

# Nonparametric Logistic Regression with Deep Neural Networks

Atsutomo Yara<sup>1</sup>, Yoshikazu Terada<sup>1, 2</sup>

<sup>1</sup> Graduate School of Engineering Science, Osaka University

<sup>2</sup> Center for Advanced Intelligence Project (AIP), RIKEN

Sep 30th

## Contents

### 1. Introduction

- 1.1. Classification and Estimation of Conditional Probability
- 1.2. Classification with Deep Learning
- 1.3. Semiparametric Estimation and Estimation of Conditional Probability

### 2. Related works

- 2.1. Theoretical Study of Classification Problem with Deep Learning
- 2.2. Remained Problems
- 2.3. Summary of this work

### 3. Main Results

- 3.1. Notations and Assumptions
- 3.2. Oracle Inequality
- 3.3. Rate of Convergence for Deep Learning

# Introduction

## Classification and Estimation of Conditional Probability<sup>4</sup>

### ► Setting of Classification Problem

- Consider the  $K$  class classification problem.
  - $\mathcal{X} \subset \mathbb{R}^d$ : input space,  $\mathcal{Y} = \{1, \dots, K\}$ : output space.
  - $P$ : distribution on  $\mathcal{X} \times \mathcal{Y}$ ,  $P_{\mathbf{X}}$ : marginal distribution of  $X$ .
  - $(\mathbf{X}, Y), (\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n) \sim_{i.i.d.} P$
  - $\mathcal{D}_n := \{(\mathbf{X}_i, Y_i)\}_{i=1}^n$ : observed data
  - $\eta_k(\mathbf{x}) := P(Y = k \mid \mathbf{X} = \mathbf{x})$ ,  $\boldsymbol{\eta}(\mathbf{x}) = (\eta_1(\mathbf{x}), \dots, \eta_K(\mathbf{x}))^T$
  - If  $\mathbf{X}$  given,  $Y$  follows the following multinomial distribution.
$$Y|\mathbf{X} = \mathbf{x} \sim \text{Multi}(\boldsymbol{\eta}(\mathbf{x}))$$

### ► The objective of the classification problem:

Finding a function  $f : \mathcal{X} \rightarrow \{1, \dots, K\}$  (decision function) that predicts  $Y$  when given new data  $\mathbf{X}$ .

## Classification and Estimation of Conditional Probability<sup>5</sup>

### ► Learning Method

- “Learning” = obtaining decision function  $f$  from observed data  $\mathcal{D}_n$ .
- Empirical Risk Minimization (ERM) is commonly used.
- ERM = Minimizing the loss function for observed data.
- Let  $l$  be a loss function and  $\mathcal{F}_n$  be

a hypothesis space (candidates of the estimator). Then, ERM  $\hat{f}_n$  is

$$\hat{f}_n \in \operatorname{argmin}_{f \in \mathcal{F}_n} \frac{1}{n} \sum_{i=1}^n l(Y_i, f(\mathbf{X}_i))$$

- Predicting  $Y$  when  $\mathbf{X}$  is given based on the estimated decision function.
- Various loss functions are employed.
  - Hinge loss:  $l(y, f(x)) := 0 \vee (1 - yf(x))$
  - Logistic loss:  $l(y, f(x)) := \log(1 + e^{-yf(x)})$

## Classification and Estimation of Conditional Probability<sup>6</sup>

### ► Estimation of Conditional Probability

- By using the hinge loss, we can only estimate the label.
- Conditional probabilities are required in many applications.
  - e.g., AI-based disease detection
- Conditional probabilities are helpful for human decision-making!
- Estimation of conditional probabilities  $\Rightarrow$  Logistic regression!

#### Nonparametric Maximum Likelihood Estimator (NPMLE)

$\hat{\mathbf{p}}_n(\mathbf{x}) = (\hat{p}_1(\mathbf{x}), \dots, \hat{p}_K(\mathbf{x}))$  is defined as

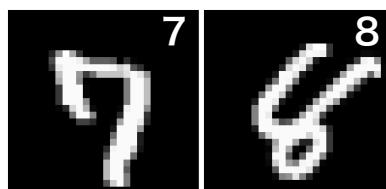
$$\hat{\mathbf{p}}_n \in \operatorname{arg min}_{\mathbf{p} \in \mathcal{F}_n} l(\mathbf{p}), \quad l(\mathbf{p}) := -\frac{1}{n} \sum_{i=1}^n \mathbf{Y}_i^\top \log \mathbf{p}(\mathbf{X}_i).$$

