

深層学習の数理

鈴木大慈

東京大学大学院情報理工学系研究科数理情報学専攻

理研AIP



2023年(令和5年)7月18日

東北大学集中講義

所属

- 東京大学大学院情報理工学系研究科数理情報学専攻・准教授
- 東大次世代知能科学研究センター研究部門研究者（研究知能部門）
- 理化学研究所 革新知能統合研究センター 深層学習理論チーム チームリーダー



鈴木大慈
情報理工学系研究科

専門

- 機械学習, 数理統計学, 統計的学習理論

主な研究内容

- 深層学習を含む様々な学習機構について理論的側面から研究を進めています。学習理論を通じて各種学習手法の汎化性能や学習アルゴリズムの収束性能を解明し複雑な学習過程の本質への理解を深め、理論をもとに新しい機械学習手法の構築や応用への還元を行っています。また、確率的最適化などの方法により大規模かつ複雑な機械学習問題を効率的に解く手法の開発も行っています。

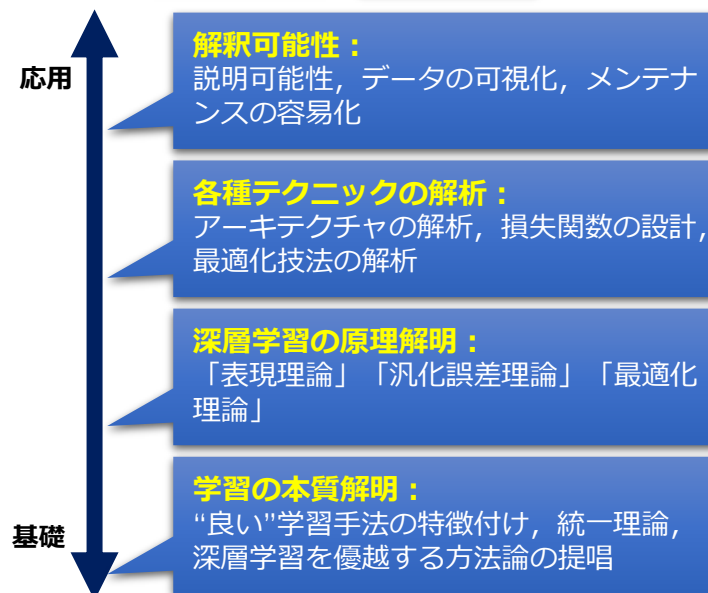


著書/授賞

- 『確率的最適化（機械学習プロフェッショナルシリーズ）』講談社，2015年8月8日。
- 金森敬文，鈴木大慈，竹内一郎，佐藤一誠：『機械学習のための連続最適化（機械学習プロフェッショナルシリーズ）』講談社，2016年12月7日。
- 文部科学大臣表彰・若手科学者賞「深層学習の原理解明に向けた統計的学習理論の研究」。文部科学省，2020年4月7日。
- 第11回日本統計学会研究業績賞（2017年度）。2017年9月5日。
- Satoshi Hayakawa and Taiji Suzuki: 日本神経回路学会論文賞。日本神経回路学会，2021年9月23日。
- 日本応用数理学会，ベストオーサー賞（論文部門）。2019年9月4日。

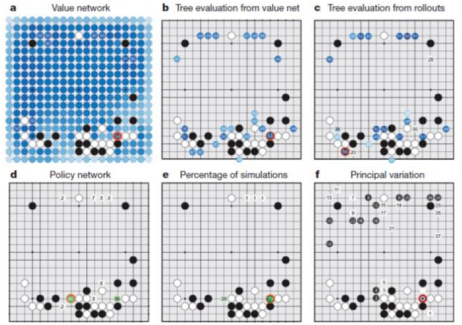
研究室URLとメール連絡先

- <http://ibis.t.u-tokyo.ac.jp/suzuki/>
- taiji@mist.i.u-tokyo.ac.jp



様々なタスクで高い精度
なぜ？

AlphaGo/Zero



[Silver et al. (Google Deep Mind): Mastering the game of Go with deep neural networks and tree search, Nature, 529, 484—489, 2016]

Image recognition



[He, Gkioxari, Dollár, Girshick: Mask R-CNN, ICCV2017]

Large language model

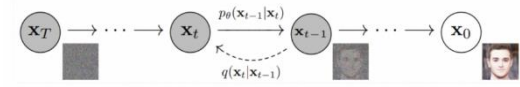
The diagram shows a sequence of tokens: "a robot must obey the orders given it". Below this, a "Transformer Decoder" is shown with 96 layers. A text input "what is the fermat's little theorem" is fed into the decoder. To the right, a screenshot of a ChatGPT response is shown, providing a detailed explanation of Fermat's Little Theorem. A vertical label "Learning efficiency of few:" is on the right side of the screenshot.

[Alammar: How Animations, <https://jalammar.github.io/how-gpt-works-visualizations-animations/>]

[ChatGPT. OpenAI2022]

[Brown et al. "Language Models are Few-Shot Learners", NeurIPS2020]

Generative models (diffusion models)



[Ho, Jain, Abbeel: Denoising Diffusion Probabilistic Models. 2020]



Stable diffusion, 2022.



Jason Allen "Théâtre D'opéra Spatial" generated by **Midjourney**. Colorado State Fair's fine art competition, 1st prize in digital art category

なぜ深層学習はうまくいくのか？

- 「〇〇法が良い」という様々な仮説の氾濫.
- 世界的課題
- 原理解明
- どうすれば“良い”学習が実現できるか？→新手法の開発

学会の問題意識



Ali Rahimi's talk at NIPS(NIPS 2017 Test-of-time award presentation)



Ali Rahimi's talk at NIPS(NIPS 2017 Test-of-time award presentation)

Ali Rahimi's talk at NIPS2017 (test of time award).
“Random features for large-scale kernel methods.”

“錬金術”という批判

民間の問題意識

- 中で何が行われているか分からないものは用いたくない.
- 企業の説明責任. 深層学習のホワイトボックス化.

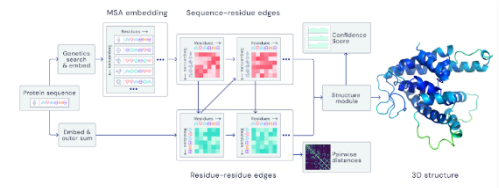


理論の必要性

深層学習(AI)の研究

応用

応用：AI手法の各種問題への応用
画像生成，パターン認識，たんぱく質構造予測



物理の対応物

半導体，新材料の開発

方法論：各種機械学習手法の開発

損失関数の設計，正則化法の開発，学習アルゴリズムの開発

各種方法論の定式化，学習アルゴリズムの開発

理論：統計的学習理論，最適化理論

深層学習の理論，収束レート解析，最適化アルゴリズム

各種機械学習手法の原理解明，最適性の理論的保証，
アルゴリズムの数理研究

素粒子論など

我々の得意分野

基礎

深層学習の理論概観



解釈可能性 :

説明可能性, データの可視化, メンテナンスの容易化

各種テクニックの解析 :

アーキテクチャの解析, 損失関数の設計, 最適化技法の解析

深層学習の原理解明 :

表現理論, 汎化誤差理論, 最適化の収束理論

学習の本質解明 :

“良い”学習手法の特徴付け, 統一理論, 深層学習を優越する方法論の提唱

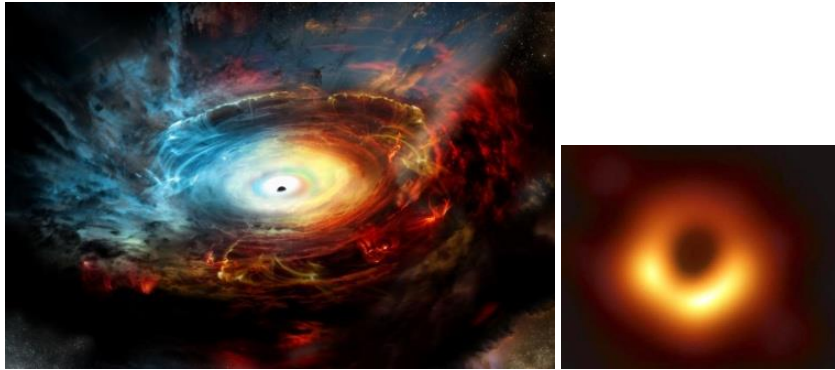
理論を通して深層学習の不可思議な挙動を理解したい.

- 説明責任
- 可能性と限界の把握
- 学習手法設計の指針

応用から基礎まで広い範囲に“理論”は遍在する.

今日の範囲

物理学



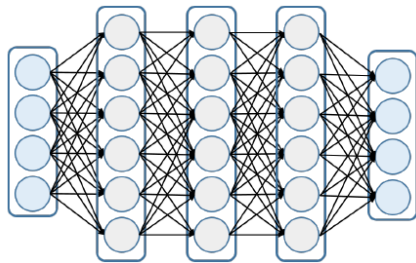
現象

$$R_{\mu\nu} - \frac{1}{2}g_{\mu\nu}R = \frac{8\pi G}{c^4}T_{\mu\nu}$$

- 相対性理論
- 量子力学
- リーマン幾何
- 関数解析

数学

機械学習 (情報学)



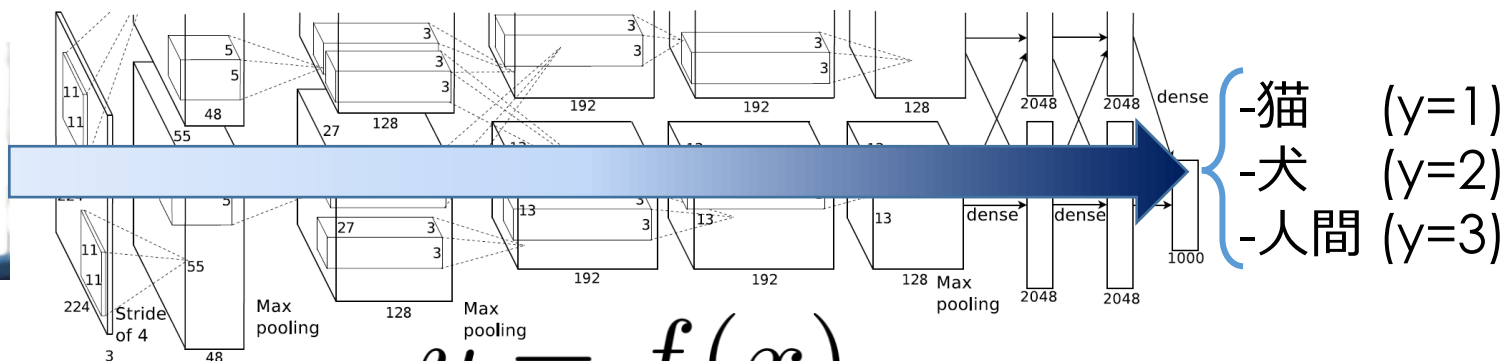
深層学習

数学者・物理学者も参入

- 確率論
- 関数解析
- Wasserstein幾何
- 熱拡散方程式
- 統計学
- 最適化理論
- 数値計算

数学

画像



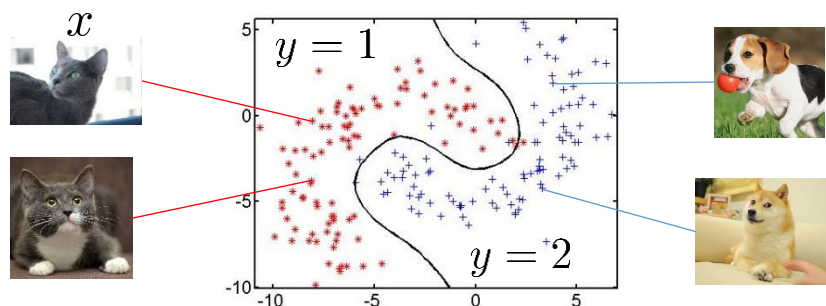
x

ベクトル $\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$

$$y = f(x)$$

y

ベクトル $\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_M \end{bmatrix}$



学習：「関数」をデータに当てはめる

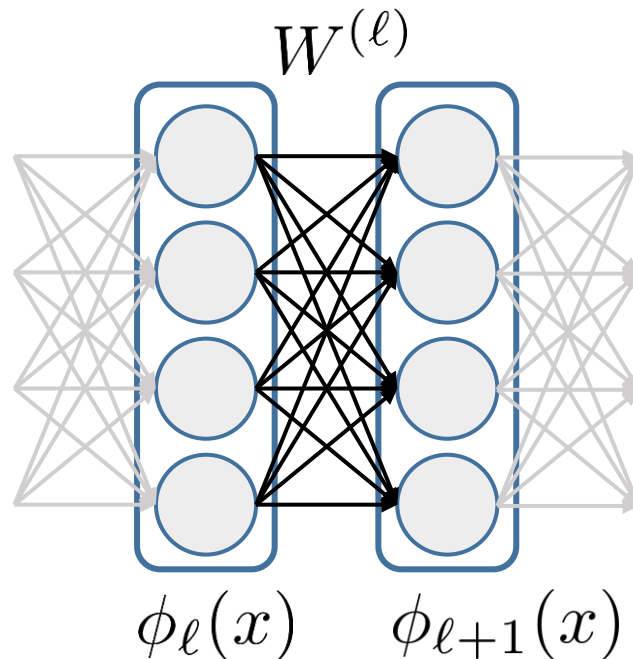
モデル：関数の集合（例：深層NNの表せる関数の集合）

アフィン変換 + 活性化関数

$$\phi_{\ell+1}(x) = \eta(W^{(\ell)}\phi_{\ell}(x) + b^{(\ell)})$$

$$W^{(\ell)} \in \mathbb{R}^{m_{\ell+1} \times m_{\ell}}$$

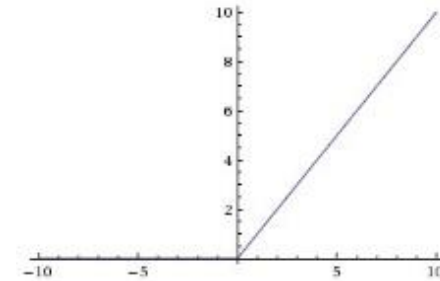
$$b^{(\ell)} \in \mathbb{R}^{m_{\ell+1}}$$



活性化関数

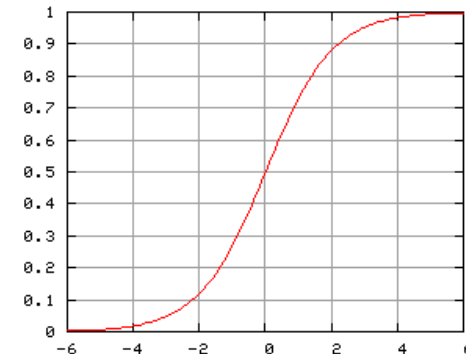
- ReLU (Rectified Linear Unit)

$$\eta(u) = \max\{u, 0\}$$



- シグモイド関数

$$\eta(u) = \frac{1}{1 + e^{-u}}$$



訓練誤差と汎化誤差

パラメータ θ : ネットワークの構造を表す変数

損失関数 $\ell(Y, f(X, \theta))$: パラメータ θ がデータをどれだけ説明しているか

予測誤差 : 損失の期待値

$$E[\ell(Y, f(X, \theta))]$$

$$L(\theta)$$

本当に最小化したいもの.

訓練誤差 : 有限個のデータで代用

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i, \theta))$$

$$\hat{L}(\theta)$$

代わりに最小化するもの.

(訓練データはテストデータと同じ分布に従っていると仮定)

この二つには大きなギャップがある.

[過学習]

※クラスタリング等, 教師なし学習も尤度を使ってこのように書ける.

学習理論の設定

- 汎化ギャップ(汎化誤差)と余剰誤差

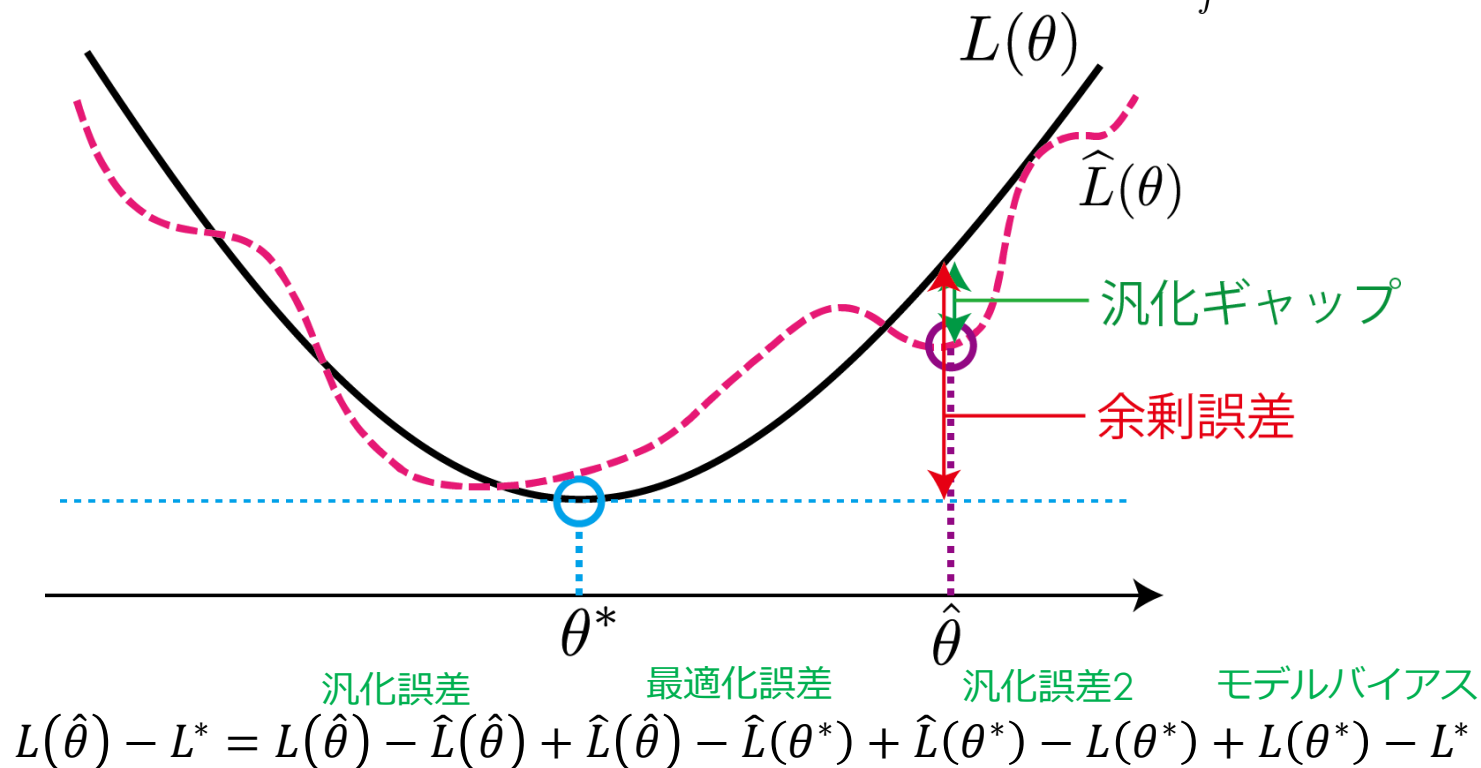
Generalization gap

Excess risk

汎化誤差: $L(\hat{\theta}) - \hat{L}(\hat{\theta})$

余剰誤差: $L(\hat{\theta}) - \inf_{\theta} L(\theta)$

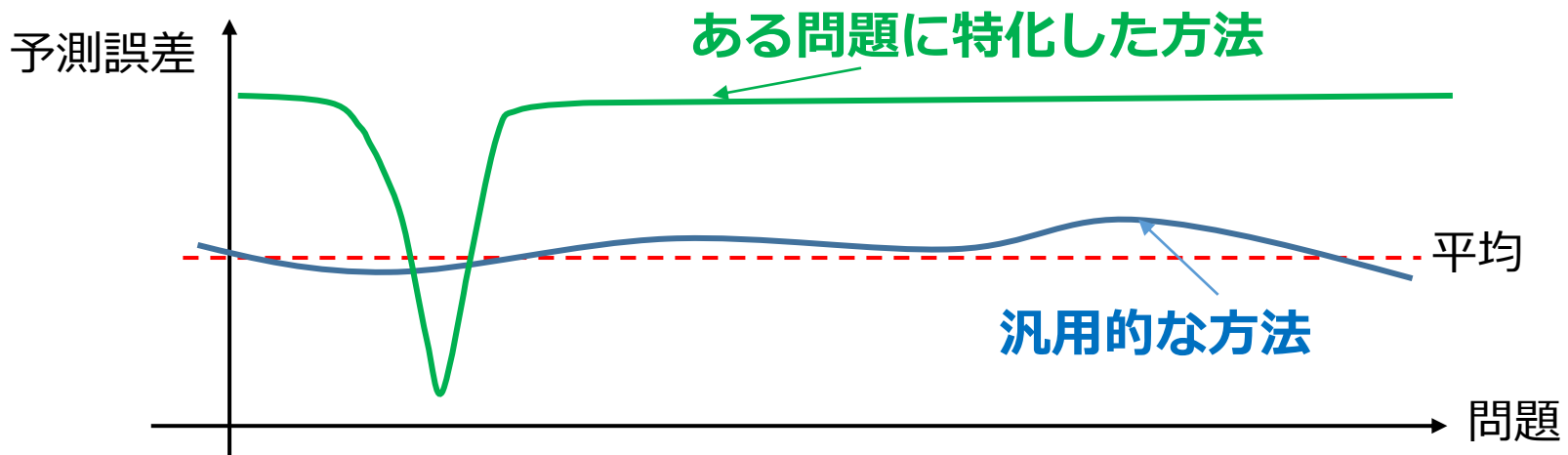
もしくは $L(\hat{\theta}) - \inf_f \mathbb{E}[\ell(Y, f(X))]$



学習機の複雑さと学習能力

- No free lunch theorem

「あらゆる問題で最高の性能を出す汎用的学習機は実現不可能であり，ある問題に特殊化された手法に勝てない」



学習手法は「どこかを“最良”する必要がある」
→ モデリングの重要性 (オッカムの剃刀)

DLも例外ではない: CNN, Transformer, ResNet, ...

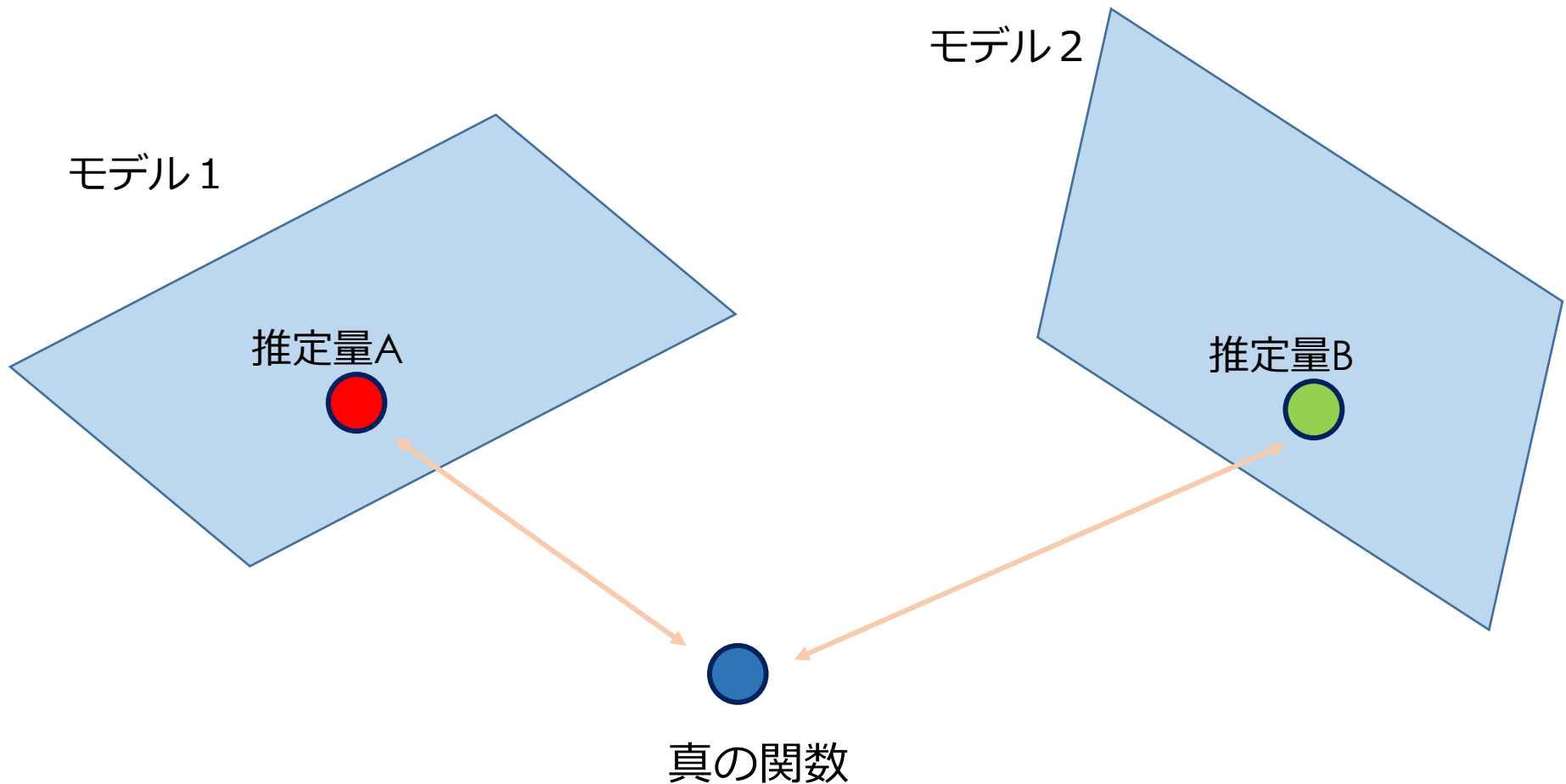
William of Ockham : 1285-1347. スコラ学の神学者, 哲学者.

No free lunch theorem: [D.H.Wolpert and W.G. Macready: 1995,1997][Y.C. Ho and D.L. Pepyne: 2002]

「なんでもできる手法は個別問題で負ける」

「必要以上に複雑なモデルを当てはめると失敗する」

(注：すでに学習済みのモデルに関する話ではない。大量のデータで学習した汎用モデルが個別問題でも勝つことはある。あくまでゼロ知識の状態から学習を始めた場合の精度比較である。)



各推定量は対応するモデル内に値を取るし、正則化や事前分布の影響でモデル内でも取りやすい値と取りにくい値がある。
→ そういった場所を最良しているということ。

リスクの概念

・ 真によらず常に予測誤差を最小化する推定量は存在しない。

「汎化誤差を最小化する」みたいな文言はそれ自体厳密な意味はない。

➤ 「最適」な推定量をどう特徴づけるか？

\mathcal{P} : 真の分布のモデル

$P^* \in \mathcal{P}$: 真の分布

$D^n = (x_i, y_i)_{i=1}^n$: 訓練データ

$\hat{\theta}, \tilde{\theta}$: 推定量

「真の分布のモデル」は本当に実際のデータ生成過程がそのモデルに入っている必要はなく仮想的なもので良い。あくまで、その推定量がどの設定で最適であるかを特徴づけるための仮想的モデルと考えてよい。

■ ミニマックス最適性

$$\sup_{P^* \in \mathcal{P}} \mathbb{E}_{D^n \sim P^*} [L(\hat{\theta})] = \inf_{\tilde{\theta}: \text{Estimator}} \sup_{P^* \in \mathcal{P}} \mathbb{E}_{D^n \sim P^*} [L(\tilde{\theta})]$$

(最悪誤差が最小)

■ 許容性 つぎのような $\tilde{\theta}$ が存在しない:

$$\mathbb{E}_{D^n \sim P^*} [L(\tilde{\theta})] \leq \mathbb{E}_{D^n \sim P^*} [L(\hat{\theta})] \quad (\forall P^* \in \mathcal{P})$$

$$\mathbb{E}_{D^n \sim P^*} [L(\tilde{\theta})] < \mathbb{E}_{D^n \sim P^*} [L(\hat{\theta})] \quad (\exists P^* \in \mathcal{P})$$

(一様に改善されることはない)

■ ベイズ最適性

π_0 : 事前分布

ベイズリスク $\int \mathbb{E}_{D^n \sim P^*} [L(\hat{\theta})] d\pi_0(P^*)$ を最小にする推定法 $\hat{\theta}$.
(\rightarrow ベイズ推定量)

(事前分布で重みづけたリスクを最小化)

深層学習の理論

表現能力

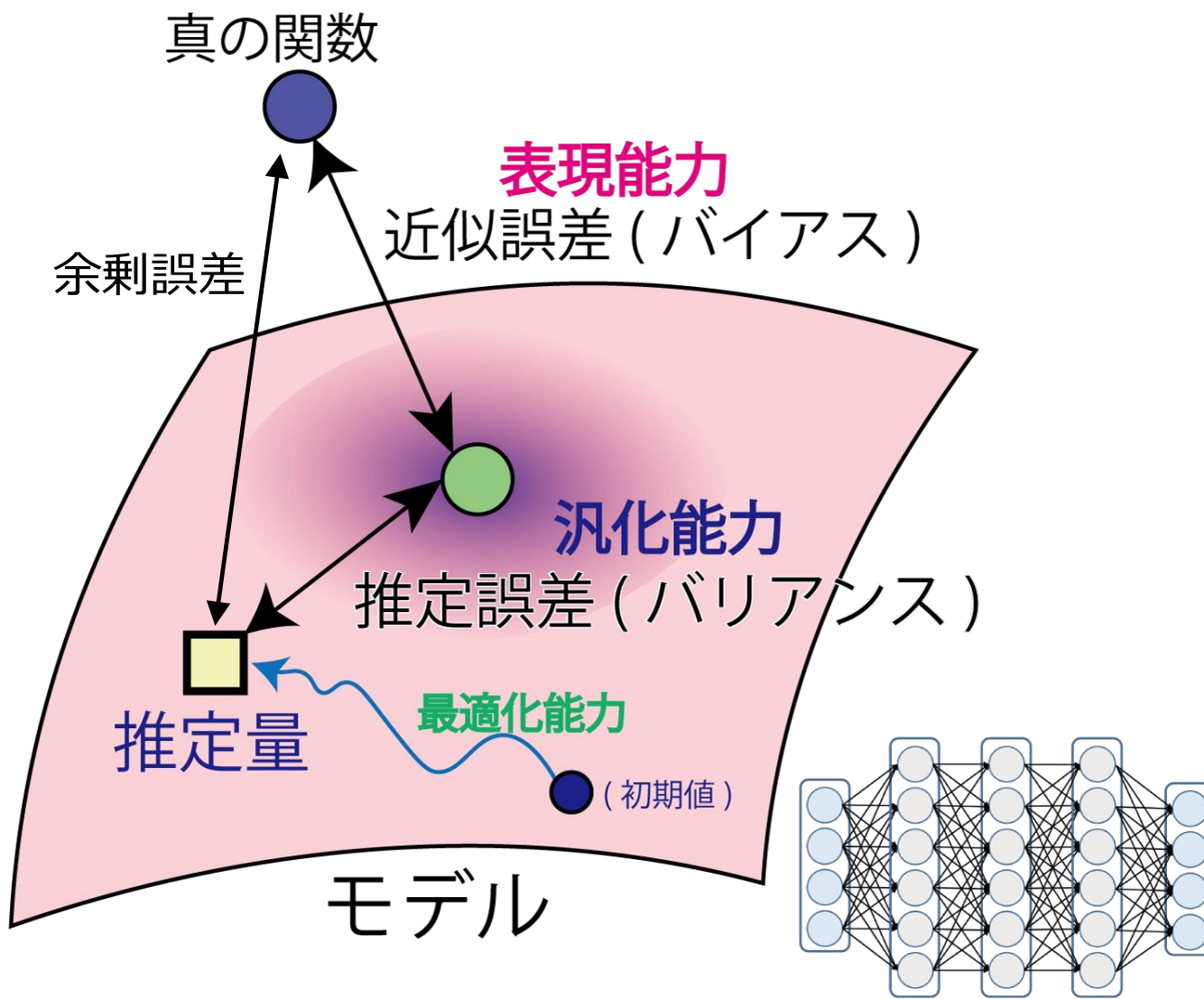
どれだけ難しい問題まで学習できるようにになるか？

汎化能力

有限個のデータで学習した時、どれだけ正しく初見のデータを正解できるようにになるか？

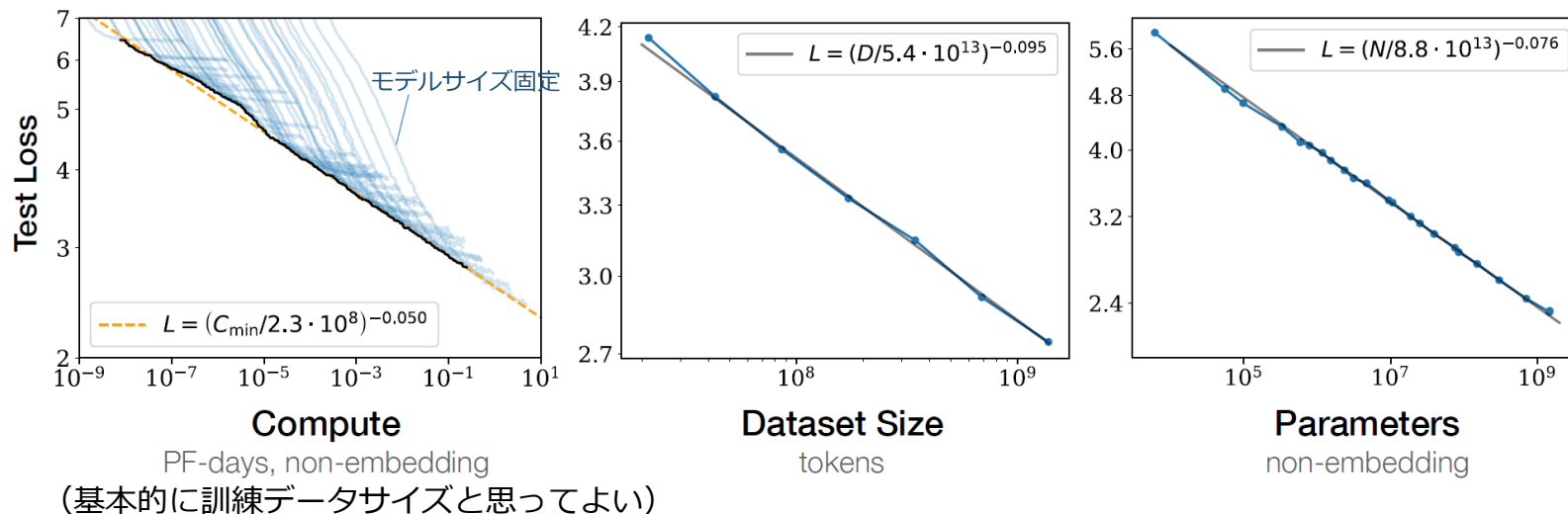
最適化能力

最適な重みを高速に計算機で求めることが可能か？

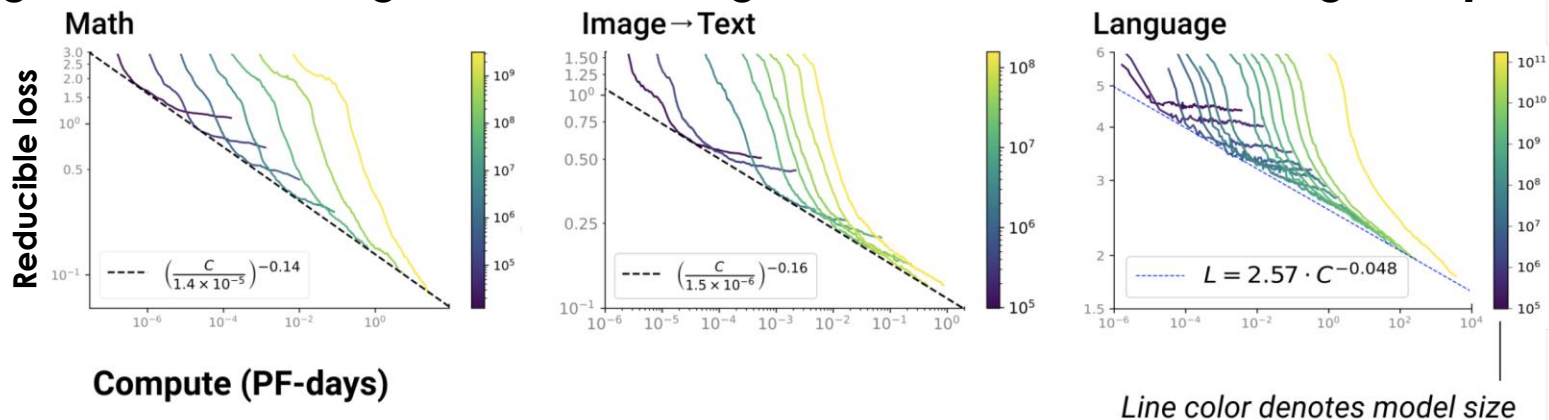


スケーリング則

[Kaplan et al.: Scaling Laws for Neural Language Models, 2020]



[Henighan et al.: Scaling Laws for Autoregressive Generative Modeling, 2020]



$$\log(\text{予測精度}) = -\alpha \log(n) + \log(C)$$

[Brown et al.: Language Models are Few-Shot Learners, 2020] (GPT-3モデルの解析)

基本的考え方

- スケーリング則は古典的な学習理論でも現れる.

真のモデル

$$f^\circ(x) = \sum_{j=1}^{\infty} \alpha_j \varphi_j(x) \quad (\text{正規直交系 in } L_2)$$

$$\text{観測データ: } y_i = f^\circ(x_i) + \epsilon_i$$

ただし

$$\mu_j \sim j^{-a} \text{ を用いて} \\ \sum_{j=1}^{\infty} \alpha_j^2 / \mu_j < \infty$$

学習モデル

正則化学習法 (カーネル法)

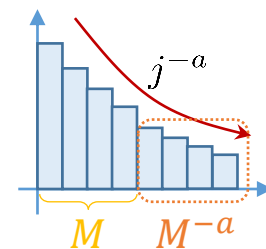
$$\hat{f}(x) = \sum_{j=1}^M \hat{\alpha}_j \varphi_j(x) \quad \leftarrow \min_{(\hat{\alpha}_j)_{j=1}^M} \frac{1}{n} \sum_{i=1}^n \left(y_i - \sum_{j=1}^M \hat{\alpha}_j \varphi_j(x_i) \right)^2 + \lambda \sum_{j=1}^M \hat{\alpha}_j^2 / \mu_j$$

予測誤差

$$\|f^\circ - \hat{f}\|_{L_2(P_X)}^2 \leq C \left(\overbrace{\frac{M}{n}}^{\text{バリエンス}} + \overbrace{M^{-a}}^{\text{バイアス}} \right)$$

最適なモデルサイズ

$$C n^{-\frac{a}{1+a}}$$



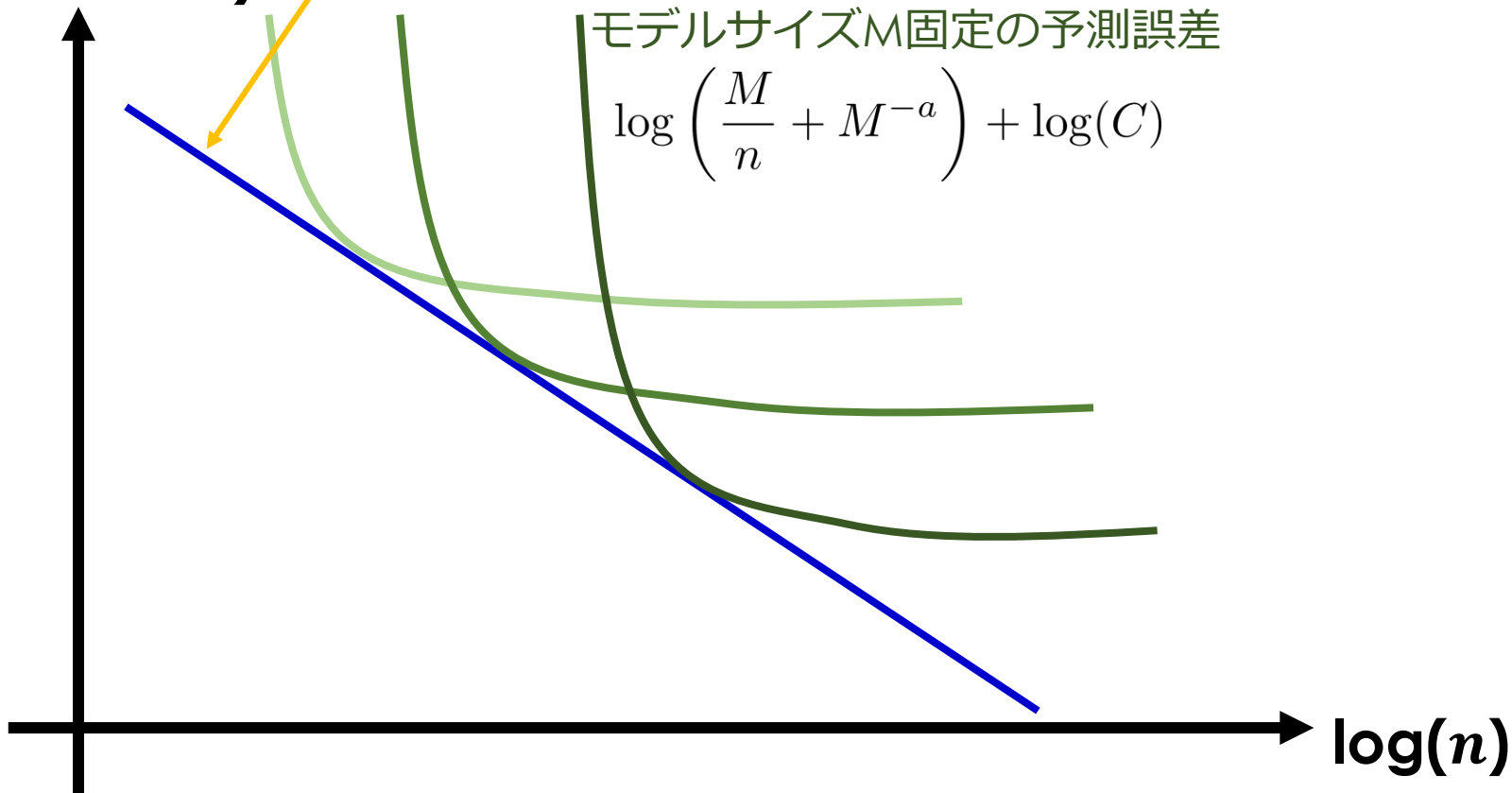
$$\log(\text{予測誤差}) = -\frac{a}{1+a} \log(n) + \log(C)$$

バリエンス=モデルの次元/n
バイアス=切り捨てた係数の二乗和

最適なモデルサイズの予測誤差

$$\log(\text{予測誤差}) = -\frac{a}{1+a}\log(n) + \log(C)$$

$\log(\text{予測誤差})$



カーネル法の学習理論

- Caponnetto and De Vito. Optimal Rates for the Regularized Least-Squares Algorithm. *Foundations of Computational Mathematics*, volume 7, pp.331–368 (2007).
- Steinwart and Christmann. *Support Vector Machines*. 2008.

関連する最近の論文

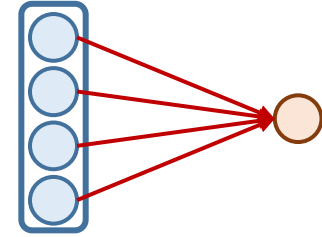
- Mei, Misiakiewicz, Montanari. Generalization error of random features and kernel methods: hypercontractivity and kernel matrix concentration. arXiv:2101.10588.
- Bordelon, Canatar, Pehlevan. Spectrum Dependent Learning Curves in Kernel Regression and Wide Neural Networks. arXiv:2002.02561.
- Canatar, Bordelon, Pehlevan. Spectral Bias and Task-Model Alignment Explain Generalization in Kernel Regression and Infinitely Wide Neural Networks. arXiv:2006.13198.

- 線形モデル

非線形化



$$f(x) = \sum_{j=1}^d \alpha_j x_j$$

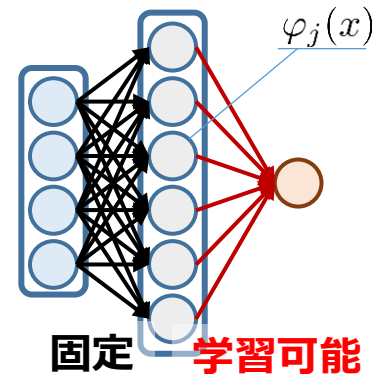


- カーネルモデル

可変基底化

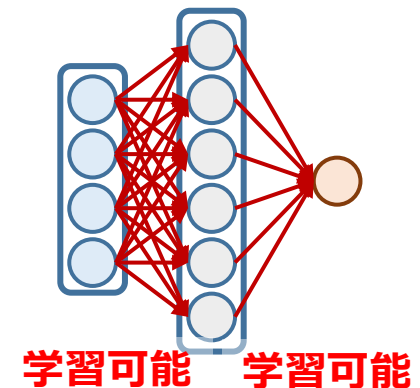


$$f(x) = \sum_{j=1}^M \alpha_j \underbrace{\varphi_j(x)}_{\text{固定}}$$



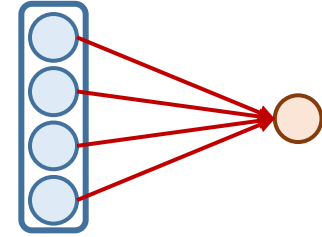
- 深層モデル

$$f(x) = \sum_{j=1}^M \alpha_j \underbrace{\varphi_j(x; \theta)}_{\text{学習可能}}$$



- 線形モデル

$$f(x) = \sum_{j=1}^d \alpha_j x_j$$

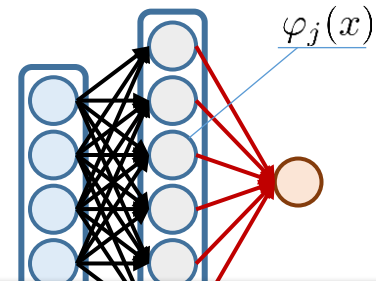


非線形化



- カーネルモデル

$$f(x) = \sum_{j=1}^M \alpha_j \varphi_j(x)$$



可変基底化

問題意識

- 基底を学習可能にすることで何が良くなるか？
- 逆に過学習を起こさないか？
- 最適化可能か？

学習可能

学習可能 学習可能

Deep learning theory lecture notes

by Matus Telgarsky

- <https://mjt.cs.illinois.edu/dlt/>

その他, 学習理論の教科書

- 一様バウンドを含む経験過程および無限次元統計モデルの教科書:
Giné and Nickl, *Mathematical foundations of infinite-dimensional statistical models*. Cambridge University Press, 2016.
- 一様バウンドや Fast learning rate が網羅的に収録された教科書:
Wainwright, *High-dimensional statistics: A non-asymptotic viewpoint*. Cambridge University Press, 2019.
- 経験過程の教科書: van der Vaart and Wellner, *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer, 1996.
- 統計的漸近論の教科書: van der Vaart, *Asymptotic statistics*. Cambridge University Press, 2000.
- 学習理論の教科書 (和書): 金森敬文, 統計的学習理論 (MLP シリーズ). 講談社, 2015.
- 学習理論の教科書:
 - Mohri et al., *Foundations of machine learning*. MIT press, 2018.
 - Shalev-Shwartz and Ben-David, *Understanding machine learning: From theory to algorithms*. Cambridge University Press, 2014.

