k	—	—	—	—	—	0.56(0.02)	0.62(0.01)
	5k	—	—	—	—	—	0.70(0.02)	0.65(0.02)
E-II	1k	—	—	—	—	—	0.61(0.03)	0.79(0.01)
	2k	—	—	—	—	—	0.65(0.02)	0.81(0.01)
	5k	—	—	—	—	—	0.70(0.01)	0.82(0.01)
E-III	1k	—	—	—	—	—	0.29(0.07)	0.24(0.03)
	2k	—	—	—	—	—	0.45(0.03)	0.28(0.03)
	5k	—	—	—	—	—	0.62(0.02)	0.34(0.04)



Simulation of $d = 2$

Table: Distance correlation with true predictor when $p = 50$, $d = 2$ for GMDDNet-F and GSIR.

(X, Y)	sample size	1	2	3	4	5	6	GSIR
F-I	1k	0.76(0.03)	0.76(0.03)	0.72(0.03)	0.63(0.05)	0.77(0.03)	0.77(0.02)	0.79(0.01)
	2k	0.80(0.02)	0.77(0.03)	0.81(0.01)	0.75(0.04)	0.82(0.02)	0.79(0.03)	0.80(0.01)
	5k	0.86(0.01)	0.85(0.02)	0.87(0.01)	0.79(0.02)	0.89(0.01)	0.86(0.02)	0.80(0.01)
F-II	1k	0.81(0.02)	0.77(0.03)	0.72(0.03)	0.57(0.03)	0.80(0.02)	0.79(0.04)	0.83(0.01)
	2k	0.84(0.02)	0.79(0.03)	0.79(0.02)	0.65(0.02)	0.85(0.02)	0.82(0.02)	0.84(0.01)
	5k	0.91(0.01)	0.85(0.04)	0.85(0.01)	0.71(0.03)	0.90(0.01)	0.86(0.02)	0.84(0.01)
F-III	1k	0.77(0.04)	0.66(0.06)	0.72(0.03)	0.62(0.03)	0.76(0.03)	0.69(0.05)	0.64(0.03)
	2k	0.79(0.06)	0.61(0.11)	0.73(0.03)	0.64(0.09)	0.80(0.02)	0.68(0.12)	0.65(0.02)
	5k	0.85(0.02)	0.65(0.04)	0.79(0.02)	0.68(0.03)	0.85(0.01)	0.72(0.03)	0.63(0.01)



Dimension determination

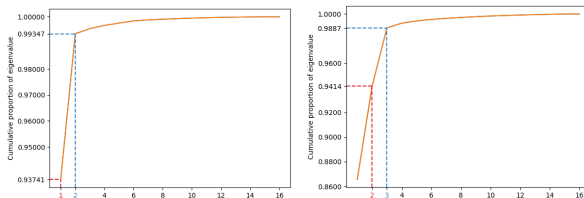


Figure: The eigenvalue ratio plot for for setting A-I with $n = 2000, p = 50, d = 1$ and setting E-I with $n = 2000, p = 50, d = 2$.



Outline

- 1 Introduction
- 2 Method
- 3 Theoretical Analysis
- 4 Simulation Studies
- 5 Real data**
- 6 Summary



The dataset called “Communities and Crime ” consists of 1994 instances with 127 features and one response variable, “violent crimes per capita”.

Table: Mean square error of different methods on the regression task.

Method	Indirect methods			Direct methods	
	GMDDNet-S	GMDDNet-F	GSIR	DNN	Linear regression
MSE	0.0170	0.0161	0.0174	0.0195	8.9×10^{17}



Classification

The MNIST database and the Fashion-MNIST both comprise 70,000 28×28 images belonging to ten categories, with 60,000 images in the training set and 10,000 in the test set. Target dimension $d = 9$ can be observed below.

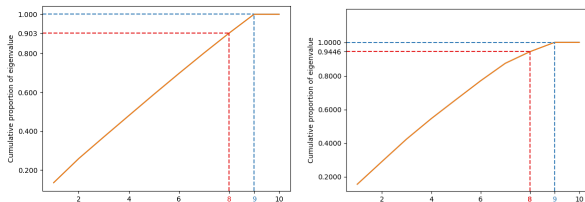


Figure: The eigenvalue ratio plot for MNIST and FashionMNIST



We apply the logistic regression to the reduced features of all SDR methods. Accuracy and visualization are shown for comparison.

Table: Accuracy of different methods on MNIST with sample size 5000.

