
Nonparametric estimation of conditional class probabilities using deep neural networks

Atsutomo Yara¹, Yoshikazu Terada^{1, 2}

¹Graduate School of Engineering Science, Osaka University,
²Center for Advanced Intelligence Project (AIP), RIKEN

0. Table of Contents

1. Introduction

- 1.1. Deep learning
- 1.2. Generalization error analysis of deep learning
- 1.3. Classification and conditional probability
- 1.4. Classification using deep learning
- 1.5. Semi-parametric estimation and conditional probability

2. Related works

- 2.1. Convergence of conditional probability
- 2.2. Overview of this study

3. Main results

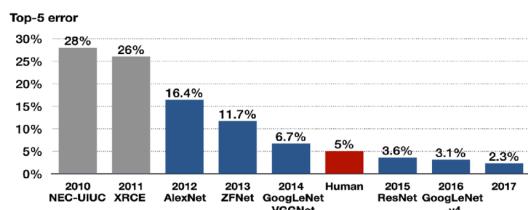
- 3.1. Overview of main results
- 3.2. Notations and assumptions
- 3.3. Non-parametric maximum likelihood estimation of conditional probability
- 3.4. Logistic regression using deep learning
- 3.5. Summary

1. Introduction

1.1 Deep learning

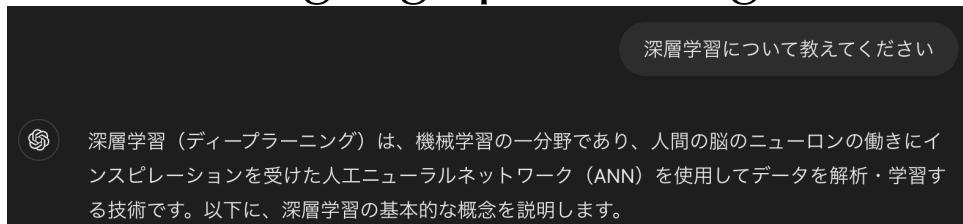
► Development of deep learning

- ✿ Deep learning has shown extremely high performance in various real-world problems:
 - Image recognition



Comparison of performance in ILSVRC.
(Gray) Pre-deep learning methods,
(Blue) Deep learning, (Red) Human, from Kang et al. (2020)

- Natural language processing



ChatGPT, OpenAI

- Image generation



1.1 Deep learning

► Importance of theoretical study

📌 A simple question arises

Q. Can deep learning solve any problem?

⇒ In fact, theoretical studies have provided **negative answers** to this question! (**No-free-lunch theorem, Slow rates of convergence**)

📌 **Claim of the No-free-lunch theorem (informal) :**

There are problems for which any learning algorithm performs poorly, meaning there is no universal learning method that works well for all the problems

⇒ **It is important to develop various methods and investigate when each method works well. This highlights the significance of theoretical research**

1.1 Deep learning

► Explanation of deep learning performance

- ➲ Theoretical studies to explain the performance of deep learning are also actively being pursued
- ➲ Major perspectives in theoretical research on deep learning:
 - Generalization ability
 - ✓ How well it performs on new data
 - Approximation ability
 - ✓ How efficiently it can be approximate functions
 - Optimization
 - ✓ Whether the parameters that fit the data can be obtained through optimization

1.2 Generalization error analysis in deep learning

► Statistical learning theory

- The goal of statistical learning theory is to provide a theoretical justification for machine learning methods.
- Evaluation of justification = Evaluation of generalization error
 - Small generalization error \Rightarrow "Good" estimator
- Points of view for evaluating generalization error
 - Does the generalization error converge to zero as the amount of training data increases?
 - If it does converge, how fast does it converge? (Convergence rate)
 - Optimality of convergence rate

1.2 Generalization error analysis in deep learning

► Empirical Risk Minimization, ERM

- 📌 Generalization error (Excess risk):

$$\mathbb{E}_{(X,Y)}[\ell(Y, \hat{f}(X))] - \inf_{f:\text{measurable}} \mathbb{E}_{(X,Y)}[\ell(Y, f(X))]$$