第1部

深層学習の表現能力

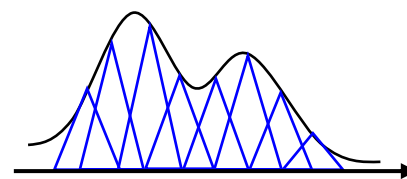
表現能力「万能近似能力」

二層ニューラルネットワークは
どんな関数も任意の精度で近似できる。

理論的にはデータが無限にあり、素子数が無限にあるニューラルネットワークを用いればどんな問題でも学習できる。

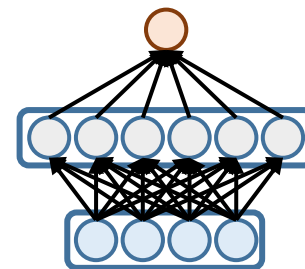
[Hecht-Nielsen, 1987][Cybenko, 1989]

「関数近似理論」

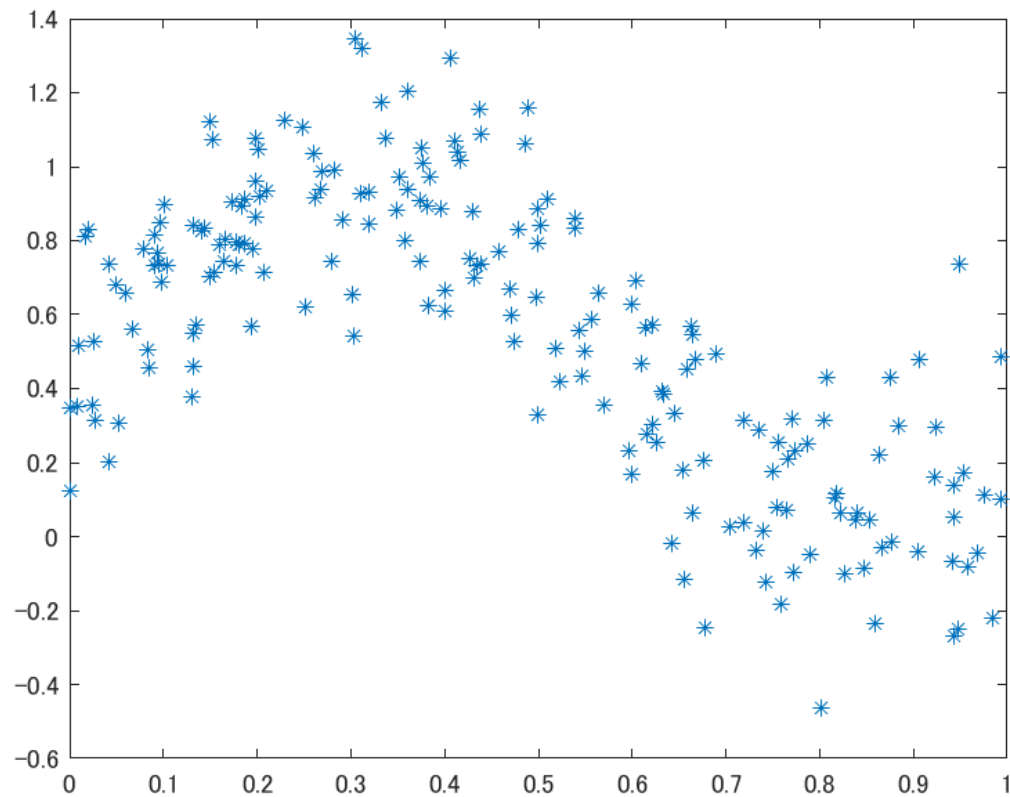


年		基底関数	空間
1987	Hecht-Nielsen	対象毎に構成	$C(R^d)$
1988	Gallant & White	Cos	$L_2(K)$
	Irie & Miyake	integrable	$L_2(R^d)$
1989	Carroll & Dickinson	Continuous sigmoidal	$L_2(K)$
	Cybenko	Continuous sigmoidal	$C(K)$
	Funahashi	Monotone & bounded	$C(K)$
1993	Mhaskar + Micchelli	Polynomial growth	$C(K)$
2015	Sonoda + Murata	Unbounded , admissible	$L_1(R^d)/L_2(R^d)$

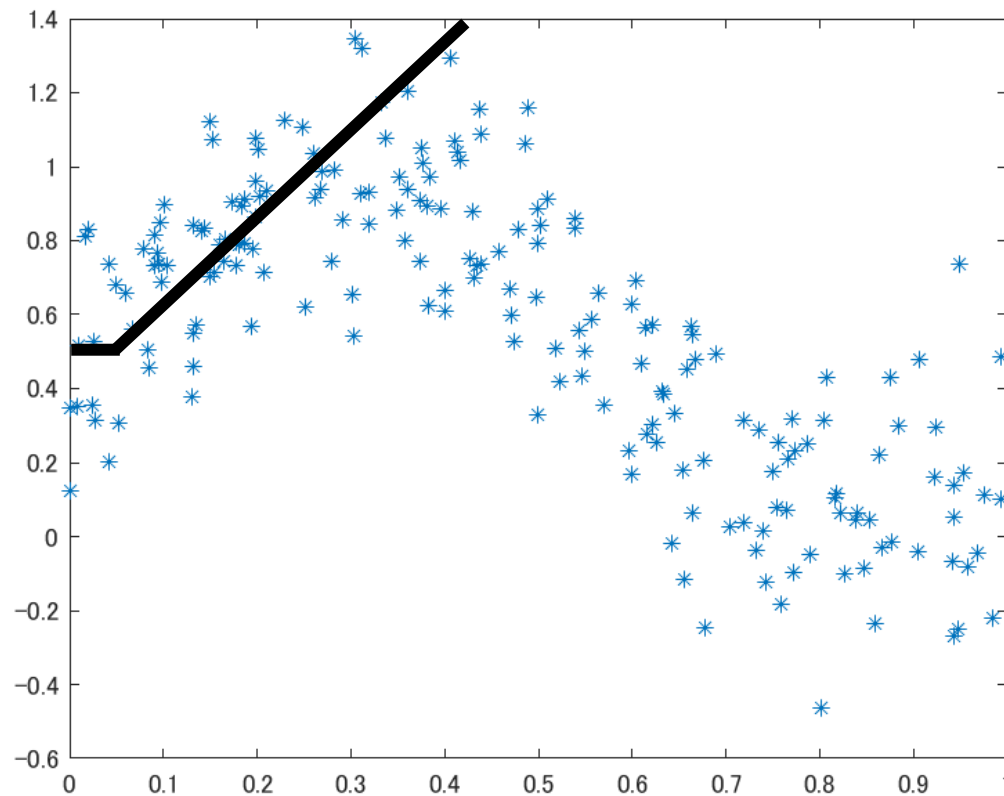
$$\hat{f}(x) = \sum_{j=1}^m v_j \eta(w_j^T x + b_j)$$



関数近似の様子

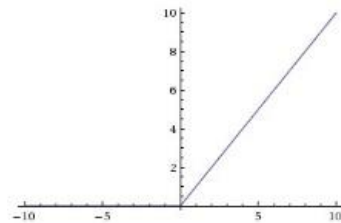


関数近似の様子



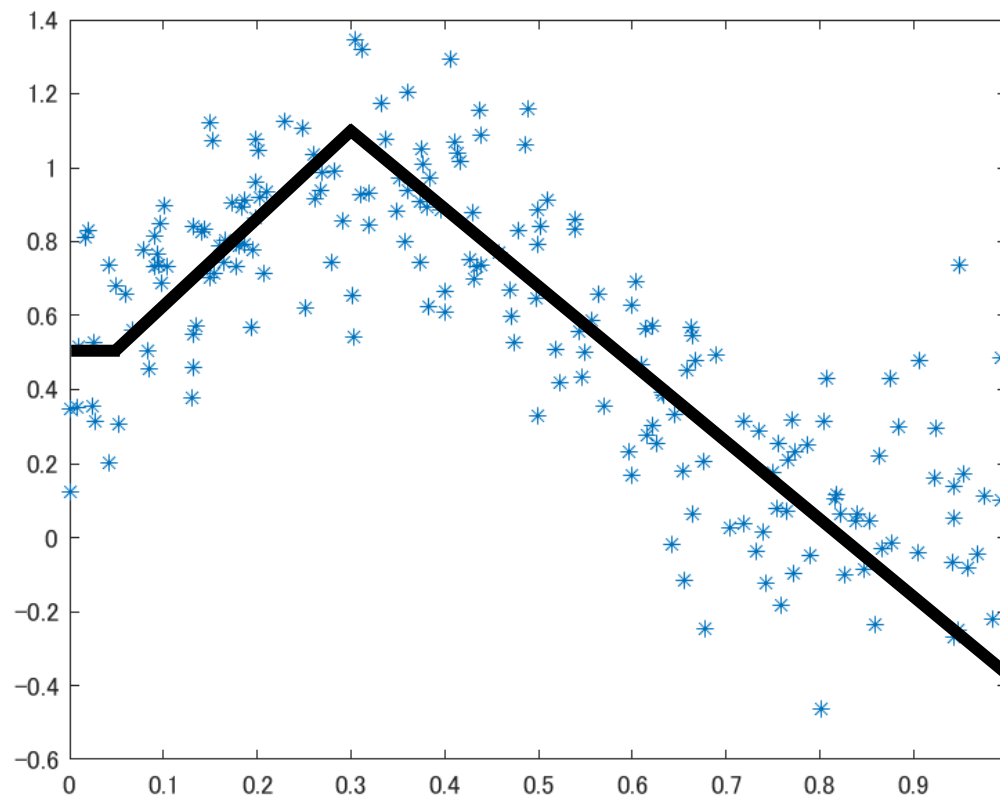
$$a_1 \eta(b_1 x + c_1)$$

ReLU活性化関数

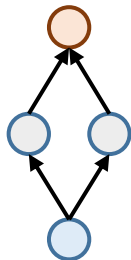


$$\eta(u) = \max\{u, 0\}$$

関数近似の様子

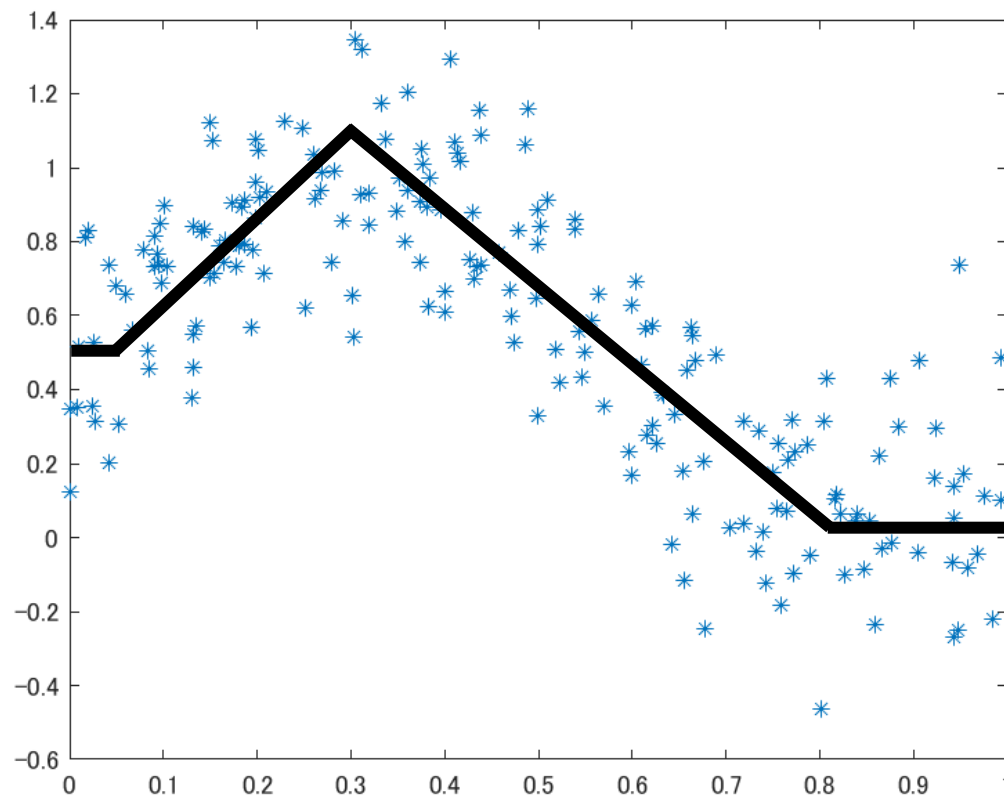


$$a_1 \eta(b_1 x + c_1)$$

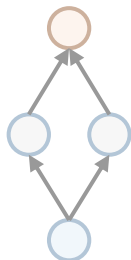


$$a_1 \eta(b_1 x + c_1) \\ + a_2 \eta(b_2 x + c_2)$$

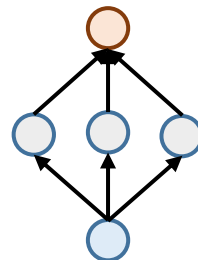
関数近似の様子



$$a_1 \eta(b_1 x + c_1)$$

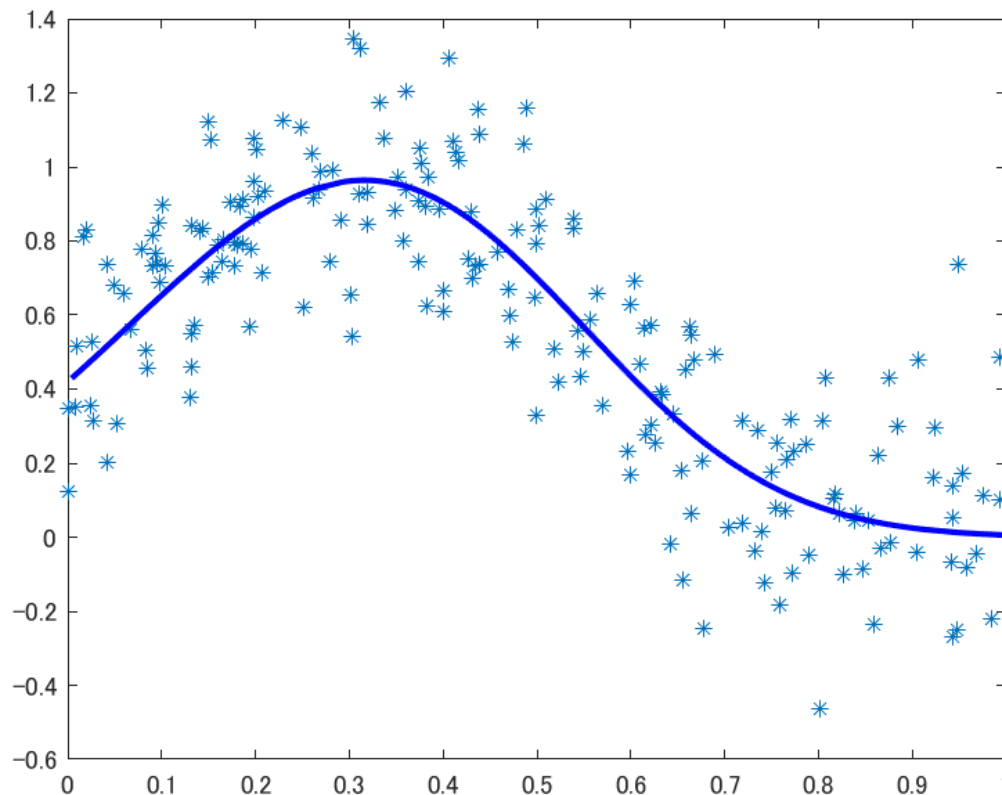


$$a_1 \eta(b_1 x + c_1) + a_2 \eta(b_2 x + c_2)$$

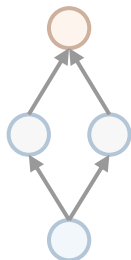


$$\sum_{j=1}^3 a_j \eta(b_j x + c_j)$$

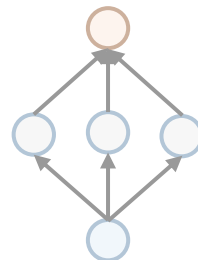
関数近似の様子



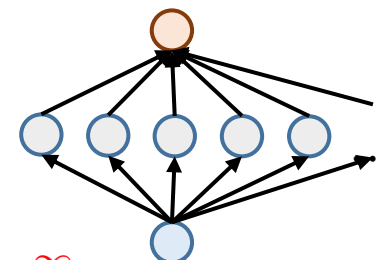
$$a_1 \eta(b_1 x + c_1)$$



$$a_1 \eta(b_1 x + c_1) + a_2 \eta(b_2 x + c_2)$$



$$\sum_{j=1}^3 a_j \eta(b_j x + c_j)$$



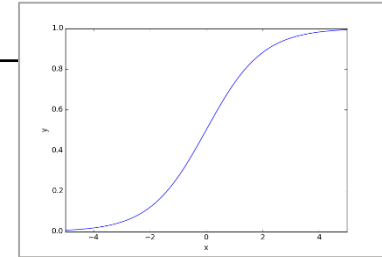
$$\sum_{j=1}^{\infty} a_j \eta(b_j x + c_j)$$

• Cybenkoの理論

[Cybenko: Approximation by superpositions of a sigmoidal function.
Mathematics of control, signals and systems, 2(4): 303-314, 1989]

定義

活性化関数 η がシグモイド的 $\Leftrightarrow \eta(x) \rightarrow \begin{cases} 1 & (x \rightarrow \infty) \\ 0 & (x \rightarrow -\infty) \end{cases}$



定理

活性化関数 η が連続なシグモイド的関数なら, 任意の $f \in C([0,1]^d)$ と, 任意の $\epsilon > 0$ に対して, ある $g(x) = \sum_{i=1}^N \alpha_i \eta(a_i x_i + b_i)$ が存在して,

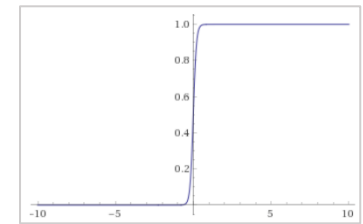
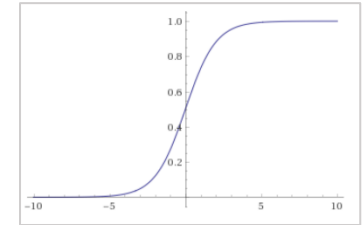
$$\sup_{x \in [0,1]^d} |f(x) - g(x)| \leq \epsilon$$

とできる.

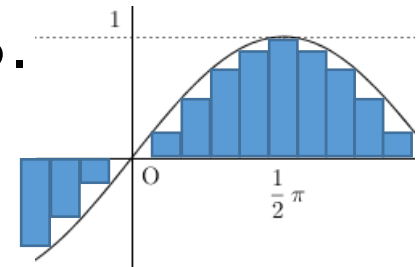
- シグモイド型の関数に対し,

$$h(a(\alpha^\top x + \beta) + \theta) \xrightarrow{a \rightarrow \infty} \begin{cases} 1 & (\alpha^\top x + \beta > 0) \\ h(\theta) & (\alpha^\top x + \beta = 0) \\ 0 & (\alpha^\top x + \beta < 0) \end{cases}$$

が成り立つ。つまり、スケールを適切に選べば、階段関数をいくらでもよく近似できる。



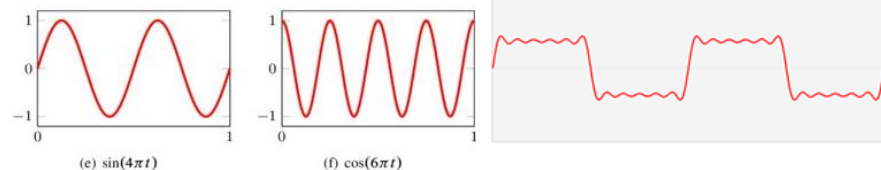
- 階段関数を近似できれば、それらを足し引きすることで、 $\cos(\alpha^\top x + \beta)$ や $\sin(\alpha^\top x + \beta)$ をいくらでもよく近似できる。
- \cos , \sin が実現できるなら Fourier(逆)変換もできる。
- 任意の連続関数が近似できる。



積分表現 (Ridgelet変換)

• Fourier変換

$$f(x) = \int_{\omega \in \mathbb{R}^d} F(\omega) e^{i\omega^\top x} d\omega$$

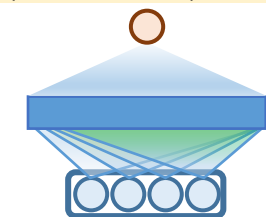


[Wikipedia, フーリエ変換]

• Ridgelet変換

NNはFourier変換におけるsin, cosの代わりに非線形ノードの足し合わせで関数を表現.

$$f(x) = \int T(a, b) \eta(a^\top x - b) db da$$



ウェーブレット変換 + ラドン変換

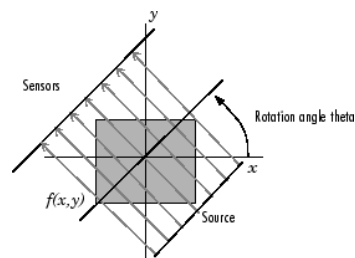
ある $\psi: \mathbb{R} \rightarrow \mathbb{R}$ が以下の「許容条件」を満たすとする :

$$K_{\psi, \eta} := (2\pi)^{d-1} \int_{\mathbb{R} \setminus \{0\}} \frac{\hat{\psi}(\xi) \hat{\eta}(\xi)}{|\xi|^d} d\xi < \infty \quad (\hat{\psi}, \hat{\eta} \text{ はFourier変換})$$

$$\mathcal{R}_\psi f(a, b) = \int_{\mathbb{R}^d} f(x) \overline{\psi(x^\top a - b)} \|a\| dx \quad (\text{Ridgelet変換})$$

$$\mathcal{R}_\eta^\dagger T(x) = \int_{a \in \mathbb{R}^d} \int_{\mathbb{R}} T(a, b) \eta(a^\top x - b) \|a\|^{-1} db da \quad (\text{双対Ridgelet変換})$$

CTスキャン



定理

$$(\mathcal{R}_\eta^\dagger \mathcal{R}_\psi f)(x) = K_{\psi, \eta} f(x)$$

(再構成定理)

カーネル法

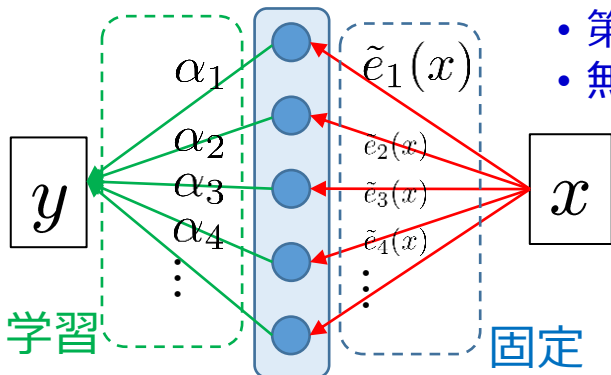
- 浅い手法の代表格.
- これも万能近似能力がある.

第1層目を固定した横幅無限の2層ニューラルネットワーク

似た手法

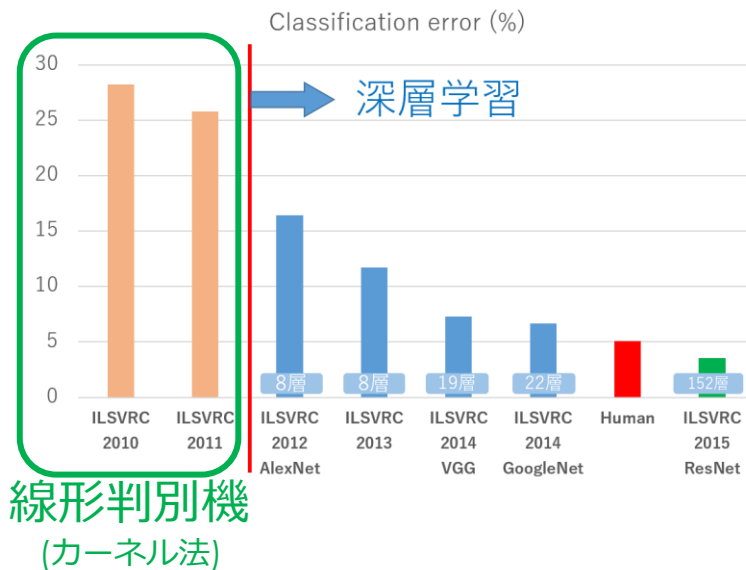
- スプライン法
- 局所多項式回帰
- シリーズ推定量

(線形推定量と呼ばれるクラス)



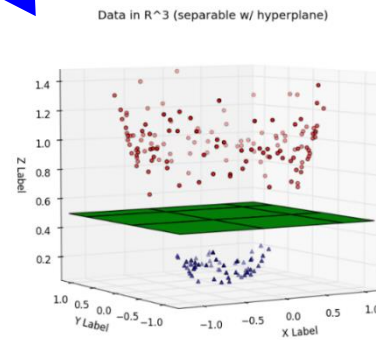
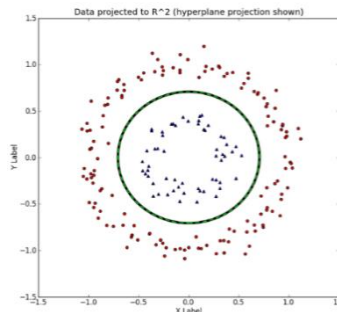
- 第一層目固定
- 無限個の素子

$$k(x, x') = \sum_{m=1}^{\infty} \tilde{e}_m(x) \tilde{e}_m(x') : \text{カーネル関数}$$



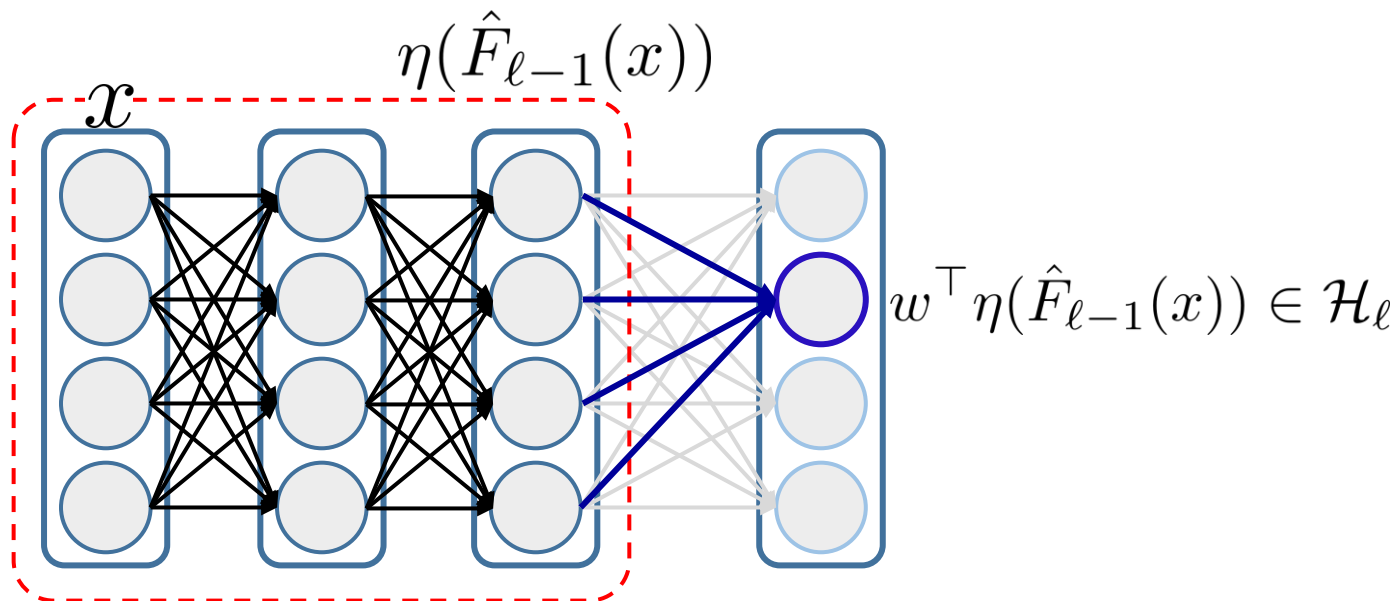
再生核ヒルベルト空間の理論

非線形写像 ϕ_x



深層NNとカーネル関数

深層NNは「カーネル関数をデータに合わせて学習する方法」と言える



$$\hat{k}_\ell(x, x') = \eta(\hat{F}_{\ell-1}(x))^\top \eta(\hat{F}_{\ell-1}(x'))$$

\hat{k}_ℓ に対応した再生核ヒルベルト空間

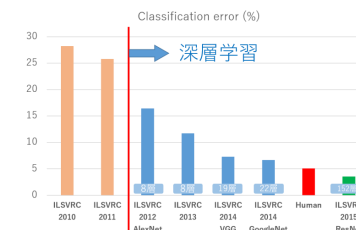
$$\mathcal{H}_\ell = \{f(x) = w^\top \eta(\hat{F}_{\ell-1}(x)) \mid w \in \mathbb{R}^{m_\ell}\},$$

$$\|f\|_{\mathcal{H}_\ell} = \|w\|.$$

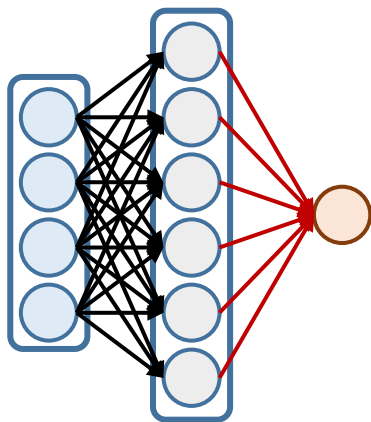
(より正確には同じ f を与える w の中で \inf を取る)

これまででわかったこと

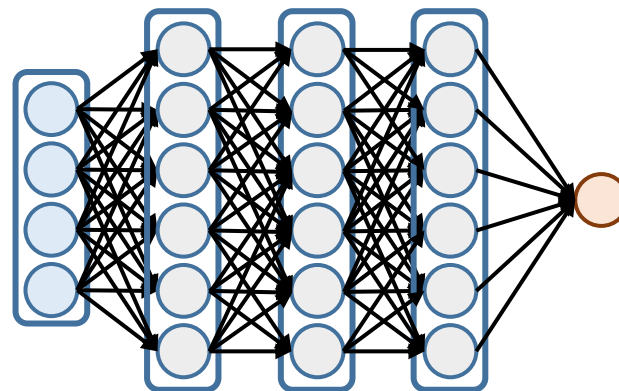
- **[理論]** 万能近似能力という意味では浅層で十分.
- **[実際]** 実際は多層を使うことが多い.
→ この差はどう埋める？



カーネル法
浅層



多層ニューラルネット
深層学習



→ 「表現力」を比べてみる.

(1) ランダム特徴量の近似能力

ランダム特徴：
$$\hat{f}(x) = \sum_{j=1}^M r_j \sigma(u_j^\top x + b_j)$$

これは近似対象に合わせて自由に設定

ランダムに生成
(球面上の一様分布) (固定)

疑問： ある対象の関数 f^* を近似するのに必要なランダム特徴量の数 M はどれくらい？

➡ 実は、一つのニューロンを近似するのにも $M = \exp(\Omega(d))$ 必要。 (次元の呪い)

つまり、ランダムに特徴量を生成する方法は非常に効率が悪い。
深層学習のように第一層目のパラメータも学習した方が効率的。

※ ランダム特徴量の方法は、カーネルの低ランク近似法とみなせる。 (特徴学習の必要性)

定理 (Yehudai and Shamir 2019)

ある $u^* \in \mathbb{R}^d$, $b^* \in \mathbb{R}$ s.t. $\|u^*\| \leq d^2$ が存在して、

$$\mathbb{E}_X \left[(\hat{f}(X) - \sigma(u^{*\top} X + b^*))^2 \right] \leq 0.1$$

であるためには、 $M \cdot \max_j |r_j| \geq \exp(\Omega(d))$ が必要。

(2-a) Barron classとランダム特徴量

π : probability measure

- π で決まる再生核ヒルベルト空間 : $(\sigma: \text{ReLU})$

$$\mathcal{H}_\pi = \left\{ \int_{\mathbb{S}^{d-1}} a(w) \sigma(w^\top x) \pi(dw) \mid \int_{\mathbb{S}^{d-1}} a(w)^2 \pi(dw) < \infty \right\}$$

$$\|f\|_{\mathcal{H}_\pi}^2 := \mathbb{E}_{w \sim \pi} [a(w)^2] \quad \text{ただし } f = \int a(w) \sigma(w^\top x) \pi(dw)$$

(実際, これは再生核を $k(x, x') = \int_{\mathbb{S}^{d-1}} \sigma(w^\top x) \sigma(w^\top x') \pi(dw)$ とするRKHSになっている.)

- M 個のランダム特徴量で張られるモデル :

$$\mathcal{H}_{\text{rand}}(M) = \left\{ \hat{f}(x) = \sum_{j=1}^M r_j \sigma(u_j^\top x) \mid r_j \in \mathbb{R} \right\}$$

ランダムに生成して固定
(近似対象によって変えない)

定理 (Barron classの近似定理; E, Ma, Wu, 2019)

ある確率測度 π が存在して, その π によって決まる再生核ヒルベルト空間の元 $f \in \mathcal{H}_\pi$ ($\|f\|_{\mathcal{H}_\pi} \leq 1$) を適切に取ってくると, 以下が成り立つ:

$$\inf_{\hat{f} \in \mathcal{H}_{\text{rand}}(M)} \|f - \hat{f}\|_{L_2(P_X)} \gtrsim \frac{1}{dM^{1/d}}$$



次元の呪い

(誤差 ϵ を達成するには $M = \epsilon^{-\Omega(d)}$ 必要)

(2-b) NNによる次元の呪いの解消

一方で、近似対象の関数 f に対して、適切に一層目のパラメータ $(u_j)_{j=1}^M$ を設定すると、次元の呪いを回避できる。

$$\mathcal{H}_{\text{NN}}(M) = \left\{ \hat{f}(x) = \sum_{j=1}^M r_j \sigma(u_j^\top x) \mid r_j \in \mathbb{R}, u_j \in \mathbb{S}^{d-1} \right\} : \text{NNの集合}$$

$$\inf_{\hat{f} \in \mathcal{H}_{\text{NN}}} \|f - \hat{f}\|_{L_2(P_X)}^2 = O\left(\frac{1}{M}\right)$$

NNによって関数近似の効率が改善されている。

$$M^{-1/d} \Rightarrow M^{-1}$$

(一層目固定)

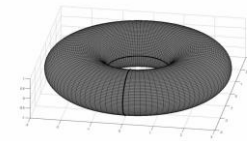
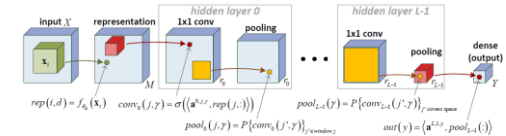
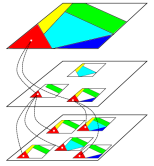
(一層目可変)

直感：ランダム特徴量だと $(d-1)$ 次元単位球面全体を覆うのに、 $1/\epsilon^d$ 個のニューロンが必要。



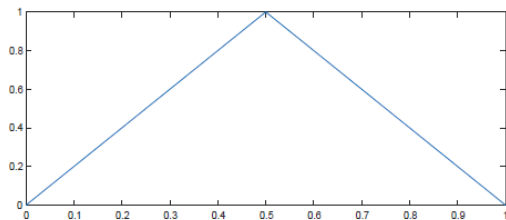
深さに対して指数関数的に“表現力”が上がる。

- **超平面アレンジメント** [Montufar et al., 2014]
空間の領域分割数
- **多項式展開, テンソル解析** [Cohen et al., 2016; Cohen & Shashua, 2016]
単項式の次数
- **代数トポロジー** [Bianchini & Scarselli, 2014]
ベッチ数(Pfaffian)
- **リーマン幾何 + 平均場理論** [Poole et al., 2016]
埋め込み曲率

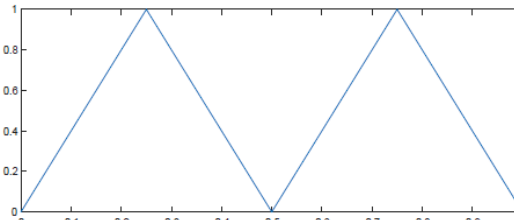


特に, 対称性の高い関数の近似は深層NNが有利

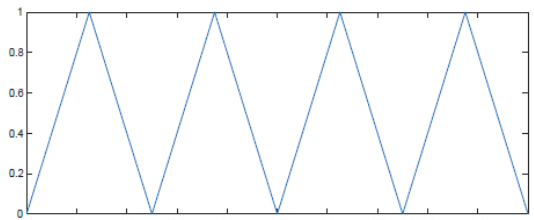
$$h(x) = \begin{cases} 2x & (0 \leq x \leq 1/2) \\ 2(1-x) & (1/2 \leq x \leq 1) \\ 0 & (\text{otherwise}). \end{cases}$$



$h(x)$



$h \circ h(x)$



$h \circ h \circ h(x)$

多層が得する例 (1): 領域分割数

NNの“表現力”：領域を何個の多面体に分けられるか？

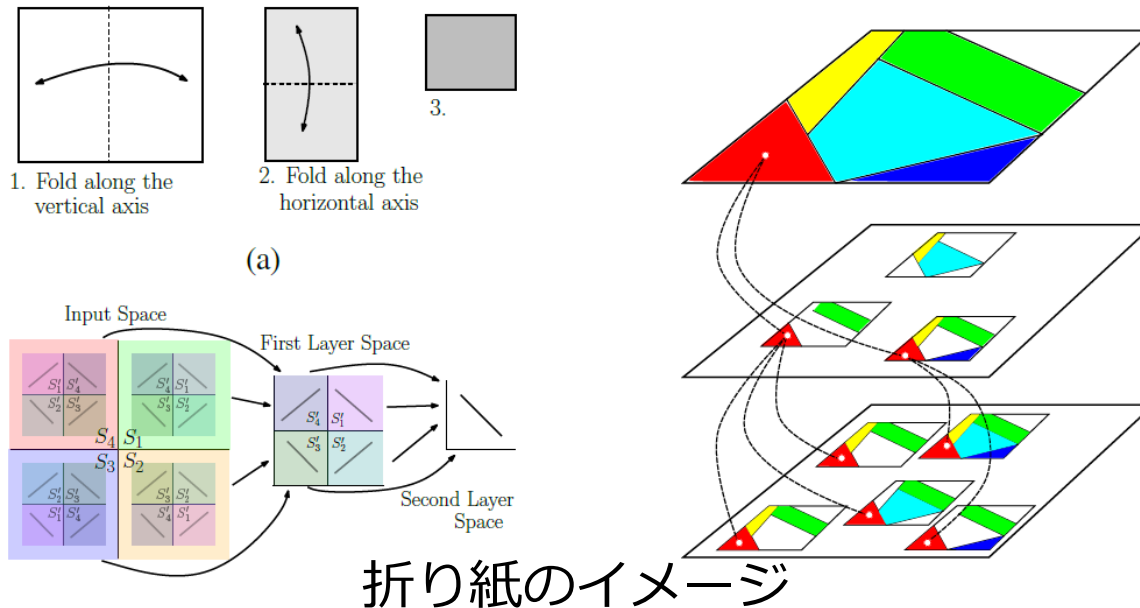
- 層の数に対して表現力は指數的に上がる。

$$\left(\frac{n}{n_0}\right)^{L-1} \sum_{j=0}^{n_0} \binom{n}{j}$$

- 中間層のユニット数 (横幅) に対しては多項式的。

$$\sum_{j=0}^{n_0} \binom{n}{j}$$

L : 層の数
 n : 中間層の横幅
 n_0 : 入力の次元

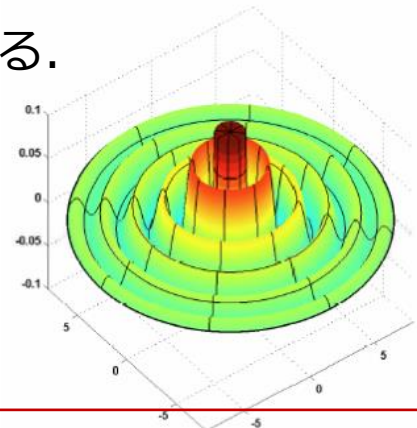
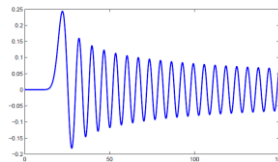


多層が得する例 (2): 対称な関数 参考

深層NNは特徴量を適切に抽出することで次元の呪いを回避できる。

$$g(\|x\|^2) = g(x_1^2 + x_2^2 + \dots + x_{d_x}^2)$$

g はBessel関数を元に構成



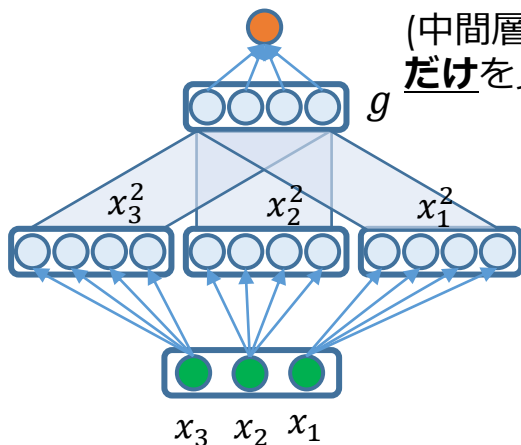
- 三層 (中間層二層) : $O(\text{poly}(d_x))$ ノードで十分
- 二層 (中間層一層) : $\Omega(\exp(d_x))$ ノードが必要 (次元の呪い)

(Eldan&Shamir, 2016)

三層

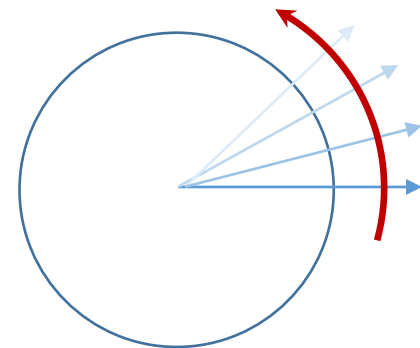
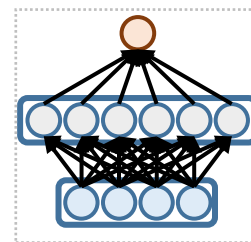
まず二乗和 $x_1^2 + \dots + x_{d_x}^2$ を作ってから g を作用。

(中間層で座標軸方向
だけを見ればよい)



二層

全方向をケアする必要がある
(座標軸方向だけではダメ)

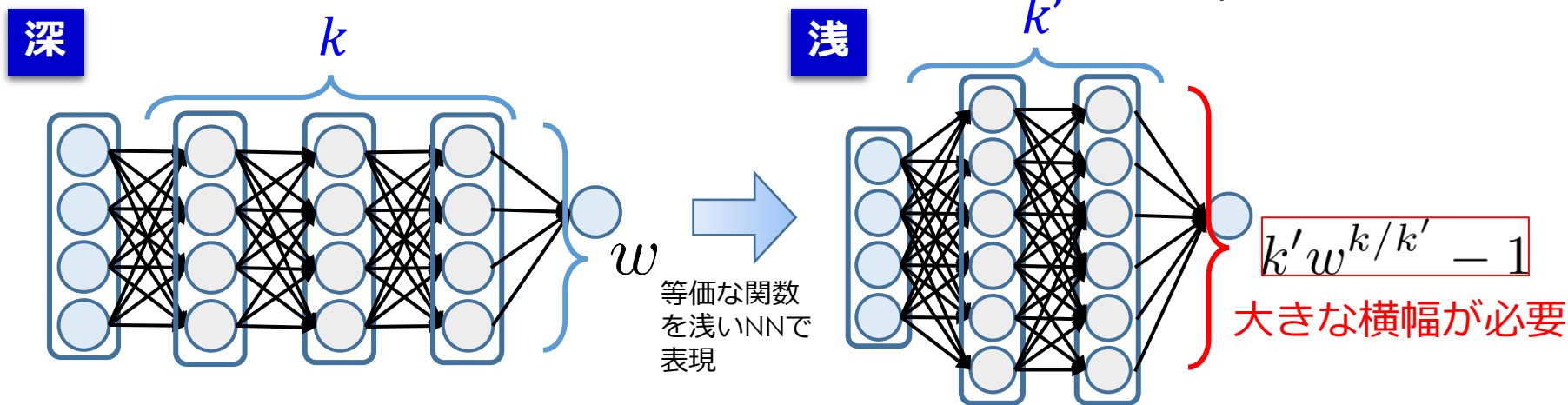


多層が得する例 (3): 区分線形関数の表現⁴⁵

参考

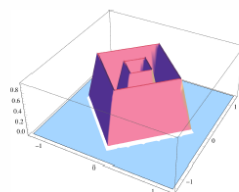
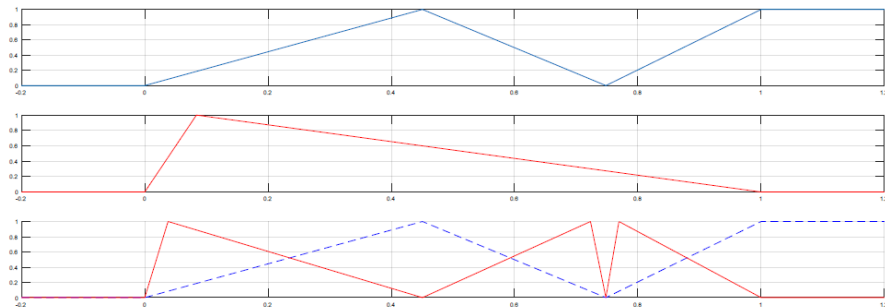
- 任意の区分線形関数($\mathbb{R}^d \rightarrow \mathbb{R}$)は深さ $\lceil \log_2(d + 1) \rceil$ のReLU-DNNで表現可能
- ある横幅 w , 縦幅 k のReLU-DNNが存在して, それを縦幅 $k' (< k)$ のネットワークで表現するには横幅 $k'w^{k/k'} - 1$ が必要.