- $\mathbf{Y}_i$  is a one-hot representation of  $Y_i$   
i.e.,  $\mathbf{Y}_i = (0, \dots, \underset{k}{1}, \dots, 0)$  if  $Y_i = k$

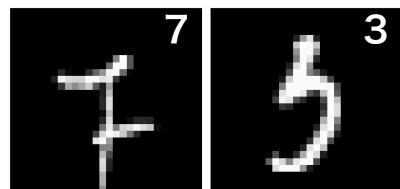
## Classification and Estimation of Conditional Probability<sup>7</sup>

### ► Advantages of Deep Learning (Deep Neural Networks)

- In recent years, it has been shown that, in the regression problem, the curse of dimensionality can be mitigated by using DNN models.  
e.g. Shmidt-Hieber (2020), Suzuki (2019)
- We can expect this benefit of DNN even for the classification problem!
- Example: MNIST handwriting data (DNN vs kNN)
  - For images that are difficult to distinguish,  
DNN provides small conditional probabilities.



DNN   4: 10%, 7: 90%   4: 10%, 8: 90%  
kNN   7: 80%, 9: 20%   4: 30%, 8: 20%, 9: 50%



7: 60%, 8: 40%   3: 60%, 4: 10%, 5: 30%  
1: 100%   3: 40%, 5: 50%, 9: 10%

## Semiparametric Estimation and Estimation of Conditional Probability<sup>8</sup>

- In the context of semiparametric estimation, we sometimes need to estimate conditional probabilities.
- In double machine learning (DML; Chernozhukov et al., 2018), we need to show the specific convergence rate ( $o_p(n^{-1/4})$ ) for the nonparametric estimator of the nuisance parameters.
- When constructing plug-in estimator using an estimator of conditional probability, it is also important that recover the theoretical property of the estimator of conditional probability.

# Related Works

## Theoretical Study of Classification Problem with Deep Learning<sup>10</sup>

### ► Results about the convergence of misclassification rate

- ✿ misclassification rate = the percentage of times a classifier is incorrect  
 $\mathbb{P}(Y \neq \hat{f}_n(\mathbf{X}))$

✿ misclassification rate is the expectation of 0-1 loss

- 0-1 loss :  $l(y, f(\mathbf{x})) = \mathbb{1}(y \neq f(\mathbf{x}))$

⇒ we want to minimize 0-1 loss.

✿ 0-1 loss is discontinuous, making optimization challenging.

⇒ convex surrogate losses are used instead

- e.g., hinge loss, logistic loss

✿ Research is being conducted on the convergence of misclassification rate when using deep learning models with various surrogate loss functions.

(Hu et al., 2020; Kim et al., 2021)

## Theoretical Study of Classification Problem with Deep Learning<sup>12</sup>

### ► Results about the convergence of conditional probability

✿ In classification problems using DNNs, you can output the conditional probabilities for each class by using the softmax function in the final layer.

✿ Natural evaluation metric : The difference between the expectation of negative log-likelihood of the estimator and the true conditional probability

$$\mathbb{E}_{\mathbf{X}} [\text{KL}(\boldsymbol{\eta}(\mathbf{X}) \| \hat{\mathbf{p}}(\mathbf{X}))] \quad (*)$$

- KL : Kullback-Leibler divergence

✿ The KL divergence can easily diverge, making it difficult to evaluate (\*).

- If  $\mathcal{F}_n$  contains all piecewise constant conditional class probabilities with at most two pieces or all piecewise linear conditional class probabilities with at most three pieces.  
⇒ (\*) diverge (Lemma 2.1 in Bos & Schmidt-Hieber, 2022)

✿ Various approaches have been proposed to address this issue.

(Ohn & Kim, 2022; Bos & Schmidt-Hieber, 2022)

## Theoretical Study of Classification Problem with Deep Learning<sup>11</sup>

### ► Results regarding the convergence of misclassification rate

► Hu et al. (2020, arXiv:2001.06892)

- **evaluation metric** : misclassification rate

- **loss function** : 0-1 loss, hinge loss

- **assumptions about population distribution** :

1. The conditional density function of  $\mathbf{X}$  given  $Y$  can be represented by DNN
2. The probability of data concentrating near the decision boundary is sufficiently low.