- 📌 The expected loss of new data is unknown \Rightarrow **ERM**

- 📌 ERM: Minimizes the loss over observed data \mathcal{D}_n

- The estimator \hat{f} is searched within a hypothesis space (model) \mathcal{F}
- In this study, the model \mathcal{F} is Deep Neural Networks (DNNs)

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(X_i))$$

1.2 Generalization error analysis in deep learning

► Bias-variance decomposition of generalization error

⌚ $L(f) = \mathbb{E}_{(X,Y)}[\ell(Y, f(X))]$: Expected risk

$$\begin{aligned} L(\hat{f}) - \inf_{f:\text{measurable}} L(f) \\ = L(\hat{f}) - \underbrace{\inf_{f \in \mathcal{F}} L(f)}_{\text{Variance}} + \underbrace{\inf_{f \in \mathcal{F}} L(f) - \inf_{f:\text{measurable}} L(f)}_{\text{Bias}} \end{aligned}$$

Variance

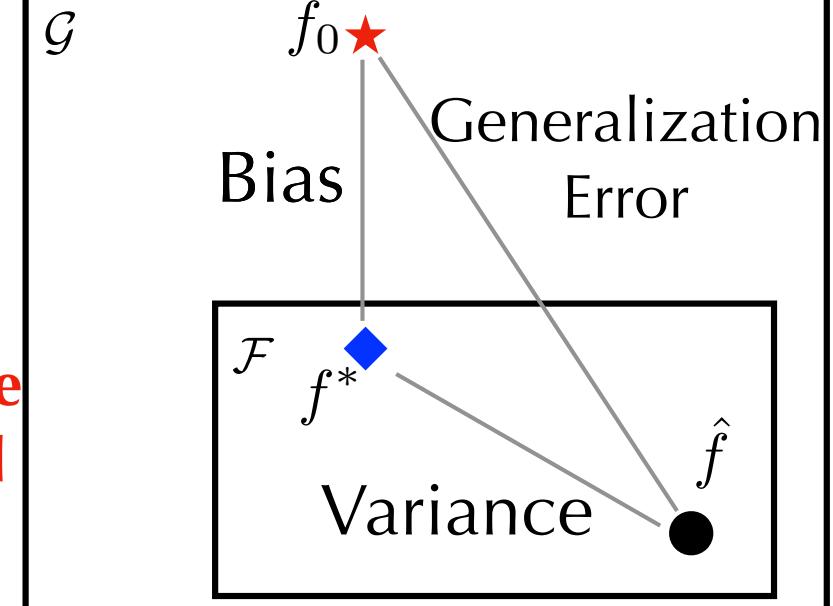
\mathcal{G} : Measurable functions

\mathcal{F} : Model

★ $f_0 = \min_{f \in \mathcal{G}} L(f)$

◆ $f^* = \min_{f \in \mathcal{F}} L(f)$

Bias



There is a trade-off between bias and variance
that depends on the complexity of the model