(Arora et al., 2018)

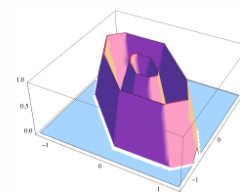


やはり層の深さに対し指数関数的に表現力が増加

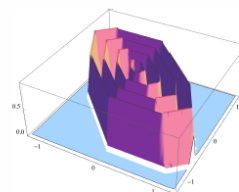
上記のネットワークの例



(a) $H_{\frac{1}{2}, \frac{1}{2}} \circ N_{\epsilon_1}$



(b) $H_{\frac{1}{2}, \frac{1}{2}} \circ \gamma_Z(b^1, b^2, b^3, b^4)$



(c) $H_{\frac{1}{2}, \frac{1}{2}, \frac{1}{2}} \circ \gamma_Z(b^1, b^2, b^3, b^4)$

多層が得する例 (4): 有理関数の近似

- 有理関数をReLU-DNNで近似

$$p : [0, 1]^d \rightarrow [-1, +1], \quad q : [0, 1]^d \rightarrow [2^{-k}, +1] \quad : \text{r次多項式}$$

p/q をReLU-DNNで近似したい

あるReLU-DNN f が存在してノード数と近似誤差が次のように抑えられる :

ノード数

$$O(\text{poly}(k, r, d) \text{poly}(\log(1/\epsilon)))$$

近似誤差

$$\sup_{x \in [0, 1]^d} \left| f(x) - \frac{p(x)}{q(x)} \right| \leq \epsilon$$

- ReLU-DNNを有理関数で近似

k -層で各層のノード数 m の任意のReLU-DNN f に対しては, 次数と近似誤差が以下で抑えられる有理関数 p/q が存在 :

次数 (分母 q と分子 p の次数の最大値)

$$O(\log(k/\epsilon)^k m^k)$$

深さに対して指数的に増大

近似誤差

$$\sup_{x \in [0, 1]^d} \left| f(x) - \frac{p(x)}{q(x)} \right| \leq \epsilon$$

- ReLU-DNNを多項式で近似 : $\Omega(\text{poly}(1/\epsilon))$ の次数が必要
→有理関数に比べて表現力が低い

多層が得する例 (4): 有理関数の近似

参考

- 有理関数をReLU-DNNで近似

$$p : [0, 1]^d \rightarrow [-1, +1], \quad q : [0, 1]^d \rightarrow [2^{-k}, +1] \quad : \text{r次多項式}$$

p/q をReLU-DNNで近似したい

あるReLU-DNN f が存在してノード数と近似誤差が次のように抑えられる :

ノード数

$$O(\text{poly}(k, r, d) \text{poly}(\log(1/\epsilon)))$$

近似誤差

$$\sup_{x \in [0, 1]^d} \left| f(x) - \frac{p(x)}{q(x)} \right| \leq \epsilon$$

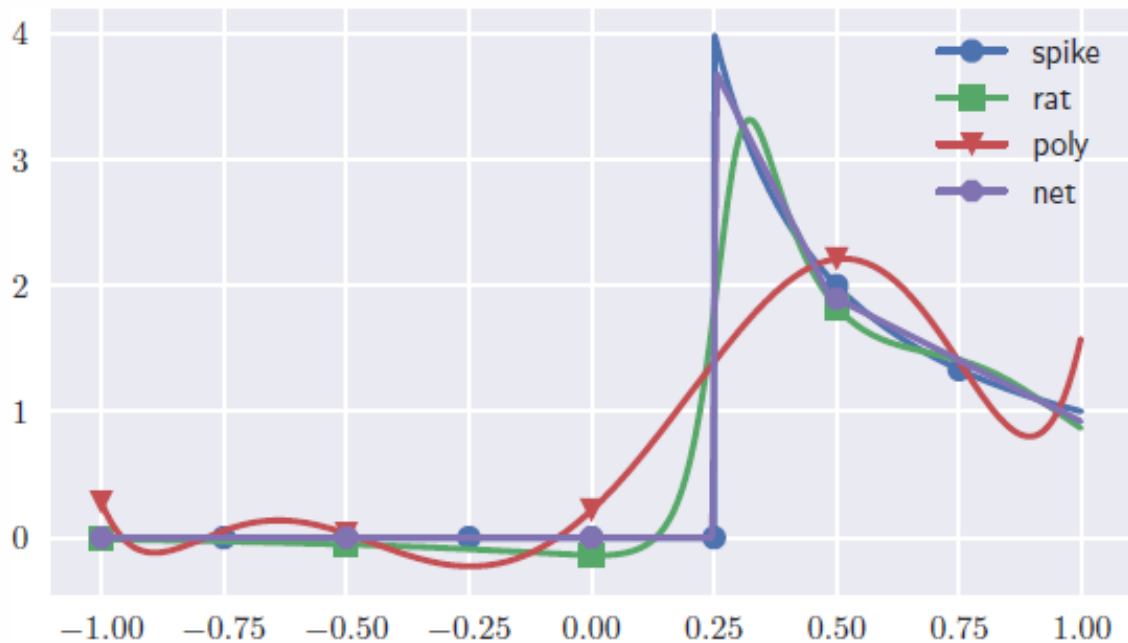
- ReLU-DNN

k -層で各
数と近似

次数 (分母)

$$O(k)$$

- ReLU-DNN



次

$$\left| \frac{p(x)}{q(x)} \right| \leq \epsilon$$

要

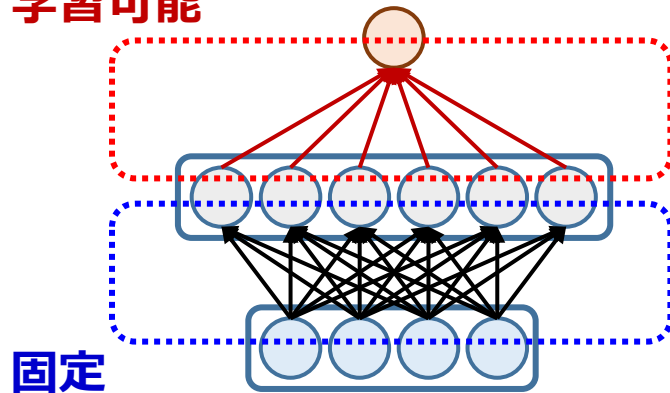
表現力が低い

深層学習の適応的推定能力 -ノンパラメトリック回帰理論-

ニューラルネットワークは「なぜ良いのか？」

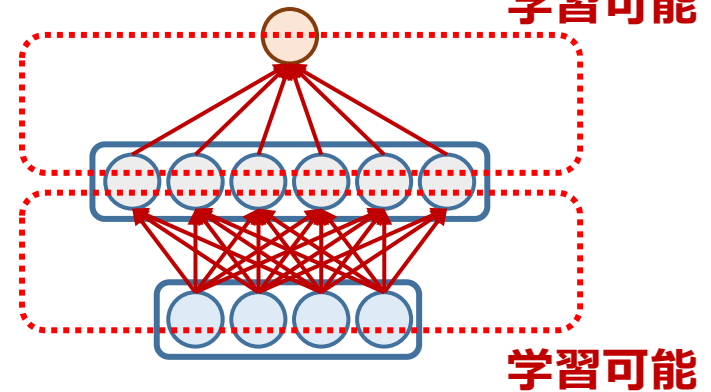
カーネル法 (線形モデル)

学習可能



深層モデル

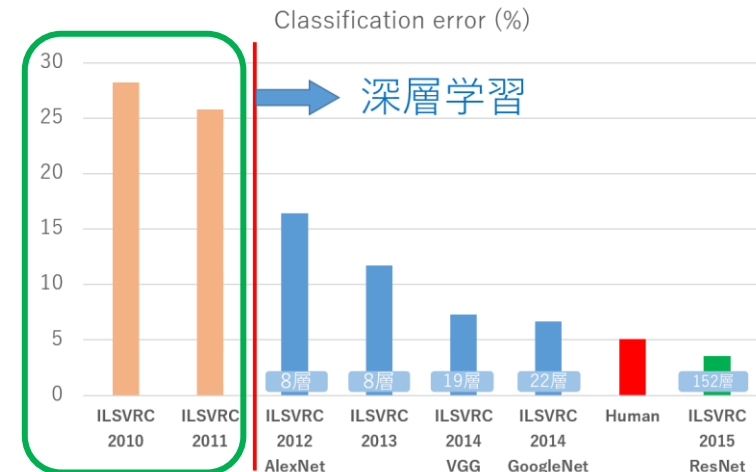
学習可能



解析対象

- 統計的効率性
- 最適化の効率性
(トレードオフの関係)

- 深層とカーネル法の性能を理論的に比較
- 深層学習の最適化を理論的に解析
- 深層学習の「適応能力」は本質的



カーネル法

適応能力とは？

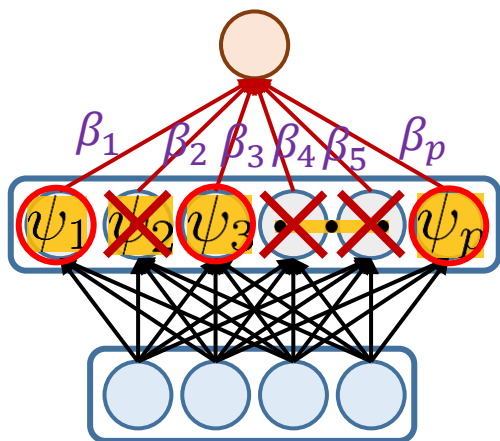
代表的な適応的推定法：**スパース推定**

やりたいこと：入出力関係を表す関数 $y = f(x)$ を推定したい。

- 事前にたくさんの基底を用意
- それらの線形結合で真の関数を表現

$$\hat{f}(x) = \beta_1\psi_1(x) + \beta_2\psi_2(x) + \cdots + \beta_p\psi_p(x)$$

- 基底は冗長性があるように用意しておき、
データに合わせて必要な基底のみ選択して使う。
→スパース推定
- データに合わせて基底を選ぶので適応的と言う。



$$\beta_2 = 0$$

$$\beta_4 = 0$$

$$\beta_5 = 0$$

多くの係数が0
⇒ スパース

スパース推定の恩恵

- 推定精度が上がる。
- 解釈性が向上する。

David Donoho



Gauss prize (2018)
Sparse estimation,
wavelet-shrinkage,
compressive sensing, ...

Wavelet shrinkage (Donoho&Johnstone, 1992)

→ Besov空間, スパース推定の適応能力

Robert Tibshirani



Lasso (1996)

→ L1正則化によるスパース推定

深層学習との関係

- スパース推定：
基底を沢山用意してその中から**選択**.
- 深層学習：
基底を事前に用意せず，データに合わせて**構築**.

$$\hat{f}(x) = \sum_{j=1}^p \beta_j \psi_j(x)$$

基底関数

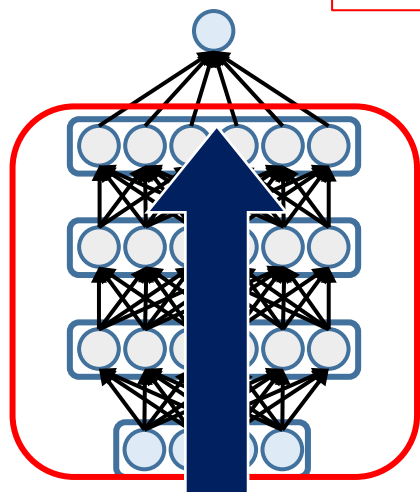
非適応的手法：

➤ データによらず全部使う：リッジ回帰

適応的手法：

➤ 沢山の候補の中から選ぶ：スパース推定

➤ データに合わせて構築：深層学習



中間層で基底を構成

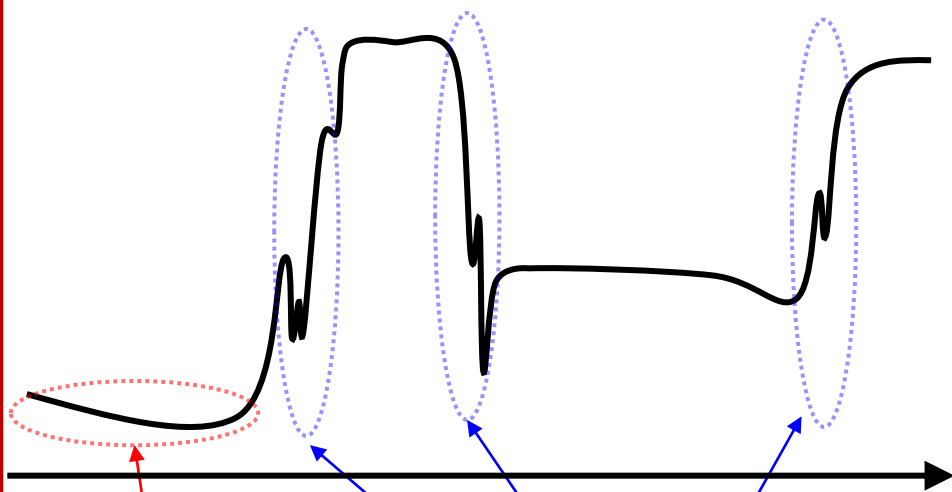


適応的推定能力

典型的な例

[Imaizumi&Fukumizu, 2019]
[Suzuki, 2019]

滑らかな部分と
そうでない部分が混在



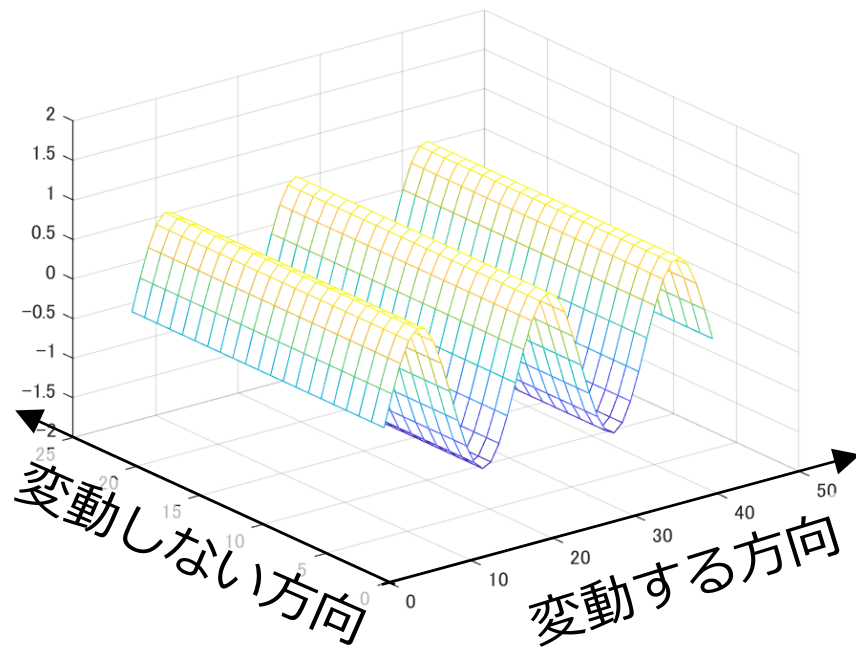
滑らか

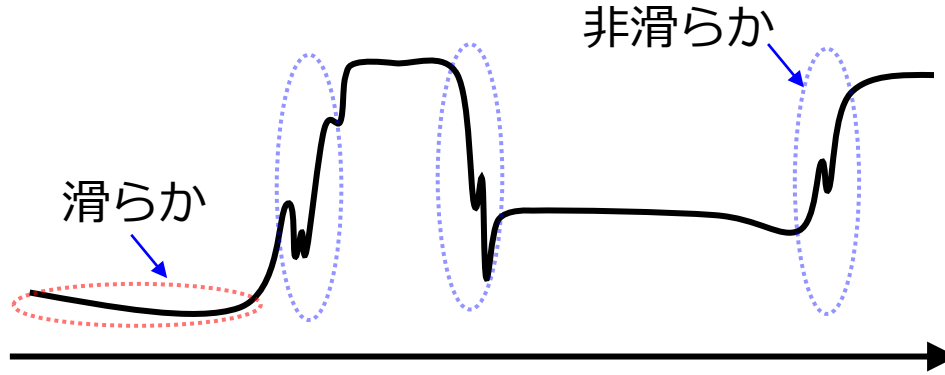
変動が大きい
(滑らかでない)

Besov空間

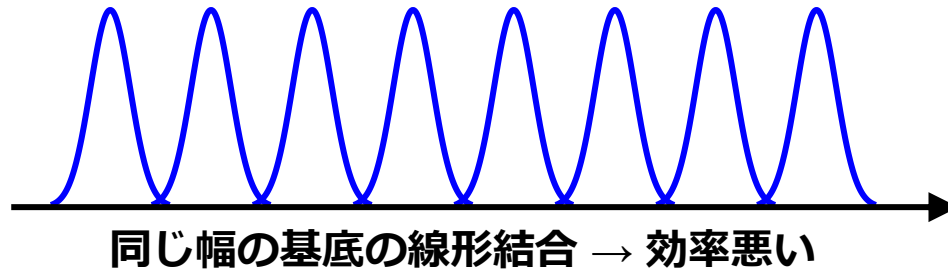
[Schmidt-Hieber, 2019] [Nakada&Imaizumi,
2019] [Chen et al., 2019] [Suzuki&Nitanda, 2019]

大きく変動する方向と
そうでない方向が混在

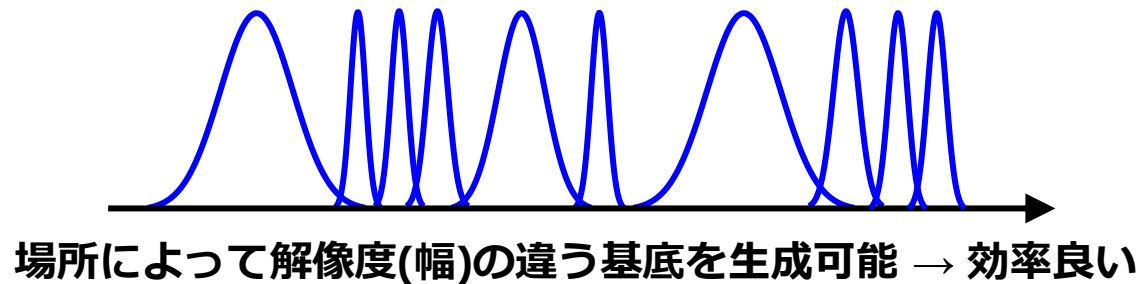




- ガウスカーネルを用いた関数近似 (カーネル法・非適応的)



- NNによる関数近似



深層 vs 浅層 の統計理論

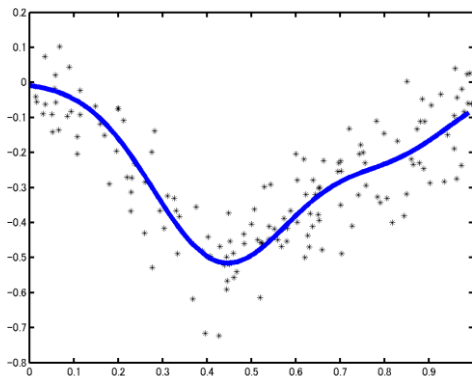
→ 「関数近似精度/推定精度」を比べてみる.

「多層」による特徴抽出と推定精度

ノンパラメトリック回帰の設定

$$y_i = f^\circ(x_i) + \xi_i \quad (i = 1, \dots, n)$$

$\xi_i \sim N(0, \sigma^2)$ は観測誤差



推定誤差 (平均二乗誤差):

$$\mathbb{E}[\|\hat{f} - f^\circ\|_{L_2(P)}^2] < ?$$

※実はこれは二乗損失の平均余剰誤差になっている.

$$\mathbb{E}[L(\hat{\theta}) - \inf_f \mathbb{E}[\ell(Y, f(X))]]$$

Hölder, Sobolev, Besov空間

$$\Omega = [0, 1]^d \subset \mathbb{R}^d$$

- Hölder space ($C^\beta(\Omega)$)

直観的意味

$$\|f\|_{C^\beta} = \max_{|\alpha| \leq m} \|\partial^\alpha f\|_\infty + \max_{|\alpha|=m} \sup_{x, y \in \Omega} \frac{|\partial^\alpha f(x) - \partial^\alpha f(y)|}{|x - y|^\beta}$$

滑らかさ

$$\|f\|_{B_{p,q}^s(\Omega)} = \|f\|_{L^p(\Omega)} + \|D^s f\|_{L^p(\Omega)}$$

空間的一様性

- Besov space ($B_{p,q}^s(\Omega)$) ($0 < p, q \leq \infty, 0 < s \leq m$)

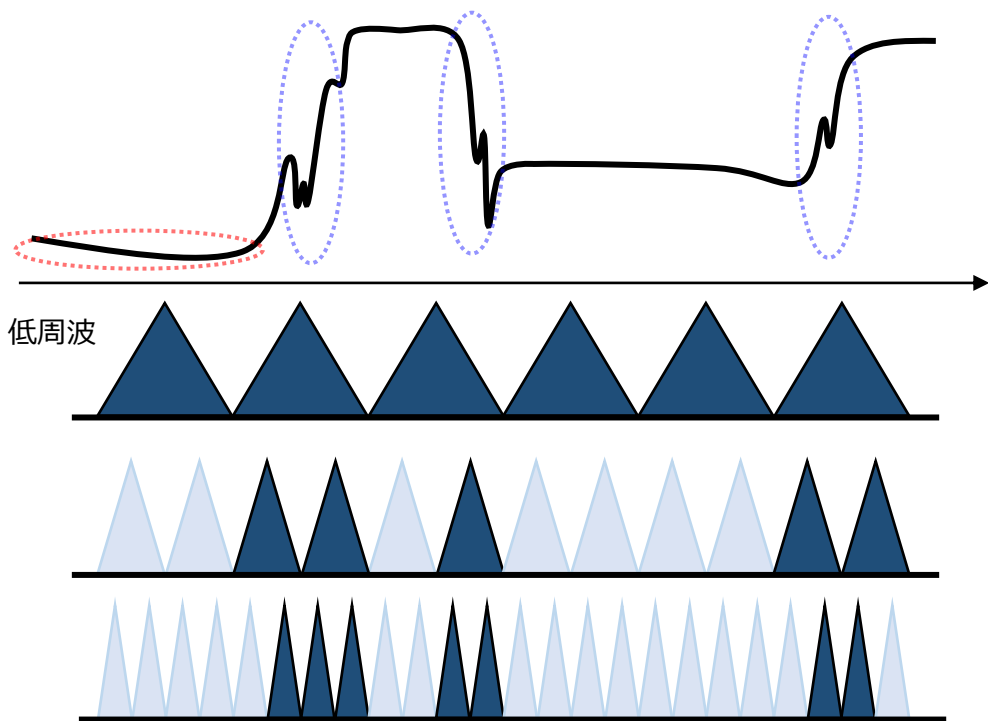
$$\omega_m(f, t)_p := \sup_{\|h\| \leq t} \left\| \sum_{j=0}^m (-1)^{m-j} \binom{m}{j} f(\cdot + jh) \right\|_{L^p(\Omega)}$$

空間的非一様性

$$\|f\|_{B_{p,q}^s(\Omega)} = \|f\|_{L^p(\Omega)} + \left(\int_0^\infty [t^{-s} \omega_m(f, t)_p]^q \frac{dt}{t} \right)^{1/q}$$

滑らかさの度合い

Besov空間とスパース性との関係

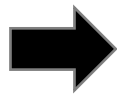


場所によって滑らかさが違うのでウェーブレット基底のスパースな線形結合が有効

$$f = \sum_{k \in \mathbb{N}_+} \alpha_k \phi_k$$

Wavelet基底による展開

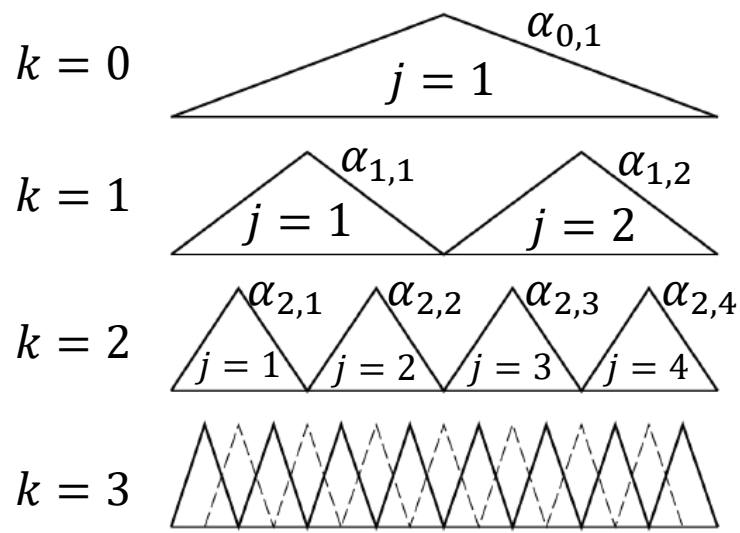
小さな p = スパースな係数



空間的な滑らかさの非一様性

Wavelet基底

解像度



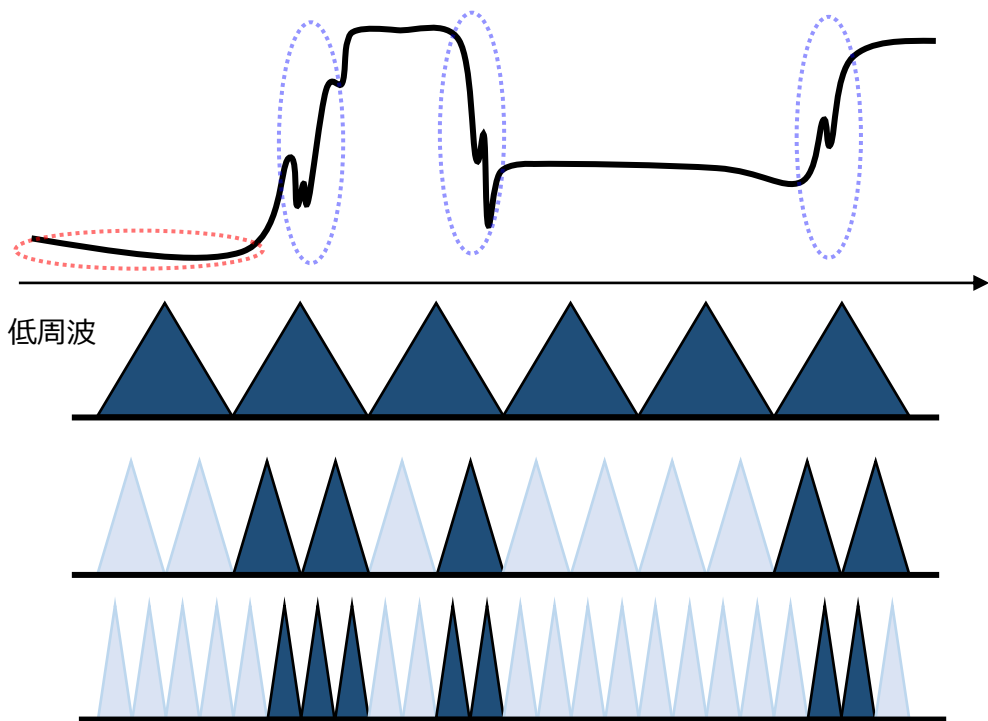
Multiresolution expansion

$$\|f\|_{B_{p,q}^s} = \left(\sum_{k \in \mathbb{N}_+} |\alpha_k|^p \right)^{1/p} \quad (0 < p)$$

(informal)

$$\|f\|_{B_{p,q}^s} \simeq \left[\sum_{k=0}^{\infty} \{2^{sk} (2^{-kd} \sum_{j \in J(k)} |\alpha_{k,j}|^p)^{1/p}\}^q \right]^{1/q}$$

Besov空間とスパース性との関係



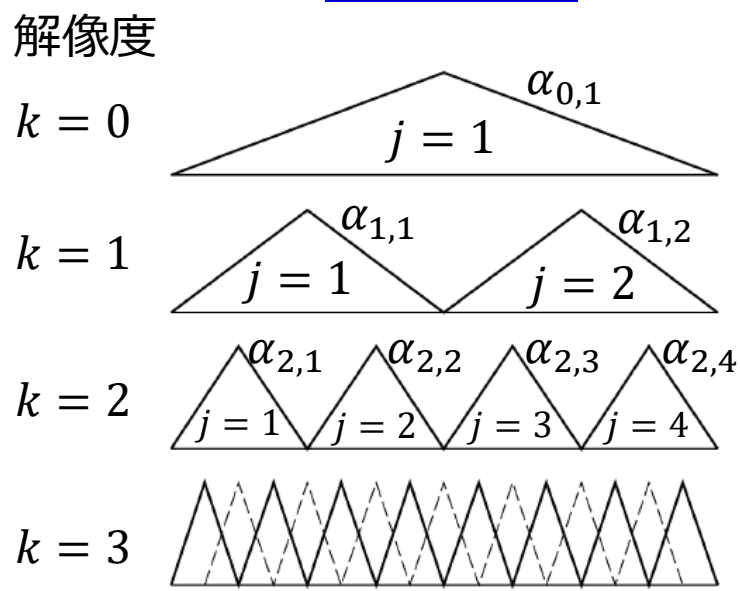
場所によって滑らかさが違うのでウェーブレット基底のスパースな線形結合が有効

$$f = \sum_{k \in \mathbb{N}_+} \alpha_k \phi_k$$

Wavelet基底による展開

小さな p = スパースな係数

Wavelet基底

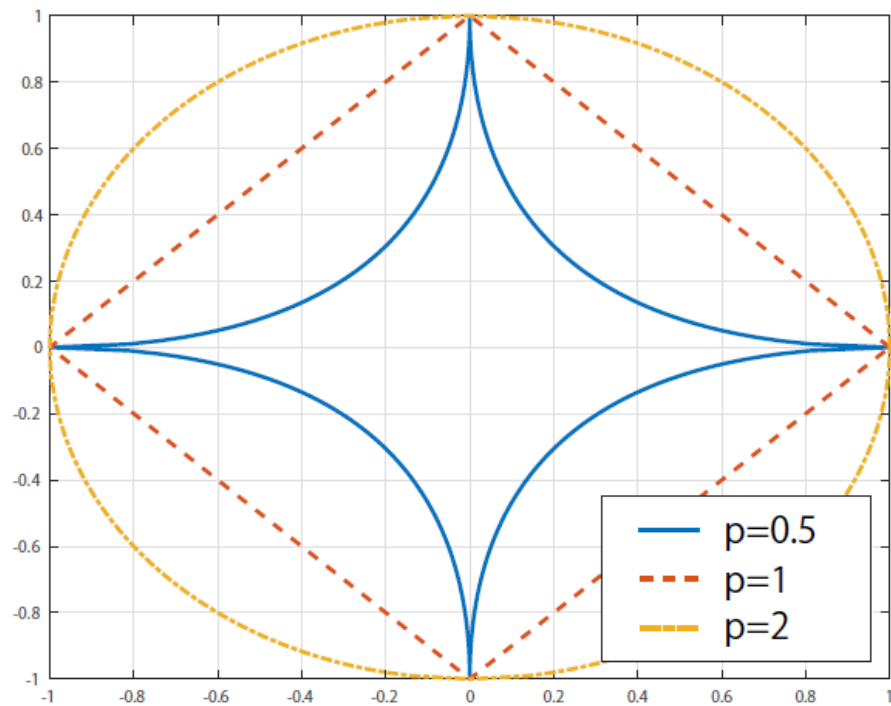


solution expansion

$$\|f\|_{B_{p,q}^s} = \left(\sum_{k \in \mathbb{N}_+} |\alpha_k|^p \right)^{1/p} \quad (0 < p)$$

$$\|f\|_{B_{p,q}^s} \simeq \left[\sum_{k=0}^{\infty} \{2^{sk} (2^{-kd} \sum_{j \in J(k)} |\alpha_{k,j}|^p)^{1/p}\}^q \right]^{1/q}$$

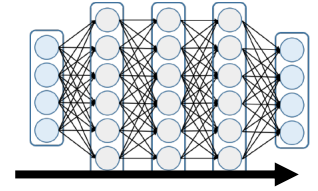
Wavelet-shrinkage
空間的な滑らかさの
非一様性



L^p ノルムのボール

「浅い」学習との比較

- 深層学習は場所によって解像度を変える適応力がある
→学習効率が良い
- 浅い学習は様々な関数を表現できる基底をあらかじめ十分用意して“待ち構える”必要がある。
→学習効率が悪い



仮定 $f^\circ \in B_{p,q}^s([0,1]^d)$: 真が“Besov空間”に入っている。

[Suzuki, ICLR2019]

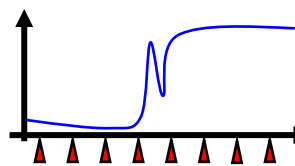
線形推定量 (非適応的手法)

カーネルリッジ回帰等：

$$n^{-\frac{2s - 2d(1/p - 1/2)_+}{2s + d - 2d(1/p - 1/2)_+}}$$

最適ではない

(n : sample size, p : uniformity of smoothness, s : smoothness)

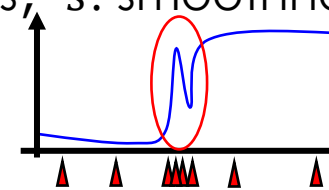


一様な解像度

深層学習

$$n^{-\frac{2s}{2s + d}}$$

最適



適応的解像度

平均二乗誤差 $E[\|\hat{f} - f^*\|^2]$ がサンプルサイズが増えるにつれ減少するレート
ミニマックス最適性の意味で理論上これ以上改善
できない精度を達成できている。

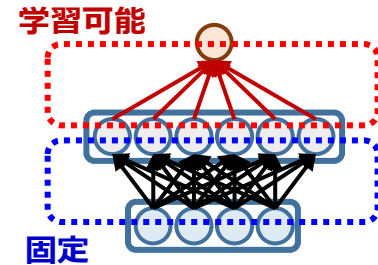
- Wavelet shrinkageより弱い条件
- 基底を用意せず最適化するだけでOK

“浅い” 学習法

Kernel ridge regression:

正則化付き最小二乗推定量

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^\infty} \sum_{i=1}^n (y_i - \beta^\top \psi(x_i))^2 + \lambda \beta^\top \beta$$



$\psi : \mathcal{X} \rightarrow \mathbb{R}^\infty$ (特徴マップ)
固定

$K_{X,X} = (\psi(x_i)^\top \psi(x_j))_{i,j=1}^{n,n}$
グラム行列 (カーネル関数)

$$\hat{f}(x) = K_{x,X} (K_{X,X} + \lambda I)^{-1} \underline{Y}$$

(see also [Imaizumi&Fukumizu, 2019])

線形推定量: 観測値 $Y = (y_i)_{i=1}^n$ に対して線形な推定量。

$$X_n = (x_1, \dots, x_n)$$

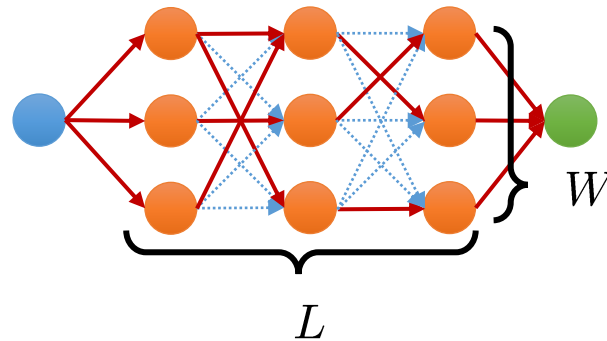
$$\hat{f}(x) = \sum_{i=1}^n \varphi_i(x; X_n) \underline{y_i}$$

線形

例

- Kernel ridge estimator
- Sieve estimator
- Nadaraya-Watson estimator
- k-NN estimator

証明: (0) ノーテーション

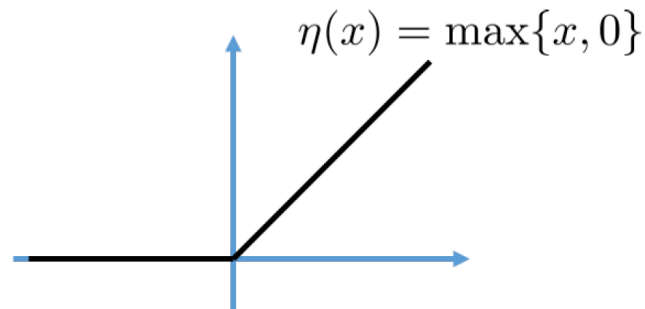


$$f(x) = (W^{(L)}\eta(\cdot) + b^{(L)}) \circ (W^{(L-1)}\eta(\cdot) + b^{(L-1)}) \circ \dots \circ (W^{(1)}x + b^{(1)})$$

$$\mathcal{F}(L, W, S, B) \left\{ \begin{array}{l} \bullet \text{ 縦幅: } L \\ \bullet \text{ 横幅: } W \\ \bullet \text{ 枝の数: } S \\ \bullet \text{ 各パラメータの上限: } B \end{array} \right.$$

の深層NNモデルの集合

- 活性化関数はReLUを仮定



証明: (1) 近似誤差の評価

- $0 < p, q, r \leq \infty$ と $0 < s < \infty$ が以下を満たすとする:

$$s > d(1/p - 1/r)_+ \quad (L^r\text{-可積分性})$$

- m を $s < \min\{m, m - 1 + 1/p\}$ を満たす整数とする.

深層ニューラルネットワークの近似誤差

ある自然数 N と用いて深さ L , 横幅 W , 枝の数 S , ノルム上界 B を以下のように定める:

$$\begin{aligned} L &= O(\log(N)), & W &= O(N), \\ S &= O(N \log(N)), & B &= O(N^{(d/p-s)_+}), \end{aligned}$$

すると, 深層NNは以下の誤差でBesov空間の元を近似できる: 大体パラメータ数

$$\sup_{f^o \in U(B_{p,q}^s([0,1]^d))} \inf_{\check{f} \in \mathcal{F}(L,W,S,B)} \|f^o - \check{f}\|_{L^r([0,1]^d)} \lesssim N^{-s/d}.$$

Pinkus (1999), Mhaskar (1996): $p = r$ かつ $1 \leq p$, ReLU活性化関数ではない.
 Petrushev (1998): $p = r = 2$, ReLU活性化関数ではない ($s \leq k + 1 + (d - 1)/2$).

証明: (1) 近似誤差の評価の導出

- **Step 1:** Besov空間の基底展開

$$f^\circ \in \mathcal{F} \quad \Rightarrow \quad f^\circ(x) = \sum_{i=1}^{\infty} \alpha_i \psi_i(x)$$

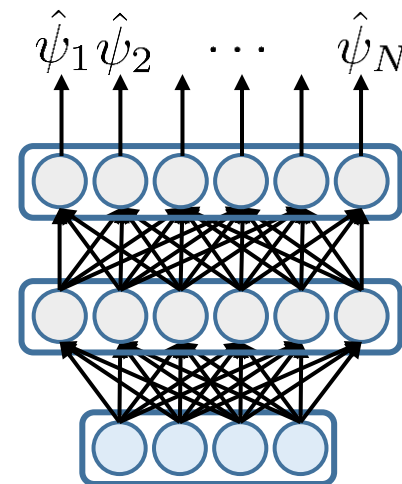
$$f^\circ = \sum_{i=1}^N \alpha_i \psi_i + \underbrace{\sum_{i=N+1}^{\infty} \alpha_i \psi_i}_{\|\cdot\|_{L^r} \leq N^{-s/d}}$$

∴ B-Splineによる適応的近似
[DeVore & Popov, 1988; Dung, 2011]

- **Step 2:** 各基底をDNNで近似.

$\psi_i \simeq \hat{\psi}_i$: DNNによる近似.

$$\Rightarrow \quad \check{f} = \sum_{i=1}^N \alpha_i \hat{\psi}_i \quad \text{: 線形結合}$$



- **Step 3:** 二つの評価を統合

$$\|f^\circ - \check{f}\|_{L^r} \leq \sum_{i=1}^N |\alpha_i| \underbrace{\|\psi_i - \hat{\psi}_i\|_{L^r}}_{\leq O(e^{-L})} + \underbrace{\left\| \sum_{i=N+1}^{\infty} \alpha_i \psi_i \right\|_{L^r}}_{\leq N^{-s/d}}$$

証明: (2) バイアス-バリエンス分解

$$\begin{aligned} & \mathbb{E}[\|f^\circ - \hat{f}\|_{L^2(P_X)}^2] \\ & \lesssim \underbrace{\frac{S[L \log(BW) + \log(Ln)]}{n}}_{\text{Variance}} + \underbrace{\inf_{f \in \mathcal{F}(L, W, S, B)} \|f - f^\circ\|_{L^2(P_X)}^2}_{\text{Bias}} \end{aligned}$$

(局所Rademacher complexityを用いて証明)

古典的なノンパラ回帰の方法でOK. DNNに関する評価は[Schmidt-Hieber, 2019; Hayakawa&Suzuki,2020]

深さ

横幅

スパース性
(非零パラメータ数)

各パラメータの絶対値の上界

$$L = O(\log(N)), W = O(N), S = O(N \log(N)), B = O(N^{(d/p-s)_+})$$

なら

$$\text{Bias} = N^{-s/d}$$

$$\text{Variance} = \frac{N \log(N)^3}{n}$$

⇒ バイアスとバリエアンスのトレードオフをバランスすればよい。

証明: (3) 推定精度の導出

- 最小二乗解 (訓練誤差最小化)

$$\hat{f} = \arg \min_{\bar{f}: f \in \mathcal{F}(L, W, S, B)} \sum_{i=1}^n (y_i - \bar{f}(x_i))^2$$

ただし, $\bar{f} = \min\{\max\{f, -F\}, F\}$ (clipping).

定理 (推定精度)

$\|f^0\|_{B_{p,q}^s} \leq 1, \|f^0\|_{\infty} \leq 1$ かつ $0 < p, q \leq \infty, s > d(1/p - 1/2)_+$ のとき,
 $N \asymp n^{\frac{d}{2s+d}}$ とすることで,

$$\|f^0 - \hat{f}\|_{L^2(P_X)}^2 \leq n^{-\frac{2s}{2s+d}} \log(n)^3.$$

$p = q = \infty$ のとき, Schmidt-Hieber (2017) に帰着.

スパース推定との繋がり



Wavelet shrinkage



Lasso

圧縮センシング

- Donoho
- Candes&Tao

第三次AIブーム

- 深層学習
- 産業応用

ビッグデータ
ブーム

スパース推定の流行

機械学習の興隆

ILSVRC
Supervision 優勝

深層学習の理論

- 適応的推定理論
- Besov空間を用いた解析



他にも様々な理論が

- 真の関数 f° の形状によって深層が有利になる

縮小ランク回帰

特徴空間の次元
が低い状況は深
層学習が得意

$$Y_i = U V X_i$$

深層

$$\frac{r(M+N)}{n}$$

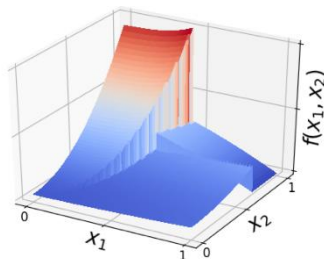
カーネル

$$\frac{MN}{n}$$

区分滑らかな関数

[Imaizumi&Fukumizu, 2019]

不連続な関数の
推定は深層学習
が得意



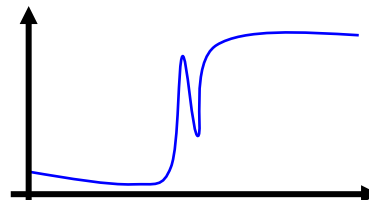
$$n^{-\frac{2s}{2s+d}} \vee n^{-\frac{\alpha}{\alpha+D-1}}$$

$$\frac{1}{\sqrt{n}}$$

Besov空間

[Suzuki, 2019]

滑らかさが非一
様な関数の推定
は深層学習が得
意



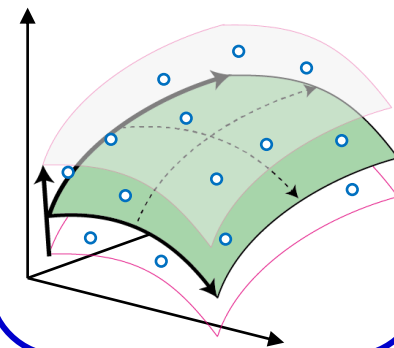
$$n^{-\frac{2s}{2s+d}}$$

$$n^{-\frac{2s-2d(1/p-1/2)_+}{2s+d-2d(1/p-1/2)_+}}$$

低次元データ

[Schmidt-Hieber, 2019] [Nakada&Imaizumi, 2019][Chen et al., 2019][Suzuki&Nitanda, 2019]

データが低次元
部分空間上に分
布していたら深
層学習が有利



$$n^{-\frac{2s}{2s+D}}$$

$$n^{-\frac{2(s-D/p+d/2)}{2(s-D/p+d/2)+d}} \vee n^{-\frac{2s}{2s+D}}$$

推定精度

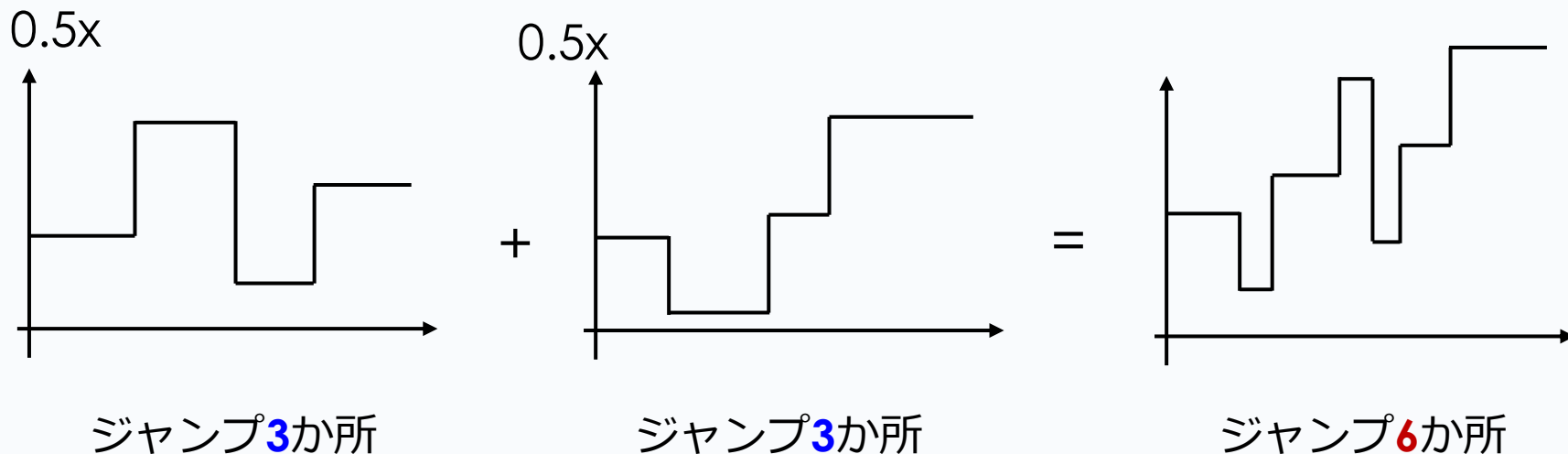
数学的に一般化

「滑らかさの非一様性」「不連続性」「データの低次元性」
凸結合を取って崩れる性質をもった関数の学習は深層学習が強い

→ 様々な性質を“凸性”で統一的に説明できる

例：ジャンプが3か所の区分定数関数

深層: $1/n$
カーネル: $1/\sqrt{n}$ \sqrt{n} 倍の違い

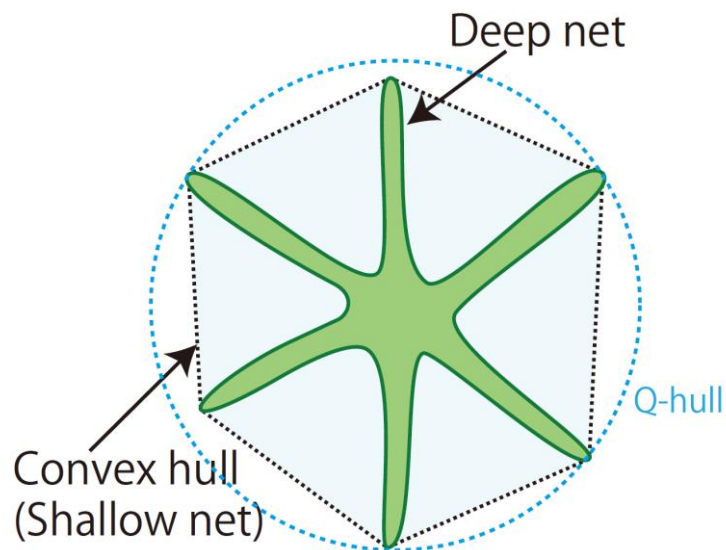


→ 「スパース性」と「非凸性」

線形推定量の最悪誤差

線形推定量： $\hat{f}(x) = \sum_{i=1}^n y_i \varphi_i(x_1, \dots, x_n; x) + b$ と書ける**任意の推定量**

例: カーネルリッジ回帰 $\hat{f}(x) = K_{x,X}(K_{X,X} + \lambda I)^{-1}Y$ (“浅い”学習法とみなす)



- 凸包を取って大きく膨らむような関数クラスにおいては深層NNを使うメリットがある.
- 特徴量(基底関数)を適応的に作ることで、ターゲットを「狙い撃ち」して近似.
- 一方、線形推定量は広めにモデルを取って「待ち構えて」いる必要がある.

$$\inf_{\hat{f}: \text{Linear}} \sup_{f^\circ \in \mathcal{F}} \mathbb{E}[\|\hat{f} - f^\circ\|_{L_2(P)}^2] = \inf_{\hat{f}: \text{Linear}} \sup_{f^\circ \in \text{conv}(\mathcal{F})} \mathbb{E}[\|\hat{f} - f^\circ\|_{L_2(P)}^2]$$

さらに条件を仮定すれば「Q-hull」まで拡張できる.