Method	DNN			LeNet			GSIR
	GMDDNet-S	GMDDNet-F	Direct	GMDDNet-S	GMDDNet-F	Direct	
accuracy	0.9583	0.949	0.9534	0.9731	0.9719	0.9708	0.8313



2D visualization

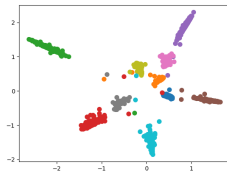


Figure: DNN with norm multiplier

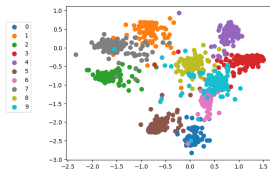


Figure: LeNet with norm multiplier

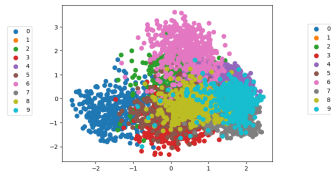


Figure: GSIR

Figure: 2D visualization of MNIST using DNN, LeNet, GSIR with sample size 5k



2D visualization

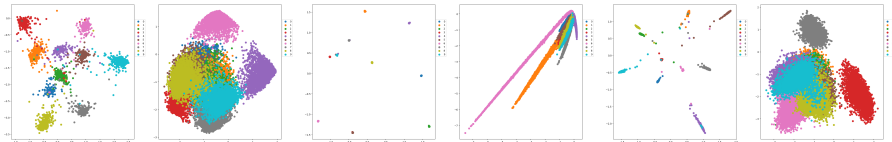


Figure: LeNet

Figure: ResNet

Figure: DenseNet

Figure: 2D-visualization of MNIST using LeNet, ResNet and DenseNet with sample size 60k (Left panel: GMDDNet-F with convolutional neural networks; Right panel: direct convolutional neural networks).



Table: Accuracy of different CNNs on MNIST with sample size 60000.

Method	LeNet		ResNet18		DenseNet	
	GMDDNet-F	Direct	GMDDNet-F	Direct	GMDDNet-F	Direct
Accuracy	0.9909	0.988	0.9948	0.9939	0.996	0.996

Table: Accuracy of different CNNs on FashionMNIST with sample size 60000.

Method	LeNet		ResNet18		DenseNet	
	GMDDNet-F	Direct	GMDDNet-F	Direct	GMDDNet-F	Direct
Accuracy	0.9143	0.9087	0.9157	0.9148	0.9345	0.9325



Outline

- 1 Introduction
- 2 Method
- 3 Theoretical Analysis
- 4 Simulation Studies
- 5 Real data
- 6 Summary**



Summary

- We propose a novel nonlinear SDR optimization objective from the perspective of GMDD, which is suitable for multidimensional Y and frees us from the restrictive step of slicing response variables in the classical SDR process.
- The nonlinear dimensional reduction function is estimated by a deep ReLU neural network.
- Our method is developed under the classical theoretical framework of SDR and the most basic unbiasedness can be theoretically guaranteed without any stringent assumptions.
- The excess risk nearly attains the minimax rate.
- The validity of our deep nonlinear SDR is demonstrated through simulations and real data.



Thank You!





Bartlett, P. L., Harvey, N., Liaw, C., and Mehrabian, A. (2019).

Nearly-tight vc-dimension and pseudodimension bounds for piecewise linear neural networks.

Journal of Machine Learning Research, 20(63):1–17.



Bauer, B. and Kohler, M. (2019).

On deep learning as a remedy for the curse of dimensionality in non-parametric regression.

The Annals of Statistics, 47(4):2261–2285.



Böttcher, B., Keller-Ressel, M., and Schilling, R. L. (2018).

Detecting independence of random vectors: generalized distance covariance and gaussian covariance.

Modern Stochastics: Theory and Applications, 5(3):353–383.



Farrell, M. H., Liang, T., and Misra, S. (2021).

Deep neural networks for estimation and inference.

Econometrica, 89(1):181–213.



Fukumizu, K., Bach, F. R., and Jordan, M. I. (2009).

Kernel dimension reduction in regression.



Györfi, L., Kohler, M., Krzyzak, A., Walk, H., et al. (2002).
A distribution-free theory of nonparametric regression, volume 1.
Springer.



Jiao, Y., Shen, G., Lin, Y., and Huang, J. (2021).
Deep nonparametric regression on approximately low-dimensional manifolds.
arXiv preprint arXiv:2104.06708.



Kohler, M., Krzyzak, A., and Langer, S. (2022).
Estimation of a function of low local dimensionality by deep neural networks.
IEEE Transactions on Information Theory.



Lee, K.-Y., Li, B., and Chiaromonte, F. (2013).
A general theory for nonlinear sufficient dimension reduction: Formulation and estimation.
The Annals of Statistics, 41(1):221 – 249.



Schmidt-Hieber, J. (2020).



Nonparametric regression using deep neural networks with relu activation function.

The Annals of Statistics, 48(4):1875–1897.



Shen, Z. (2020).

Deep network approximation characterized by number of neurons.

Communications in Computational Physics, 28(5):1768–1811.



Vaart, A. W. and Wellner, J. A. (1996).

Weak convergence.

In *Weak convergence and empirical processes*, pages 16–28. Springer.