1.3 Classification and conditional probability

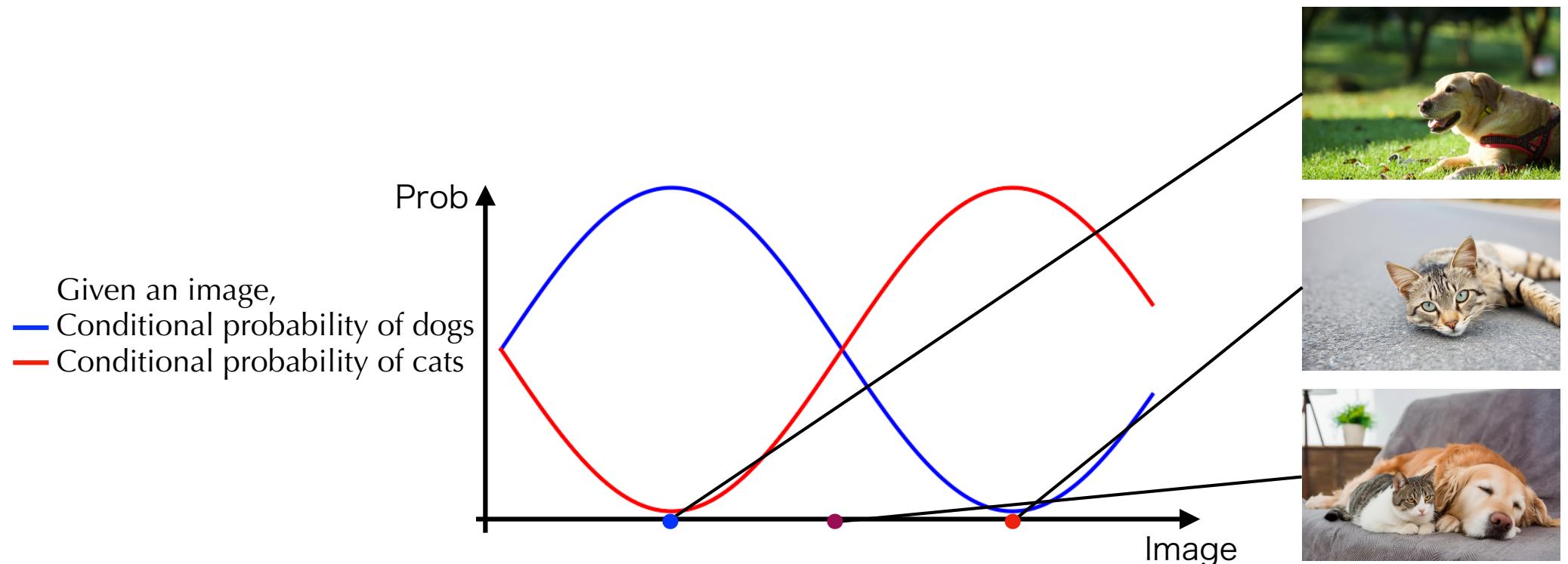
► Set-up of the classification problem

- We want to predict its label Y (e.g. whether it is ill or not) from input \mathbf{X} (e.g. medical images)
- True conditional class probability:

$$\eta_k(\mathbf{x}) := \mathbb{P}(Y = k \mid \mathbf{X} = \mathbf{x}), \boldsymbol{\eta}(\mathbf{x}) := (\eta_1(\mathbf{x}), \dots, \eta_K(\mathbf{x}))^\top$$