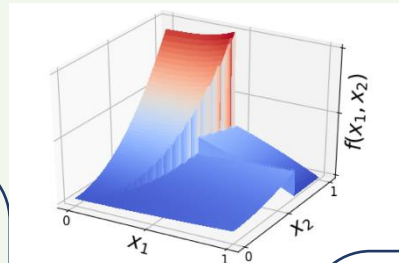
[Hayakawa&Suzuki: 2019][Donoho & Johnstone, 1994]



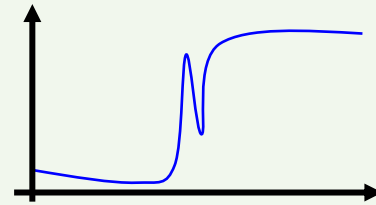
縮小ランク回帰

$$Y_i = U V X_i$$

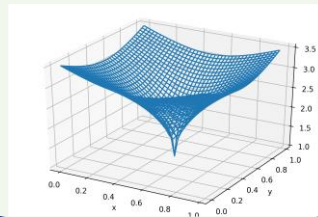
区分滑らかな関数



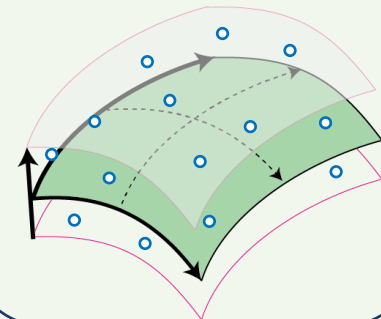
Besov空間



変動指数
Besov空間



低次元データ

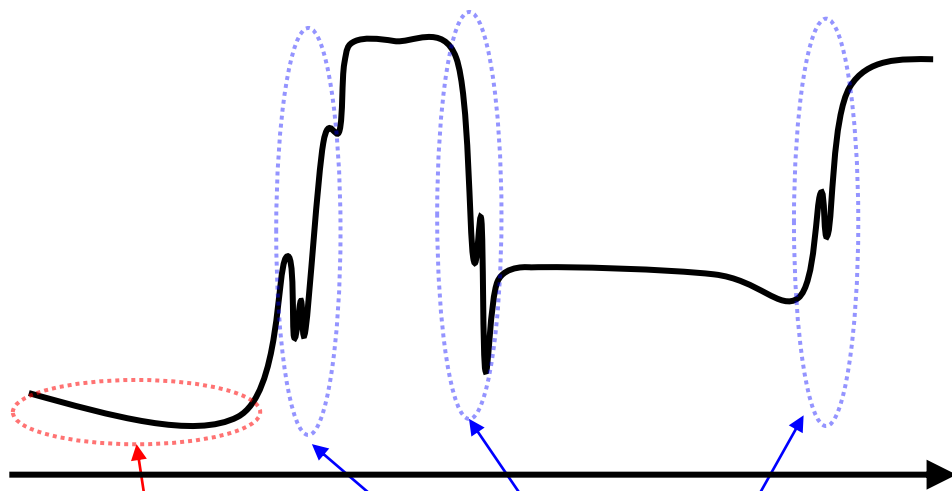


非凸性
スパース性

適応的推定法

- (勾配)ブースティング
- スパース推定
- 深層学習

滑らかな部分と
そうでない部分が混在

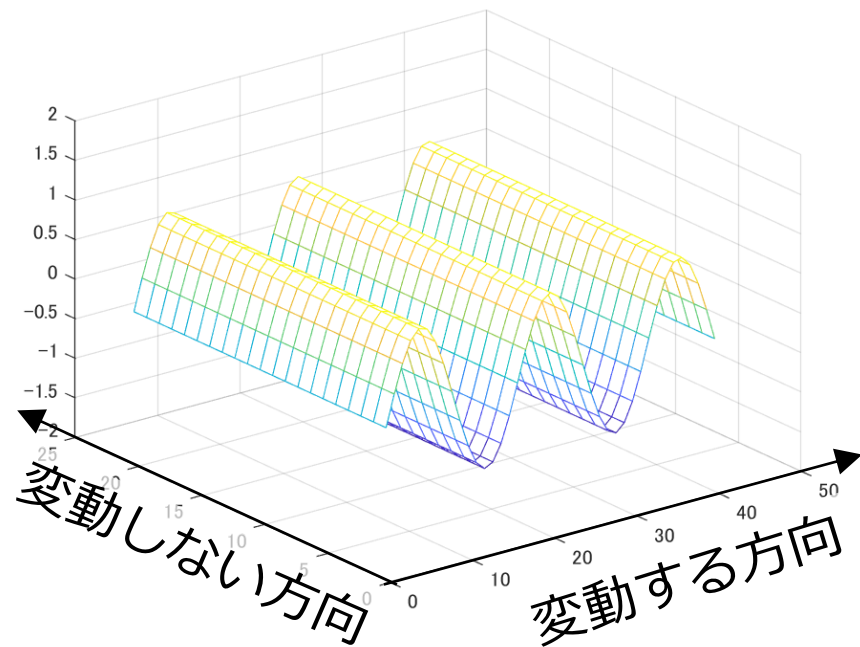


滑らか

変動が大きい
(滑らかでない)

Besov空間

大きく変動する方向と
そうでない方向が混在



深層ニューラルネットは特徴抽出に長けている。

▶ 合成関数を推定できる。

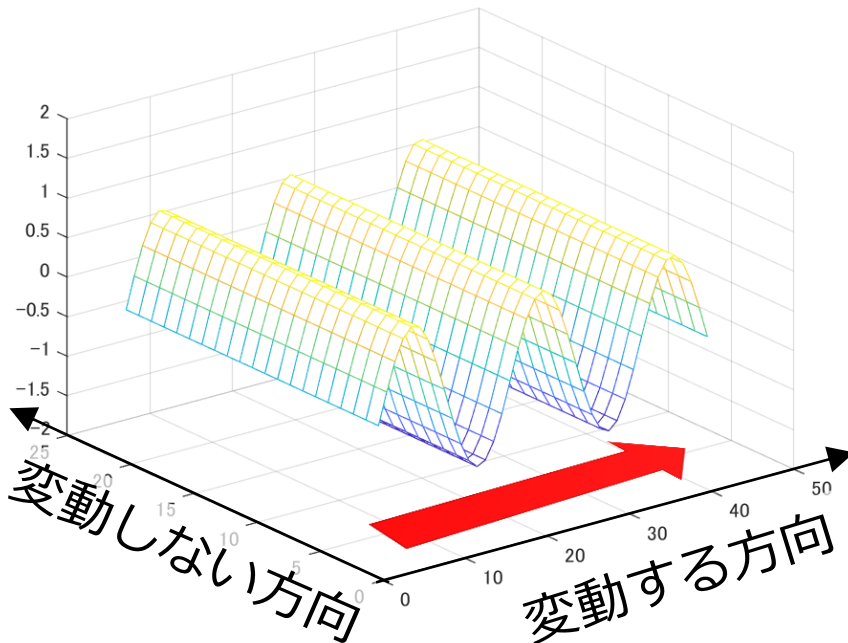
(合成される各関数 h_ℓ が“単純”であれば層に分割することで得をする)

$$f^\circ(x) = h_H \circ \underbrace{h_{H-1} \circ \cdots \circ h_1}_{\text{特徴抽出}}(x)$$

特徴抽出

(無駄な情報を削除)

典型例



• **ガウスクアーネルを用いた関数近似**

変動する方向を特定できない
→ 次元の呪い

• **NNによる関数近似**

変動する方向を表す特徴量を
中間層で抽出
→ 次元の呪いを回避

推定誤差のバウンド：

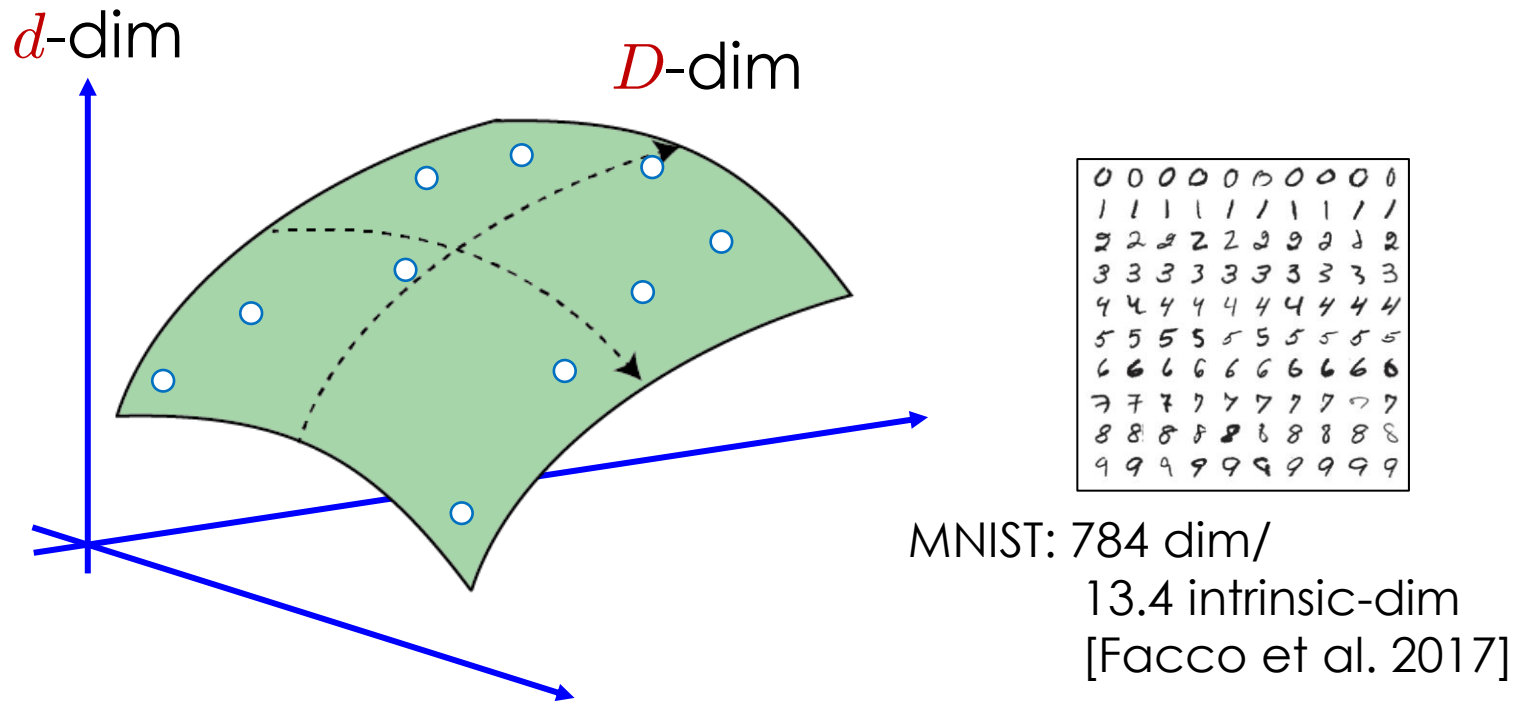
$$n^{-\frac{2s}{2s+d}}$$

近似誤差のバウンド：

$$N^{-\frac{s}{d}}$$

→ 次元の呪い

アプローチ (1): 多様体回帰



- **Classic nonparametric method:** Bickel & Li (2007); Yang & Tokdar (2015); Yang & Dunson (2016).
- **Deep learning:** Nakada & Imaizumi (2019); Schmidt-Hieber (2019); Bauer & Kohler (2019); Chen et al. (2019a,b); Liu et al. (2021).

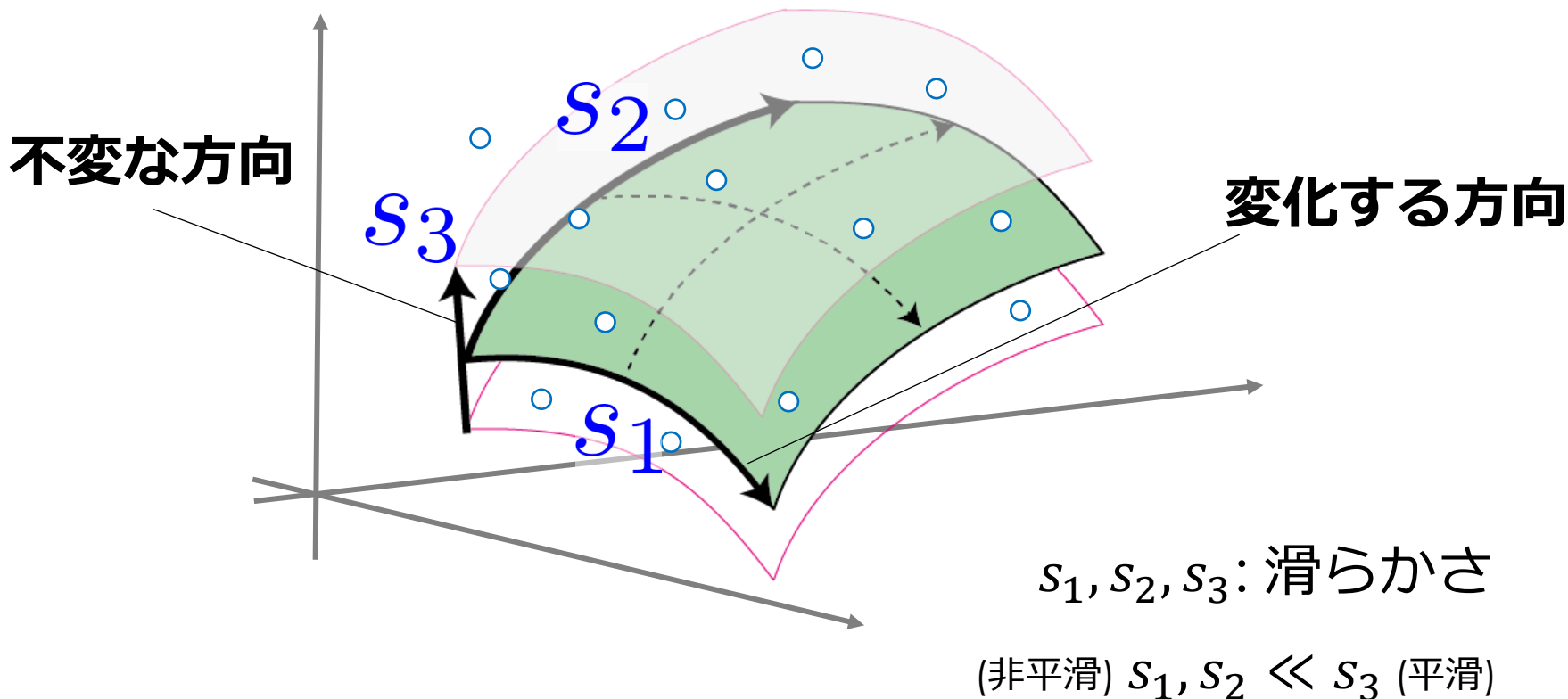
$$n^{-\frac{2s}{2s+D}}$$

データが低次元多様体*に分布していれば次元の呪いを回避できる！

* Nakada&Imaizumi (2019) では非滑らかな低次元構造も許容 (Hausdorff次元が小さい)

アプローチ (2): 関数の平滑性の非等方性⁷⁶

[Suzuki&Nitanda: Deep learning is adaptive to intrinsic dimensionality of model smoothness in anisotropic Besov space. NeurIPS2021.]



データが低次元多様体からはみ出る場合：

- 真の関数の滑らかさが方向に依存.
- 多様体に直交する方向にはほぼ定数 (滑らかさ大)

(超)高次元入力NNの学習理論

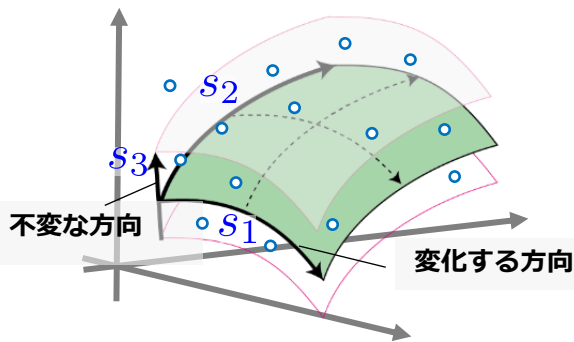
Suzuki&Nitanda: Deep learning is adaptive to intrinsic dimensionality of model smoothness in anisotropic Besov space. NeurIPS2021, spotlight.

- 真の関数が方向によって異なる滑らかさを持つ状況ではDNNは重要な方向を見つけ、次元の呪いを回避する。
- 一方で、浅い学習法は次元の呪いを受ける。

関連研究：

- 教師生徒設定における大域的最適化と次元の呪いの回避
Suzuki and Akiyama: ICLR2021, spotlight.
- 深層学習の浅層学習への優位性
Hayakawa and Suzuki: Neural Networks 2020, 日本神経回路学会論文賞.

s_1, s_2, s_3 : 滑らかさ



真の関数の滑らかさが方向によって大きく異なる状況滑らかでない方向が少なければ次元の呪いを回避できる。

→ 「**非等方的Besov空間**」を用いた理論。

真の関数のモデル：

非等方的Besov空間の元の合成関数

$$f^\circ(x) = h_H \circ \dots \circ h_1(x)$$

$h_\ell : \mathbb{R}^{m_\ell} \rightarrow \mathbb{R}^{m_{\ell+1}}$: **非等方的Besov空間** ($B_{p,q}^{s(\ell)}$).

- 滑らかさが方向によって異なる関数空間
- 合成することで様々な形状を実現

(例：多様体上の関数: 一層目で座標を抽出、二層目はその座標上の関数)

Def. (非等方的Besov空間)

$$\Delta_h^r(f)(x) := \Delta_h^{r-1}(f)(x+h) - \Delta_h^{r-1}(f)(x),$$
$$\Delta_h^0(f)(x) := f(x) \quad (h \in \mathbb{R}^d)$$

$$w_{r,p}(f,t) = \sup_{h \in \mathbb{R}^d: |h_i| \leq t_i} \|\Delta_h^r(f)\|_p$$

$$s = (s_1, \dots, s_d) \in \mathbb{R}_{++}^d$$

$$|f|_{B_{p,q}^s} := \begin{cases} \left(\sum_{k=0}^{\infty} [2^k w_{r,p}(f, (2^{-k/s_1}, \dots, 2^{-k/s_d}))]^q \right)^{1/q} & (q < \infty), \\ \sup_{k \geq 0} 2^k w_{r,p}(f, (2^{-k/s_1}, \dots, 2^{-k/s_d})) & (q = \infty). \end{cases}$$

$$\|f\|_{B_{p,q}^s} := \|f\|_p + |f|_{B_{p,q}^s}$$

(超)高次元入力NNの学習理論

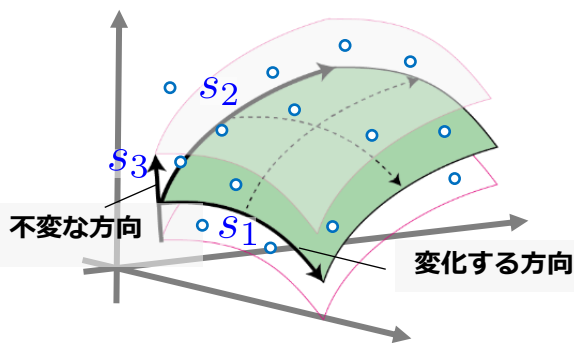
Suzuki&Nitanda: Deep learning is adaptive to intrinsic dimensionality of model smoothness in anisotropic Besov space. NeurIPS2021, spotlight.

- 真の関数が方向によって異なる滑らかさを持つ状況ではDNNは重要な方向を見つけ、次元の呪いを回避する。
- 一方で、浅い学習法は次元の呪いを受ける。

関連研究：

- 教師生徒設定における大域的最適化と次元の呪いの回避
Suzuki and Akiyama: ICLR2021, spotlight.
- 深層学習の浅層学習への優位性
Hayakawa and Suzuki: Neural Networks 2020, 日本神経回路学会論文賞.

s_1, s_2, s_3 : 滑らかさ



真の関数の滑らかさが方向によって大きく異なる状況滑らかでない方向が少なければ次元の呪いを回避できる。

→ 「**非等方的Besov空間**」を用いた理論。

真の関数のモデル：

非等方的Besov空間の元の合成関数

$$f^\circ(x) = h_H \circ \dots \circ h_1(x)$$

$h_\ell : \mathbb{R}^{m_\ell} \rightarrow \mathbb{R}^{m_{\ell+1}}$: **非等方的Besov空間** ($B_{p,q}^{s_\ell}$)

- 滑らかさが方向によって異なる関数空間
- 合成することで様々な形状を実現

(例：多様体上の関数: 一層目で座標を抽出、二層目はその座標上の関数)

Def. (非等方的Besov空間)

直感

$$\|f\|_{B_{p,q}^s} = \|f\|_{L^p} + \sum_{i=1}^d \left\| \frac{\partial^{s_i} f}{\partial x_i^{s_i}} \right\|_{L^p}$$

$$|f|_{B_{p,q}^s} := \begin{cases} \left(\sum_{k=0}^{\infty} [2^k w_{r,p}(f, (2^{-k/s_1}, \dots, 2^{-k/s_d}))]^q \right)^{1/q} & (q < \infty), \\ \sup_{k \geq 0} 2^k w_{r,p}(f, (2^{-k/s_1}, \dots, 2^{-k/s_d})) & (q = \infty). \end{cases}$$

$$\|f\|_{B_{p,q}^s} := \|f\|_p + |f|_{B_{p,q}^s}$$

推定誤差の評価

$$f^\circ(x) = h_H \circ \cdots \circ h_1(x) \quad h_\ell \in B_{p,q}^{(s_1^{(\ell)}, \dots, s_{m_\ell}^{(\ell)})}([0, 1]^{m_\ell}) \quad h_\ell : [0, 1]^{m_\ell} \rightarrow [0, 1]^{m_{\ell+1}}$$

$$\hat{f} = \arg \min_{f: \text{deep neural-net}} \sum_{i=1}^n (y_i - f(x_i))^2 \quad (\text{最小二乗推定量})$$

※今回は最適化手法に関しては議論せず，最適化はできるものと仮定する。

主結果

$$\text{Let } \tilde{s}^{(\ell)} := \left(\frac{1}{s_1^{(\ell)}} + \cdots + \frac{1}{s_{m_\ell}^{(\ell)}} \right)^{-1}, \quad \tilde{s}^{*(\ell)} := \tilde{s}^{(\ell)} \prod_{k=\ell+1}^H [(\min_j s_j^{(k)} - 1/p) \wedge 1]$$

$$\mathbb{E}[\|\hat{f} - f^\circ\|_{L^2(P_X)}^2] \lesssim \max_{\ell \in [H]} n^{-\frac{2\tilde{s}^{*(\ell)}}{2\tilde{s}^{*(\ell)}+1}} \log(n)^3$$

各方向への滑らかさの調和平均が収束レートを定める。

例: $H = 1$ の時

$$n^{-\frac{2\tilde{s}}{2\tilde{s}+1}}$$

$$\tilde{s} = (s_1^{-1} + \cdots + s_d^{-1})^{-1}$$

少ない数の方向において s_i が小さく (滑らかでない), その他の方向には s_i が大きい (滑らか) であるとき, 次元の呪いを回避できる。

浅い学習方法との比較 (informal) :

深層

$$n^{-\frac{2\tilde{s}}{2\tilde{s}+1}}$$

$$\tilde{s} = (s_1^{-1} + s_2^{-1} + s_3^{-1})^{-1}$$



浅層

$$n^{-\frac{2s_1}{2s_1+d}}$$

(次元の呪いを受ける)



- 特徴抽出能力の重要性を理論的に正当化
- 浅い学習方法は一番滑らかでない方向の滑らかさ (s_1) が支配的で, 次元の呪いを受ける。
- 証明には「凸法の議論」を用いる。

線形推定量との比較 (より正確なステートメント)

$$\mathcal{F} = \{f^\circ = g(Wx + b) \mid g \in U(B_{p,q}^s([0,1]^D)), W \in \mathbb{R}^{D \times d}, b \in \mathbb{R}^D\}$$

(s.t. $Wx + b \in [0,1]^D$ for any $x \in [0,1]^d$)

f° は D -次元部分空間にのみ依存

$$\text{If } s > \frac{D}{d-D} \left(\frac{d}{2} - \frac{D}{p} + c \right)$$

(非適応的)

深層

$$n^{-\frac{2s}{2s+D}}$$

$$(n^{-\frac{2s}{2s+1}} \text{ when } D = 1)$$

\ll

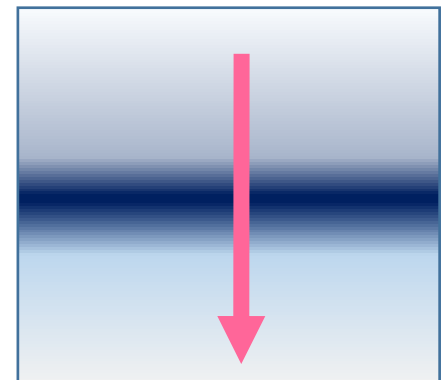
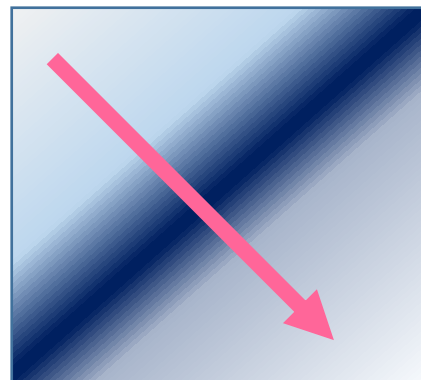
線形推定量

$$n^{-\frac{2(s-D/p+d/2+c)}{2(s-D/p+d/2+c)+d}}$$

$$c = 1 \text{ if } D < d/2, c = 0 \text{ if } D \geq d/2.$$

$$(n^{-\frac{2s+d}{2s+2d}} \text{ when } D = 1 \text{ and } p = 1)$$

深層にすることで次元の呪いを回避できている。



無限次元入力NN

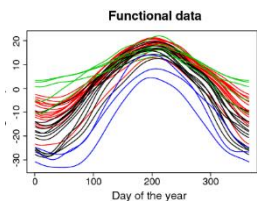
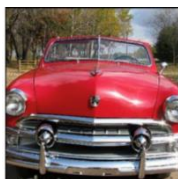
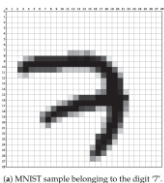
無限次元入力

(画像, 音声信号, 自然言語,...)

無限 (高) 次元データ

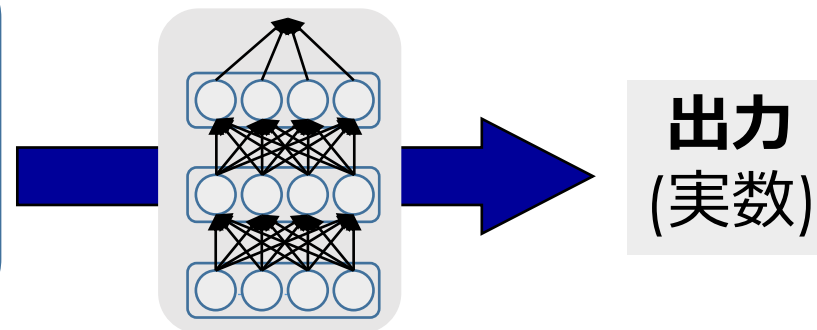
画像データ

関数データ



- 音声
- 文章
-

[Okumoto&Suzuki: Learnability of convolutional neural networks for infinite dimensional input via mixed and anisotropic smoothness. ICLR2022.]



[Ramsay, J., Hooker, Giles, & Graves, Spencer. (2009). Functional data analysis with R and MATLAB (Use R!). Dordrecht: Springer.]

さきほどの議論は入力の有限次元性を使っていた。
→ 実は**無限次元**まで拡張できる。

無限次元入力に対する深層学習の統計理論

- 次元に依存しないバウンド (有限次元の拡張)
- 畳み込みNNによる特徴量の抽出

$$\mathbb{E}[\|\hat{f} - f^\circ\|_{L_2(P_X)}^2] \lesssim n^{-\frac{2(\tilde{\alpha}-v)}{2(\tilde{\alpha}-v)+1}} (\log n)^{\frac{2}{q}+2} \max\{(\log n)^{4/q}, (\log n)^4\}$$

拡散モデルの統計理論

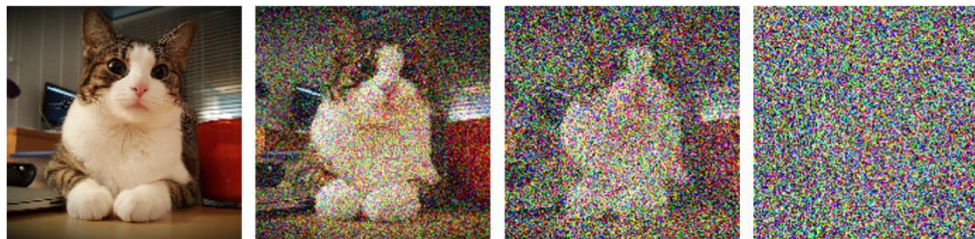
[Kazusato Oko, Shunta Akiyama, Taiji Suzuki: Diffusion Models are Minimax Optimal Distribution Estimators. ICML2023]



Stable diffusion, 2022.

$$dX_t = -X_t dt + \sqrt{2} dB_t$$

Forward process



Backward process

$$dY_t = (Y_t + 2\nabla \log(p_{\bar{T}-t}(Y_t)))dt + \sqrt{2}dB_t$$

$(Y_t \sim X_{\bar{T}-t})$

経験スコアマッチング推定量:

$$\hat{s} = \arg \min_{s \in \text{DNN}} \frac{1}{n} \sum_{i=1}^n \int_{t=\underline{T}}^{\bar{T}} \mathbb{E}_{X_t|X_0=x_{0,i}} [\|s(X_t, t) - \nabla \log p_t(X_t|x_{0,i})\|^2] dt$$

定理

Let \hat{Y} be the r.v. generated by the backward process w.r.t. \hat{s} , then

$$\mathbb{E}_{D_n} \left[\text{TV}(\hat{Y}, X_0) \right] \lesssim n^{-\frac{s}{2s+d}} \log^9(n), \quad (s: \text{密度関数の滑らかさ})$$

$$\mathbb{E}_{D_n} \left[W_1(\hat{Y}, X_0) \right] \lesssim n^{-\frac{s+1-\delta}{2s+d'}} \quad (\text{for any } \delta > 0).$$

どちらも (ほぼ) **ミニマックス最適** [Yang & Barron, 1999; Niles-Weed & Berthet, 2022].

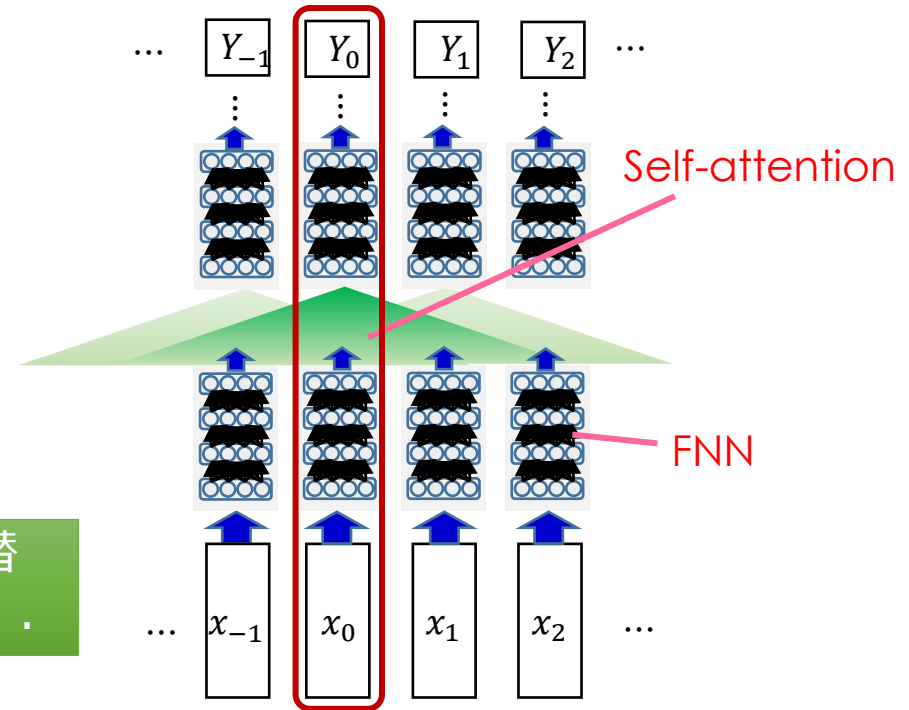
(Estimator for W_1 distance requires some modification)

[Shokichi Takakura, Taiji Suzuki: Approximation and Estimation Ability of Transformers for Sequence-to-Sequence Functions with Infinite Dimensional Input. ICML2023]

Transformerの性質

- かなり広いトークン幅から重要なトークンを選ぶ。
→ 次元の呪い？
- 入力に依存して重要なトークンを選択できる。
→ 次元の呪いを回避！

入力に依存して重要なトークンを切り替えることで、関数を「切り替えている」。



定理 (推定誤差)

$$\frac{1}{r-l+1} \sum_{j=l}^r \mathbb{E}[\|\hat{F}_j - F_j^\circ\|_{L_2(P_X)}^2] \lesssim n^{-\frac{2a^\dagger}{2a^\dagger+1}} (\log n)^{2/\alpha+2+\max\{4/\alpha, 4\}}$$

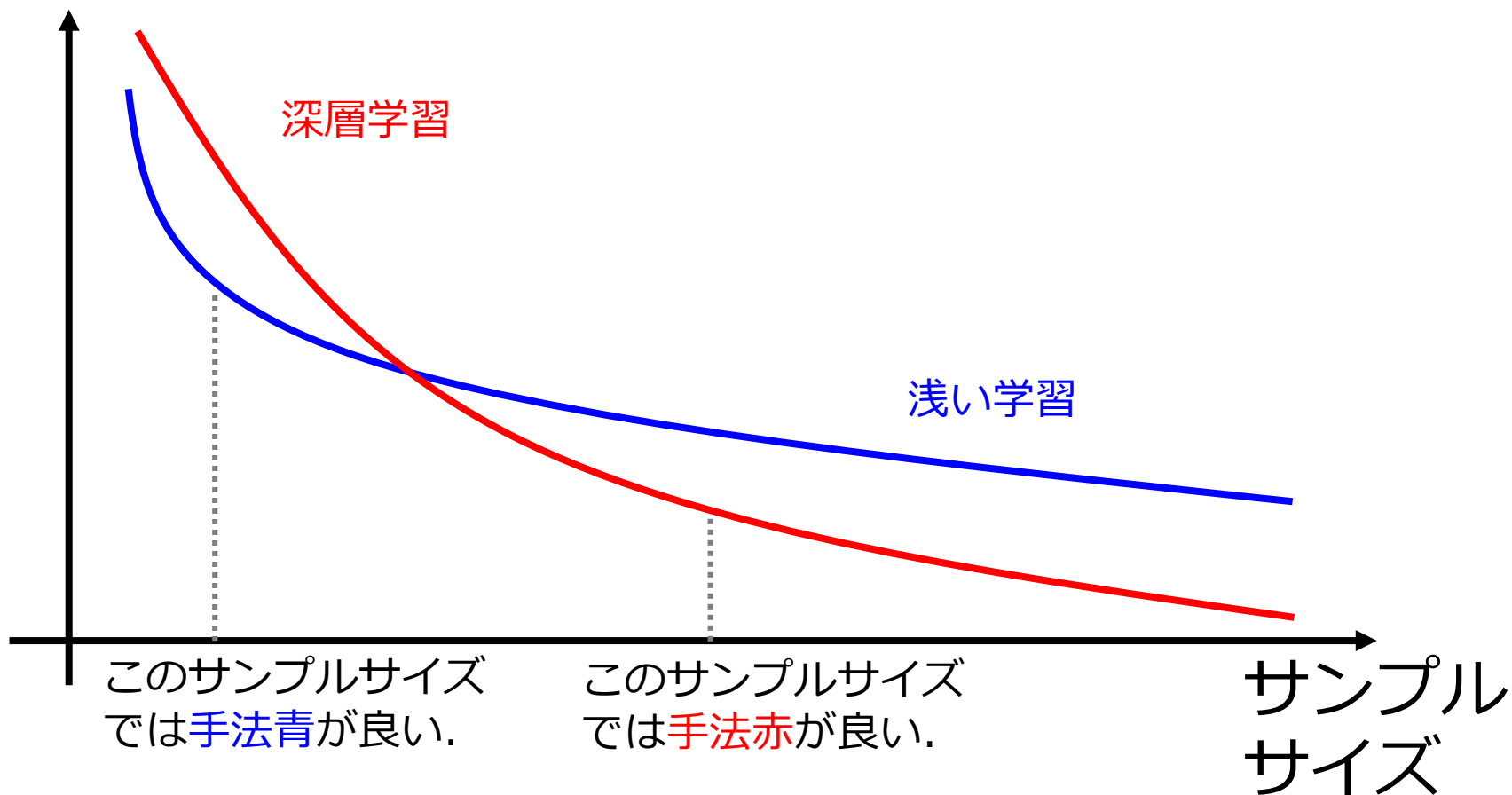
(ほぼミニマックス最適)

- 入力が無限次元でも多項式オーダーの収束レート。

収束レートに関する注意

注意：収束レートが速いからといって，その手法が常に良いとは限らない。

推定誤差



第2部

深層学習の汎化誤差

-Overparametrization-

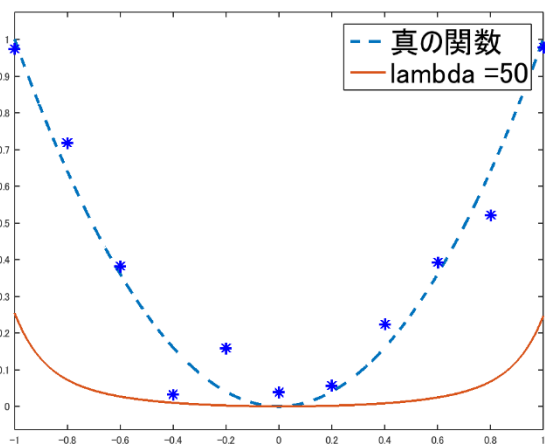
- これまでの議論は、実は問題に合わせて「適切なサイズのネットワーク」を用いた場合の議論であった。
- 実際は、かなりサイズの大きなネットワークを用いる。

→ **Overparameterization**
(過剰パラメータ化)

過学習

「なんでも表現できる方法」が最適とは限らない
 少しのノイズにも鋭敏に反応してしまう

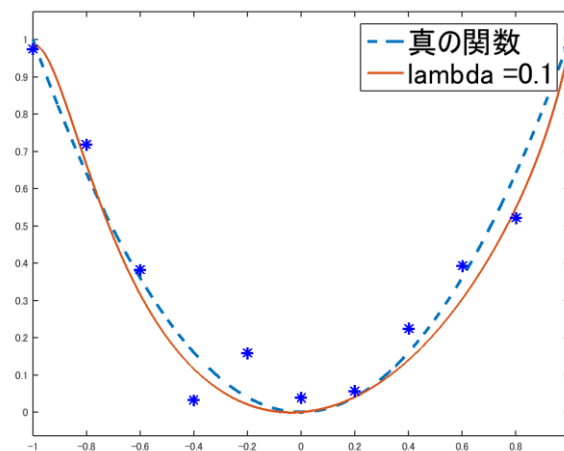
悪い学習結果



過小学習

説明力が低すぎる

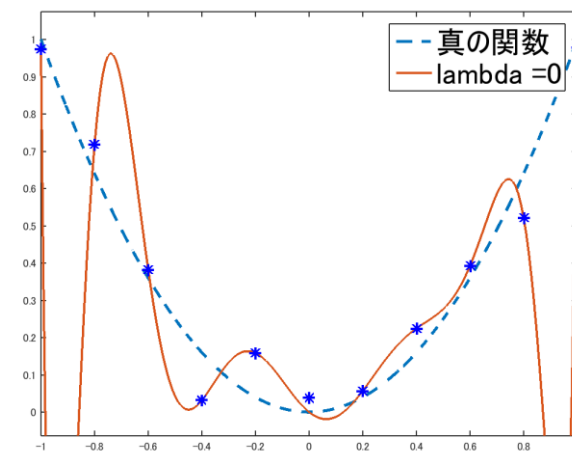
良い学習結果



適切な学習

説明力が適切

悪い学習結果



過学習

説明力が高すぎる
 (複雑すぎる)

学習に用いるデータには誤りも含まれる

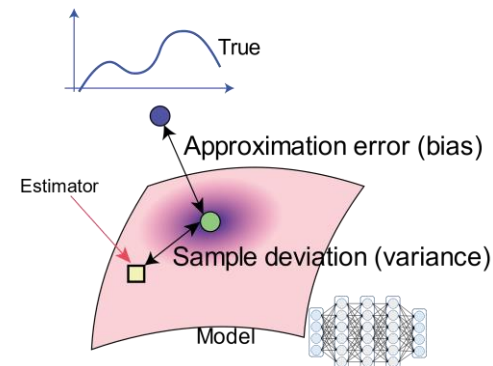
一見当てはまりが良いので危険

通常の学習理論

モデル \mathcal{F} : d -次元パラメータ

n : データサイズ

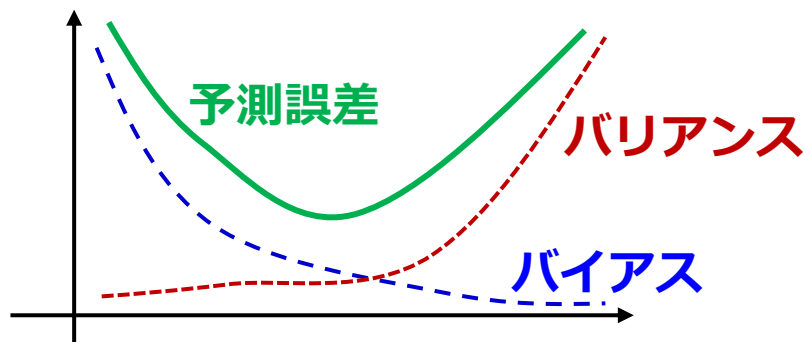
$$\hat{f} \leftarrow \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$



$$\text{予測誤差} = \underbrace{\frac{d}{n}}_{\text{Variance}} + \underbrace{\inf_{f' \in \mathcal{F}} \|f - f'\|_{L^2(P_X)}^2}_{\text{Bias (近似誤差)}}$$

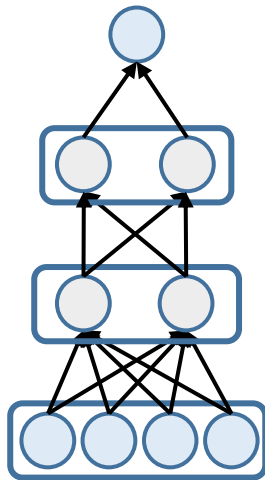
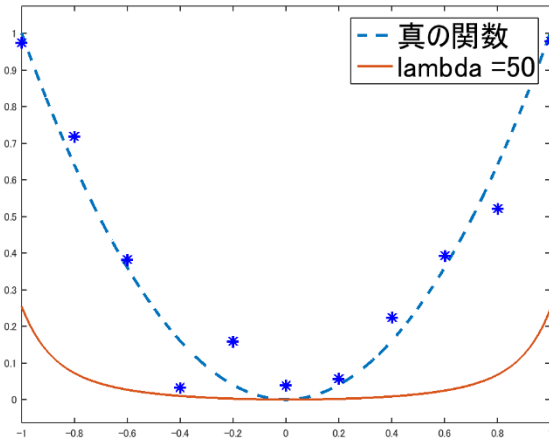
**データサイズ (n) \ll パラメータ数 (d)
の場合, 通常の学習理論をナイーブに当てはめると失敗する.**

※カーネル法のような無限次元モデルもバイアスバリエンス分解で説明できるので, あくまで「ナイーブに」当てはめた場合の話.

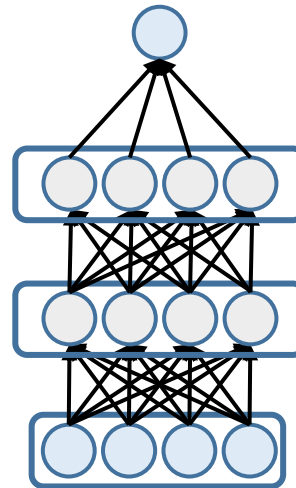
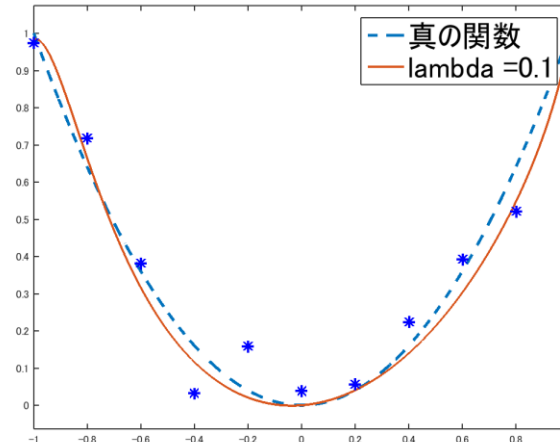


ニューラルネットワークでは？

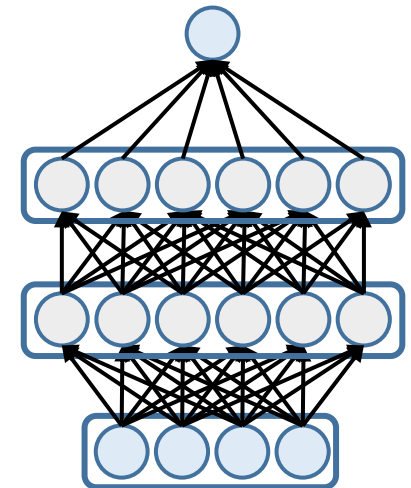
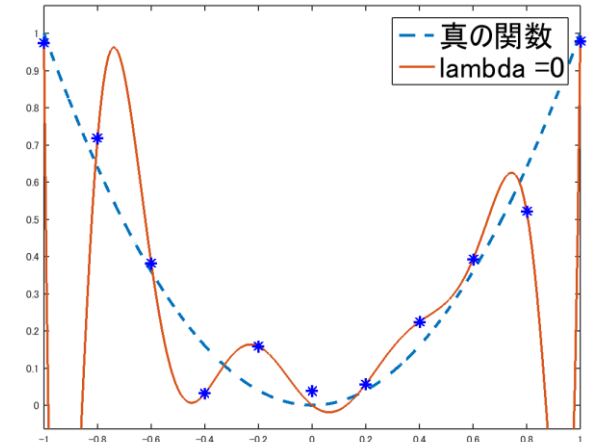
過小学習



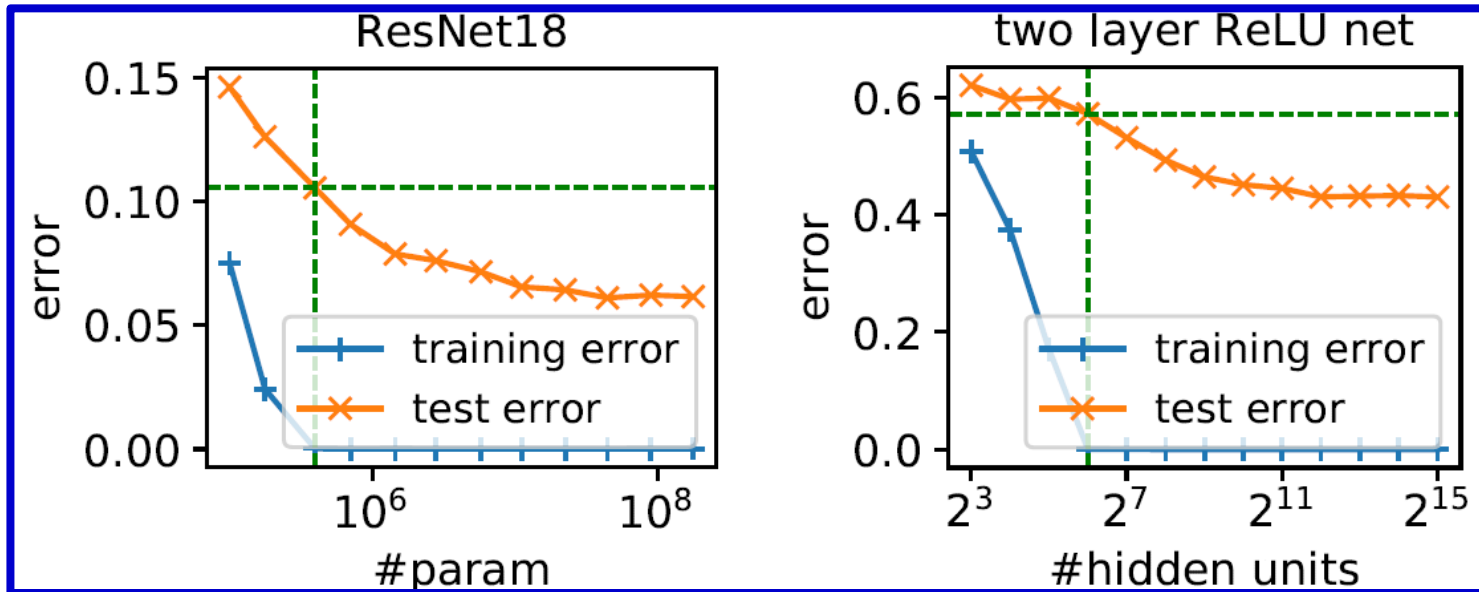
適切な学習



過学習



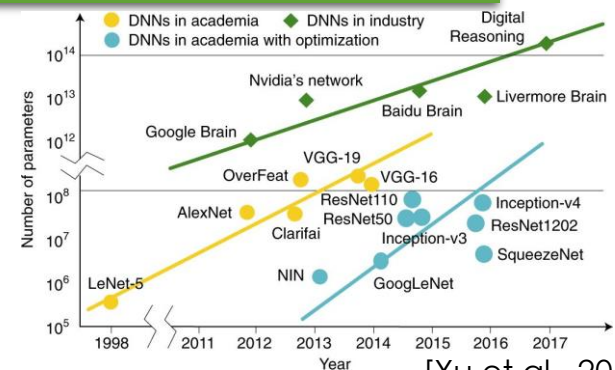
実際は...