► Kim et al. (2021, Neural Networks)

- **evaluation metrics** : misclassification rate

- **loss function** : hinge loss, logistic loss

- **assumptions about population distribution** :

1. smooth decision boundary
2. smooth true conditional probability
3. margin condition: The density of input  $\mathbf{X}$  is low near the decision boundary.

## Theoretical Study of Classification Problem with Deep Learning<sup>13</sup>

### ► Results about the convergence of conditional probability

- Ohn and Kim (2022, Neural computation)
  - **evaluation metric** : KL divergence
  - **loss function** : logistic loss
  - **assumptions about population distribution** :  
The true conditional probability is bounded away from 0 and 1.
- Bos & Schmidt-Hieber (2022, Electronic Journal of Statistics)
  - **evaluation metric** : truncated KL divergence
  - **loss function** : negative log-likelihood
  - **assumptions about population distribution** : The true conditional probability belongs to the class  $\mathcal{H}$ , which is  $\alpha$ -SVB.  
$$\exists C > 0, \forall p \in \mathcal{H}; \mathbb{P}_{\mathbf{X}}(p_k(\mathbf{X}) \geq t) \leq Ct^\alpha$$
for all  $t \in (0, 1]$  and all  $k \in \{1, \dots, K\}$

## Remained Problems

- There is limited discussion on the theoretical properties of non-parametric estimation of conditional probabilities when the estimator lacks an explicit representation.
- Ohn & Kim (2022) : Results in limited situations.
- Bos & Schmidt-Hieber (2022) : By evaluating the truncated KL divergence instead of KL divergence, assumptions were removed to demonstrate a slightly weaker form of convergence.
  - There has been no discussion on the optimality of the derived rates.
- This work
  - Theoretical investigation of the properties of the NPMLE in classification problems.
  - Deriving the rate of convergence of the NPMLE in deep learning.
  - Prooving minimax optimality.

## Summary of the work

- Evaluating KL divergence can be challenging in NPMLE.  
e.g., in Ohn & Kim (2022), strong assumptions are required.  
⇒ Resolved the issue by employing the approach of van de Geer (2000) to evaluate the Hellinger distance directly. (Theorem 1, 2)
  - No strong assumptions are required.
  - There is a good insight into the proof.
- Application of the theoretical analysis of logistic regression with deep learning. (Theorem 3)
  - Avoiding the curse of dimensionality.
- Proved that the derived convergence rate is minimax optimal. (Theorem 4)
  - It is understood that classification problems are more challenging than regression problems.

## Main Results

## Notations and Assumptions

17

- Notations

- $\bar{\mathcal{F}}_n^{1/2}(\tilde{\mathbf{p}}, \delta) = \left\{ \frac{\mathbf{p} + \tilde{\mathbf{p}}}{2} : R\left(\frac{\mathbf{p} + \tilde{\mathbf{p}}}{2}, \tilde{\mathbf{p}}\right) \leq \delta^2, \mathbf{p} \in \mathcal{F}_n \right\}$

- $N_{p,B}(\delta, \mathcal{F}, Q)$ :

$\delta$ -bracketing number of function class  $\mathcal{F}$  for  $L^p(Q)$  metric

- Representing a certain "size" or "complexity" of a function class  $\mathcal{F}$ .

- $J_B(\delta, \bar{\mathcal{F}}_n^{1/2}(\tilde{\mathbf{p}}_n, \delta), \mu) = \int_{\delta^2/(2^{13}c_0)}^{\delta} \sqrt{\log N_{2,B}(u, \bar{\mathcal{F}}_n^{1/2}(\tilde{\mathbf{p}}_n, \delta), \mu)} du \vee \delta$

- $\mu$  is the product measure of a counting measure on  $\{1, \dots, K\}$  and the  $P_X$ .

- $J_B(\delta, \bar{\mathcal{F}}_n^{1/2}(\tilde{\mathbf{p}}_n, \delta), \mu)$  represents local complexity around  $\tilde{\mathbf{p}}_n$ .

⇒ Using local complexity allows for deriving tight inequalities.