- Given \mathbf{X}, Y follows the following multinomial distribution

$$Y | \mathbf{X} = \mathbf{x} \sim \text{Multi}(\boldsymbol{\eta}(\mathbf{x}))$$

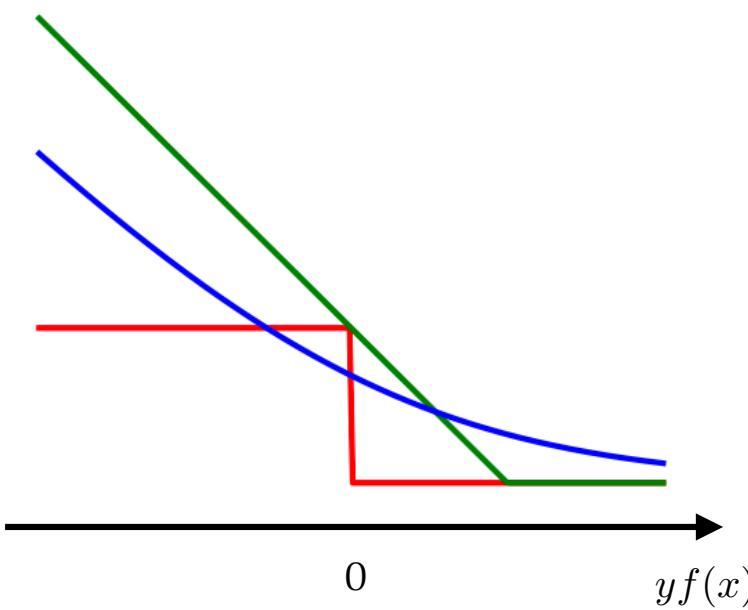


1.3 Classification and conditional probability

► Learning method

- Estimators are obtained by ERM
- There are various loss functions
 - 0-1 loss ← Most intuitive loss, but difficult to optimise
 $\ell(y, f(x)) := \mathbb{1}(yf(x) > 0)$

- Hinge loss
 $\ell(y, f(x)) := 0 \vee (1 - yf(x))$
- Logistic loss
 $\ell(y, f(x)) := \log(1 + e^{-yf(x)})$



1.3 Classification and conditional probability

► Advantage of conditional probability

- ✿ If hinge loss is used, only the classifier itself can be obtained
 - i.e. only label predictions can be made.
- ✿ It can be more convenient in some situations to estimate conditional probabilities
 - e.g. AI-based disease diagnosis
- ✿ If conditional probabilities can be estimated, they can be used to support human decision-making!
- ✿ Method for estimating conditional probabilities:
⇒ **Use the negative log-likelihood of a multinomial distribution as the loss function (Maximum Likelihood Estimation)**

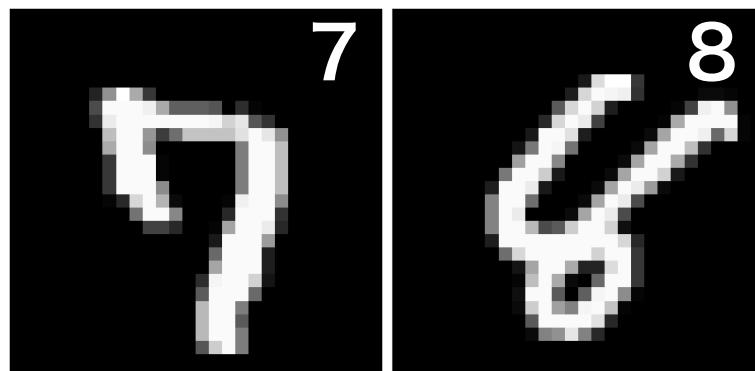
1.4 Classification using deep learning

► Classification using deep learning

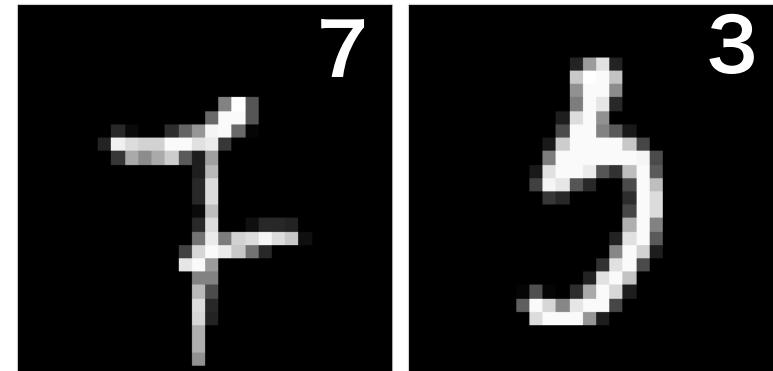
- In classification problems using deep learning, it is possible to estimate conditional probabilities by using a DNN with a softmax function in the output layer

$$\text{softmax}_k(x_1, \dots, x_K) := \frac{e^{x_k}}{\sum_{j=1}^K e^j}$$

- An example of classification in MNIST (comparison with kNN)
 - Conditional probabilities for digits that are difficult to distinguish are estimated to be smaller



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%

1.5 Semi-parametric estimation and conditional probability

► Relationship with Double/Debiased Machine Learning (DML)

- ✿ In the context of causal inference, the propensity score (i.e., the conditional probability of assignment given covariates) plays an important role as a nuisance parameter in estimating causal effects
- ✿ When the nuisance parameter is complex, machine learning (ML) methods can be used to provide a natural estimator
- ✿ However, the bias and overfitting that can result from using ML may negatively impact the estimator
 - The slow convergence rate (the $n^{-1/2}$ -consistency does not hold)
⇒ DML (Chernozhukov et al., 2018) can solve this problem!