[Neyshabur et al., ICLR2019]

ネットワークのサイズを大きくしても過学習しない

データサイズ : 120万
モデルパラメータサイズ : 10億



[Xu et al., 2018]

「Overparameterization」

パラメータサイズがデータサイズを超えている状況での汎化性能を説明したい。

パラメータ数 >> **データサイズ** >> **実質的自由度**

数十億

数百万

数十万

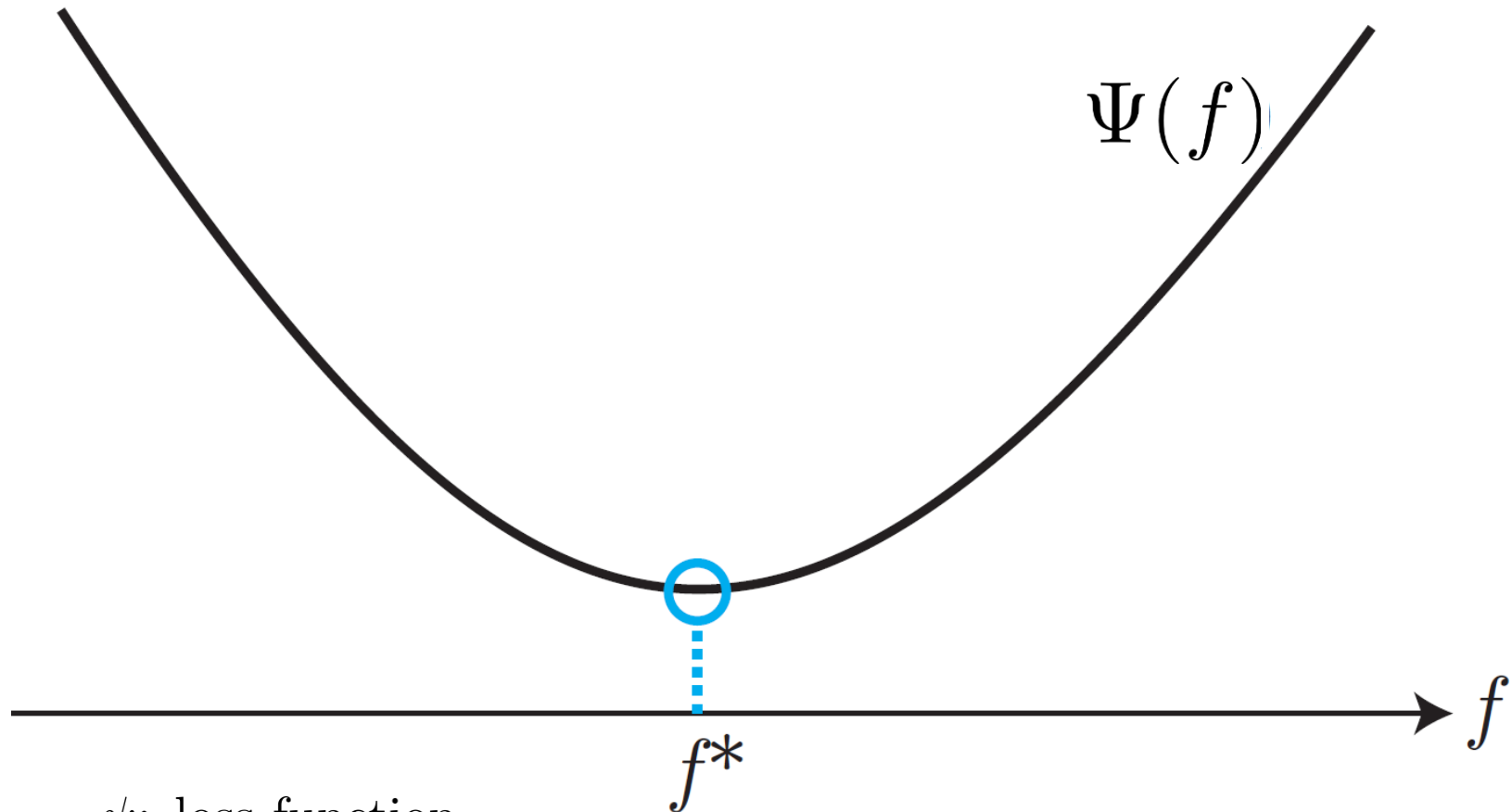
[仮説] 見かけの大きさ (パラメータ数) よりも
実質的な大きさ (自由度) はかなり小さいはず。

“実質的自由度”を調べる研究：

- ノルム型バウンド
- 圧縮型バウンド

「実質的自由度」として何が適切かを見つけることが理論上問題になる。

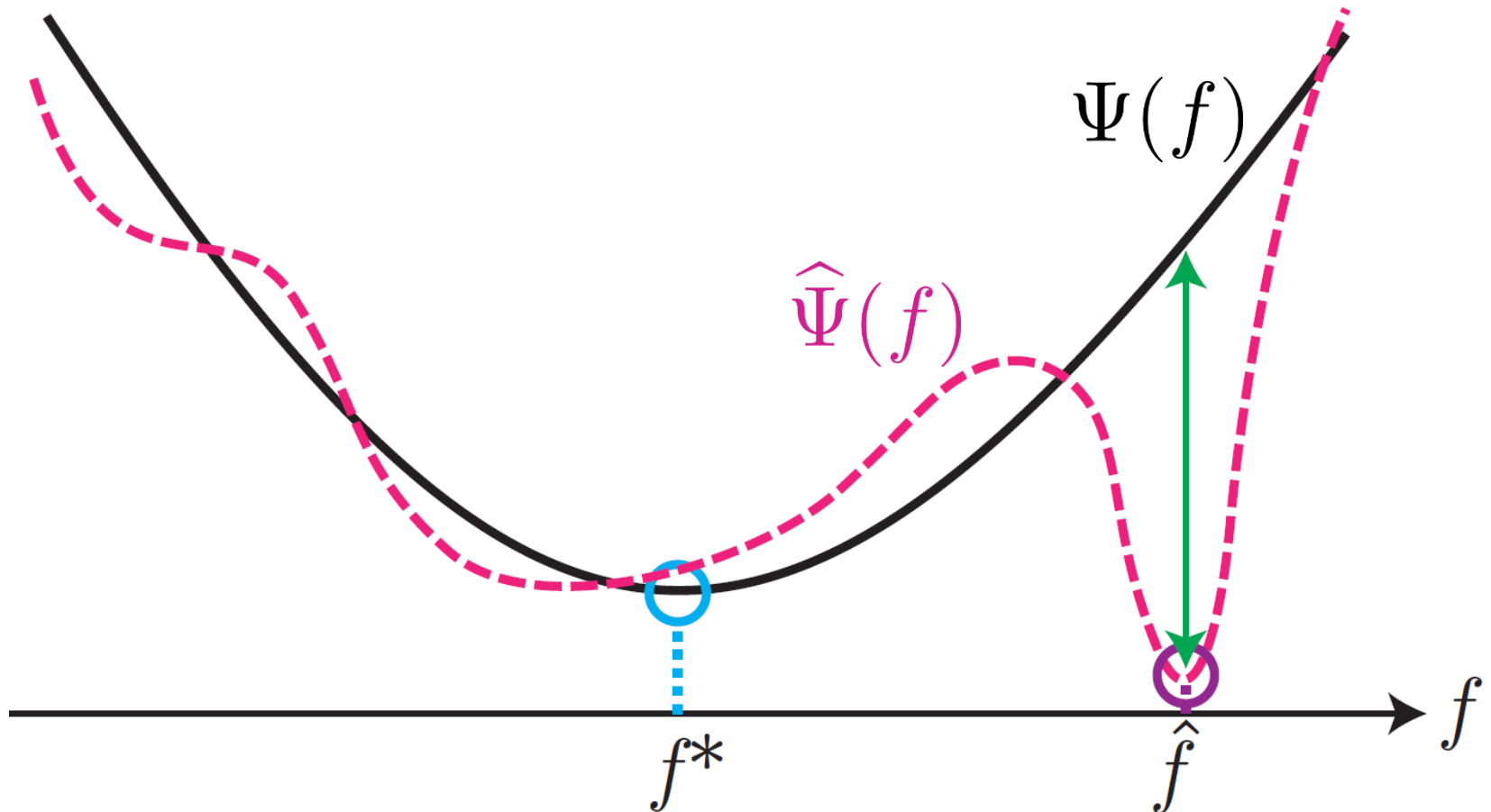
Uniform bound



ψ : loss function

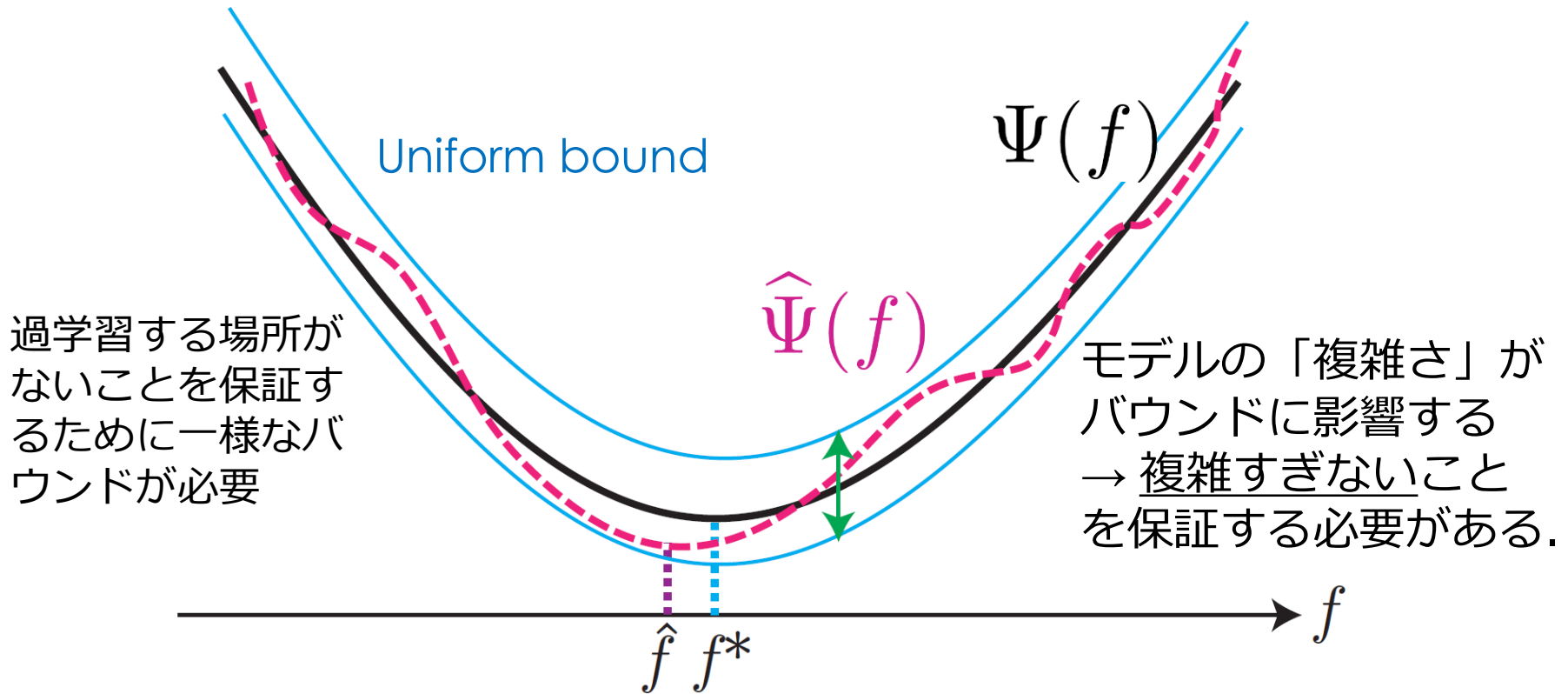
$$\hat{\Psi}(f) := \frac{1}{n} \sum_{i=1}^n \psi(y_i, f(x_i))$$

$$\Psi(f) := \mathbb{E}[\psi(Y, f(X))]$$



“運良くデータに強く当てはまる”場所があるかもしれない。
→ 過学習

Uniform bound



Uniform bound

$$\Psi(\hat{f}) - \hat{\Psi}(\hat{f}) \leq \sup_{f \in \mathcal{F}} \left\{ \Psi(f) - \hat{\Psi}(f) \right\}$$

Rademacher complexity

$$\bar{R}_n = \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i \psi(y_i, f(x_i)) \right]$$

Dudley integral (covering num.)

$$\frac{1}{\sqrt{n}} \int_{1/n}^{\infty} \sqrt{\log(N(\epsilon, \|\cdot\|_{\infty}, \mathcal{F}))} d\epsilon$$

深層学習の汎化誤差バウンド (抜粋)

ノルム型バウンド

Author	Rate	Bound type
Neyshabur et al. (2015)	$\frac{2^L \prod_{\ell=1}^L R_{\ell,F}}{\sqrt{n}}$	Norm base
Bartlett et al. (2017)	$\frac{\prod_{\ell=1}^L R_{\ell,2}}{\sqrt{n}} \left(\frac{R_{\ell,2 \rightarrow 1}^{2/3}}{R_{\ell,2}^{2/3}} \right)^{3/2}$	Norm base
Neyshabur et al. (2017)	$\frac{\prod_{\ell=1}^L R_{\ell,2}}{\sqrt{n}} \sqrt{L^2 W \sum_{\ell=1}^L \frac{R_{\ell,F}^2}{R_{\ell,2}^2}}$	Norm base
Golowich et al. (2018)	$\prod_{\ell=1}^L R_{\ell,F} \min \left\{ \frac{1}{n^{1/4}}, \sqrt{\frac{L}{n}} \right\}$	Norm base
Li et al. (2018) Harvey et al. (2017)	$\frac{\prod_{\ell=1}^L R_{\ell,2} \sqrt{L^2 W^2}}{\sqrt{n}}$	VC-dim Naïve bound
Arora et al. (2018)	$\sqrt{\frac{L^2 \max_{1 \leq i \leq n} \hat{f}(x_i) ^2 \sum_{\ell=1}^L \frac{1}{\mu_{\ell}^2 \mu_{\ell \rightarrow}^2}}{n}}$	Compression
Baykal et al. (2018)	$\sqrt{\frac{L^2 \max_{1 \leq i \leq n} \hat{f}(x_i) ^2 \sum_{\ell=2}^L (\hat{\Delta}^{\ell \rightarrow})^2 \sum_{i=1}^W S_i^{\ell}}{n}}$	Compression
Suzuki et al. (2018)	$\sum_{\ell=2}^L \sqrt{\lambda_{\ell}} + \sqrt{\frac{\sum_{\ell=1}^L m_{\ell+1}^{\#} m_{\ell}^{\#}}{n}}$	Compression

圧縮型バウンド

L : depth

$W = \max_{\ell} m_{\ell}$: width

R_F : Frobenius norm, R_2 : operator norm

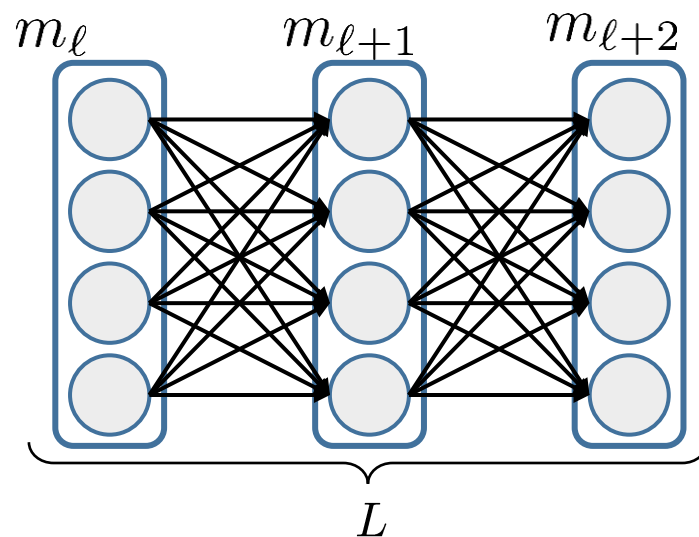
Naïve bound (VC-次元)

$$\Psi(\hat{f}) \leq \hat{\Psi}(\hat{f}) + \boxed{?}$$

VC-次元 [Harvey et al.2017]

$$\sqrt{L \frac{\sum_{\ell=1}^L m_{\ell} m_{\ell+1}}{n} \log(n)}$$

$$\left(\sqrt{L \frac{\# \text{ of parameters}}{n}} \right)$$



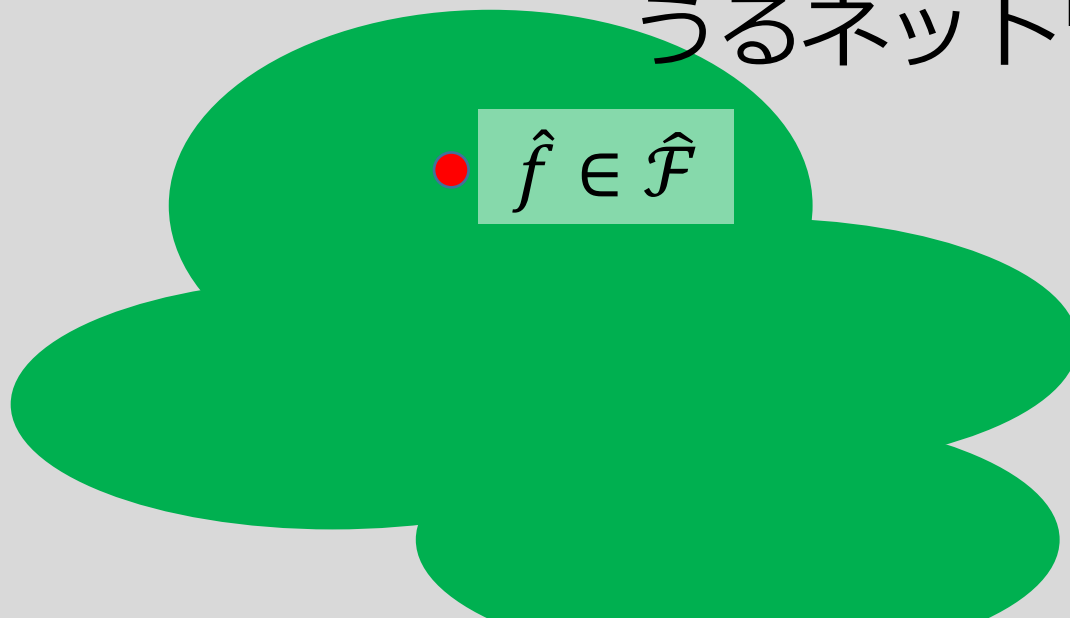
- m_{ℓ} : width of ℓ -th layer
- L : depth

- ⊗ パラメータ数 $\sum_{\ell=1}^L m_{\ell} m_{\ell+1}$ がそのままバウンドに現れてしまう.
- ⊗ パラメータ数 \gg データサイズ の状況を説明できていない.

どうやって改善するか？

\mathcal{F} : ネットワーク全体

$\hat{\mathcal{F}}$: 学習済みモデルが入り
うるネットワークの集合



※ $\hat{\mathcal{F}}$ は \mathcal{F} よりもはるかに小さいと考えられる.

“典型的な学習済みネットワーク”の集合を解析する.

ノルム型バウンド

NN model: $f(x) = (W^{(L)}\eta(\cdot)) \circ (W^{(L-1)}\eta(\cdot)) \circ \dots \circ (W^{(1)}x)$

- Bartlett et al. (2017): Spectrally-normalized margin bound

$$\frac{1}{\sqrt{n}} \prod_{\ell=1}^L R_{\ell,2} \left(\sum_{\ell=1}^L \frac{R_{\ell,2 \rightarrow 1}^{2/3}}{R_{\ell,2}^{2/3}} \right)^{3/2}$$

$$R_{\ell,2} := \|W^{(\ell)}\| = \sigma_1(W^{(\ell)}) \quad R_{\ell,2 \rightarrow 1} := \sum_j \|W_{j,:}^{(\ell)}\|$$

(最大特異値)

- $\prod_{\ell=1}^L R_{\ell,2}$ は中間層から最終層にどう誤差が伝搬するかを表すリップシッツ定数.
- $R_{\ell,2 \rightarrow 1}/R_{\ell,2}$ は行列 $W^{(\ell)}$ のスパース性を表す. もし $W^{(\ell)}$ が小さなノルムの行をたくさん含むなら, この値は小さくなる.

☺ 横幅に依存しない → 過剰パラメータ化されたネットワークにも使える!

ノルム型バウンド (2)

NN model: $f(x) = (W^{(L)}\eta(\cdot)) \circ (W^{(L-1)}\eta(\cdot)) \circ \dots \circ (W^{(1)}x)$

- Neyshabur et al. (2017): PAC-Bayes bound

$$\prod_{\ell=1}^L R_{\ell,2} \sqrt{\frac{L^2 \max_{\ell} m_{\ell} \sum_{\ell=1}^L \frac{R_{\ell,F}}{R_{\ell,2}}}{n}}$$

$$R_{\ell,2} := \|W^{(\ell)}\| = \sigma_1(W^{(\ell)})$$

(maximum singular value)

$$R_{\ell,F} := \|W^{(\ell)}\|_F : \text{Frobenius norm}$$

(smaller than $R_{\ell,2 \rightarrow 1}$)

- $R_{\ell,2 \rightarrow 1}$ の代わりに $R_{\ell,F}$ (which is smaller) 出ている。
- その代わりに, 横幅 m_{ℓ} も出てきている。
- 重み行列がスパースではない状況では, こちらの方がタイト。

[B. Neyshabur, S. Bhojanapalli, D. McAllester, and N. Srebro. A PAC-Bayesian approach to spectrally-normalized margin bounds for neural networks. ICLR2018.]

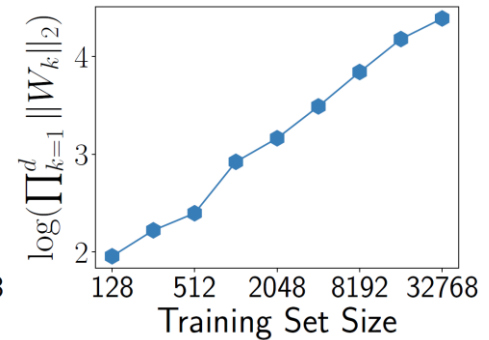
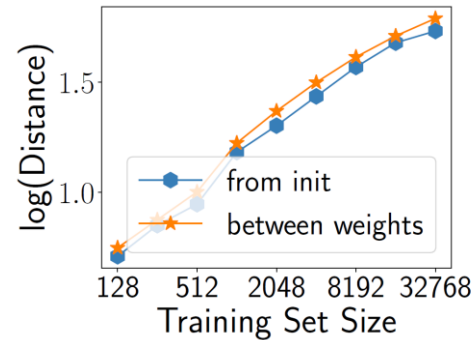
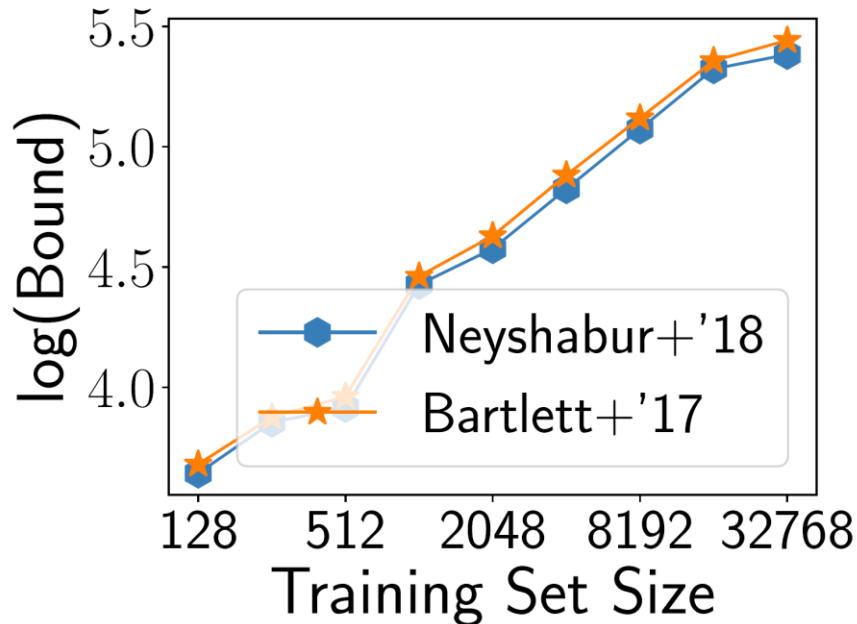
- Golowich et al. (2018)

$$\prod_{\ell=1}^L R_{\ell,F} \min \left\{ \frac{1}{n^{1/4}}, \sqrt{\frac{L}{n}} \right\}$$

- 完全にFrobeniusノルムだけで特徴づけられている。
- 横幅が出てこない。

[N. Golowich, A. Rakhlin, and O. Shamir. Size-independent sample complexity of neural networks. COLT2018, pp. 297–299.]

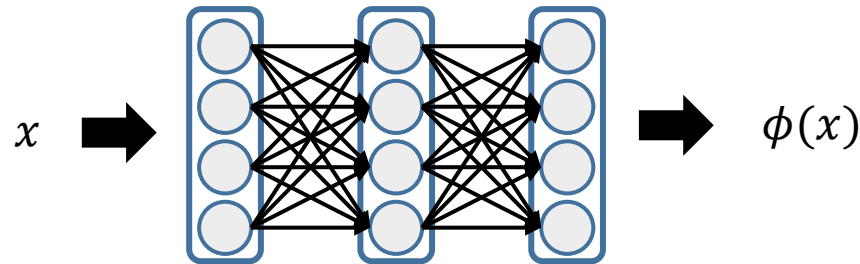
ノルム型バウンドは大きすぎる。



[Nagarajan and Kolter: Uniform convergence may be unable to explain generalization in deep learning. NeurIPS2019]

実験結果：訓練データサイズが増えるほどバウンドが大きくなる。

- これはリップシツツ定数 ($\prod_{\ell=1}^L \|W_{\ell}\|_2 = \prod_{\ell=1}^L R_{\ell,2}$) が大きくなることによる。
- この量はよりデータ依存な量に置換できる (Wei&Ma (2019), Arora et al. (2018)). → かなりタイトになる。



κ : 「データ依存な」リプシッツ連続性:

$$\|\phi(x) - \phi(x')\| \leq \kappa \|x - x'\|$$

for any training data point x and its neighborhood point x' .

$$\frac{1}{\sqrt{n}} \left(\sum_{\ell=1}^L \kappa^{2/3} R_{\ell, 2 \rightarrow 1}^{2/3} \right)^{3/2}$$

(informal)

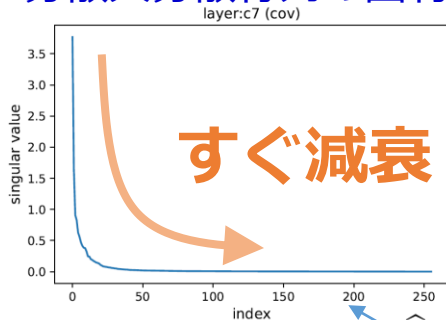
[Wei and Ma: Data-dependent Sample Complexity of Deep Neural Networks via Lipschitz Augmentation. NeurIPS 2019]

➤ このバウンドは先ほどのバウンドよりかなりタイトになる。

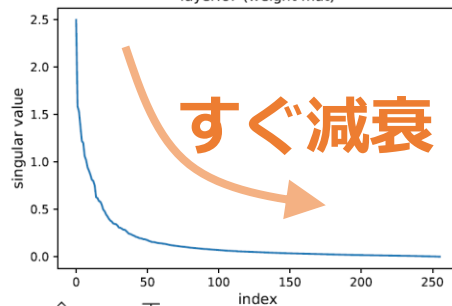
圧縮型バウンド

- 中間層の分散共分散行列の固有値分布
 - 中間層の重み行列の特異値分布
- が速く減衰するなら圧縮しやすい。

分散共分散行列の固有値



重み行列の特異値



$$\hat{\Sigma}_\ell = \frac{1}{n} \hat{\phi}_\ell(x_i) \hat{\phi}_\ell(x_i)^\top \quad \hat{W}(\ell)$$

分散共分散行列も重み行列も
特異値が速く減衰
→ **小さい統計的自由度**

(AIC, Mallows' Cp)

$$\Psi(\hat{f}) \leq \hat{\Psi}(\hat{f}) + O\left(\left(\frac{\sum_{\ell=1}^L m_\ell}{n} \log(n)\right)^{\frac{2\alpha}{1+2\alpha}} + \sqrt{L^{1+\delta} \left(\frac{\sum_{\ell=1}^L m_\ell}{n}\right)^{\frac{4/\beta}{4/\beta+2(1-1/2\alpha)}} \log(n)^3}\right)$$

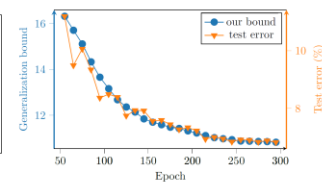
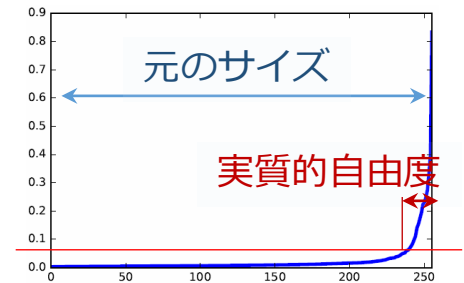
カーネル法の理論
(そもそもカーネルは無次元モデル)

(次ページに詳細)

[実験的観察] 実際に学習したネットワークは圧縮しやすい。

Layer	元サイズ Original	圧縮可能 サイズ Our bound
1	1,728	1,013
4	147,456	84,499
6	589,824	270,216
9	1,179,648	50,768
12	2,359,296	4,583
15	2,359,296	3,886

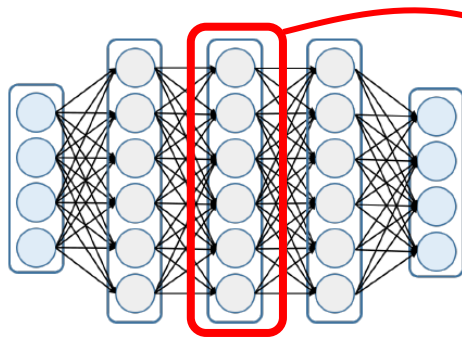
大 **小**



(a) Bound comparison

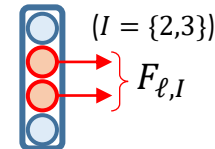
(b) Generalization bound

[Suzuki: Fast generalization error bound of deep learning from a kernel perspective. AISTATS2018]
 [Li, Sun, Liu, Suzuki and Huang: Understanding of Generalization in Deep Learning via Tensor Methods. AISTATS2020]
 [Suzuki, Abe, Nishimura: Compression based bound for non-compressed network: unified generalization error analysis of large compressible deep neural network, ICLR2020]
 [Suzuki et al.: Spectral pruning: Compressing deep neural networks via spectral analysis and its generalization error. IJCAI-PRICAI 2020]



$F_\ell(x) \in \mathbb{R}^{m_\ell}$: 中間層の特徴量 (m_ℓ 次元)

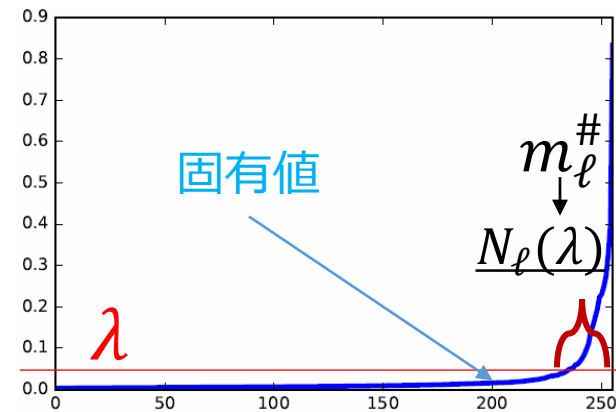
$F_{\ell,I}(x) \in \mathbb{R}^{|I|}$: 部分ベクトル



対応する共分散行列: $\Sigma_\ell := \mathbb{E}[F_\ell F_\ell^\top]$

その固有値: $\mu_j^{(\ell)}$ ($j = 1, \dots, m_\ell$)

統計的自由度: $N_\ell(\lambda) = \sum_{j=1}^{m_\ell} \frac{\mu_j^{(\ell)}}{\mu_j^{(\ell)} + \lambda}$



定理 (Bach, 2017)

任意の $\lambda > 0$ に対して, $m_\ell^\# = \Omega(N_\ell(\lambda) \log(N_\ell(\lambda)))$ 個の特徴量の集合 $\hat{I}_\ell \subset [m_\ell]$ ($|\hat{I}_\ell| = m_\ell^\#$) が存在して,

$$\sup_{w \in \mathbb{R}^{m_\ell} : \|w\| \leq 1} \inf_{\hat{w} \in \mathbb{R}^{m_\ell^\#}} \mathbb{E}[(w^\top F_\ell(X) - \hat{w}^\top F_{\ell, \hat{I}_\ell}(X))^2] \lesssim \lambda$$

(\hat{I}_ℓ に対応した部分ベクトル)

つまり, 中間層を大体 $m_\ell^\#$ 個の特徴量で代替できる.

分散共分散行列と重み行列の低ランク性¹⁰⁴

$$f(x) = (W^{(L)}\eta(\cdot)) \circ (W^{(L-1)}\eta(\cdot)) \circ \dots \circ (W^{(1)}x)$$

- 近似的に低ランクな重み行列と分散共分散行列:

$$\triangleright \sigma_j(\widehat{W}^{(\ell)}) \lesssim j^{-\alpha}$$

特異値が早く減衰

$$\triangleright \sigma_j(\widehat{\Sigma}^{(\ell)}) \lesssim j^{-\beta}$$

($\sigma_j(\cdot)$: j -th largest singular-value)

+ Other boundedness condition.

定理 (Suzuki, Abe, Nishimura, ICLR2020)

$$\Psi(\widehat{f}) \leq \widehat{\Psi}(\widehat{f})$$

$$+ O \left(\left(L \frac{\sum_{\ell=1}^L m_{\ell}}{n} \log(n) \right)^{\frac{2\alpha}{1+2\alpha}} + \sqrt{L^{1+\delta} \frac{\left(\sum_{\ell=1}^L m_{\ell} \right)^{\frac{4/\beta}{4/\beta+2(1-1/2\alpha)}}}{n} \log(n)^3} \right)$$

$$\text{where } \delta = \frac{\beta}{\frac{4\alpha}{(2\alpha-1)} + \beta} .$$

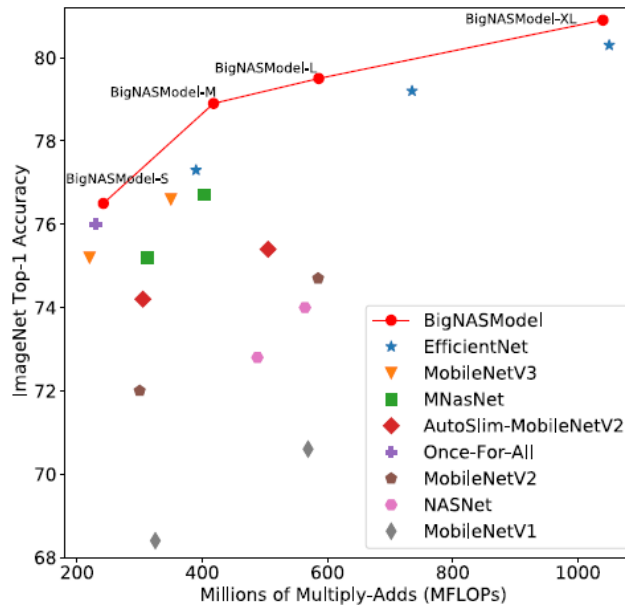
VC-次元によるバウンドより大きく改善:

$$\Psi(\widehat{f}) \leq \widehat{\Psi}(\widehat{f}) + \sqrt{L \frac{\sum_{\ell=1}^L \overset{\text{横幅の二乗}}{m_{\ell} m_{\ell+1}}}{n} \log(n)}$$

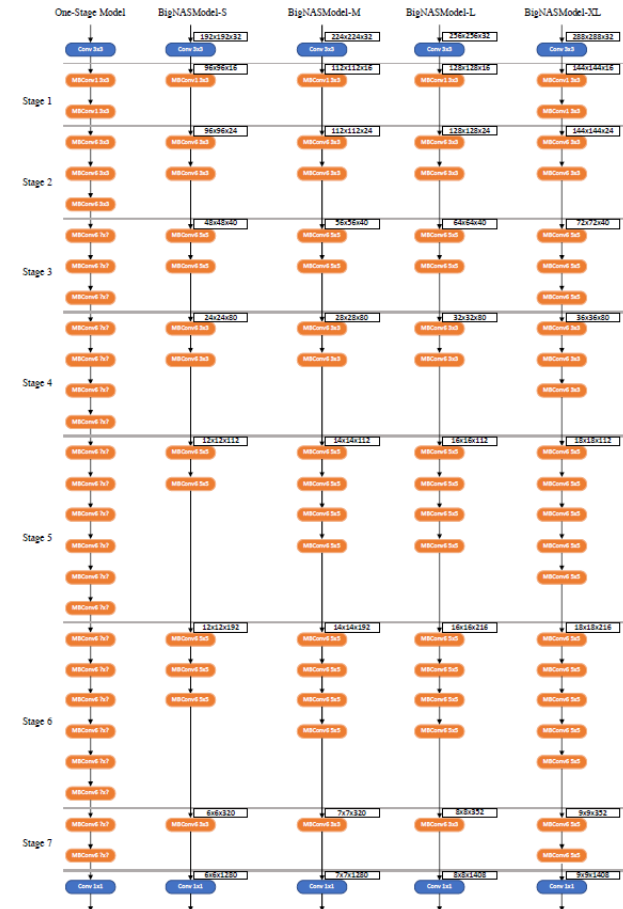
[Yu et al.: BigNAS: Scaling Up Neural Architecture Search with Big Single-Stage Models. ECCV2020]

(理論と関係あるNAS手法)

- 学習後のネットワークが圧縮できるように学習
- 大きなネットワークから小さなネットワークを生成できる
- EfficientNetを上回る効率性を実現



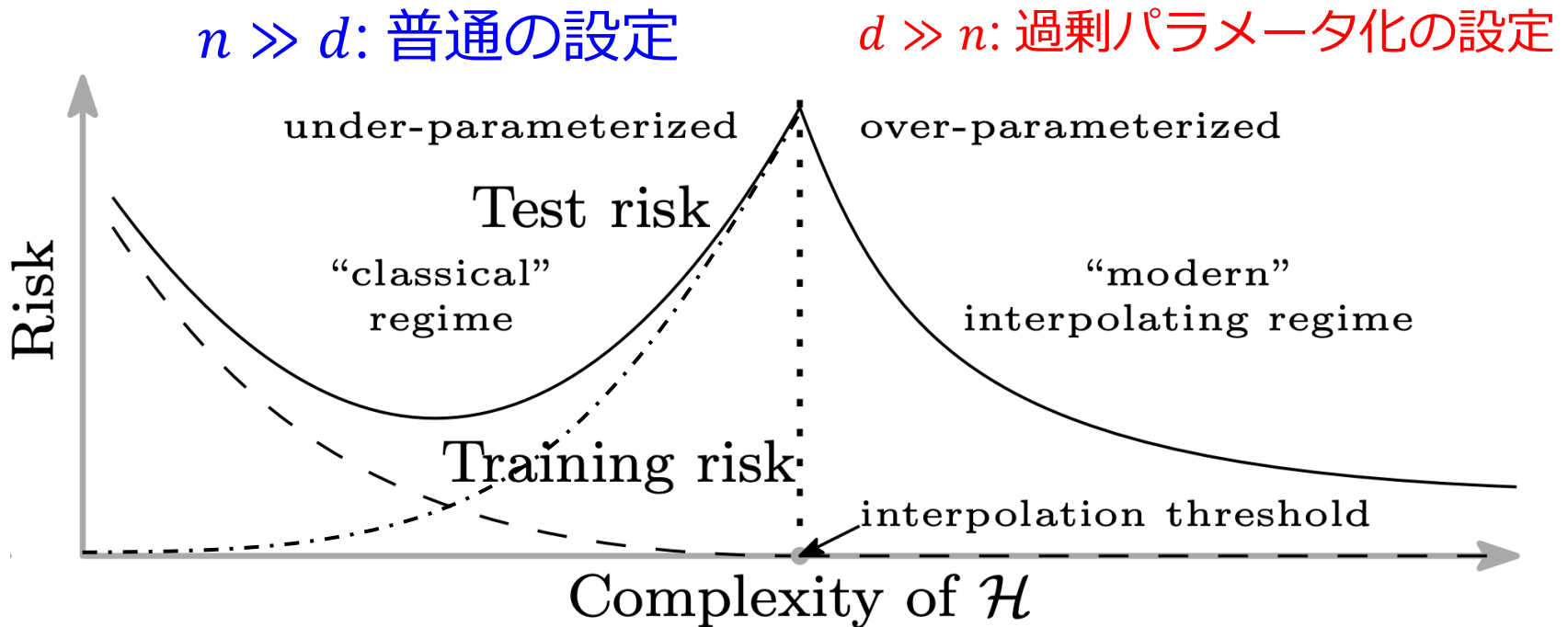
Group	Model Family	Params	FLOPs	Top-1
200M FLOPs	MobileNetV1 $0.5\times$	1.3M	150M	63.3
	MobileNetV2 $0.75\times$	2.6M	209M	69.8
	AutoSlim-MobileNetV2	4.1M	207M	73.0
	MobileNetV3 $1.0\times$	5.4M	219M	75.2
	MNASNet A_1	3.9M	315M	75.2
	Once-For-All	4.4M	230M	76.0
	Once-For-All $finetuned$	4.4M	230M	76.4
	BigNASModel-S	4.5M	242M	76.5
400M FLOPs	NASNet B	5.3M	488M	72.8
	MobileNetV2 $1.3\times$	5.3M	509M	74.4
	MobileNetV3 $1.25\times$	8.1M	350M	76.6
	MNASNet A_3	5.2M	403M	76.7
	EfficientNet B_0	5.3M	390M	77.3
	BigNASModel-M	5.5M	418M	78.9
600M FLOPs	MobileNetV1 $1.0\times$	4.2M	569M	70.9
	NASNet A	5.3M	564M	64.0
	DARTS	4.9M	595M	73.1
	EfficientNet B_1	7.8M	734M	79.2
	BigNASModel-L	6.4M	586M	79.5
1000M FLOPs	EfficientNet B_2	9.2M	1050M	80.3
	BigNASModel-XL	9.5M	1040M	80.9



圧縮できるように学習するとスクラッチ学習より性能が向上することもある。

Overparametrizeされた ネットワークの統計学

Double-descent (二重降下)



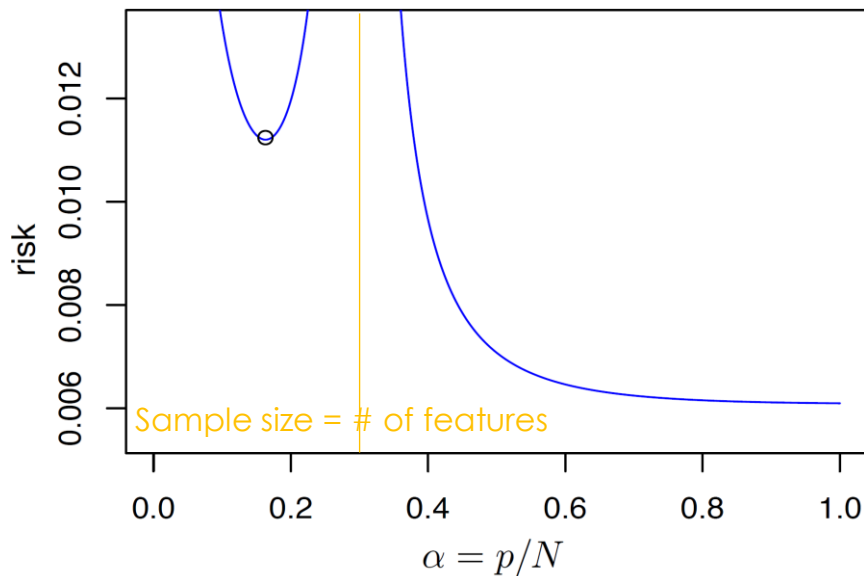
[Belkin et al.: Reconciling modern machine learning practice and the bias-variance trade-off. 2018]

- モデルがある複雑度 (サンプルサイズ) を超えた後, 第二の降下が始まる.
- モデルサイズがデータより多いと推定量のバリエーションがむしろ減る.

※設定によるので注意が必要.