(Bartlett, 2005)

## Notations and Assumptions

18

- Notations

- $H(P, Q)^2 := \frac{1}{2} \int \left( \sqrt{dP} - \sqrt{dQ} \right)^2$ : squared Hellinger distance

- $R(\boldsymbol{\eta}(\mathbf{X}), \mathbf{p}(\mathbf{X})) := \mathbb{E}_{\mathbf{X}} [H^2(\boldsymbol{\eta}(\mathbf{X}), \hat{\mathbf{p}}(\mathbf{X}))]$

• In this study, we discuss the convergence of  $R(\boldsymbol{\eta}(\mathbf{X}), \hat{\mathbf{p}}(\mathbf{X}))$  instead of

$$\mathbb{E}_{\mathbf{X}} [\text{KL}(\boldsymbol{\eta}(\mathbf{X}) \| \hat{\mathbf{p}}(\mathbf{X}))].$$

•  $R(\boldsymbol{\eta}(\mathbf{X}), \hat{\mathbf{p}}_n(\mathbf{X})) \leq \mathbb{E}_{\mathbf{X}} [\text{KL}(\boldsymbol{\eta}(\mathbf{X}) \| \hat{\mathbf{p}}_n(\mathbf{X}))]$

- Assumptions

- $\exists \tilde{\mathbf{p}}_n \in \mathcal{F}_n, \exists c_0, \forall \mathbf{x} \in \mathcal{X}; \frac{\eta_k(\mathbf{x})}{\tilde{p}_k(\mathbf{x})} \leq c_0^2, k = 1, \dots, K$

• The only assumption imposed in this study.

• If  $\exists c > 0, \forall \mathbf{p} \in \mathcal{F}_n; p_k \geq c, k = 1, \dots, K$  then the assumption is satisfied.

• Automatically satisfied with appropriately chosen network parameters in deep learning.

## Oracle Inequality

19

### ★ Oracle Inequality

#### Theorem 1 (Evaluation of Variance)

Consider the  $K$  class classification problem. Let  $\hat{\mathbf{p}}_n$  be a NPMLE. Take  $\Psi(\delta) \geq J_B(\delta, \bar{\mathcal{F}}_n^{1/2}(\tilde{\mathbf{p}}_n, \delta), \mu)$  in such a way  $\Psi(\delta)/\delta^2$  is a non-increasing function of  $\delta$ . Then for universal constant  $c$  and for

$$\sqrt{n}\delta_n^2 \geq c\Psi(\delta_n), \quad (1)$$

we have for all  $\delta \geq \delta_n$ ,

$$\mathbb{P}(R(\hat{\mathbf{p}}_n, \boldsymbol{\eta}) > 514(1 + c_0^2)(\delta^2 + R(\tilde{\mathbf{p}}_n, \boldsymbol{\eta}))) \leq c \exp\left(-\frac{n\delta^2}{c^2}\right)$$

• According to Theorem 1, if the hypothesis space  $\mathcal{F}_n$  is defined, and its local complexity  $J_B$  can be computed, then by solving the inequality (1) (the critical inequality), we can find  $\delta_n$ .

• Assuming a function class to which  $\boldsymbol{\eta}$  belongs, and evaluating the approximation error  $R(\tilde{\mathbf{p}}_n, \boldsymbol{\eta})$ , we can determine rate of convergence.

## Oracle Inequality

20

### ★ Oracle Inequality

#### Theorem 2 (oracle inequality)

Under the conditions of Theorem 1, we have for some universal constant  $c$ ,

$$\mathbb{E}(R(\hat{\mathbf{p}}_n, \boldsymbol{\eta})) \leq 514(1 + c_0^2)(\delta_n^2 + R(\tilde{\mathbf{p}}_n, \boldsymbol{\eta})) + \frac{c}{n}$$

• In Theorem 2,  $\delta_n$  represents the variance of the estimator, while  $R(\tilde{\mathbf{p}}_n, \boldsymbol{\eta})$  represents the bias of the estimator.

• This inequality is the expected value version of Theorem 1's inequality.