1.5 Semi-parametric estimation and conditional probability

► Relationship with Double/Debiased Machine Learning (DML)

- ✿ In DML, it is necessary to demonstrate a $o_p(n^{-1/4})$ convergence rate for the nonparametric estimation of nuisance parameters
⇒ For example, when estimating propensity scores using deep learning, it is essential to analyze the generalization error of the estimator with deep learning
- ✿ In this study, we derive the convergence rate for the estimation of conditional probabilities using deep learning

2. Related works

2.2 Convergence of conditional probability

- The difficulty of theoretical analysis in maximum likelihood estimation of conditional probabilities

- 💡 Generalization error:

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

- KL: Kullback-Leibler divergence

- 💡 Since the KL divergence can easily diverge, directly evaluating (*) is challenging.

- 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
⇒ (*) diverges. (Lemma 2.1 in Bos & Schmidt-Hieber, 2022)

- 💡 Various approaches have been proposed to address this problem (Ohn & Kim, 2022; Bos & Schmidt-Hieber, 2022; Bilodeau et al., 2023)

2.4 Overview of this study

- In nonparametric maximum likelihood estimation, evaluating the KL divergence is challenging
 - In studies such as Ohn and Kim (2022), strong assumptions are required...
⇒ By using the approach of van de Geer (2000) and directly evaluating the Hellinger distance, the issue can be resolved (Theorem 1)
 - No strong assumptions are needed, and the proof is more straightforward
 - Remark: In Bos & Schmidt-Hieber (2022), the issue is avoided by considering a truncated KL divergence
 - Remark: In Bilodeau et al. (2023), the issue is avoided by using a transformed estimator of the MLE.
- Application to the theoretical analysis of logistic regression using deep learning (Theorem 2)
 - It suggests the possibility of avoiding the curse of dimensionality even in classification problems
 - Results have been demonstrated with non-sparse DNNs
- It has been proven that deep learning is minimax optimal (Theorem 3)

2.4 Overview of this study

► Comparison with related works

	Metric	Loss functions	Model	Notes
Kim et al.	Misclassification rate	Hinge loss Logistic loss	Sparse DNN	Only classifier Minimax optimal
Ohn & Kim	KL divergence	Logistic loss	Sparse DNN	Strong assumption With regularization
Bos & Schmidt-Hieber	Truncated KL divergence	Negative log likelihood	Sparse DNN	Weak assumption
Bildeau et al.	KL divergence	Negative log likelihood	Well-specified	Transformed MLE
This work	Hellinger distance	Negative log likelihood	Non-sparse DNN	Weak assumption Minimax optimal

3. Main results

3.2 Notations and assumptions

► Notations

- $H(P, Q)^2 := \frac{1}{2} \int \left(\sqrt{dP} - \sqrt{dQ} \right)^2$: The squared Hellinger distance
- $R(\boldsymbol{\eta}(\mathbf{X}), \mathbf{p}(\mathbf{X})) := \mathbb{E}_{\mathbf{X}} [H^2(\boldsymbol{\eta}(\mathbf{X}), \mathbf{p}(\mathbf{X}))]$

✿ In this study, we discuss convergence with respect to

$$R(\boldsymbol{\eta}(\mathbf{X}), \hat{\mathbf{p}}(\mathbf{X})) \text{ instead of } \mathbb{E}_{\mathbf{X}} [\text{KL}(\boldsymbol{\eta}(\mathbf{X}) \parallel \hat{\mathbf{p}}(\mathbf{X}))]$$