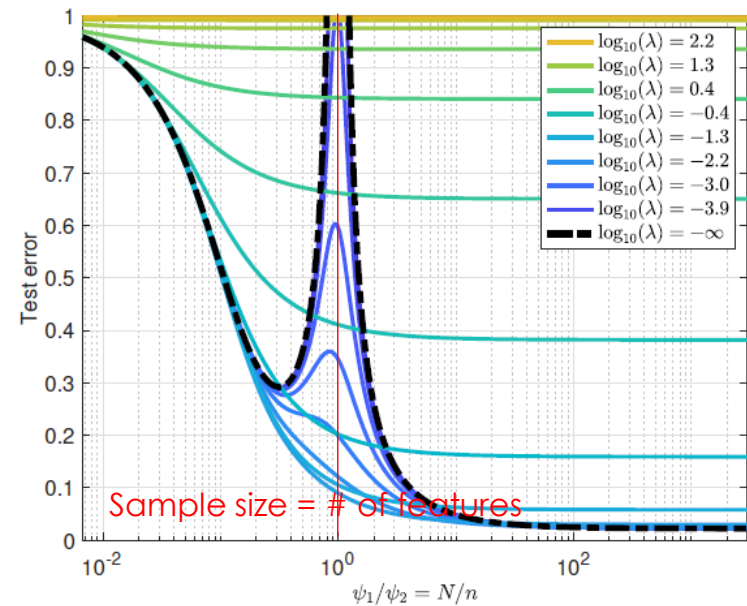
Principal component regression

$$\kappa = 2$$



(いくつかの主成分を用いたか)
Populationの分散共分散を知っているとして、
その主成分を利用

2-layer neural network



(いくつかのニューロンを用いたか)

[Xu and Hsu: On the number of variables to use in principal component regression. NeurIPS2019.]

[Mei and Montanari. "The generalization error of random features regression: Precise asymptotics and double descent curve." arXiv preprint arXiv:1908.05355 (2019)]

典型的なアプローチ (抜粋)

• ランダム行列理論

- $d/n \rightarrow \gamma > 0$ という漸近的設定で, 厳密なリスクの導出
- Marchenko–Pastur則, Stieltjes変換

$$\tilde{m}(z) := \frac{1}{d} \text{Tr} \left[(X^\top X/n - zI)^{-1} \right] \xrightarrow{a.s.} m(z)$$

$$m(z) = \int_0^\infty \frac{1}{\tau(1 - \gamma(1 + zm(z))) - z} dH(\tau) \quad (H(\tau): \text{Spectral measure of } \Sigma_x)$$

- Dobriban&Wager: High-dimensional asymptotics of prediction: Ridge regression and classification. The Annals of Statistics, 46(1):247–279, 2018.
- Hastie et al.: Surprises in High-Dimensional Ridgeless Least Squares Interpolation, arXiv:1903.08560.
- Song&Montanari. The generalization error of random features regression: Precise asymptotics and double descent curve. Communications on Pure and Applied Mathematics. arXiv:1908.05355 (2019).

• 集中不等式による評価

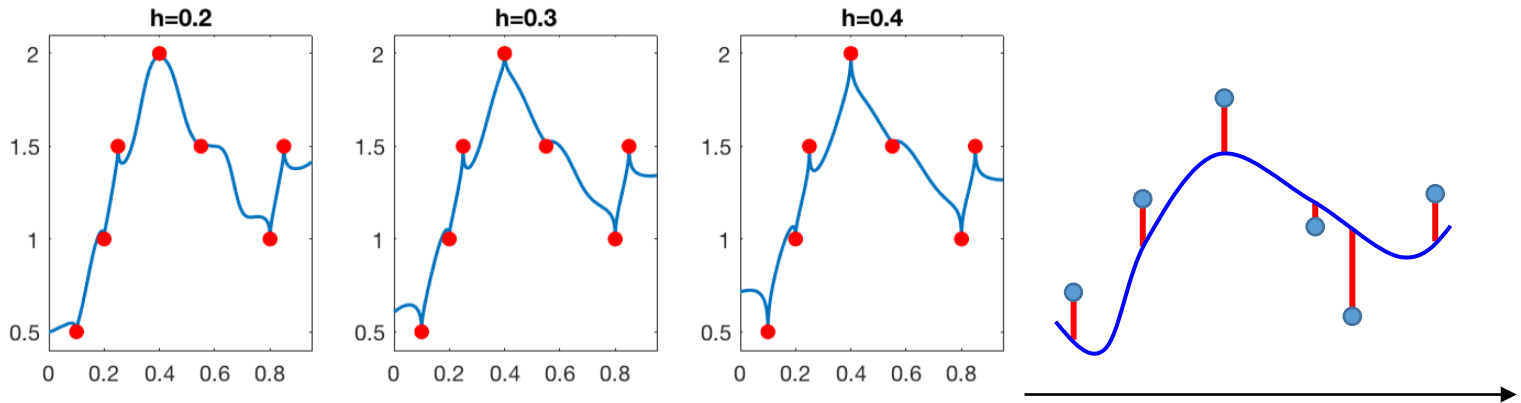
- 有限サンプルサイズにおける予測誤差の上界評価 ($n < \infty$)
- 収束レートが評価できる. ($n \rightarrow \infty$ を取る前の振る舞いを評価)

- Belkin, Rakhlin&Tsybakov: Does data interpolation contradict statistical optimality? AISTATS2019.
- Bartlett, Long, Lugosi&Tsigler: Benign Overfitting in Linear Regression. PNAS, 117(48):30063-30070, 2020.
- Liang&Rakhlin: Just interpolate: Kernel “Ridgeless” regression can generalize. The Annals of Statistics, 48(3):1329–1347, 2020.

• CGMT (Convex Gaussian min-max Theorem)

- Thrampoulidis, Oymak & Hassibi: Regularized linear regression: A precise analysis of the estimation error. COLT2015.
- Thrampoulidis, Abbasi & Hassibi: Precise error analysis of regularized m-estimators in high dimensions. IEEE Transactions on Information Theory, vol. 64, no. 8, pp. 5592–5628, 2018.

[Belkin, Rakhlin&Tsybakov: Does data interpolation contradict statistical optimality? AISTATS2019]



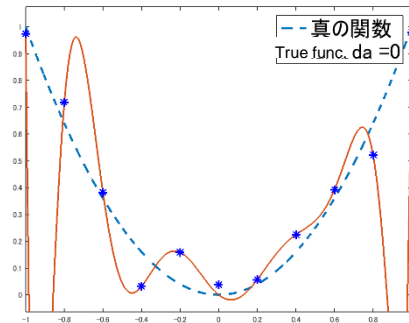
- 過剰パラメータ化されたモデルは“スパイク”成分を持つ.
- スパイク成分がノイズを説明.
- 大まかな関数形はモデルの主成分が説明.

d-次元空間

訓練データ x_i ($i = 1, \dots, n$)

x (テスト時の入力)

- 「スパイク成分」とはほぼ直交.
- 直交するには高次元性が必要.
- 高次元空間で2つのランダムベクトルはほぼ直交.

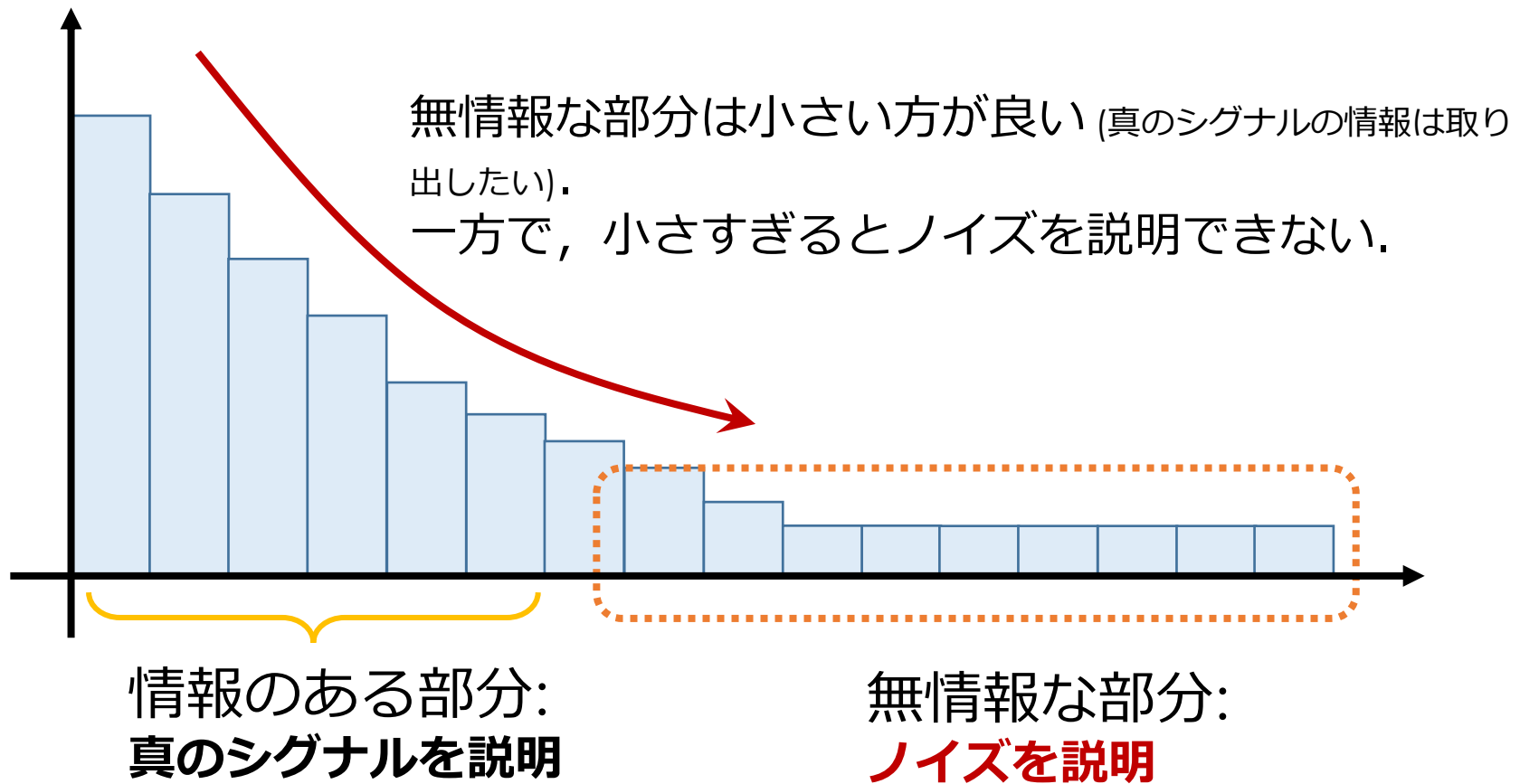


悪性過学習
(モデルの説明力が中途半端なので、無理してノイズも説明)

線形モデル: 分散共分散行列での直感

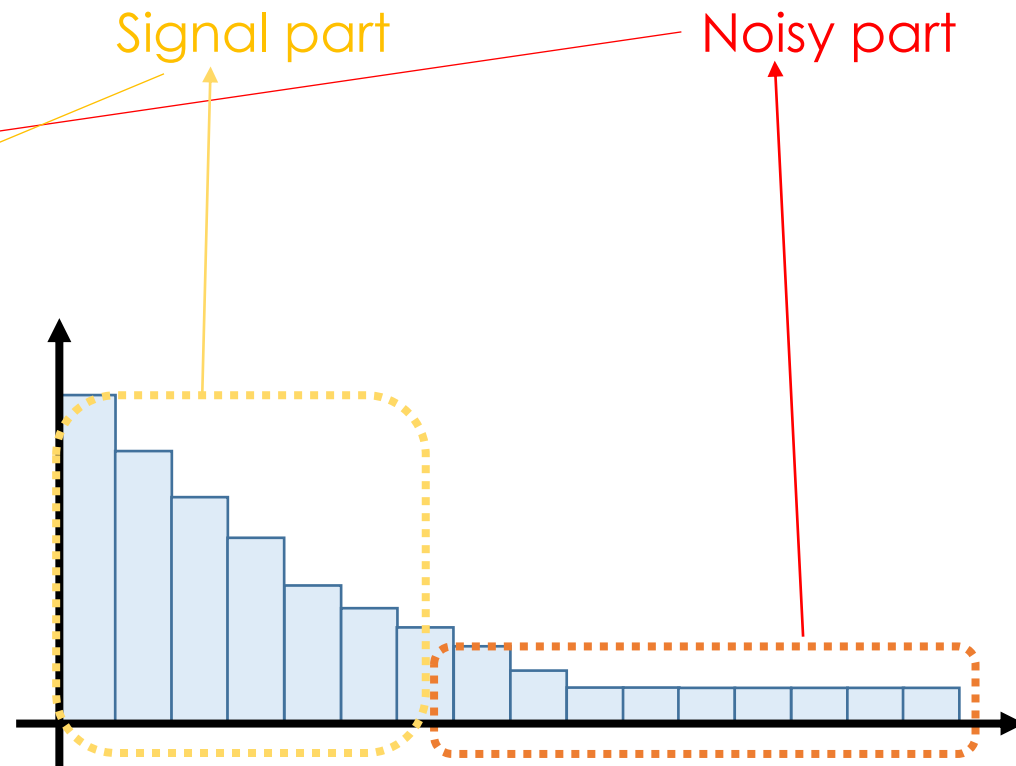
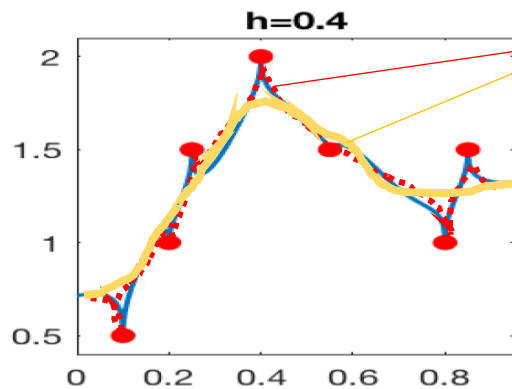
$$\Sigma_X = \mathbb{E}[xx^\top] \in \mathbb{R}^{d \times d}$$

Eigenvalues of Σ_X (spectrum)



※サンプルサイズに応じてどこから無情報になるかは変わる。あくまで「そのサンプルサイズの解像度では無情報に見える」ということ。

$$\hat{f}(x) = \hat{\beta}^\top x = \underbrace{x^\top (X^\top X)^+ X^\top X \beta^*}_{\text{Signal part}} + \underbrace{x^\top (X^\top X)^+ X^\top \epsilon}_{\text{Noisy part}}$$



問題設定: 補間推定量

$$y_i = x_i^\top \beta + \epsilon_i \quad (d > n)$$

overparameterization

$$\mathbb{E}[\epsilon_i] = 0, \text{Var}(\epsilon_i) = \sigma^2, \text{Cov}(x_i) = \Sigma$$

$$(X = [x_1, \dots, x_n]^\top, Y = [y_1, \dots, y_n]^\top)$$

- **最小ノルム補間推定量 (Minimum-norm interpolator):**

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^d} \|\beta\|$$

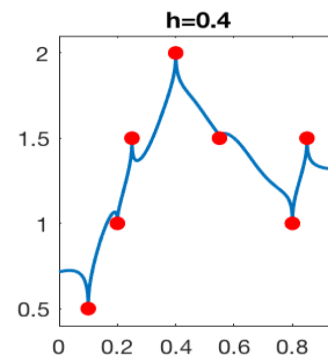
過学習しているものの
中で一番単純な推定量

$$\text{s.t. } y_i = x_i^\top \beta \quad (\forall i = 1, \dots, n) \quad : \text{interpolation}$$

$$\Rightarrow \hat{\beta} = (X^\top X)^+ X^\top Y$$

- **予測誤差 (predictive risk, excess risk):**

$$\mathbb{E}_\epsilon [\|\hat{\beta} - \beta\|_\Sigma^2 | X]$$



$$\hat{\beta}_\lambda := \arg \min_{\beta \in \mathbb{R}^d} \|X\beta - Y\|^2 + \lambda \|\beta\|^2 \quad \Rightarrow \quad \hat{\beta} = \lim_{\lambda \searrow 0} \hat{\beta}_\lambda$$

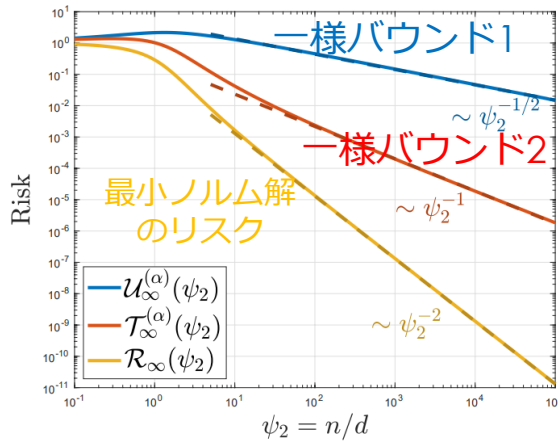
Ridge regression

一様バウンドではダメ

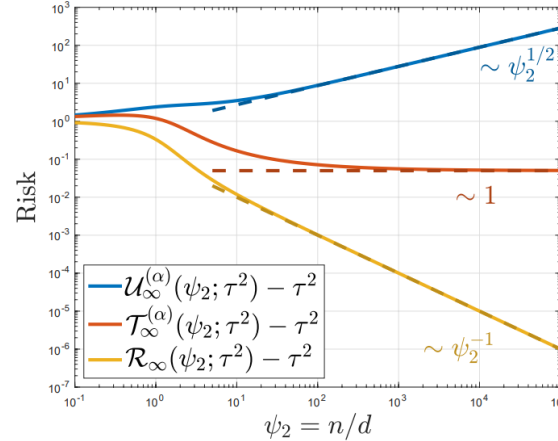
[Yang, Bai, Mei: Exact Gap between Generalization Error and Uniform Convergence in Random Feature Models. ICML2021.]

[Bartlett&Long, JMLR, 22(204):1-15, 2021]も参照

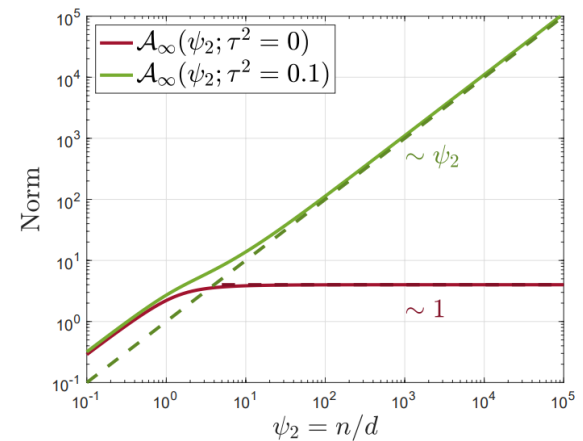
$\alpha = 1.5$



(a) Noiseless response ($\tau^2 = 0$)



(b) Noisy response ($\tau^2 = 0.1$)



(c) Minimum norm $\mathcal{A}_\infty(\psi_2)$

以下の量の $n, N, d \rightarrow \infty$ の極限を取る。ただし, $N/d \rightarrow \infty, n/d \rightarrow \psi_2$.

$U(A, N, n, d) \equiv \sup_{(N/d) \|\mathbf{a}\|_2^2 \leq A} (R(\mathbf{a}) - \widehat{R}_n(\mathbf{a}))$, 一様バウンド1 (Aに最小ノルム解のノルム $\times \alpha$ を代入)
経験誤差=0とは限らない解の中で一様バウンドを取る。

$T(A, N, n, d) \equiv \sup_{(N/d) \|\mathbf{a}\|_2^2 \leq A, \widehat{R}_n(\mathbf{a})=0} R(\mathbf{a})$. 一様バウンド2 (Aに最小ノルム解のノルム $\times \alpha$ を代入)
経験誤差=0の中で一様バウンドを取る。 $\alpha > 1$ の時は最小ノルム解以外も含まれる。

ランダム特徴モデル

$$\mathcal{F}_{\text{RF}}(\Theta) \equiv \left\{ f(\mathbf{x}) = \sum_{j=1}^N a_j \sigma(\langle \mathbf{x}, \boldsymbol{\theta}_j \rangle / \sqrt{d}) : \mathbf{a} \in \mathbb{R}^N \right\} \quad N: \text{特徴数}, d: \text{入力の次元}, n: \text{サンプルサイズ}$$

ランダム行列のアプローチ

- Hastie et al.: Surprises in High-Dimensional Ridgeless Least Squares Interpolation, Ann. Statist. 50(2): 949--986.

$$y_i = x_i^\top \beta^* + \epsilon_i \quad \hat{\beta} = (X^\top X)^+ X^\top Y : \text{最小ノルム補間推定量}$$

proportional limit $n, d \rightarrow \infty, d/n \rightarrow \gamma$ における漸近リスク

$$R(\gamma) := \lim_{n \rightarrow \infty} \mathbb{E}_\epsilon[\|\hat{\beta} - \beta^*\|_\Sigma^2 | X] = ?$$

Theorem

$$r^2 := \|\beta^*\|^2, \sigma^2 := \mathbb{E}[\epsilon_i^2]$$

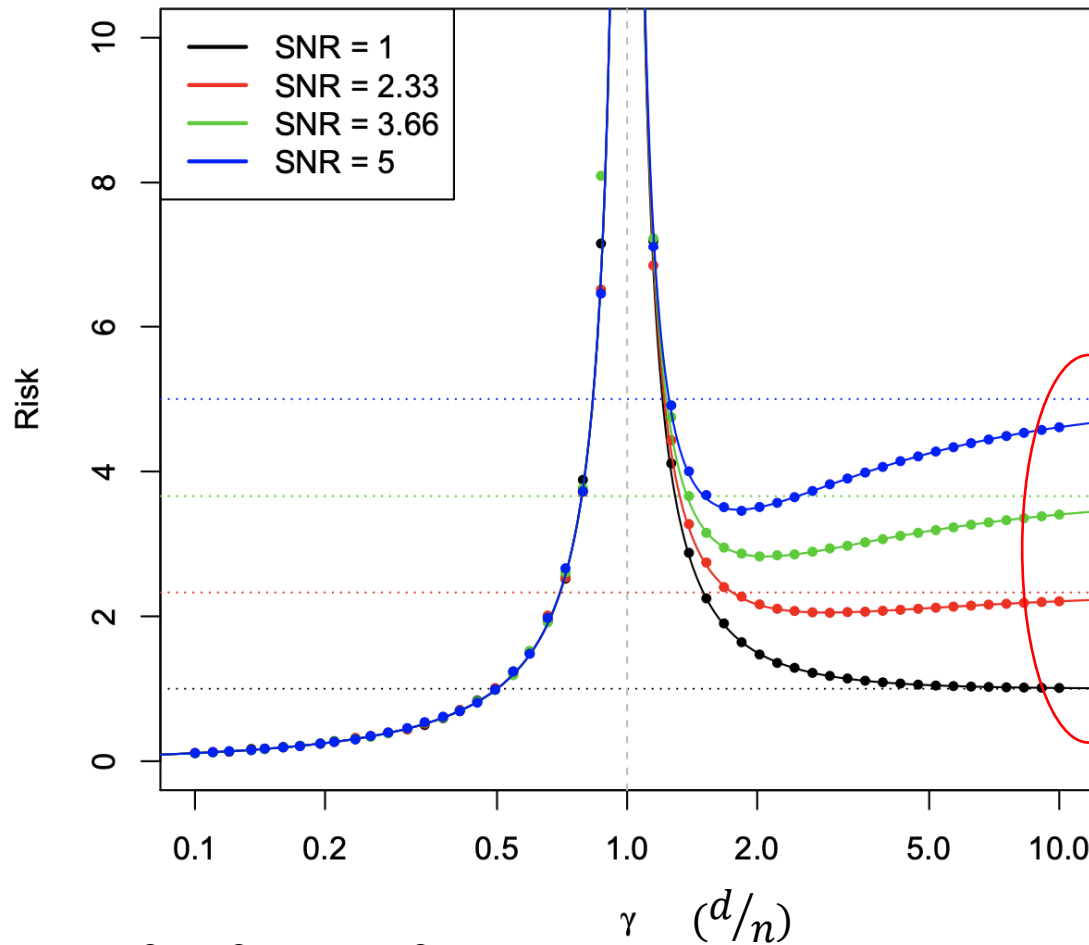
$n, d \rightarrow \infty$ かつ $d/n \rightarrow \gamma \in (0, \infty)$ (proportional limit)
の極限で期待予測誤差は以下の値に概収束する:

$$R(\gamma) = \begin{cases} \sigma^2 \frac{\gamma}{1-\gamma} & (\gamma < 1) \\ r^2 \left(1 - \frac{1}{\gamma}\right) + \sigma^2 \frac{1}{\gamma-1} & (\gamma > 1) \end{cases}$$

Bias
Variance

$$\mathbb{E}_{x, D_n}[(x^\top \beta^* - x^\top \hat{\beta})^2] = \underbrace{\mathbb{E}_x[(x^\top \beta^* - x^\top \mathbb{E}[\hat{\beta}])^2]}_{B: \text{Bias}} + \underbrace{\text{tr}[\text{Cov}(\hat{\beta})\Sigma_x]}_{V: \text{Variance}}$$

Numerical Experiment



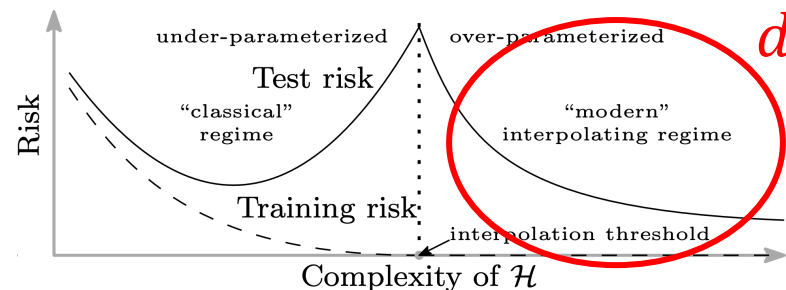
Bias = $\|\beta^*\|^2$, Variance = 0.
※ここでは $\hat{\beta} = 0$ である. これをもつて「推定している」とは言いにくいので注意.

$$r^2 := \|\beta^*\|^2, \sigma^2 := \mathbb{E}[\epsilon_i^2]$$

$$R(\gamma) = \begin{cases} \sigma^2 \frac{\gamma}{1-\gamma} & (\gamma < 1) \\ \underbrace{r^2 \left(1 - \frac{1}{\gamma}\right)}_{\text{Bias}} + \underbrace{\sigma^2 \frac{1}{\gamma-1}}_{\text{Variance}} & (\gamma > 1) \end{cases}$$

前処理付き勾配法

Amari, Ba, Grosse, Li, Nitanda, Suzuki, Wu, Xu: When Does Preconditioning Help or Hurt Generalization? ICLR2021.



$d \gg n$: overparameterized regime

$$y_i = x_i^\top \beta^* + \epsilon_i$$

$$\min_{\beta \in \mathbb{R}^d} \underbrace{\|\beta\|_{P^{-1}}^2}_{=\beta^\top P^{-1} \beta} \quad \text{s.t.} \quad y_i = x_i^\top \beta \quad (\text{interpolation})$$

$$(\forall i \in [n])$$



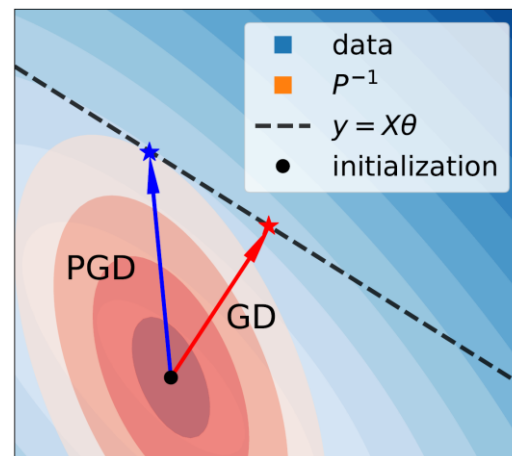
Q: P によって予測性能がどのように影響受けるか？

$$\frac{d\beta(t)}{dt} = -P X^\top (Y - X\beta(t)) / n$$

Preconditioned Gradient Descent

$P = I$: Gradient descent (GD)

$P = \Sigma_x^{-1}$: Natural Gradient descent (NGD)
(真の分布で期待値取ったFisher情報行列)



最適な前処理行列

Bias-variance分解

$$\mathbb{E}_{x, D_n} [(x^\top \beta^* - x^\top \hat{\beta})^2] = \underbrace{\mathbb{E}_x [(x^\top \beta^* - x^\top \mathbb{E}[\hat{\beta}])^2]}_{B: \text{Bias}} + \underbrace{\text{tr}[\text{Cov}(\hat{\beta}) \Sigma_x]}_{V: \text{Variance}}$$

定理 (informal)

[Amari, Ba, Grosse, Li, Nitanda, Suzuki, Wu, Xu: When Does Preconditioning Help or Hurt Generalization? ICLR2021]

$d/n \rightarrow \gamma > 1$ ($n \rightarrow \infty$)の極限における予測誤差の漸近値を厳密に導出

1. バリアンス:

$P = \Sigma_x^{-1}$ (Fisher情報行列) がバリエンスを最小化.

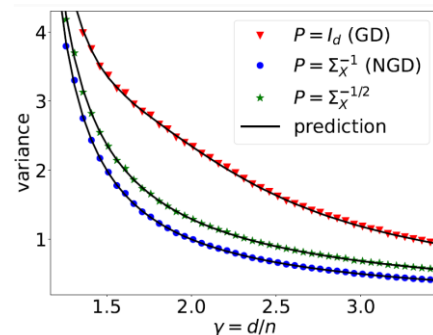
NGDがバリエンスの意味で最適.

2. バイアス:

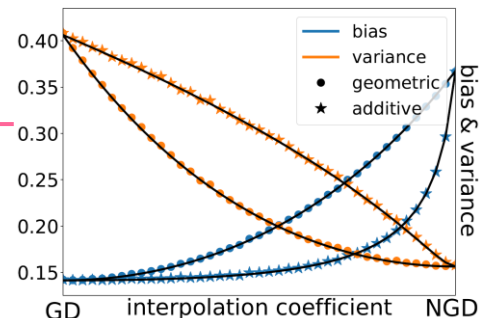
No free-lunch: 事前に最適なPは決定できない:

- 真が等方的分布 $\Sigma_{\beta^*} = I$ に従っていれば**GD**が良い.
- 真が非等方的分布 $\Sigma_{\beta^*} = \Sigma_x^{-1}$ に従っていれば**NGD**が良い.

(バイズの設定: β^* の実現値に関する予測誤差の期待値を比較 $\mathbb{E}[\beta^* \beta^{*\top}] = \Sigma_{\beta^*}$)



(dが増えるほど
バリエンスが減る)



➡ GDとNGDの間が良い

$$\begin{cases} \text{Additive: } P = (\alpha \Sigma_x + (1 - \alpha) I_d)^{-1} \\ \text{Geometric: } P = \Sigma_x^{-\alpha} \end{cases} \quad \alpha \in [0, 1]$$

より詳細な結果

(A2) $\Sigma_{XP} := P^{1/2}\Sigma P^{1/2}$ のスペクトル分布が H_{XP} に弱収束すると仮定.

• $m(z)$ を自己整合条件を満たす関数とする: $\frac{1}{m(z)} = -z + \gamma \int \frac{\tau}{1 + \tau m(z)} dH_{XP}(\tau)$

→ $\frac{1}{n}XPX^\top$ のスペクトルの漸近分布を表現.

1. バリانس:

$$V \xrightarrow{P} \sigma^2 \left(\lim_{\lambda \rightarrow +0} m'(-\lambda)m^{-2}(-\lambda) - 1 \right)$$

$V \geq \sigma^2(\gamma - 1)^{-1}$ and equality holds by $P = \Sigma^{-1}$

(A3) P と Σ が同じ固有ベクトル U を共有.

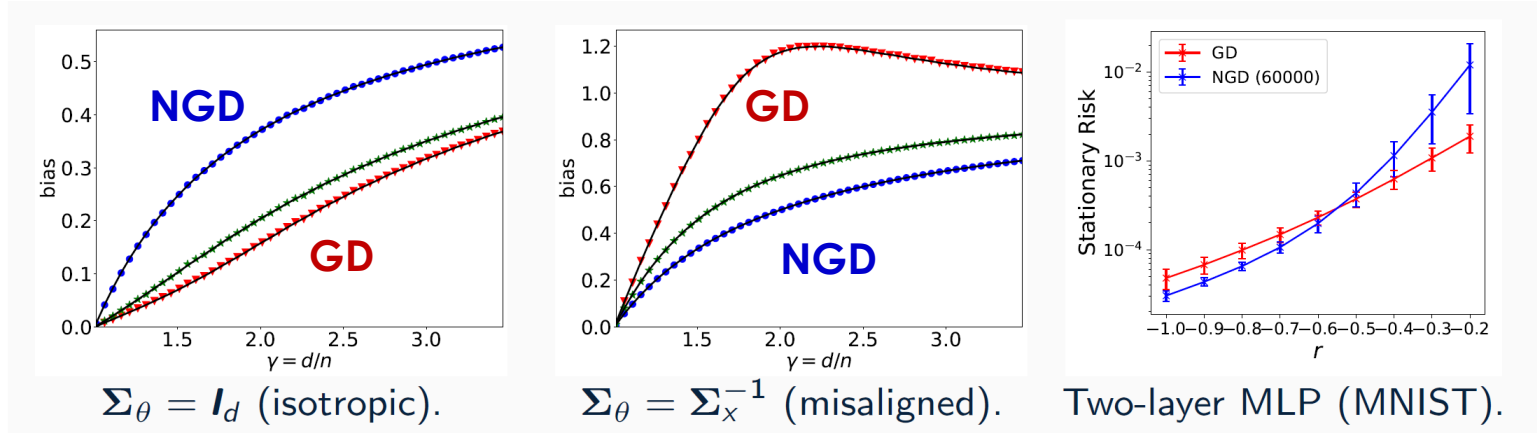
2. バイアス:

$$\mathbb{E}_{\beta^*}[B] \xrightarrow{P} \lim_{\lambda \rightarrow +0} m'(-\lambda)m^{-2}(-\lambda)\mathbb{E}[v_x v_\theta (1 + v_{xp}m(-\lambda))^{-2}]$$

ただし, (e_x, e_θ, e_{xp}) は $\Sigma, \Sigma_{XP}, \text{diag}(U^\top \Sigma_{\beta^*} U)$ の固有値で (v_x, v_θ, v_{xp}) に弱収束するものとする.

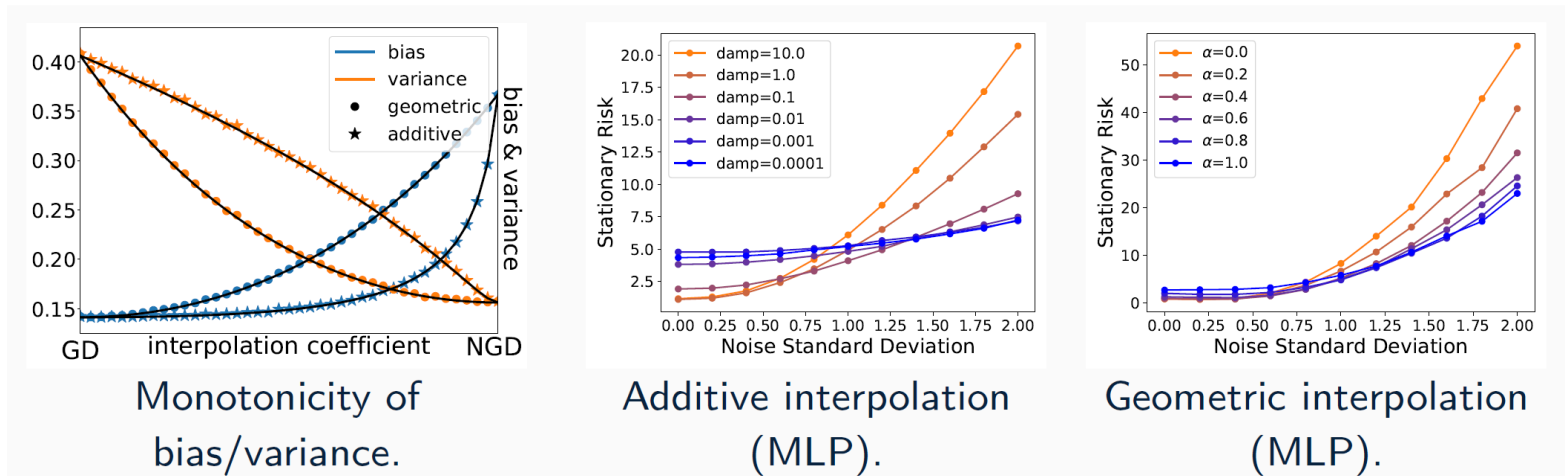
Experiments

Bias:



Bias/Variance trade-off:

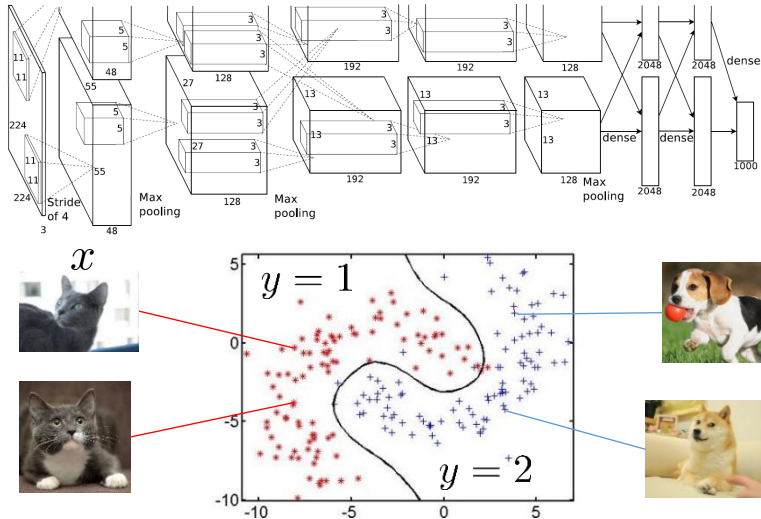
$$\begin{cases} \text{Additive: } P = (\alpha \Sigma_x + (1 - \alpha) I_d)^{-1} \\ \text{Geometric: } P = \Sigma_x^{-\alpha} \end{cases}$$



第3部

深層学習の最適化

-非凸性の影響-



深層ニューラルネットワークをデータにフィットさせるとは？

$$L(W) = \frac{1}{n} \sum_{i=1}^n \ell_i(W)$$

W : パラメータ

i 番目のデータで正解していれば小さく、間違っていれば大きく

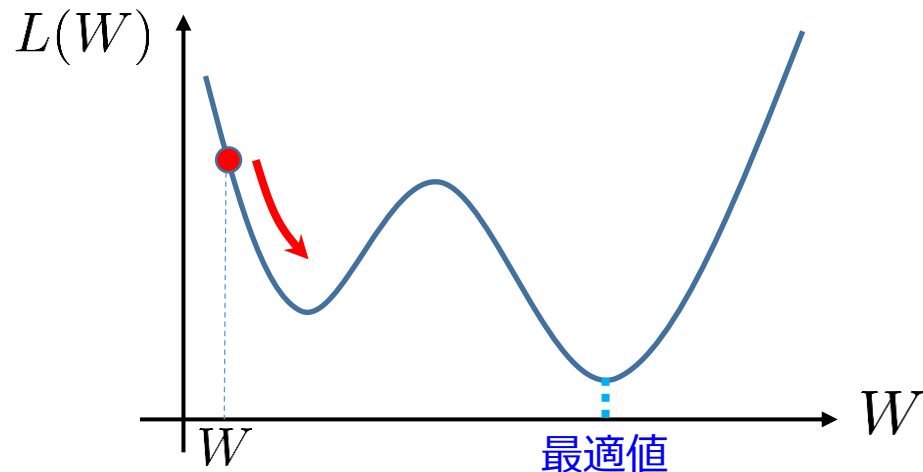
損失関数：データへの当てはまり度合い

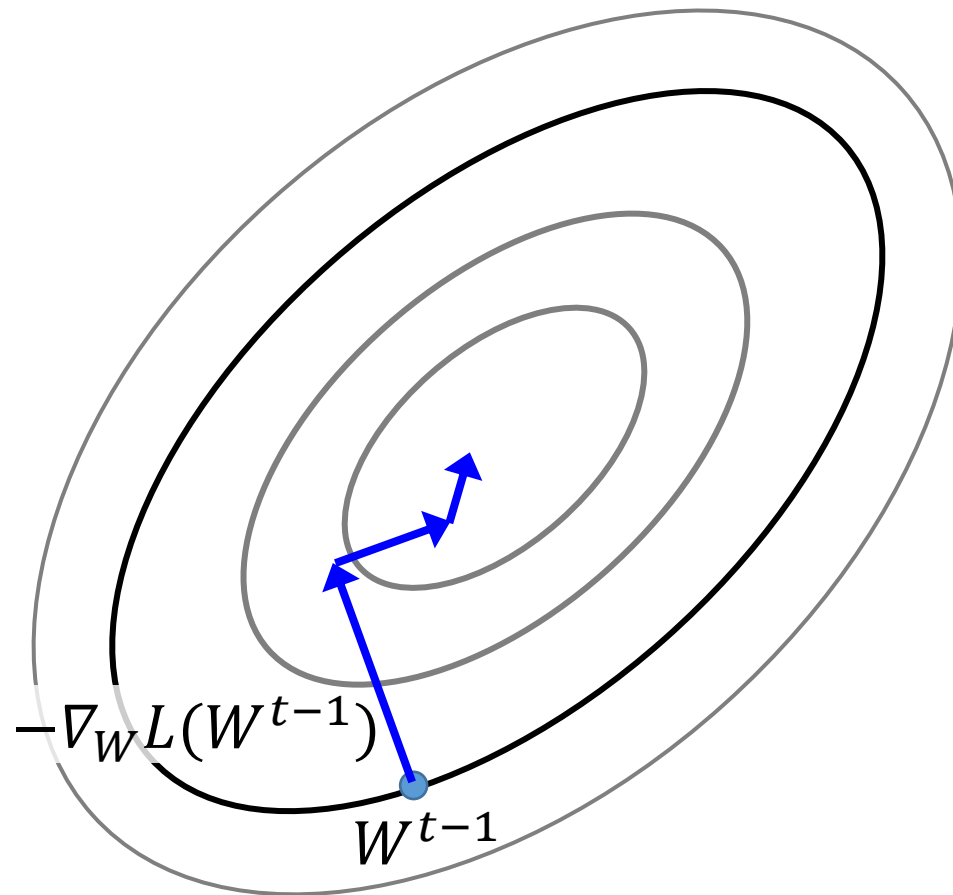
損失関数最小化

$$\min_W L(W)$$

(W は数十億次元)

通常、**確率的勾配降下法**で最適化



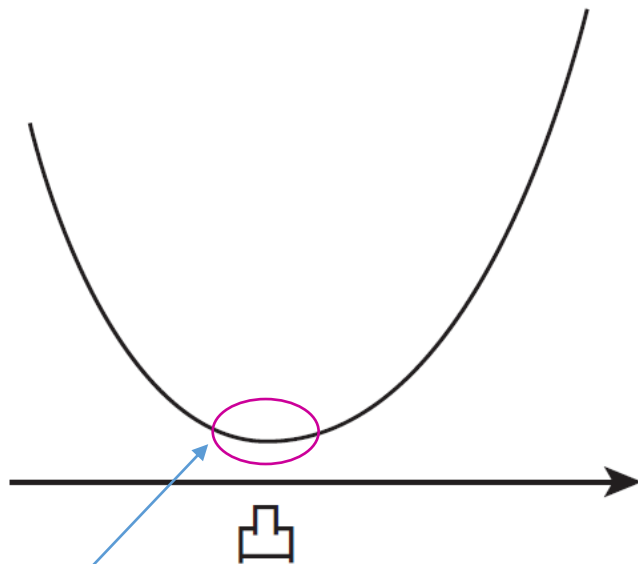


$$W^t = W^{t-1} - \eta \nabla_W L(W^{t-1})$$

問題点

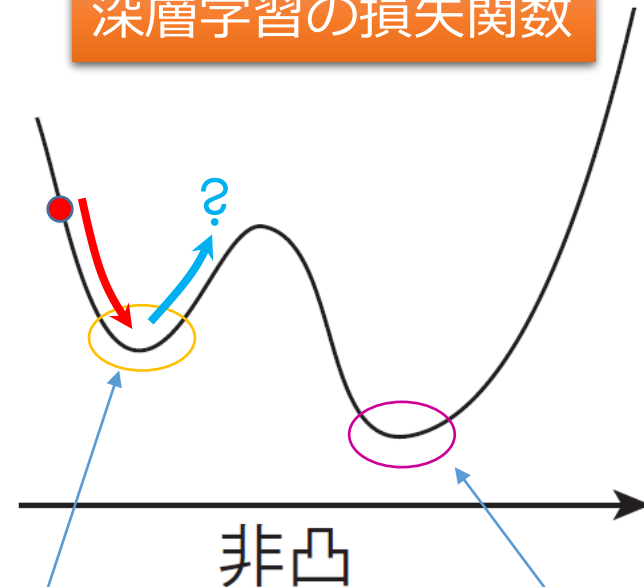
目的関数が非凸関数

凸関数 $\theta f(x) + (1 - \theta)f(y) \geq f(\theta x + (1 - \theta)y) \quad (\forall x, y \in \mathbb{R}^p, \theta \in [0, 1])$



局所最適解 = 大域的最適解

深層学習の損失関数



局所最適解

大域的最適解

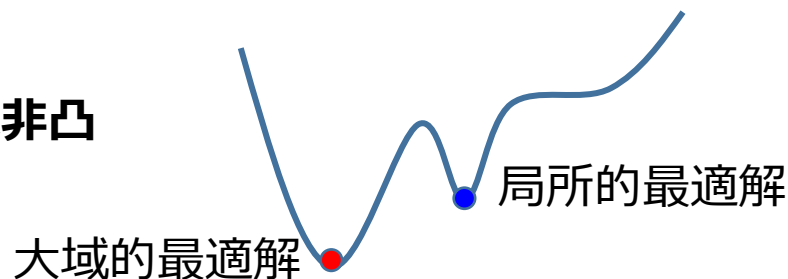
局所最適解や鞍点にはまる可能性あり

“狭い”ネットワークの学習はNP-完全:

- Judd (1988), Neural Network Design and the Complexity of Learning.
- Blum&Rivest (1992), Training a 3-node neural network is NP-complete.

局所最適性

深層学習の目的関数は非凸



- 線形深層NNの局所的最適解は全て大域的最適解：
Kawaguchi, 2016; Lu&Kawaguchi, 2017.

※ただし対象は線形NNのみ.

→ 臨界点が大域的最適解であること条件も出されている
(Yun, Sra&Jadbabaie, 2018)

- 低ランク行列補完の局所的最適解は全て大域的最適解：
Ge, Lee&Ma, 2016; Bhojanapalli, Neyshabur&Srebro, 2016.

$$\min_{U \in \mathbb{R}^{M \times k}} \sum_{(i,j) \in E} (Y_{i,j} - (UU^T)_{i,j})^2$$

Loss landscape

- 横幅の広いNNの訓練誤差には孤立した局所最適解がない。(局所最適解は大域的最適解とつながっている) ※とはいえ、勾配法で大域的最適解に到達可能かは別問題.

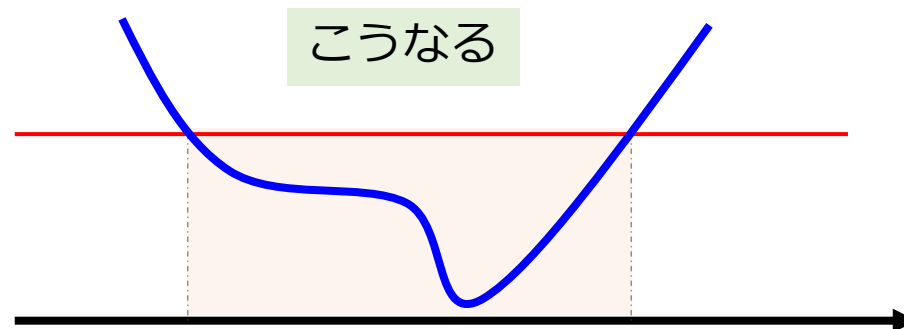
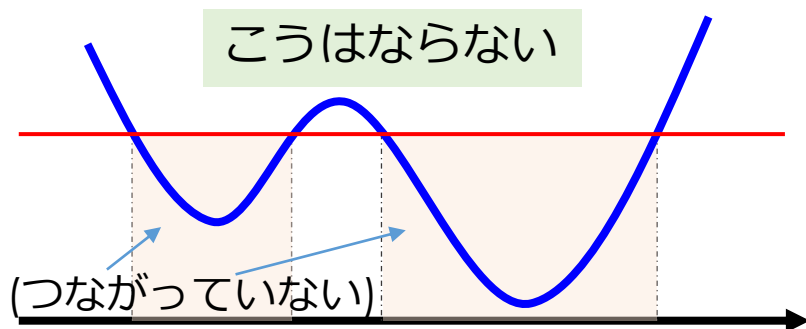
定理

n 個の訓練データ $(x_i, y_i)_{i=1}^n$ が与えられているとする. 損失関数 ℓ は凸関数とする.