## Rate of Convergence in Deep Learning

21

- Settings of hypothesis space

- $\sigma(x) := \max(0, x)$  : activation function,  $\mathbf{v} = (v_1, \dots, v_r) \in \mathbb{R}$  : bias
  - $\sigma_{\mathbf{v}}(y_1, \dots, y_r) = (\sigma(y_1 - v_1), \dots, \sigma(y_r - v_r))$  : activation with bias
  - $L$  : depth of network,  $\mathbf{m} = (m_0, \dots, m_{L+1})$  : width of network
  - $W_j \in \mathbb{R}^{m_{j+1} \times m_j}$  : weight matrix,  $\mathbf{v}_i \in \mathbb{R}^{m_i}$  : bias vector
  - $\Phi$  : softmax function
  - Feedforward neural network (FFNN) :
- $$f : \mathbb{R}^{m_0} \rightarrow \mathbb{R}^{m_{L+1}}, \mathbf{x} \mapsto f(\mathbf{x}) = \Phi W_L \sigma_{\mathbf{v}_L} W_{L-1} \sigma_{\mathbf{v}_{L-1}} \cdots W_1 \sigma_{\mathbf{v}_1} W_0 \mathbf{x} \quad (\star)$$
- $\mathcal{F}(L, \mathbf{m}) := \left\{ f \text{ of the form } (\star) : \max_{j=0, \dots, L} \|W_j\|_{\infty} \vee |\mathbf{v}_j|_{\infty} \leq 1 \right\}$
- $\|W_j\|_0$  : the number of non-zero entries of  $W_j$
  - $|\mathbf{v}_j|_0$  : the number of non-zero entries of  $\mathbf{v}_j$
  - $\mathcal{F}(L, \mathbf{m}, s) := \left\{ f \in \mathcal{F}(L, \mathbf{m}) : \sum_{j=0}^L (\|W_j\|_0 + |\mathbf{v}_j|_0 \leq s) \right\}$

## Rate of Convergence in Deep Learning

23

### ★ Rate of Convergence in Deep Learning

- $\mathcal{F}_n = \mathcal{F}(L, \mathbf{m}, s)$

#### Theorem 3 (Rate of Convergence in Deep Learning)

Coonsider the  $K$  class classification problem. Let  $\hat{\mathbf{p}}_n$  be a NPMLE taking values in the network class  $\mathcal{F}(L, \mathbf{m}, s)$  satisfying

$$(i) n\phi_n \lesssim \min_{i=1, \dots, L} m_i, \quad (ii) s \asymp n\phi_n \log(n), \quad (iii) L \asymp \log(n).$$

Then, there exists constant  $C$  only depending on  $r, d, t, \beta$  such that

$$\mathbb{E}[R(\hat{\mathbf{p}}_n, \boldsymbol{\eta})] \leq C\phi_n L \log^2(n),$$

where  $\beta_i^* = \beta_i \prod_{l=i+1}^r (\beta_l \wedge 1)$ ,  $\phi_n = \max_{i=0, \dots, r} K^{\frac{\beta_i^* + 3t_i}{\beta_i^* t_i}} n^{-\frac{\beta_i^*}{\beta_i^* + t_i}}$

- The rate of convergence is  $\phi_n$  up to log-factor.
- **The rate of convergence is not dependent on dimension  $d$ . (avoiding the curse of dimensionality)**

## Rate of Convergence in Deep Learning

22

- Settings of underlying true conditional probability

- Hölder space with smoothness  $\beta$  :

$$\mathcal{C}^{\beta}(D, Q) = \left\{ f : D \subset \mathbb{R}^m \rightarrow \mathbb{R} : \right.$$

$$\sum_{\alpha: |\alpha| < \beta} \|\partial^{\alpha} f\|_{\infty} + \sum_{\alpha: |\alpha| = \lfloor \beta \rfloor} \sup_{\mathbf{x}, \mathbf{y} \in D, \mathbf{x} \neq \mathbf{y}} \frac{|\partial^{\alpha} f(\mathbf{x}) - \partial^{\alpha} f(\mathbf{y})|}{|\mathbf{x} - \mathbf{y}|_{\infty}^{|\beta| - \lfloor \beta \rfloor}} \leq Q \left. \right\}$$

- Composition Structured Functions :

$$\begin{aligned} \mathcal{G}_{\text{comp}}(r, \mathbf{d}, \mathbf{t}, \beta, Q) = \{f = \mathbf{g}_r \circ \cdots \circ \mathbf{g}_0 : \mathbf{g}_i = (g_{ij})_j : [a_i, b_i]^{d_i} \rightarrow [a_{i+1}, b_{i+1}]^{d_{i+1}}, \\ g_{ij} \in \mathcal{C}^{\beta_i}([a_i, b_i]^{t_i}, Q), \text{ for some } |a_i|, |b_i| \leq Q\} \\ \text{• } f(x_1, x_2, x_3) = g_1 \circ g_2(x_1, x_2, x_3) = g_{11}(g_{01}(x_1, x_3), g_{02}(x_1, x_2)) \end{aligned}$$

is an example when  $d_0 = 3, t_0 = 2, d_1 = t_1 = 2, d_2 = 1$ .