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

► Assumptions

- $\exists \tilde{\mathbf{p}}_n \in \mathcal{F}_n, \exists c_0 > 0; \frac{\eta_k(\mathbf{x})}{\tilde{p}_{n,k}(\mathbf{x})} \leq c_0^2, \forall n \in \mathbb{N}, k = 1, \dots, K$

✿ This is automatically satisfied when estimation is performed using DNNs

3.3 Maximum likelihood estimation of conditional probability

★ Oracle inequality

Theorem 1 (Oracle inequality)

Consider a classification problem with K classes. Let \hat{p}_n be the maximum likelihood estimator. Under some conditions, we have

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

- ⌚ δ_n corresponds to the complexity of the model \mathcal{F}_n
⇒ The more complex the model, the larger the variance becomes
- ⌚ The more complex the model, the smaller the bias becomes (as a more complex model is more likely to capture the true conditional probabilities)
- ⌚ If the complexity of the model and the bias can be determined, the convergence rate of the MLE can be derived!

3.4 Logistic regression using deep learning

► DNN model

- $\sigma(x) := \max(0, x)$: ReLU activation, $\mathbf{v} = (v_1, \dots, v_r) \in \mathbb{R}$: 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, $\mathbf{m} = (m_0, \dots, m_{L+1})$: Width of each layer,
- $W_j \in \mathbb{R}^{m_{j+1} \times m_j}$: Weight matrix,
- Φ : softmax function

$$f : \mathbb{R}^{m_0} \rightarrow \mathbb{R}^{m_{L+1}}, \quad \mathbf{x} \mapsto f(\mathbf{x}) = \Phi W_L \sigma_{\mathbf{v}_L} W_{L-1} \sigma_{\mathbf{v}_{L-1}} \cdots W_1 \sigma_{\mathbf{v}_1} W_0 \mathbf{x} \quad (\star)$$

$$\mathcal{F}(L, \mathbf{m}, B) := \left\{ \text{f of the form } (\star) : \max_{j=0, \dots, L} \|W_j\|_\infty \vee |\mathbf{v}_j|_\infty \leq B \right\}$$

3.4 Logistic regression using deep learning

- Underlying function space for the true conditional probability
 - Hölder space with smoothness β

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

- Composition structured functions

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

3.4 Logistic regression using deep learning

- ▶ **Underlying function space for the true conditional probability**

- Composition structured functions

$$\mathcal{G}_{\text{comp}}(r, \mathbf{d}, \mathbf{t}, \boldsymbol{\beta}, Q) = \{f = \mathbf{g}_r \circ \cdots \circ \mathbf{g}_0 : \mathbf{g}_i = (g_{ij})_j : [a_i, b_i]^{d_i} \rightarrow [a_{i+1}, b_{i+1}]^{d_{i+1}}, \\ g_{ij} \in \mathcal{C}^{\beta_i}([a_i, b_i]^{t_i}, Q), \text{ for some } |a_i|, |b_i| \leq Q\}$$

• $f(\underline{x_1, x_2, x_3}) = g_1 \circ \mathbf{g}_2(x_1, x_2, x_3) = g_{11}(\underline{g_{01}(x_1, x_3)}, \underline{g_{02}(x_1, x_2)})$

3 dim 2 dim

• **There is an underlying low-dimensional structure**

• It contains a reasonable example of conditional probability

$$\eta_k(x_1, \dots, x_d) = \Psi_k \left(\sum_{i=1}^d f_{ki}(x_i) \right)$$

where $\Psi_k(\mathbf{x}) = e^{x_k} / \sum_{i=1}^d e^{x_i}$

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

3.4 Logistic regression using deep learning

★ Convergence rate

Theorem 3 (Convergence rate)

Consider the family of neural networks $\mathcal{F}(L, \mathbf{m}, B)$ that satisfy the following conditions as the model.