任意の連続な活性化関数について, 横幅がデータサイズより広い

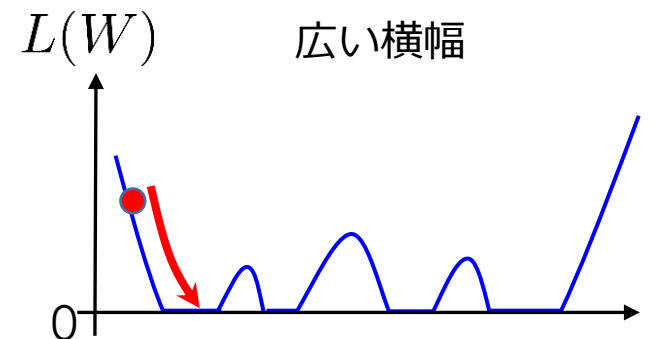
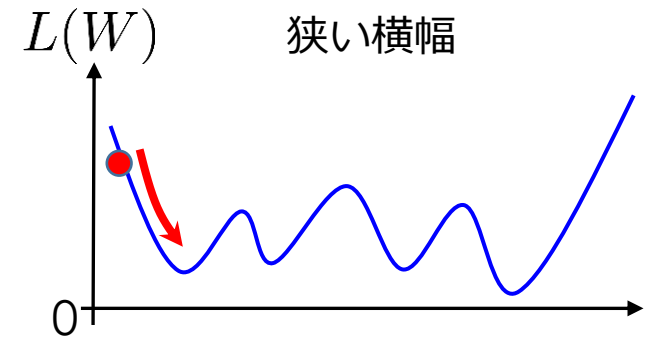
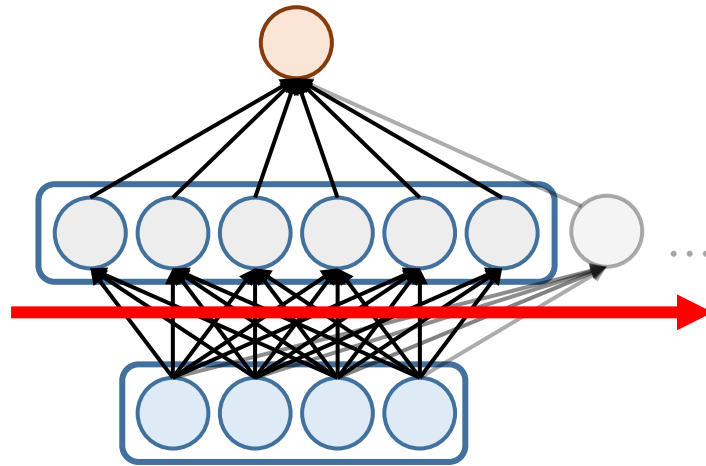
($M \geq n$) 二層NN $f_{(a,W)}(x) = \sum_{m=1}^M a_m \eta(w_m^T x)$ に対する訓練誤差 $\hat{L}(a, W) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{(a,W)}(x_i))$ の任意のレベルセットの弧状連結成分は大域的最適解を含む. 言い換えると, **任意の局所最適解は大域的最適解である.**

[Venturi, Bandeira, Bruna: Spurious Valleys in One-hidden-layer Neural Network Optimization Landscape JMLR, 20:1-34, 2019.]



オーバーパラメトライゼーション

横幅が広いと局所最適解が大域的最適解になる。



自由度が高いため、目的関数を減少させる方向が見つけやすい。

• 二種類の解析手法

- Neural Tangent Kernel (NTK)
- Mean-field analysis (平均場解析)

$$f_W(x) = \sum_{j=1}^M a_j \eta(w_j^\top x)$$

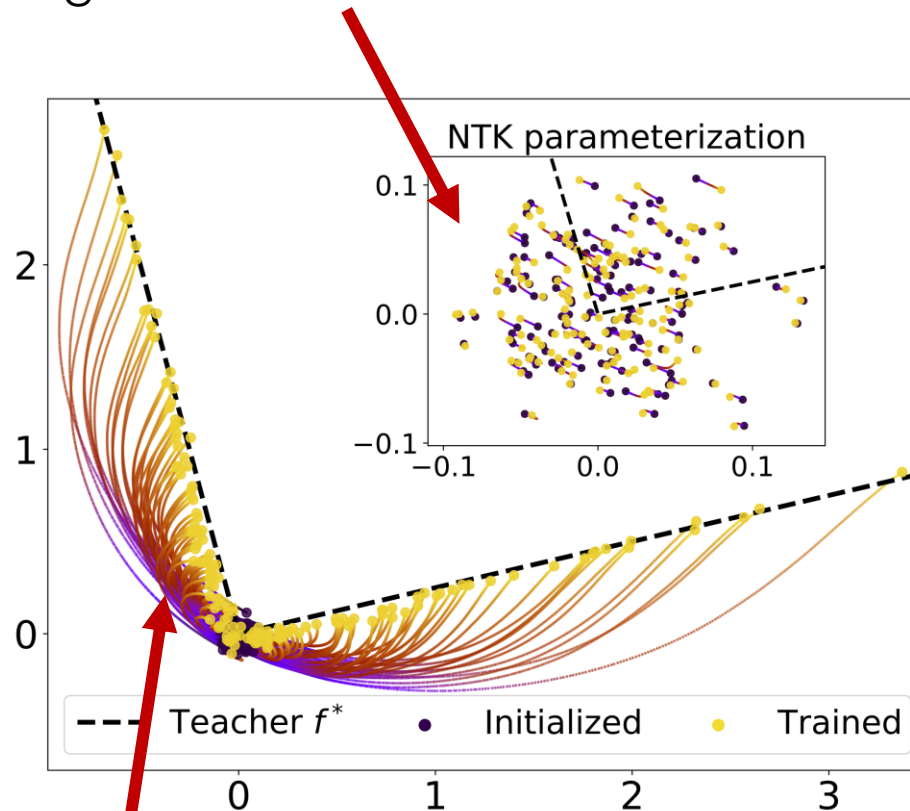
- Neural Tangent Kernelのregime (lazy learning)
 - $a_j = \mathbf{O}(1/\sqrt{M})$ [Jacot+ 2018][Du+ 2019][Arora+ 2019]
(Xavier initialization/He initialization)
- 平均場解析の設定 (mean field)
 - $a_j = \mathbf{O}(1/M)$ [Nitanda & Suzuki (2017), Chizat & Bach (2018), Mei, Montanari, & Nguyen (2018)]

初期化のスケーリングによって、初期値と比べて学習によって動く大きさの割合が変わる。
→ 学習のダイナミクス、汎化性能に影響
(解析の難しさも違う)

変数のスケールと特徴学習

$$f(x) = \frac{1}{\sqrt{M}} \sum_{j=1}^M r_j \sigma(w_j^\top x)$$

NTK: Large scale initialization \rightarrow features are (almost) freezed.



二層NNにおける一層目の
パラメータの軌跡

$$f(x) = \sum_{j=1}^M a_j \sigma(w_j^\top x)$$

[Ba et al., 2022]

Mean field: Small scale initialization \rightarrow features need to move significantly.

$$f(x) = \frac{1}{M} \sum_{j=1}^M r_j \sigma(w_j^\top x)$$

(参考) 初期化スケールと学習率の取り方¹³¹

• ABCパラメトライゼーション [Yang&Hu, 2021]

$$x^l(\xi) = \phi(h^l(\xi)) \in \mathbb{R}^n, \quad h^{l+1}(\xi) = W^{l+1}x^l(\xi) \in \mathbb{R}^n, \quad \text{for } l = 1, \dots, L-1, \quad n:\text{横幅}$$

(1) パラメータ設定

$$W^l = n^{-a_l} w^l$$

(w^l が学習パラメータ)

(2) 初期化法

$$w_{\alpha\beta}^l \sim \mathcal{N}(0, n^{-2b_l})$$

(3) 学習率のスケール

$$\eta n^{-c}$$

A

B

C

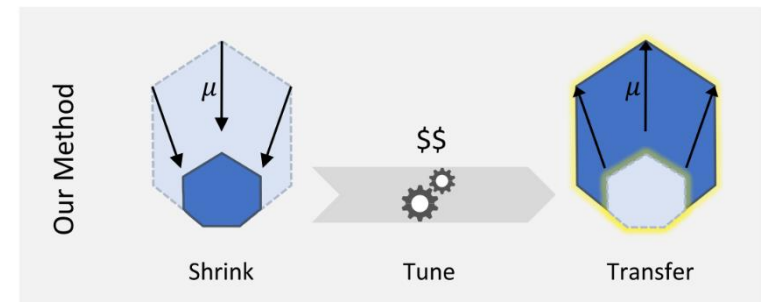
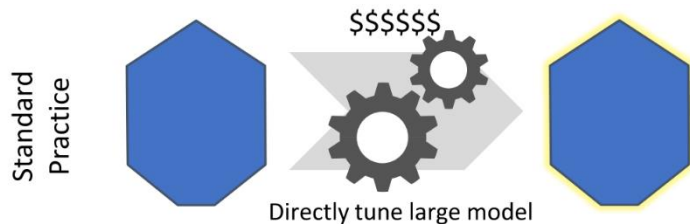
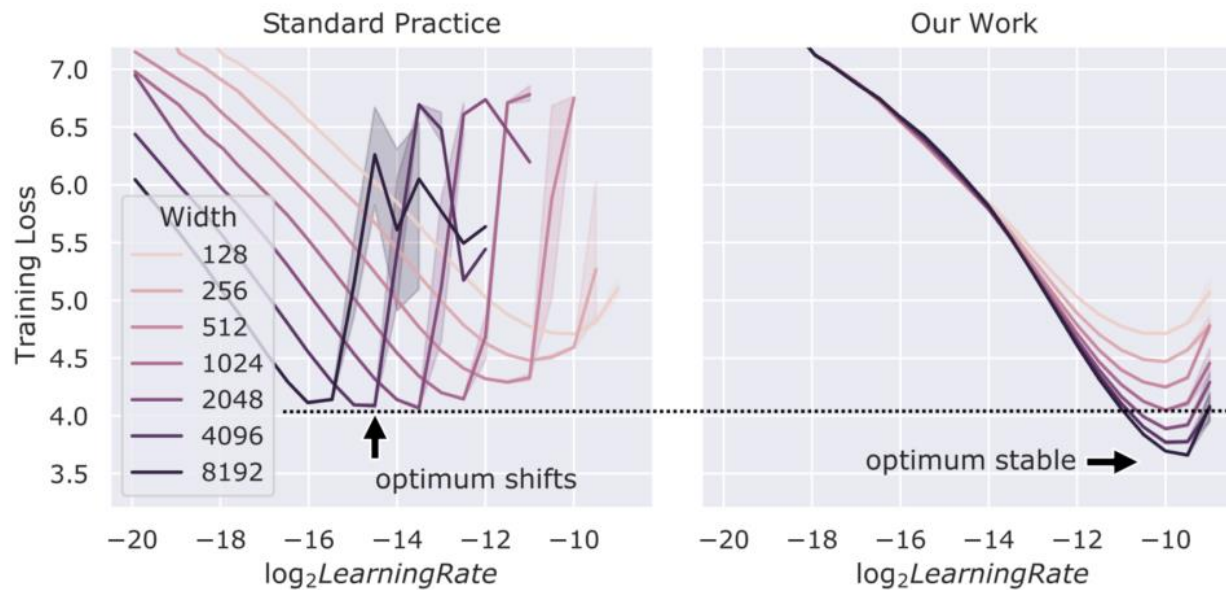
$$h^1 = W^1 \xi \in \mathbb{R}^n, x^l = \phi(h^l) \in \mathbb{R}^n, h^{l+1} = W^{l+1} x^l \in \mathbb{R}^n, f(\xi) = W^{L+1} x^L$$

	Definition	SP (w/ LR $\frac{1}{n}$)	NTP	MFP ($L = 1$)	μ P (ours)
a_l	$W^l = n^{-a_l} w^l$	0	$\begin{cases} 0 & l = 1 \\ 1/2 & l \geq 2 \end{cases}$	$\begin{cases} 0 & l = 1 \\ 1 & l = 2 \end{cases}$	$\begin{cases} -1/2 & l = 1 \\ 0 & 2 \leq l \leq L \\ 1/2 & l = L + 1 \end{cases}$
b_l	$w_{\alpha\beta}^l \sim \mathcal{N}(0, n^{-2b_l})$	$\begin{cases} 0 & l = 1 \\ 1/2 & l \geq 2 \end{cases}$	0	0	1/2
c	$LR = \eta n^{-c}$	1	0	-1	0
r	Definition 3.2	1/2	1/2	0	0
	$2a_{L+1} + c$	1	1	1	1
	$a_{L+1} + b_{L+1} + r$	1	1	1	1
	Nontrivial?	✓	✓	✓	✓
	Stable?	✓	✓	✓	✓
	Feature Learning?			✓	✓
	Kernel Regime?	✓	✓		

(適切なスケール)

[Yang et al.:Tensor Programs V: Tuning Large Neural Networks via Zero-Shot Hyperparameter Transfer. arXiv:2203.03466]

小さいモデルのハイパーパラメータを大きなモデルに転用できる。
GPTの学習に利用. 数億円の学習コストを抑えられる。



$$f_W(x) = \sum_{j=1}^M a_j \eta(w_j^\top x)$$

重要

テイラー展開

特徴写像とみなせる

$$f_W(x) \simeq (W - W^{(0)})^\top \nabla_W f_{W^{(0)}}(x) \quad (\text{線形モデル})$$

初期値のスケールが大きいため、初期値周りの線形近似でデータにフィットできてしまう。

カーネル関数は特徴写像の内積

$$\begin{aligned} \blackrightarrow \quad k_W(x, x') &= \langle \nabla_W f_{W^{(0)}}(x), \nabla_W f_{W^{(0)}}(x') \rangle \\ &= \sum_{j=1}^M \underbrace{a_j^2}_{\frac{1}{M}} (x^\top x') \eta'(w_j^\top x) \eta'(w_j^\top x') \end{aligned}$$

Neural Tangent Kernel各ニューロンの微分で得られる
特徴写像の内積を全ニューロンで平均

最適化のダイナミクスや汎化性能などは、
NTKをカーネル関数とするカーネル法としてとらえられる。

以下のように初期化する:

- $a_j \sim (\pm 1) \frac{1}{\sqrt{M}}$ (+, - is generated evenly)
- $w_j \sim N(0, I)$

$$f_W(x) = \sum_{j=1}^M a_j \eta(w_j^\top x)$$

Theorem [Arora et al., 2019]

$M = \Omega(n^2 \log(n) / \lambda_{\min})$ とすれば, 勾配法によって大域的最適解へ **線形収束** し, その汎化誤差は $\sqrt{\mathbf{y}^\top (K_{W(0)})^{-1} \mathbf{y} / n}$ で抑えられる.

$\exp(-ct)$ で収束

See also [Du et al., 2018; Allen-Zhu, Li & Song, 2018; Li & Liang, 2018]

- 訓練誤差0の解に線形収束する.
- 汎化誤差も一応抑えられている.

- データに完全にフィットさせてしまうので過学習の可能性あり.
- Early stoppingや正則化を入れれば過学習を防げる. (次ページ)

NTKの収束定理の導出

連続時間ダイナミクスを考える。

$$\text{Model : } f_W(x) = \sum_{j=1}^M a_j \eta(w_j^\top x)$$

- a_j は固定
- w_j を学習

[Jacot, Gabriel&Hongler, NeurIPS2018]

$$\frac{dw_j}{dt} = -\nabla_{w_j} \hat{L}(f_W) \quad (\text{Gradient descent, GD})$$

$$= -\frac{1}{n} \sum_{i=1}^n \ell'_i(f_W(x_i)) a_j \nabla_{w_j} \eta(w_j^\top x_i)$$

$$\nabla_{w_j} \eta(w_j^\top x_i) = x_i \eta'(w_j^\top x_i)$$

➡
(勾配法による更新)

$$\frac{df_W(x)}{dt} = \sum_{j=1}^M a_j \nabla_{w_j}^\top \eta(w_j^\top x) \frac{dw_j}{dt}$$

$O(1/M)$:
特徴写像の内積の平均

$$= -\frac{1}{n} \sum_{i=1}^n \left(\underbrace{\sum_{j=1}^M a_j^2 \nabla_{w_j}^\top \eta(w_j^\top x) \nabla_{w_j} \eta(w_j^\top x_i)}_{k_W(x, x_i)} \right) \ell'_i(f_W(x_i))$$

residual
(関数勾配)

$$k_W(x, x_i)$$

Neural Tangent Kernel

$$\begin{aligned}\frac{d\hat{L}(f_W)}{dt} &= \frac{1}{n} \sum_{i=1}^n \frac{df_W(x_i)}{dt} \ell'_i(f_W(x_i)) \\ &= -\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \ell'_i(f_W(x_i)) \underbrace{k_W(x_i, x_j)}_{(K_W)_{i,j}} \ell'_j(f_W(x_j)) \\ &= -\frac{1}{n^2} \|\nabla_f \hat{L}(f_W)\|_{K_W}^2 \\ &\leq -\lambda_{\min} \frac{1}{n^2} \|\nabla_f \hat{L}(f_W)\|^2 \quad (\lambda_{\min}: \text{グラム行列の最小固有値})\end{aligned}$$

Fact

[Du et al., 2018; Allen-Zhu, Li & Song, 2018]

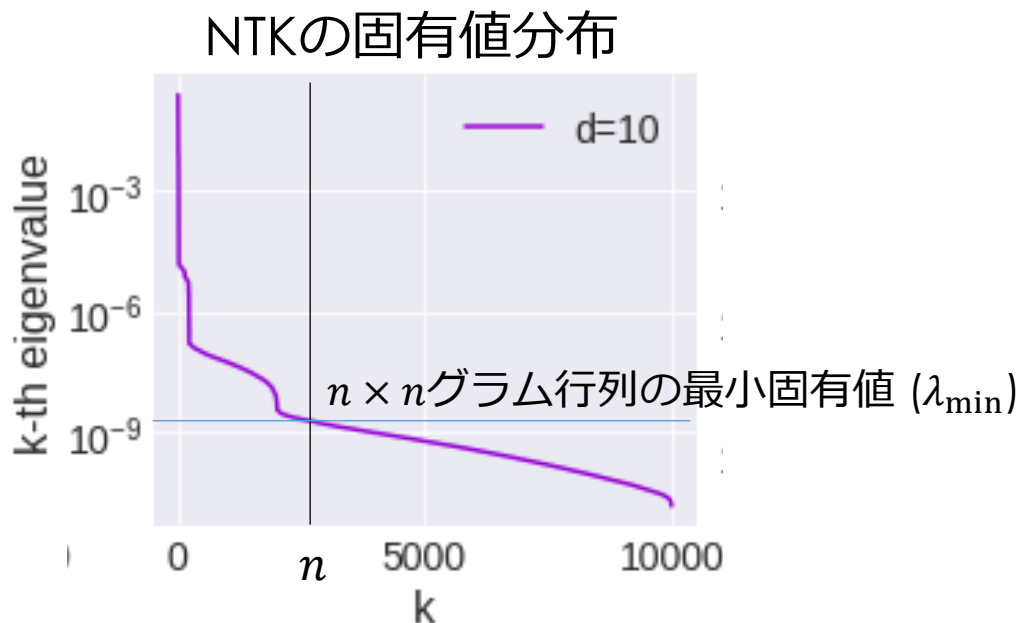
- ランダム初期化しておけば, $K_{W(0)} > \epsilon I$ が高確率で成立.
- 最適化の最中に最小固有値は正のまま ($\geq \epsilon/2$).



線形収束 ($\exp(-\lambda_{\min} t)$)

Spectral bias

- 最適化の観点からはoverparameterizationは有用に見える。
- 汎化誤差はどうであろうか？



- グラム行列の最小固有値は小さい ($1/\text{poly}(n)$).
 - 固有値の減少レートは多項式オーダー (理論+実験).
- Spectral bias: 汎化の意味では好ましい。

Kernelによる平滑化という視点

- Frechet 微分 in $L_2(P_n)$: $\nabla_f \hat{L}(f)$

$$\nabla_f \hat{L}(f) = (\ell'_i(f(x_i)))_{i=1}^n$$

$$\hat{L}(f+h) = \hat{L}(f) + \langle \nabla_f \hat{L}(f), h \rangle_{L_2(P_n)} + o(\|h\|_{L_2(P_n)}^2)$$

- **平滑化**積分作用素:

$$T_k f(x) := \int k(x, x') f(x') dP_n(x')$$

$$T_{k_W} \phi_j = \mu_j \phi_j$$

- NTKにおける勾配は関数勾配を平滑化したもの:

$$\frac{df_W}{dt} = -T_{k_W} \nabla_f \hat{L}(f_W) \leftarrow \text{勾配を平滑化!}$$

$$\left(= -\frac{1}{n} \sum_{i=1}^n k_W(\cdot, x_i) \ell'_i(f_W(x_i)) \right)$$

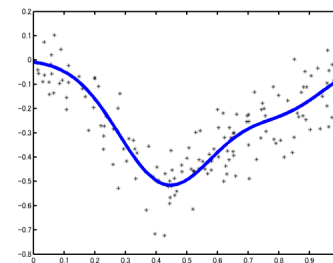
(導出は前のページ参照「NTKの収束定理の導出」)

k_W が高周波成分に小さな固有値を持てば, T_{k_W} は平滑化作用素として働く
 → 帰納的バイアス (inductive bias).

モデル:

$$f_{a,W}(x) = \frac{1}{\sqrt{M}} \sum_{j=1}^M a_j \eta(w_j^\top x)$$

(We train both of first and second layers)



目的関数:

$$L(a, W) = \mathbb{E}[(Y - f_{a,W}(X))^2] + \frac{\lambda}{2} (\|a - a^{(0)}\|^2 + \|W - W^{(0)}\|_F^2)$$

期待損失

初期値からのずれ

$$Y = f^*(X) + \epsilon \quad (\text{ノイズありの観測})$$

Averaged Stochastic Gradient Descent

for $t = 0$ **to** $T - 1$ **do**

Randomly draw a sample $(x_t, y_t) \sim \rho$

Perform SGD update for all $j \in \{1, \dots, M\}$:

$$a_j^{(t+1)} = a_j^{(t)} - \alpha_t [\nabla_a \ell(y_t, f_{a^{(t)}, W^{(t)}}(x_t)) + \lambda(a^{(t)} - a^{(0)})]$$

$$W_j^{(t+1)} = W_j^{(t)} - \alpha_t [\nabla_W \ell(y_t, f_{a^{(t)}, W^{(t)}}(x_t)) + \lambda(W^{(t)} - W^{(0)})]$$

end for

$$\text{Return } \bar{a}^{(T)} = \frac{1}{T} \sum_{t=0}^{T-1} a^{(t)}, \quad \bar{W}^{(T)} = \frac{1}{T} \sum_{t=0}^{T-1} W^{(t)}.$$

NTKにおける余剰誤差の速い収束

[Nitanda&Suzuki: Fast Convergence Rates of Averaged Stochastic Gradient Descent under Neural Tangent Kernel Regime, 2020.]

仮定：真の関数がNTKの作るRKHSに入っているとする。

NTK設定で適切な正則化を入れたSGDは“速い学習レート”を達成できる。

→ NTKによるsmoothingのおかげ。

Thm (速い収束レート)

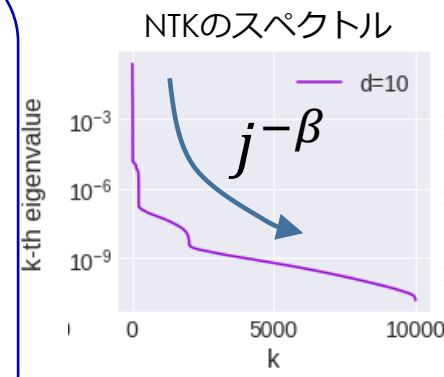
f_T : T 回更新後の解

NTKの固有値の減衰レート

$$\mathbb{E}[\|f_T - f^*\|_{L_2}^2] \leq \epsilon_M + O\left(T^{-\frac{2r\beta}{2r\beta+1}}\right)$$

$M \rightarrow \infty$ で0に収束する項

速い学習レート
($O(1/\sqrt{T})$ より速い)



→ $T^{-\frac{2r\beta}{2r\beta+1}}$ はミニマックス最適レート。
(各種パラメータの意味は次ページに詳細)

2層NNのNTK:

$$k_\infty(x, x') = \mathbb{E}_{w^{(0)}} [\eta(w^{(0)\top} x) \eta(w^{(0)\top} x')] + \mathbb{E}_{w^{(0)}} [\eta'(w^{(0)\top} x) \eta'(w^{(0)\top} x') x^\top x]$$

横幅無限における積分作用素:

$$T_{k_\infty} f(x) = \int k_\infty(x, x') f(x') dP_X$$

population

スペクトル分解: $T_{k_\infty} \phi_j = \mu_j \phi_j$, $k_\infty(x, x') = \sum_{j=1}^{\infty} \mu_j \phi_j(x) \phi_j(x')$

仮定

- $f^*(x) = \mathbb{E}[Y|X = x]$ が次のように書ける:

$$T_{k_\infty}^r h = f^*$$

for $h \in L_2(P_X)$, and $r \in [1/2, 1]$.

- 固有値減衰条件:

$$\mu_j = O(j^{-\beta}).$$

真の関数の平滑性

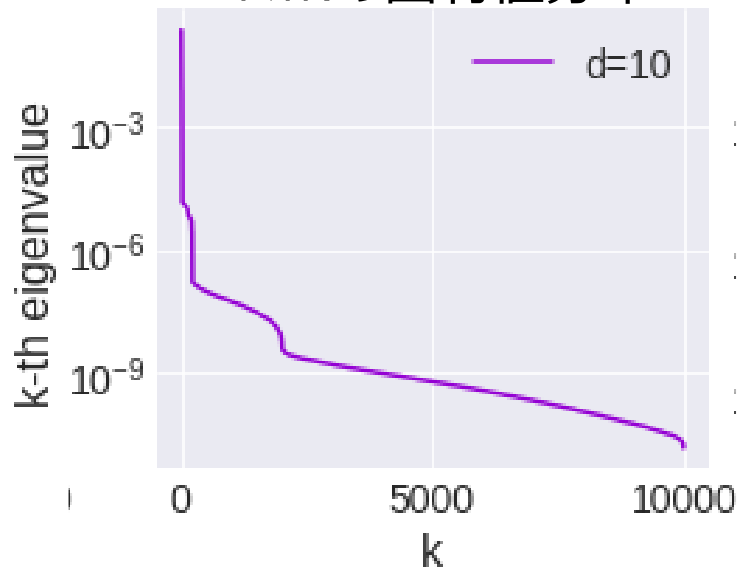
カーネル関数の
“複雑さ”

カーネルリッジ回帰の解析における標準的な仮定; see, e.g., Dieuleveut et al. (2016); Caponnetto and De Vito (2007) (r の条件はやや強め).

$$k_{\infty}(x, x') = \sum_{m=1}^{\infty} \lambda_m \phi_m(x) \phi_m(x')$$

NTKの固有値固有関数分解
 $(\phi_m)_{m=1}^{\infty}$: 固有関数. $L_2(P_X)$ 内の
 正規直交基底.

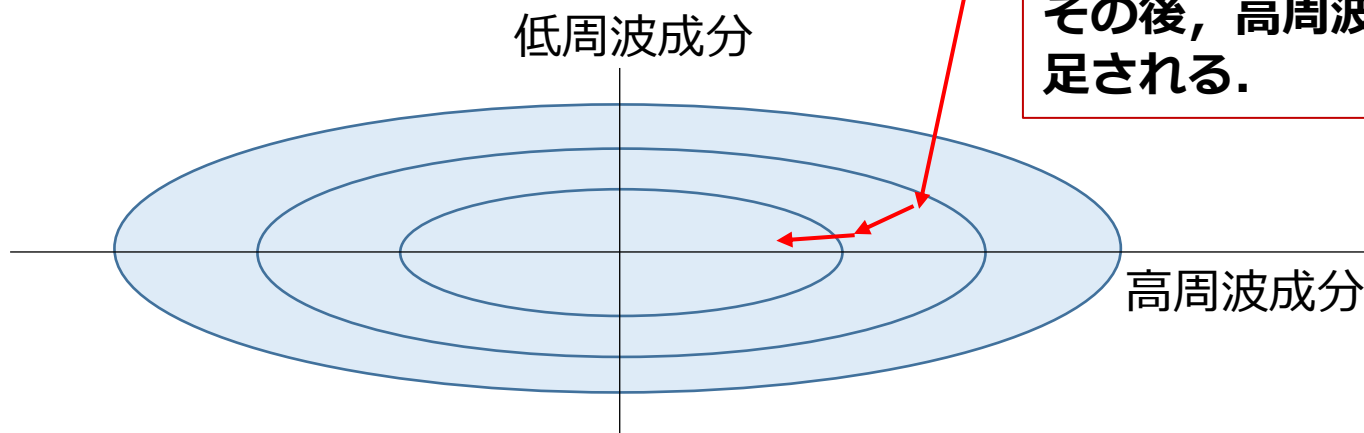
NTKの固有値分布



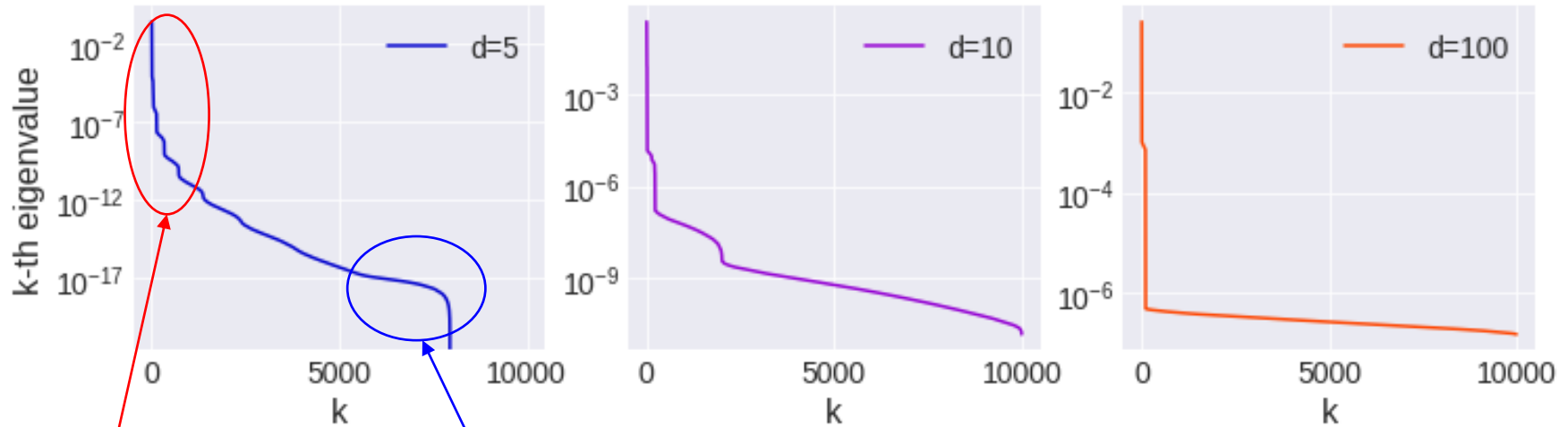
実際のNTKの固有値は多項式
 オーダーで減衰する.

[Bietti&Mairal (2019); Cao et al. (2019);
 Ronen et al. (2019)]

低周波成分が最初に補足される。
 その後、高周波成分が徐々に補
 足される。



固有値減少の数値実験による検証



High frequency components

Low frequency components

理論

ReLU, $(a, w) \sim N(0, I)$:

$$k_{\text{NTK}}(x, y) = \sum_{k=0}^{\infty} \mu_k \sum_{j=1}^{N(d,k)} Y_{k,j}(x) Y_{k,j}(y)$$

- $Y_{k,j}$: spherical harmonics functions with degree k .
- $\mu_k \sim k^{-d}$.
- $N(d, k) = \frac{2k + d - 2}{k} \binom{k + d - 3}{d - 2}$

$$f_{a,W}(x) = \frac{1}{\sqrt{M}} \sum_{j=1}^M a_j \eta(w_j^\top x)$$

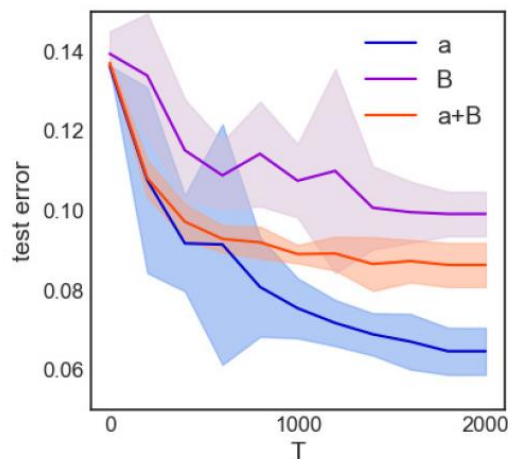
$$k_{a,W}(x, x') = k_a(x, x') + k_W(x, x')$$

二層目のNTK 一層目のNTK

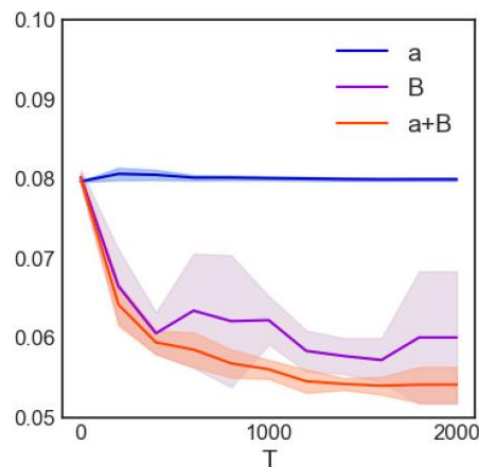
一層目と二層目のカーネルの和 : multiple kernel

仮定: f^* is in

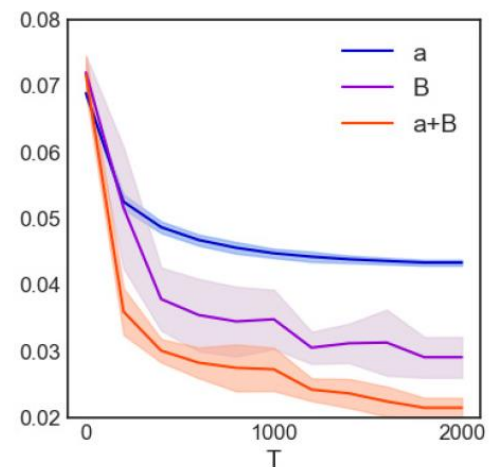
RKHS w.r.t. NTK
for the 2nd layer.



RKHS w.r.t. NTK
for the 1st layer.



RKHS w.r.t. NTK
for the both layer.



- 二層NNのNTKによる学習は, multiple kernel learningの効果がある.
- 多層NNを用いることはモデルmisspecificationに対してよりロバストになる.

問題点：NTKは解析がしやすいが，結局カーネル法の範疇なので深層学習の“良さ”が現れない。

➤ NTKをはみ出す理論の試みがいくつかなされている。

(今後発展が予想される)

- Allen-Zhu&Li (2019,2020)

Allen-Zhu&Li: What Can ResNet Learn Efficiently, Going Beyond Kernels? NIPS2019.

Allen-Zhu&Li: Backward Feature Correction: How Deep Learning Performs Deep Learning. arXiv:2001.04413.

(ResNet型ネットワークでカーネルを優越する状況)

- Li, Ma&Zhang (2019)

Li, Ma&Zhang: Learning Over-Parametrized Two-Layer ReLU Neural Networks beyond NTK. arXiv:2007.04596.

(テンソル分解の理論で深層学習がカーネルを優越することを示した)

- Bai&Lee (2020)

Bai&Lee: Beyond Linearization: On Quadratic and Higher-Order Approximation of Wide Neural Networks. ICLR2020.

(二次のテイラー展開まで使う)

2層NNの特徴学習のダイナミクス解析

- **(NTKを超えて)**勾配法によって特徴量がどう学習されるか？
 - 色々な研究がある.
- NTK近似が成り立たない領域では非線形性が強くなり, 最適化のダイナミクスの解析が難しくなる.
- 最近では妥協点として, 勾配法の最初の段階 (1ステップ or 少数ステップ) でどのような特徴量が獲得できるかを解析する研究が複数なされている.
 - 少数ステップで得られる情報は限られているが, それでも予測性能の改善が示せる.
 - 今後はより非線形性の強い特徴量学習のダイナミクスの解析が進むと思われる.

※ 計算量をあまり気にしなければ勾配ランジュバン動力学を用いた学習の解析は完全に非線形な特徴学習をとらえられる.

$$f_{\text{NN}}(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^N a_i \sigma(\langle x, w_i \rangle) = \frac{1}{\sqrt{N}} a^\top \sigma(W^\top x)$$

問：勾配法で W を更新することで，データに合った特徴量を獲得できるか？

答：NTK近似が成り立つ領域からはみ出るくらい大きなステップサイズを用いれば，一回の更新で意味のある特徴量の方向を得ることができる。

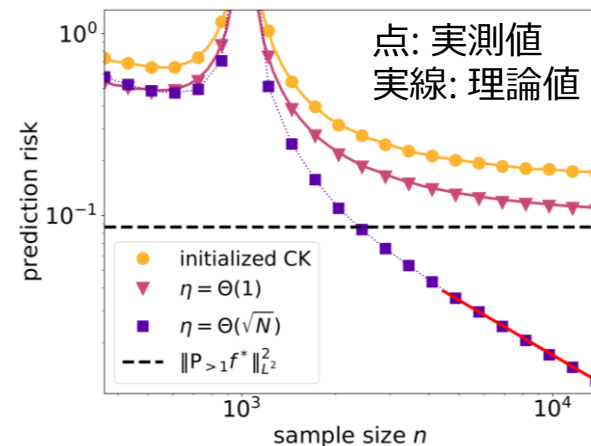
→ カーネルAlignment, 特徴量学習.

$$W_{k+1} = W_k + \eta \sqrt{N} \nabla L(f_{\text{NN}})$$

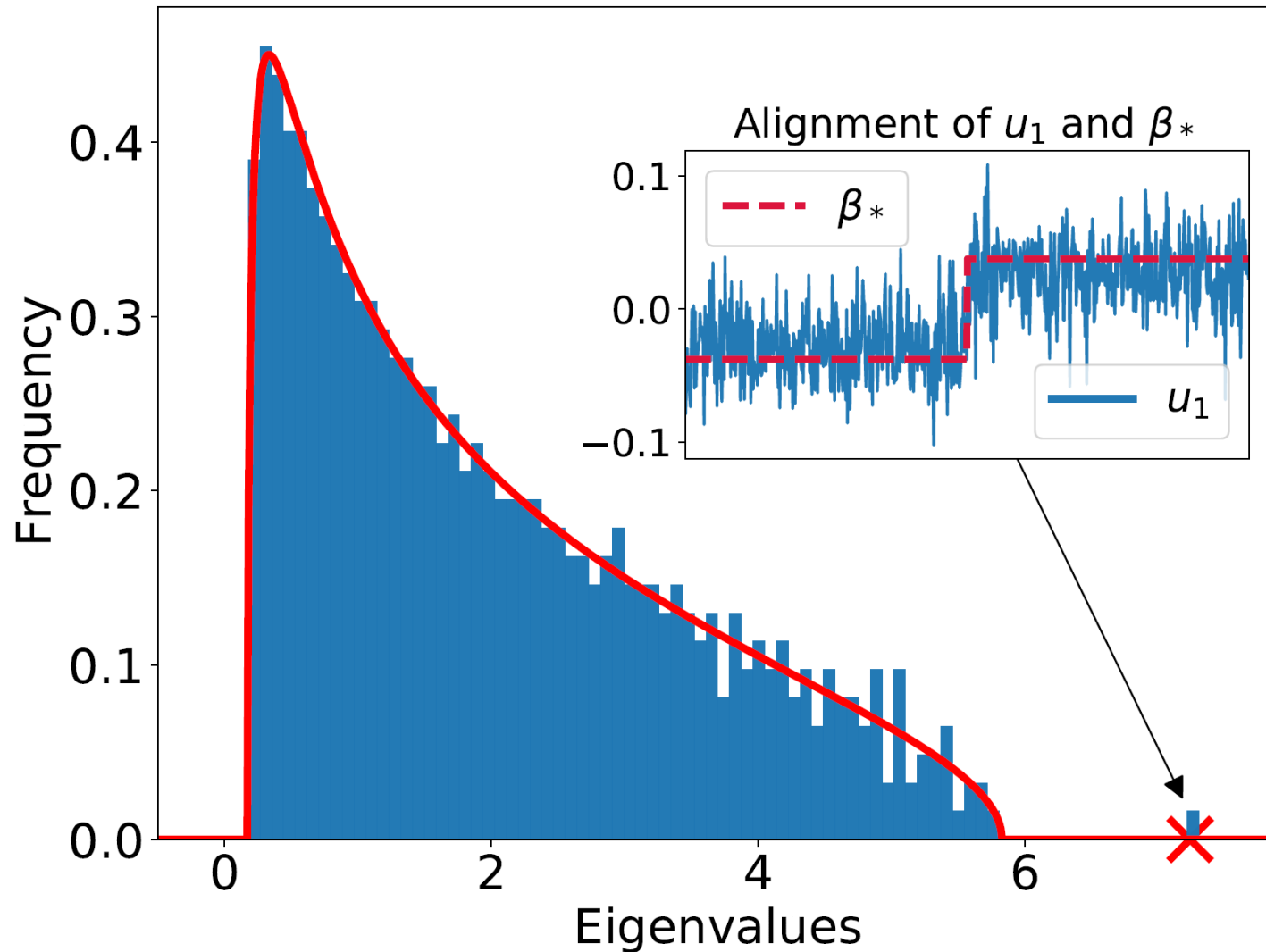
$n, d, N \rightarrow \infty$ の極限を考え，勾配法1回の更新後の予測誤差を評価してみる。

- $\eta = \sqrt{N}$ のように大きな更新を用いると，任意の線形モデルによるリッジ回帰を優越する。
- $\eta = 1$ では最適なリッジ回帰を優越しないが初期値 W は優越。
- $\eta = o(1)$ では初期値 W と同じ予測誤差 (NTK-regime)。

Gaussian equivalence property + Random matrix theory
→ Exact risk evaluation.



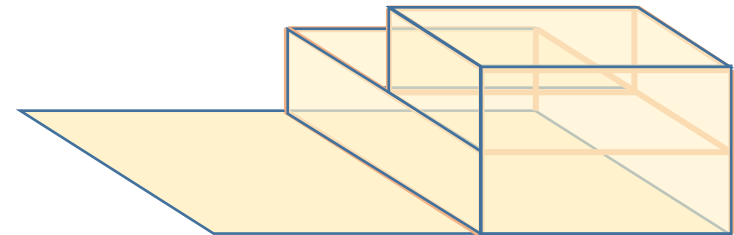
勾配法とKernel alignment



The first few step of GD with large learning rate can extract informative features.

- **Staircase function** [Abbe et al., NeurIPS2021; Abbe et al., arXiv2202.08658]

Small number of gradient descent can extract nonlinear features to estimate “staircase” function. The trained features for GD can outperform random feature model.



- **Benign overfitting with feature learning**

[Cao et al., arXiv:2202.06526; Frei et al., arXiv:2202.05928]

Gradient descent in two-layer NN can yield benign overfitting and achieves almost the Bayes error in binary classification.

問題設定

観測モデル:

$$y_i = f^*(x_i) + \epsilon_i \quad (i = 1, \dots, n)$$

where $x_i \sim N(0, I)$, $\epsilon_i \sim N(0, 1)$, and $x_i \in \mathbf{R}^d$.

➤ 平均場スケールの2-層NNを学習:

$$f_{\text{NN}}(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^N a_i \sigma(\langle x, w_i \rangle) = \overbrace{\frac{1}{\sqrt{N}} a^\top}^{\text{平均場レジーム } \Delta O(1/N)} \sigma(W^\top x) \quad (\because a_i = O_p(1/\sqrt{N}))$$

where $a_i \sim N(0, \underbrace{1/N}_{\text{var}})$ and $W_{ij} \sim N(0, \underbrace{1/d}_{\text{var}})$.

経験リスク:

$$\mathcal{L}(f) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

予測リスク:

$$\mathcal{R}(f) = \mathbb{E}[(f^*(X) - f(X))^2]$$

問: 勾配法による特徴学習で予測誤差を改善できるか?

真の関数がシングルインデックスモデルで書ける状況を解析:

$$f^*(x) = \sigma^*(\langle x, \beta^* \rangle)$$

勾配法

- 一層目のパラメータ W の更新:

$$W_{t+1} = W_t + \eta\sqrt{N} \underbrace{\left(-\nabla_W \mathcal{L}(f_{\text{NN}}^{(t)})/2\right)}_{G_t}$$

$$G_t = -\frac{1}{n} X^\top \left[\left(\frac{1}{\sqrt{N}} \left(\frac{1}{\sqrt{N}} \sigma(XW_t)a - y \right) a^\top \right) \odot \sigma'(XW_t) \right]$$

- ✓ 二層目のパラメータ a は最適化の途中で固定.
- ✓ あくまで W の特徴学習のダイナミクスに注目.

- 特に, 勾配法の“**1ステップ更新**”に注目:

$$W_1 = W_0 + \eta\sqrt{N}G_0$$

- その後, 二層目はリッジ回帰で推定: $\phi_{\text{GD}}(x) = \frac{1}{\sqrt{N}}\sigma(W_1^\top x)$

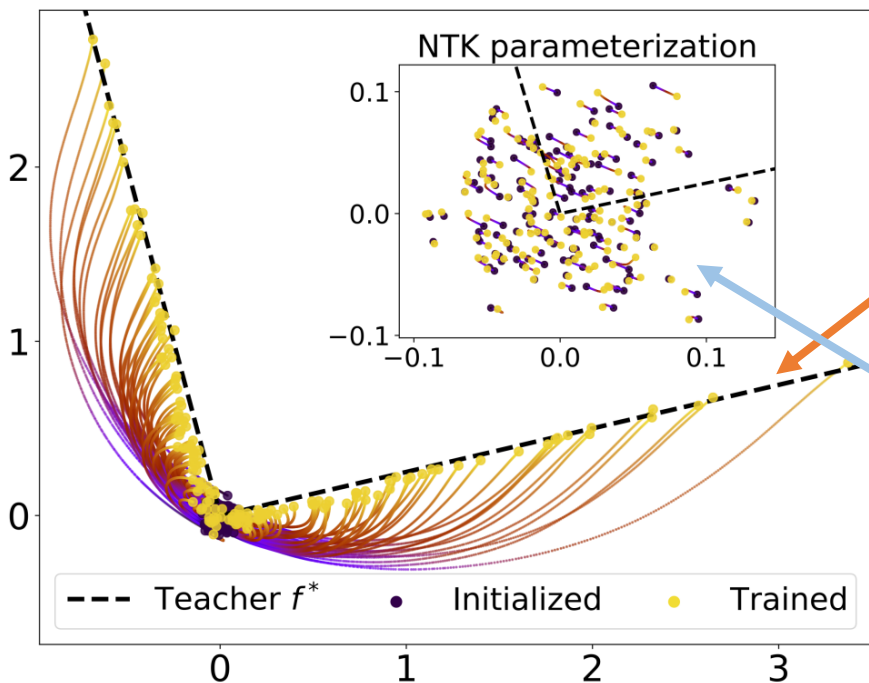
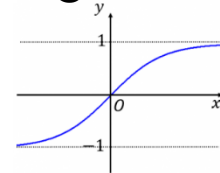
$$\hat{a}_{\text{GD}} = \arg \min_{a \in \mathbb{R}^N} \left\{ \frac{1}{n} \sum_{i=1}^n (\tilde{y}_i - \langle a, \phi_{\text{GD}}(\tilde{x}_i) \rangle)^2 + \frac{\lambda}{N} \|a\|^2 \right\} \quad (\tilde{x}_i, \tilde{y}_i)_{i=1}^n: \text{i.i.d. copy of } (x_i, y_i)_{i=1}^n$$

- W_1 (1回更新後) は, 初期値 W_0 より真の関数 f^* に「関連している」と考えられる. →より良い予測誤差.