- $\eta_k \in \mathcal{G}_{\text{comp}}(r, \mathbf{d}, \mathbf{t}, \beta, Q)$ ,  $k = 1, \dots, K$

- In cases where the dimensions become significantly smaller in each compositional step ( $t_i \leq d_i$ ), it is sufficient to consider  $t_i$  dimensions.  
⇒ avoid the curse of dimensionality.

## Rate of Convergence in Deep Learning

24

### ★ Minimax Optimality

#### Theorem 4 (minimax optimality)

Consider the  $K$  class classification problem with  $X_j$  drawn from a distribution with Lebesgue density on  $[0, 1]^d$  which is lower and upper bounded by positive constants. For any nonnegative integer  $r$ , any dimension vectors  $\mathbf{d}$  and  $\mathbf{t}$  satisfying  $t_i \leq \min(d_0, \dots, d_{i-1})$  for all  $i$ , any smoothness vector  $\beta \in \mathbb{R}^r$  and all sufficiently large constants  $Q > 0$ , there exists a positive constant  $c$  such that

$$\inf_{\hat{\mathbf{p}}_n} \sup_{\eta_k \in \mathcal{G}_{\text{comp}}(r, \mathbf{d}, \mathbf{t}, \beta, Q), k=1, \dots, K} \mathbb{E}[R(\hat{\mathbf{p}}_n, \boldsymbol{\eta})] \geq c\phi_n.$$

- According to Theorem 4, the rate of convergence derived in Theorem 3 is minimax optimal up to log-factor.

# Rate of Convergence in Deep Learning

25

## ► Comparison with regression

- Consider regression.
- The true regression function  $f$  satisfies  $f \in \mathcal{G}_{\text{comp}}(r, d, t, \beta, Q)$
- The rate of convergence in regression :  $\max_{i=1, \dots, r} n^{-\frac{2\beta_i^*}{(2\beta_i^* + t_i)}}$ 
  - minimax optimal (Schmidt-Hieber, 2020)
- The rate of convergence in classification :  $\max_{i=1, \dots, r} n^{-\frac{\beta_i^*}{\beta_i^* + t_i}}$
- Regression achieves faster rate.  
⇒ Suggesting that estimating conditional probabilities from given covariates and labels is a more challenging problem than estimating regression functions.
- Remark
  - If the true conditional probability is bounded away from 0 and 1, classification achieves the same rate as regression.
  - The true function class can be arbitrarily changed and analyzed as long as it allows for bias evaluation, not limited to this example.

# Rate of Convergence in Deep Learning

26

## • Comparison with Related Work

	evaluation metrics	Loss function	Remarks
Hu et al. (2020)	misclassification rate	0-1 loss hinge loss	assume strong condition minimax optimal
Kim et al. (2021)	misclassification rate	hinge loss logistic loss	assume strong condition minimax optimal
Ohn & Kim (2022)	estimating conditional probability KL divergence	logistic loss	assume strong condition with regularization
Bos & Schmidt-Hieber (2022)	estimating conditional probability truncated KL divergence	negative log-likelihood	general condition
This work	estimating conditional probability Hellinger distance	negative log-likelihood	general condition minimax optimal

# Reference

27

- Sara A Geer, Sara van de Geer, and D Williams. Empirical Processes in M-estimation, volume 6. Cambridge university press, 2000.
- Evarist Giné and Richard Nickl. Mathematical Foundations of Infinite-Dimensional Statistical Models, volume 40. Cambridge University Press, 2016.
- Yongdai Kim, Ilsang Ohn, and Dongha Kim. Fast convergence rates of deep neural networks for classification. *Neural Networks*, 138:179–197, 2021.
- Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with ReLU activation function. *The Annals of Statistics*, 48(4):1875 – 1897, 2020.
- Taiji Suzuki. Adaptivity of deep ReLU network for learning in Besov and mixed smooth Besov spaces: optimal rate and curse of dimensionality. In International Conference on Learning Representations, 2019.
- Martin J Wainwright. High-dimensional statistics: A non-asymptotic viewpoint, volume 48. Cambridge University Press, 2019.
- Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, James Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, Volume 21, Issue 1, 1 February 2018, Pages C1–C68

# Reference

28

- Imaiizumi, M., & Fukumizu, K. (2019). Deep neural networks learn non-smooth functions effectively. In Proceedings of the International Conference on Artificial Intelligence and Statistics.
- Ilsang Ohn, Yongdai Kim; Nonconvex Sparse Regularization for Deep Neural Networks and Its Optimality. *Neural Comput* 2022; 34 (2): 476–517.