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

Then, there exists a constant C depending only on r, d, t, β, Q, K s.t.

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

For sufficiently large n . Here,

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

- The convergence rate is ϕ_n up to log-factor
- **There is no restriction for the sparseness of DNN**
- **The convergence rate does not depend on the dimension d , thus avoiding the curse of dimensionality**
- Deep learning can effectively capture the underlying low-dimensional structure (the structure of composite functions).

3.4 Logistic regression using deep learning

► Why sparse DNN?

- Many studies assume sparsity in DNN models
 - e.g. Suzuki (2019), Schmidt-Hieber (2020)
- Typically, generalization error is bounded by

$$\frac{\text{"Complexity"}^2}{n} + \text{Bias}^2$$

- The complexity increases as the number of non-zero parameters in the DNN grows.
⇒ Imposing sparsity on the DNN leads to a tighter bound
- In this study, we derived a tighter convergence rate without imposing sparsity by using the tighter bias bound from Kohler & Langer (2021)

3.4 Logistic regression using deep learning

★ Minimax optimality

- ⌚ Is the rate derived in Theorem 2 optimal in some sense?

Theorem 3 (Minimax optimality)

Under some conditions, we have

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

where infimum is taken over all estimators.

- ⌚ The convergence rate derived in Theorem 3 is minimax optimal up to log factor

3.5 Summary

► Summary

- ★ We derived a general inequality for evaluating variance in the maximum likelihood estimation of conditional probabilities (Theorem 1)
- ★ As an application of Theorem 1, we derived the convergence rate for logistic regression in non-sparse deep learning models (Theorem 2)
- ★ We demonstrated the minimax optimality of the rate derived in Theorem 2 (Theorem 3)
- ★ The paper is available in arXiv
URL: <https://arxiv.org/abs/2401.12482>

Appendix

1.2 Generalization error analysis in deep learning

► Slow rates of convergence

📌 As an example, consider a binary classification problem:

- $\mathcal{X} \subset \mathbb{R}^d$: Input space, $\mathcal{Y} = \{0, 1\}$: Output space
- P : Probability 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

📌 **Goal: Construct a function $\hat{f}_n : \mathcal{X} \rightarrow \mathcal{Y}$ from \mathcal{D}_n that predicts Y well when new data X is given**

- $L_n = \mathbb{P}(\hat{f}_n(\mathbf{X}) \neq Y \mid \mathcal{D}_n)$: Misclassification rate (the probability of incorrect classification)
- $L^* = \inf_{f:\text{measurable}} \mathbb{P}(f(\mathbf{X}) \neq Y)$: Optimal misclassification rate

1.2 Generalization error analysis in deep learning

► Slow rates of convergence

Theorem 0 (Slow rates of convergence, Devroye, 1982)

Let (a_n) be a sequence of positive numbers converging to zero such that $1/16 \geq a_1 \geq a_2 \geq \dots$. For any learning method, there exists a distribution of (X, Y) satisfying $L^* = 0$ such that

$$\forall n \in \mathbb{N}; \mathbb{E}[L_n] \geq a_n$$

- 💡 No learning method is universally optimal; the performance of any method can deteriorate for some distributions.

3.2 記号と仮定

▶ 記号

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

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

関数クラス \mathcal{F} の $L^p(Q)$ ノルムに関する δ -bracketing number

- 関数クラス \mathcal{F} のある種の"大きさ"や"複雑度"を表す

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

- μ は $\{1, \dots, K\}$ 上の数え上げ測度と P_X の直積測度

- $J_B(\delta, \bar{\mathcal{F}}_n^{1/2}(\tilde{\mathbf{p}}_n, \delta), \mu)$ は $\tilde{\mathbf{p}}_n$ の近くで見たときの局所的な複雑度

⇒ 局所的な複雑度を用いることで tight な不等式を導く

(van de Geer, 2000; Bartlett et al., 2005)

3.3 条件付き確率のノンパラメトリック最尤推定

★オラクル不等式

Theorem 1 (オラクル不等式 確率ver.)

Kクラスの分類問題を考える. \hat{p}_n をNPMLEとする.

$\Psi(\delta) \geq J_B(\delta, \bar{\mathcal{F}}_n^{1/2}(\tilde{p}_n, \delta), \mu)$ を $\Psi(\delta)/\delta^2$ が δ の非増加関数となるようにと
る. δ_n とある普遍定数 c は

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

を満たすとする. このとき

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

3.3 条件付き確率のノンパラメトリック最尤推定

★オラクル不等式（推定量の汎化性能を評価）

Theorem 2（オラクル不等式 期待値ver.）

Kクラスの分類問題を考える. \hat{p}_n をNPMLEとする.

$\Psi(\delta) \geq J_B(\delta, \bar{\mathcal{F}}_n^{1/2}(\tilde{p}_n, \delta), \mu)$ を $\Psi(\delta)/\delta^2$ が δ の非増加関数となるようにと
る. このとき, ある普遍定数 c と

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

を満たす δ_n , 全ての $\delta \geq \delta_n$ に対して以下が成り立つ.

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

Variance **Bias**

- ⌚ δ が Variance の評価に対応
 - ⌚ J_B はモデル \mathcal{F}_n の "大きさ" = "複雑さ" を表す
- ⇒ モデルが複雑なほど Variance は大きくなってしまう

3.3 条件付き確率のノンパラメトリック最尤推定

★オラクル不等式（推定量の汎化性能を評価）

Theorem 2（オラクル不等式 期待値ver.）

Kクラスの分類問題を考える。 \hat{p}_n をNPMLEとする。

$\Psi(\delta) \geq J_B(\delta, \bar{\mathcal{F}}_n^{1/2}(\tilde{p}_n, \delta), \mu)$ を $\Psi(\delta)/\delta^2$ が δ の非増加関数となるようにと
る。このとき、ある普遍定数 c と

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

を満たす δ_n 、全ての $\delta \geq \delta_n$ に対して以下が成り立つ。

$$\mathbb{E}[R(\hat{p}_n, \eta)] \leq 514(1 + c_0^2) \left(\frac{\delta_n^2}{n} + \frac{R(\tilde{p}_n, \eta)}{\text{Variance}} \right) + \frac{c}{n},$$

Variance **Bias**

◆ モデルが複雑なほど Varianceは大きくなってしまう

Bias-Variance

◆ 一方、モデルが複雑なほど Bias $R(\tilde{p}_n, \eta)$ は小さくなる

トレードオフ

(モデルが複雑なほど真の条件付き確率を含みやすくなるため)

3.3 条件付き確率のノンパラメトリック最尤推定

★オラクル不等式（推定量の汎化性能を評価）

Theorem 2（オラクル不等式 期待値ver.）

Kクラスの分類問題を考える。 \hat{p}_n をNPMLEとする。

$\Psi(\delta) \geq J_B(\delta, \bar{\mathcal{F}}_n^{1/2}(\tilde{p}_n, \delta), \mu)$ を $\Psi(\delta)/\delta^2$ が δ の非増加関数となるようにと
る。このとき、ある普遍定数 c と

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

を満たす δ_n 、全ての $\delta \geq \delta_n$ に対して以下が成り立つ。

$$\mathbb{E}[R(\hat{p}_n, \eta)] \leq 514(1 + c_0^2) \left(\frac{\delta_n^2}{n} + \frac{R(\tilde{p}_n, \eta)}{\text{Variance Bias}} \right) + \frac{c}{n},$$

◆ モデルの複雑さ J_B とBias $R(\tilde{p}_n, \eta)$ の評価ができれば

NPMLEの収束レートがわかる！

◆ Theorem 2は深層学習に限らず一般のNPMLEに対して成立