仮定

$$f_{\text{NN}}(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^N a_i \sigma(\langle x, w_i \rangle) = \frac{1}{\sqrt{N}} a^\top \sigma(W^\top x)$$

1. **Proportional limit:** $n, d, N \rightarrow \infty$ and $n/d \rightarrow \psi_1, N/d \rightarrow \psi_2$. e.g., \tanh
2. **活性化関数の性質:**
 1. σ は有界, $\max\{\|\sigma'\|_\infty, \|\sigma''\|_\infty, \|\sigma'''\|_\infty\} \leq C < \infty$.
 2. $E[\sigma(z)] = 0, E[z\sigma(z)] \neq 0$ for $z \sim N(0,1)$. (ちょっと強い条件)
3. **教師の条件:** f^* はリプシッツ連続で $\|f^*\|_{L^2(P_X)} = \Theta_d(1)$.



W の勾配法における軌跡 ($d = 2$).
 f^* は二つの ReLU ニューロンの和.

1. 平均場設定では, 各ニューロンは初期値から大きく動いて目標となる真の関数の方向 (黒い破線; 二つのニューロン) を向く.
2. NTK のスケールでは, 一層目のパラメータはほとんど動かず特徴学習ができていない.

See also Akiyama&Suzuki (2021), Chizat (2019).

ランダム特徴量

ランダム特徴量 (特徴量の学習無し)

- Conjugate kernel at initialization:

$$\phi_{\text{CK}}(x) = \frac{1}{\sqrt{N}} \sigma(W_0^\top x)$$

正確な漸近解析がかなり研究されている (e.g., [Louart, Liao, and Couillet, 2018; Mei and Montanari, 2019])

- NTK (Neural tangent kernel):

$$\phi_{\text{NTK}}(x) = \frac{1}{\sqrt{Nd}} \text{Vec}(\sigma'(W_0^\top x) x^\top)$$

$$\hat{a}_{\text{RF}} = \arg \min_{a \in \mathbb{R}^N} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - \langle a, \phi_{\text{RF}}(x_i) \rangle)^2 + \frac{\lambda}{N} \|a\|^2 \right\}$$

RF \in {CK, NTK}

学習されたカーネル (1ステップGD):

$$\phi_{\text{GD}}(x) = \frac{1}{\sqrt{N}} \sigma(W_1^\top x)$$

よくわかってない

$$\hat{a}_{\text{GD}} = \arg \min_{a \in \mathbb{R}^N} \left\{ \frac{1}{n} \sum_{i=1}^n (\tilde{y}_i - \langle a, \phi_{\text{GD}}(\tilde{x}_i) \rangle)^2 + \frac{\lambda}{N} \|a\|^2 \right\}$$

$(\tilde{x}_i, \tilde{y}_i)_{i=1}^n$: i.i.d. copy of $(x_i, y_i)_{i=1}^n$

$$\hat{a}_{\text{RF}} \quad \text{VS.} \quad \hat{a}_{\text{GD}}$$

[El Karoui (2010); Ghorbani et al. (2019); Hu and Lu (2020); ...]

定理 (RFの予測誤差の下限)

$$\inf_{\lambda > 0} \min\{\mathcal{R}_{\text{CK}}(\lambda), \mathcal{R}_{\text{NTK}}(\lambda)\} \geq \|P_{>1}f^*\|_{L^2(P_X)}^2 + o_{p,d}(1)$$

$$P_{>1}f^* := (I - P_{\leq 1})f^*$$

ただし, $P_{\leq 1}$ は線形関数かなす空間への $L^2(P_X)$ -空間内での射影.

Remark: 同じことが “回転不変カーネル” でもなりたつ [El Karoui (2010)].

- 高次元かつ $d/n = o(1)$ の状況では, ランダム特徴量を使っている限り真の関数 f^* の線形要素しか取り出せない.
- f^* の非線形性が強ければ, 精度が出せない.
($n = d^{k+1-\epsilon}$ なら k 次多項式まで出てくる [Ghorbani et al. 2021; Mei et al. 2021])

これは, 高次元設定では以下のような中心極限定理が成り立つからである:

$$a^\top \phi_{\text{CK}}(x_i) = \frac{1}{\sqrt{N}} a^\top \sigma(W_0^\top x_i) \approx \frac{1}{\sqrt{N}} a^\top (\mu_1 W_0^\top x_i + \mu_2 z)$$

**[Gaussian等価性]
後のスライドを参照**

最初の勾配ステップはほぼランクが 1

- 勾配 G_t は, ランク 1 行列で近似できる.
 $\Rightarrow W_1$ のスペクトル分布に「スパイク」が現れる!

$$G_t = -\frac{1}{n} X^\top \left[\left(\frac{1}{\sqrt{N}} \left(\frac{1}{\sqrt{N}} \sigma(XW_t) a - y \right) a^\top \right) \odot \sigma'(XW_t) \right]$$

(σ' の非線形性より, 一般的には低ランクにならない. しかし高次元だと低ランクになる)

(Gordon-Slepian ineq.; Hanson-Wright ineq.)

定理 (勾配のランク 1 近似)

$$G_0 = \frac{1}{\eta\sqrt{N}} (W_1 - W_0) \quad (\because W_1 = W_0 + \eta\sqrt{N}G_0) \quad \text{に注意.}$$

$$\mu_1 = \mathbb{E}[z\sigma(z)], \quad \mu_2 = \sqrt{\mathbb{E}[\sigma(z)^2] - \mu_1^2}, \quad \text{where } z \sim \mathcal{N}(0, 1). \quad \text{(活性化関数の 1 次成分と 2 次成分)}$$

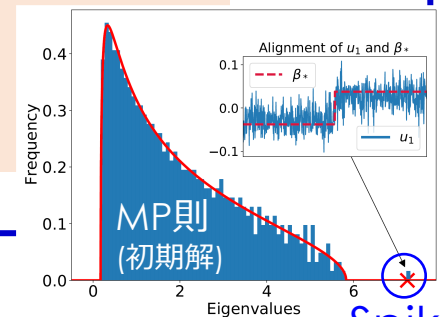
$$A := \frac{\mu_1}{n\sqrt{N}} X^\top y a^\top \quad (\text{ランク 1 行列}) \text{ とすれば, 以下を得る:}$$

$$\|G_0 - A\| \lesssim \frac{\log^2(n)}{\sqrt{n}} \cdot \|G_0\|$$

with high probability for sufficiently large n, d, N .

$$W_1 = W_0 + \eta \times (\text{rank one matrix}).$$

\Rightarrow ステップサイズ η が大きいと, スパイクが支配的.



Alignment with the Target Function

f^* can be decomposed into **first order** **higher order**

$$f^*(x) = \mu_0^* + \mu_1^* \langle x, \beta_* \rangle + P_{>1} f^*(x)$$

$$\text{where } \mu_1^* \beta_* = \mathbb{E}_{x \sim N(0, I)} [x f^*(x)].$$

Define $\bar{\mu} := \lim_{d \rightarrow \infty} \|f^*\|_{L^2(P_X)}$, and

$$\theta_1 := \sqrt{\bar{\mu}^2 \psi_1^{-1} + \mu_1^{*2}} \cdot \mu_1 \eta, \quad \theta_2 := \mu_1 \mu_1^* \eta.$$

$$n/d \rightarrow \psi_1, \quad N/d \rightarrow \psi_2$$

Theorem (Alignment to the target function)

Suppose that $\eta = \Theta(1)$ and $\sigma_1(W_1)$ is the largest singular value of W_1 and u_1 is the corresponding left singular vector, then

(i) $\theta_1 > \psi_2^{1/4}$ (large step size):

$$|\langle u_1, \beta_* \rangle|^2 \rightarrow \frac{\theta_2^2}{\theta_1^2} \left(1 - \frac{\psi_2 + \theta_1^2}{\theta_1^2 (\theta_1^2 + 1)} \right), \quad \sigma_1(W_1) \rightarrow \sqrt{\frac{(1 + \theta_1^2)(\psi_2 + \theta_1^2)}{\theta_1^2}}$$

(ii) $\theta_1 < \psi_2^{1/4}$ (small step size):

$$|\langle u_1, \beta_* \rangle|^2 \rightarrow 0, \quad \sigma_1(W_1) \rightarrow 1 + \sqrt{\psi_2}$$

In both cases, we have

$$|\langle u_i, \beta_* \rangle|^2 \rightarrow 0 \quad (\forall i \geq 2).$$

Alignment with the Target Function 158

参考

f^* can be decomposed into first order higher order

$$f^*(x) = \mu_0^* + \mu_1^* \langle x, \beta_* \rangle + P_{>1} f^*(x)$$

$$\text{where } \mu_1^* \beta_* = \mathbb{E}_{x \sim N(0, I)} [x f^*(x)].$$

Define $\bar{\mu} := \lim_{d \rightarrow \infty} \|f^*\|_{L^2(P_X)}$, and

$$\theta_1 := \sqrt{\bar{\mu}^2 \psi_1^{-1} + \mu_1^{*2}} \cdot \mu_1 \eta, \quad \theta_2 := \mu_1 \mu_1^* \eta.$$

$$n/d \rightarrow \psi_1, \quad N/d \rightarrow \psi_2$$

Theorem (Alignment to the target function)

Suppose that $\eta = \Theta(1)$ and $\sigma_1(W_1)$ is the largest singular value, u_1 is the corresponding left singular vector, then

Phase transition for the leading singular value: i.e., under sufficiently large step size, the “spike” in the weight matrix aligns with the leading term of the true function.

(i) $\theta_1 > \psi_2^{1/4}$ (large step size):

$$|\langle u_1, \beta_* \rangle|^2 \rightarrow \frac{\theta_2^2}{\theta_1^2} \left(1 - \frac{\psi_2 + \theta_1^2}{\theta_1^2 (\theta_1^2 + 1)} \right)$$

$$\sigma_1(W_1) \rightarrow \sqrt{\frac{(1 + \psi_1)(\psi_2 + \psi_1)}{\psi_2}}$$

The “alignment” becomes more significant as we (i) increase the step size; (ii) use more data (i.e., larger ψ_1)

(ii) $\theta_1 < \psi_2^{1/4}$ (small step size):

$$|\langle u_1, \beta_* \rangle|^2 \rightarrow 0,$$

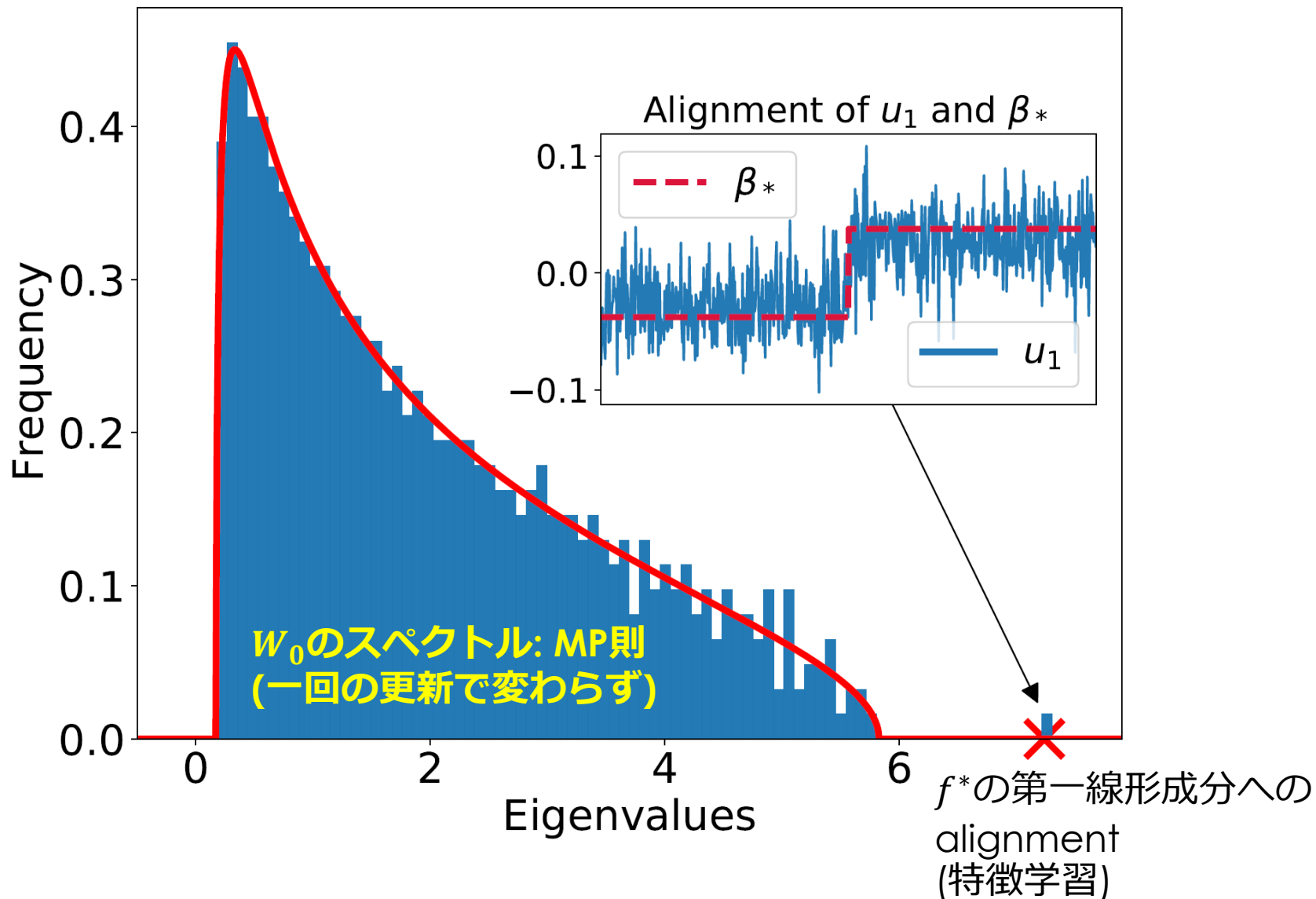
$$\sigma_1(W_1) \rightarrow 1 + \sqrt{\psi_2}$$

In both cases, we have

$$|\langle u_i, \beta_* \rangle|^2 \rightarrow 0 \quad (\forall i \geq 2).$$

Non-leading singular vectors are near orthogonal to the linear part of the true function in the high dimensional setting.

W_1 のスペクトル分布のスパイク



- t ステップ目のconjugate kernel :

$$\phi_{\text{CK}}(x) = \frac{1}{\sqrt{N}} \sigma(W_t^\top x)$$

- **Gaussian等価なモデル:**

$$\phi_{\text{GE}}(x) = \frac{1}{\sqrt{N}} (\mu_1 W_t^\top x + \mu_2 z)$$

線形成分+ガウス雑音
(非線形部分はガウス雑音になってしまう)
 $\mu_1 = \mathbb{E}[z\sigma(z)], \mu_2 = \sqrt{\mathbb{E}[\sigma(z)^2] - \mu_1^2}$,
where $z \sim \mathcal{N}(0, 1)$.

$$\hat{a}_F = \arg \min_{a \in \mathbb{R}^N} \left\{ \frac{1}{n} \sum_{i=1}^n (\tilde{y}_i - \langle a, \phi_F(\tilde{x}_i) \rangle)^2 + \frac{\lambda}{N} \|a\|^2 \right\}$$

$(\tilde{x}_i, \tilde{y}_i)_{i=1}^n$: i.i.d. copy
of $(x_i, y_i)_{i=1}^n$

where $F \in \{\text{CK}, \text{GE}\}$.

初期解に関する解析は[Hu and Lu, 2020]

定理 (Gaussian等価性)

$\eta = \Theta(1)$ として, $f^*(x) = \sigma^*(\langle x, \beta^* \rangle)$ とする. すると, 有界な t に対する最初の t ステップでは以下が成り立つ:

$$|\mathcal{R}_{\text{CK}}(\lambda) - \mathcal{R}_{\text{GE}}(\lambda)| = o_{p,d}(1)$$

$$\mathcal{R}_{\text{GE}}(\lambda) \geq \|P_{>1} f^*\|_{L^2(P_X)}^2$$

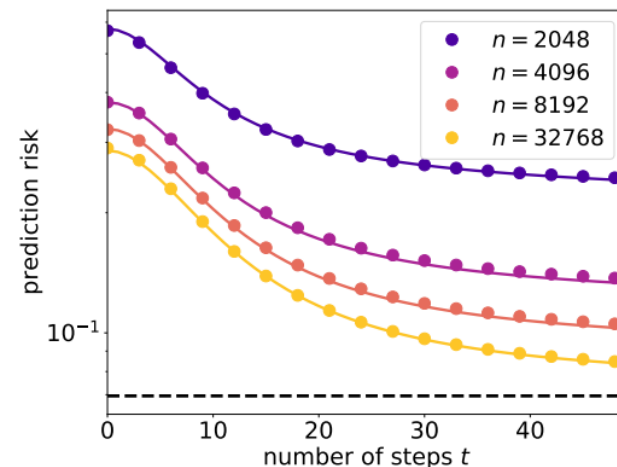
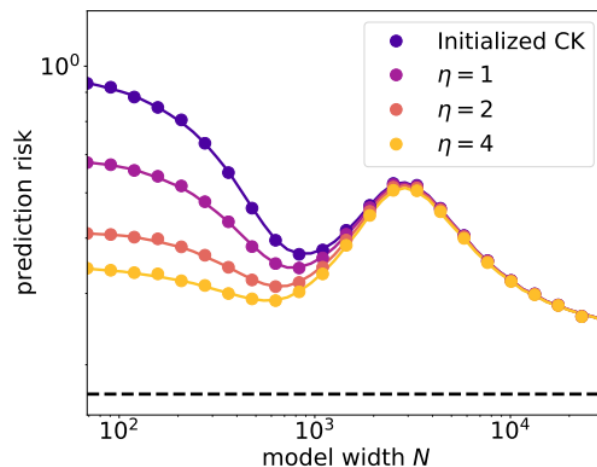
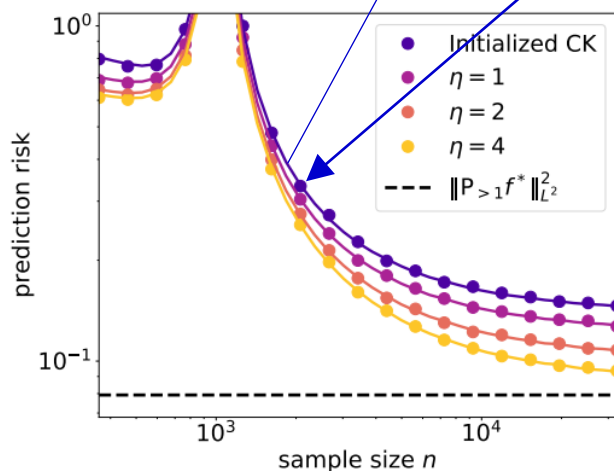
- これは特徴量の空間における“中心極限定理”と言える. (L2正則化によって a_i の全成分が満遍なく小さくなるので)
- 非線形部分はi.i.d.ガウス雑音のようにふるまう (非線形項の情報は消える).
- 小さなステップサイズで定数回勾配法で更新しても限界がある.

数値実験との整合性

1. GETは初期解(ランダム特徴量)だけでなく、**学習したモデル**でも成り立つ。
2. GETを使うことで、ランダム行列の理論から**予測誤差を厳密に計算可能**。
(なぜならカーネル関数が“線形化”されているので)。

Solid line: GEから理論的予測.

Dots: 計算機によるシミュレーション.



$\eta = \Theta(1)$ を用いて学習した特徴量を用いたリッジ回帰の予測誤差

より大きなステップサイズでより多くのステップを費やした方がより良い精度
→ 特徴学習が予測誤差を改善している！

初期値のCKからの改善

($R_W(\lambda)$ is the ridge regression estimator using W for the first layer.)

- $\eta = \Theta(1)$ (中間的な大きさのステップサイズ):

$$\mathcal{R}_{W_0}(\lambda) - \mathcal{R}_{W_1}(\lambda) \xrightarrow{p} \exists \delta > 0$$

線形な領域

- 1ステップのGDは必ず精度を改善させる.
- しかし, 大きく改善させることはない. なぜなら $\mathcal{R}_{W_1}(\lambda) \geq \|P_{>1} f^*\|_{L^2(P_X)}^2$.

- $\eta = \Theta(\sqrt{N})$ (大きな学習率): $f^*(x) = \sigma^*(\langle x, \beta^* \rangle)$ を仮定

Maximal update parameterization (μP) [Yang and Hu, 2020] として知られている.

- $\tau^* = 0$ if $\sigma = \sigma^* = \text{erf}$.
- $\tau^* \ll 1$ if $\sigma = \sigma^* = \text{tanh}$.

$$\tau^* = \inf_{\eta > 0} \mathbb{E}_{\xi_1 \sim N(0,1)} [\sigma^*(\xi_1) - \mathbb{E}_{\xi_2 \sim N(0,1)} [\sigma(\eta \xi_1 + \xi_2)]]$$

(モデルのズレを表す量)

$$n/d \rightarrow \psi_1, N/d \rightarrow \psi_2$$

$$\mathcal{R}_{W_1}(\lambda) \leq 16\tau^* + C(\sqrt{\tau^*} \psi_1^{-1/2} + \psi_1^{-1}) + o_p(1)$$

非線形領域

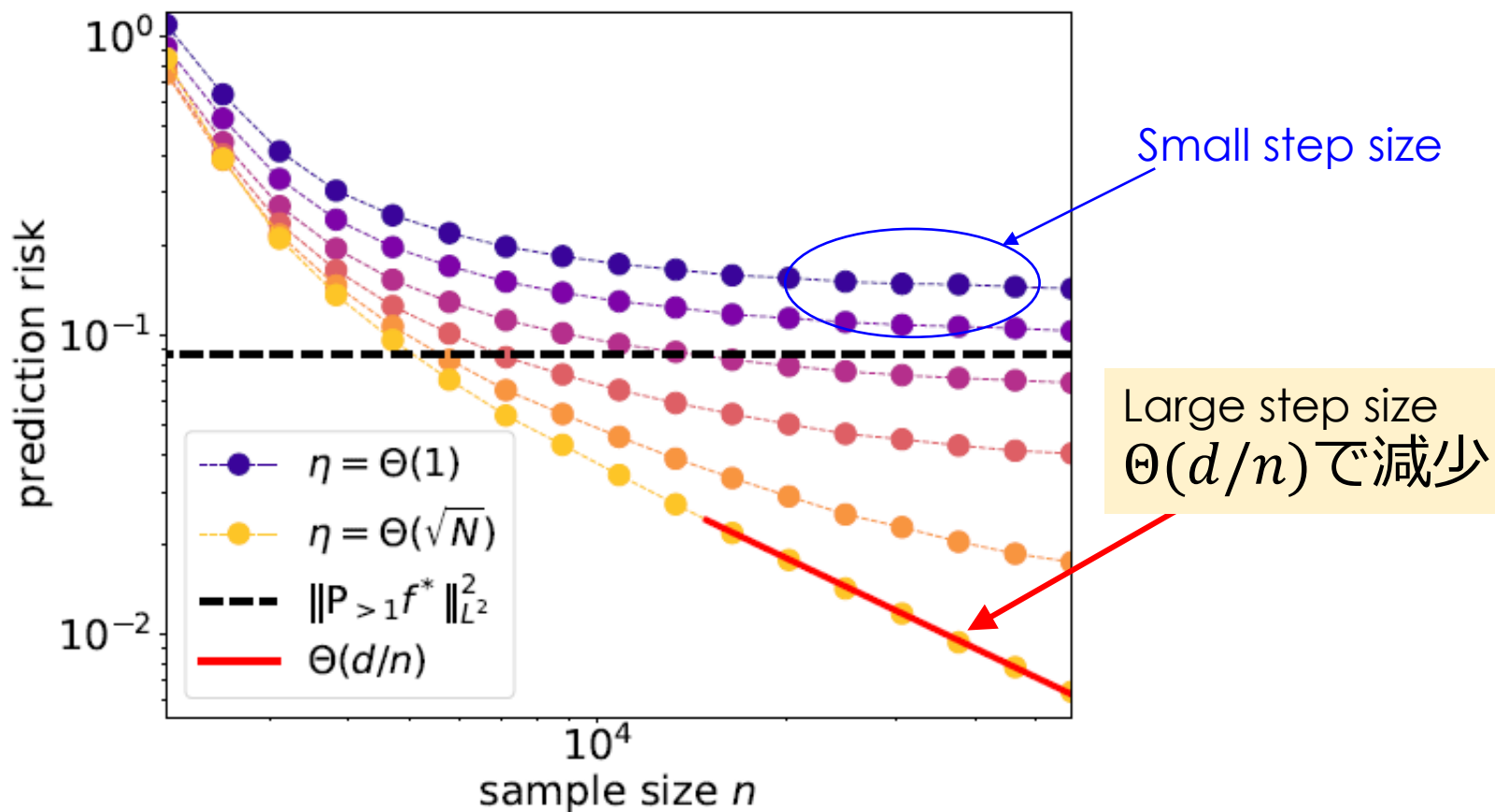
学習率を大きくすることで, 精度を大きく改善できる

($\tau^* \ll \|P_{>1} f^*\|^2$ の状況で). ※バイアス項を振りなおせば $\tau^* = 0$ とできる.

$$\|W_1 - W_0\|_F \geq \|W_0\|_F$$

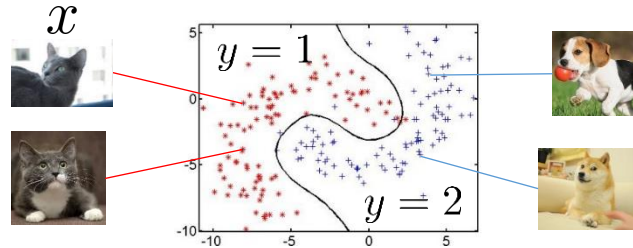
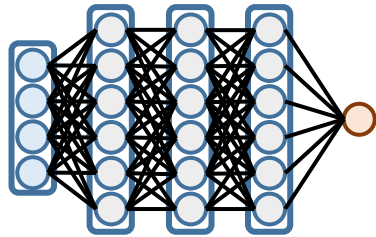
考察

Corollary もし $\sigma = \sigma^* = \text{erf}$ なら, $\tau^* = 0$.
 特に, $R_{W_1}(\lambda) = \Theta(\psi_1^{-1}) = \Theta(d/n)$ である.



Predictive risk of ridge regression on CK obtained by one step GD (empirical simulation, $d = 1024$): brighter color represents larger step size scaled as $\eta = N^\alpha$ for $\alpha \in [0, 1/2]$. We chose $\sigma = \sigma^* = \text{erf}$, $\psi_2 = 2$, $\lambda = 10^{-3}$, and $\sigma_\epsilon = 0.1$.

Learning = Fitting a model to the training data



$$f_{\theta}(x) = W_L \sigma(W_{L-1} \cdots \sigma(W_1 x))$$

Find a model that fits the data well by minimizing a loss function.

$$\min_{\theta} \sum_{i=1}^n \ell(y_i, f_{\theta}(x_i))$$

ML is intended to obtain good **prediction**.

- Lower classification error.
- Better recommendation.

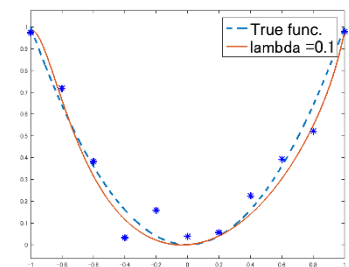
Black-box model

They do not find causality but correlation.

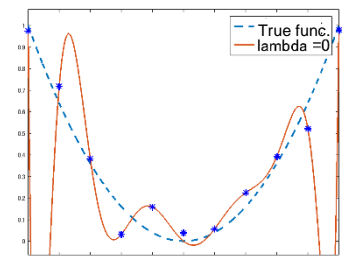
Better prediction \neq Better interpretability/explainability:

- How is a particular answer produced?
- Is the output produced by a correct reasoning?

○ Good pred.



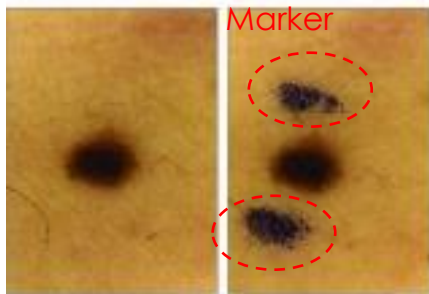
✗ Bad pred. (overfit)



Learning a Shortcut

ML models sometimes find a “**shortcut**” that lowers the training error but yields an unreasonable decision rule.

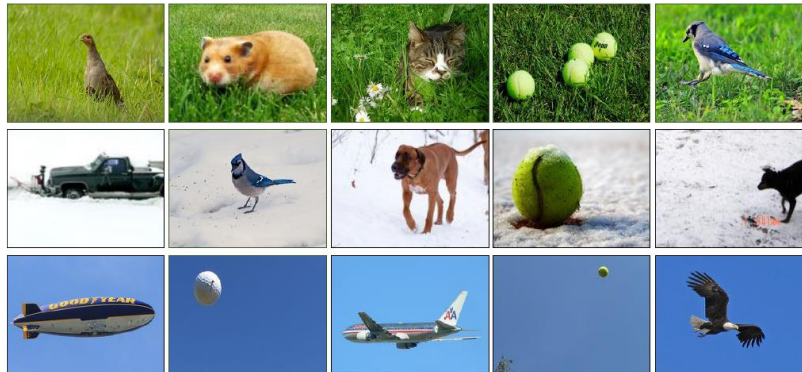
- Medical image recognition:
Marker marked by doctors produces unintentional bias.



[Winkler et al., Association Between Surgical Skin Markings in Dermoscopic Images and Diagnostic Performance of a Deep Learning Convolutional Neural Network for Melanoma Recognition. JAMA dermatology, 155(10):1135–1141, 2019]

Camera, lighting conditions, skin-colors also produce bias.

- General object recognition:



A specific object is linked to the background color.

- ML model recognizes via the BG-color.
- It fails if BG-color changes.

[Nanxuan Zhao, Zhirong Wu, Rynson W. H. Lau, and Stephen Lin. Distilling localization for self-supervised representation learning. arXiv: 2004.06638, 2020]

Interpretability is necessary for **high-risk stakes problems**.

- E.g., cancer diagnosis.

• Hallucination

TA Is there a survey paper of federated learning?

A confident response by an AI that does not seem to be justified by its training data.

Yes, there are several survey papers available on the topic of federated learning. These survey papers provide an overview of the concepts, techniques, and advancements in federated learning, along with a comprehensive analysis of existing research in the field. They often cover a wide range of topics, including privacy preservation, communication efficiency, optimization algorithms, and applications of federated learning.

While I cannot provide the full text of these survey papers, here are a few prominent ones that you can search for and access through academic databases or online research portals:

1. "A Comprehensive Survey on Federated Learning" by Yang Liu, et al. (2020).
2. "Federated Learning: Challenges, Methods, and Future Directions" by Qiang Yang, et al. (2019).

Wrong author information

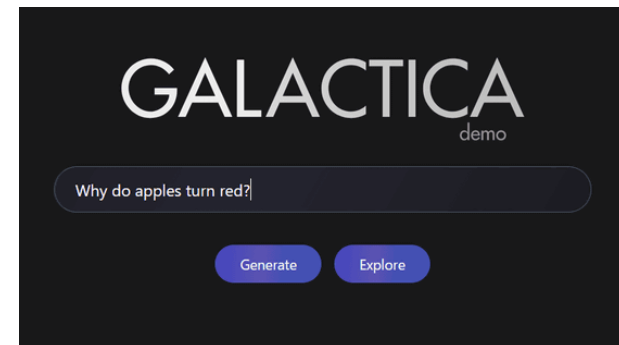
The situation could get worse because it is estimated that 33–46% of crowd workers use LLMs. → AI would be trained by AI generated data.

[Veselovsky et al: Artificial Artificial Artificial Intelligence: Crowd Workers Widely Use Large Language Models for Text Production Tasks. arXiv:2306.07899]

• Ethical issue

Galactica: LLM by Meta trained mainly by scientific papers and documents. (published on 15th/Nov/2022)

- Produced several hallucination including vicious racist responses backed up by non-existing references.
- Was closed to public in 3 days.



[Taylor et al: Galactica: A Large Language Model for Science. 2022]

- Gender bias in natural language processing

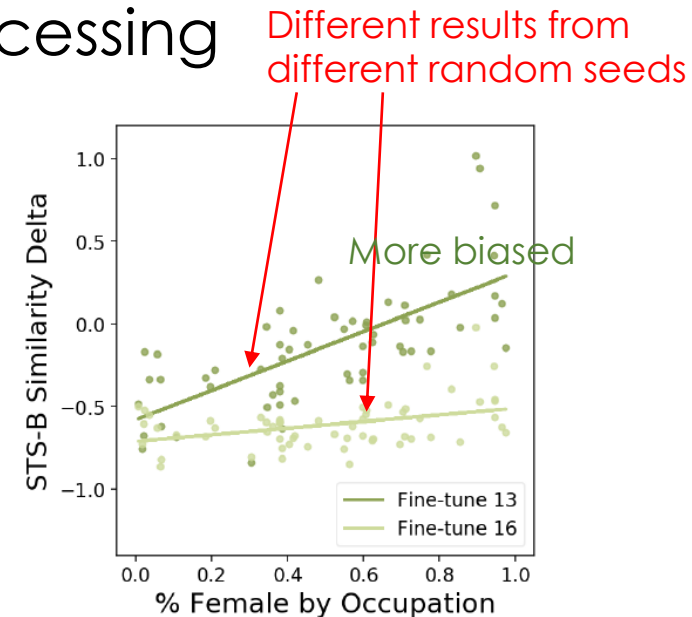
Semantic textual similarity (STS) task

$\text{sim}(\text{"a woman is walking"}, \text{"a doctor is walking"})$
 $- \text{sim}(\text{"a man is walking"}, \text{"a doctor is walking"})$.

Are the results affected by the gender ratio of the occupations?

- Random seed to fine tune the BERT model largely affected the result.
- **It is difficult to detect from the predictive accuracy.**

[D'Amour et al. (Google), Underspecification Presents Challenges for Credibility in Modern Machine Learning. JMLR, 23(226):1-61, 2022]



- Copyright issue

Article 30-4 of the Copyright Law in Japan permits the use of a copyrighted work for training ML models, but it does not include any procedure for gaining permission in advance from copyright holders.

["Copyrighted Works Get Flimsy Protection from AI Under Japanese Law," The Yomiuri Shimbun, published on April 29, 2023]

- Japanese government announced that that necessary measures will be established by summarizing the issues such as cases in which learning of copyrighted works by AI constitutes an infringement. 「知的財産推進計画2023」

G7: Hiroshima AI process

- Agreed to establish **Hiroshima AI process**.
- They are determined to work together and with others to “advance international discussions on inclusive artificial intelligence (AI) governance and interoperability to achieve our common vision and goal of **trustworthy AI**, in line with our shared democratic values.”



[HP of Ministry of Foreign Affairs of Japan.
https://www.mofa.go.jp/ecm/ec/page1e_000673.html]

[Cabinet Office report: G7 Hiroshima Summit (Session 1 (Working Lunch) "Toward an International Society of Cooperation, Not Division and Conflict/World Economy" Summary)]
[内閣府資料: G7広島サミット (セッション1 (ワーキング・ランチ) 「分断と対立ではなく協調の国際社会へ/世界経済」概要)]

The G7 Digital and Technology Ministers' Meeting was held in Takasaki City, Gunma Prefecture, prior to the G7 Summit, and adopted a joint statement that included the promotion of the development of "human-centered and reliable artificial intelligence (AI).

1. **Facilitation of Cross-Border Data Flows and Data Free Flow with Trust.**
2. **Secure and Resilient Digital Infrastructure:**
3. **Internet Governance.**
4. **Innovating the Economic Society and Enhancing Digital Skills.**
5. **Promoting Responsible AI and AI Governance**
6. **Competition Policy on Digital Market.**

1. reinforce democratic values
2. respect for human rights and fundamental freedoms
3. collective efforts to promote interoperability between AI governance frameworks

まとめ

- 深層学習はなぜうまくいくのか？ [世界的課題]
- 数学による深層学習の原理究明
 - 「表現能力」, 「汎化能力」, 「最適化」

学習

スパース推定

カーネル法

テンソル分解

特徴抽出

深層学習の理論

Besov空間

関数近似理論

確率集中不等式

数学

連続方程式

Wasserstein幾何

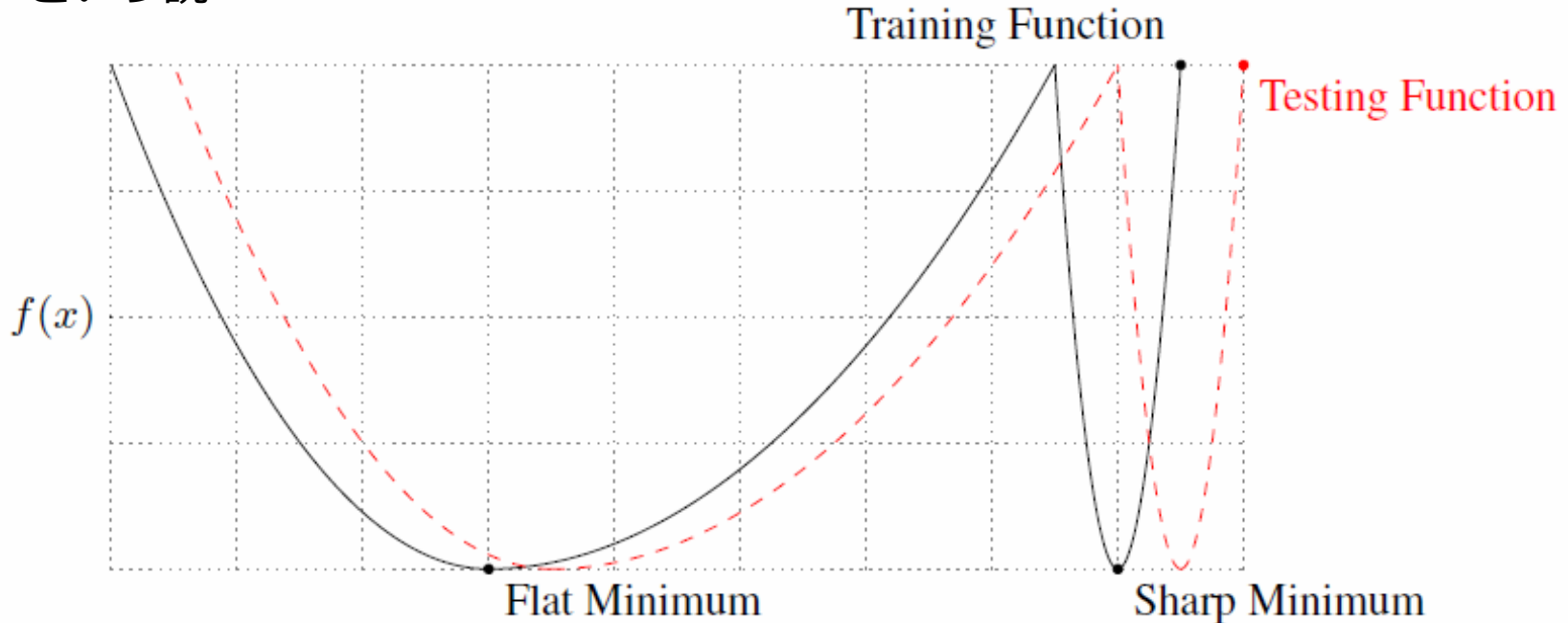
確率過程

理論により深層学習を“謎の技術”から“制御可能な技術”へ
深層学習を超える方法論の構築へ

ノイズあり勾配法と大域的最適性 (参考資料：最終日に説明予定)

Sharp minima vs flat minima

SGDは「フラットな局所最適解」に落ちやすい→良い汎化性能を示す
という説



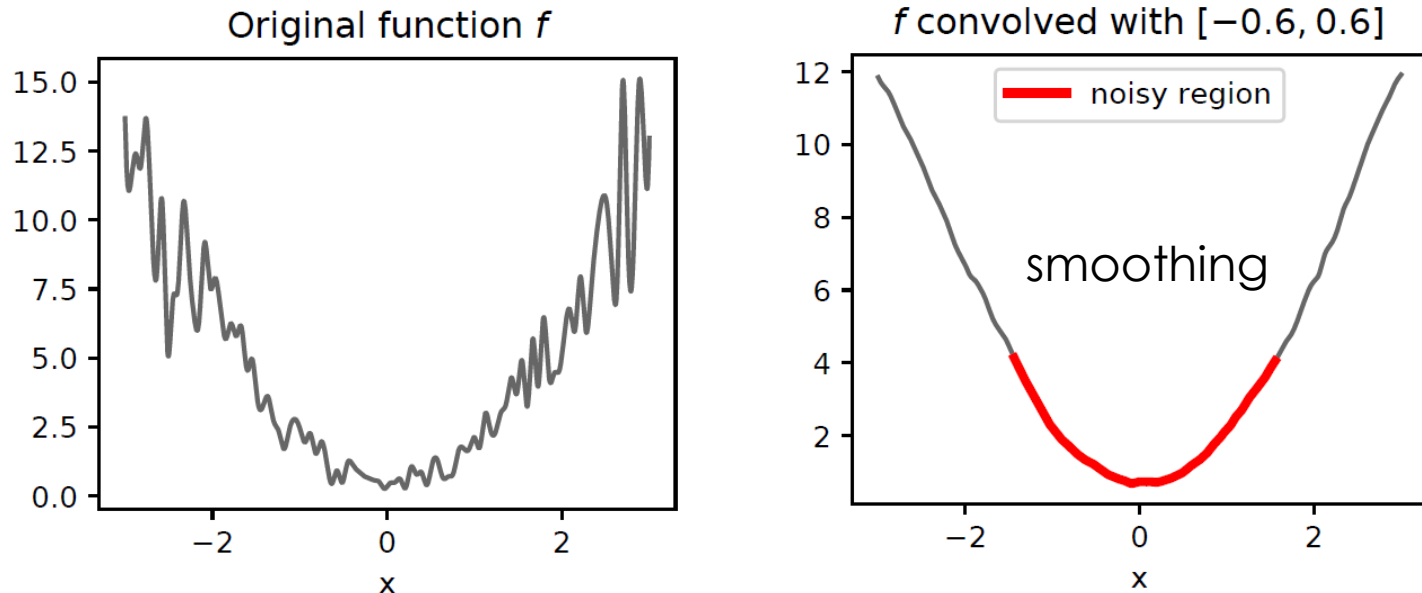
Keskar, Mudigere, Nocedal, Smelyanskiy, Tang (2017):
On large-batch training for deep learning: generalization gap and sharp minima.

$$\theta_t = \theta_{t-1} - \alpha_b \underbrace{\left(\frac{1}{b} \sum_{j=1}^b \nabla_{\theta} \ell(z_{i_j}; \theta) \right)}_{\cong \text{正規分布}}$$

→ ランダムウォークはフラットな領域にとどまりやすい

- 「フラット」という概念は座標系の取り方によるから意味がないという批判.
(Dinh et al., 2017)
- PAC-Bayesによる解析 (Dziugaite, Roy, 2017)

ノイズによる平滑化効果



[Kleinberg, Li, and Yuan, ICML2018]

確率的勾配を用いる \Rightarrow 解にノイズを乗せている \Rightarrow 目的関数の平滑化

$$x_t = x_{t-1} - \eta(\nabla L(x_{t-1}) + \xi_t) \quad (y_t = x_t + \eta\xi_t)$$

$$\Rightarrow y_t = y_{t-1} - \eta\xi_{t-1} - \eta\nabla L(y_{t-1} - \eta\xi_{t-1})$$

$$\Rightarrow \mathbb{E}_{\xi_{t-1}}[y_t] = y_{t-1} - \eta\nabla \mathbb{E}_{\xi_{t-1}}[L(y_{t-1} - \eta\xi_{t-1})]$$

ノイズを加えて平滑化した目的関数 $\bar{L}(y_t) = \mathbb{E}_{\xi_t}[L(y_t - \eta\xi_t)]$ を最適化.

- Graduated non-convexity

Blake and Zisserman: *Visual reconstruction*, volume 2. MIT press Cambridge, 1987.

- Gaussian kernelとの畳み込み

Z. Wu. The effective energy transformation scheme as a special continuation approach to global optimization with application to molecular conformation. *SIAM Journal on Optimization*, 6(3):748-768, 1996.

- Graduated optimization

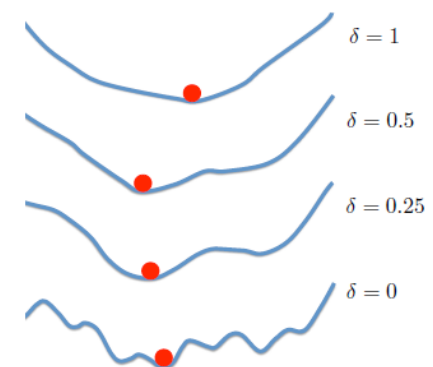
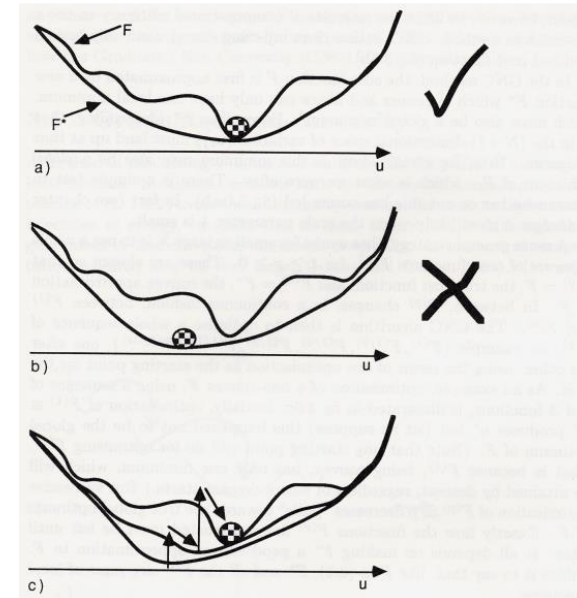
Hazan, Levy, and Shalev-Shwartz: On graduated optimization for stochastic non-convex problems. *International conference on machine learning*, pp. 1833-1841, 2016.

σ -nice性の導入. 多項式オーダーでの収束.

$$\hat{L}_\delta(x) = E_{u \sim U(B(\mathbb{R}^d))} [L(x + \delta u)]$$

Survey:

Mobahi and Fisher III. On the link between gaussian homotopy continuation and convex envelopes. *Energy Minimization Methods in Computer Vision and Pattern Recognition*, pp. 43-56, 2015.



• 確率的勾配ランジュバン動力学

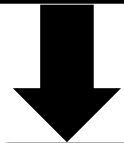
Stochastic Gradient Langevin Dynamics (SGLD)

$$\min_{x \in \mathbb{R}^d} L(x) = \min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell_i(x) \quad (\text{非凸})$$

β : 逆温度パラメータ

$$dX_t = -\nabla L(X_t)dt + \sqrt{2\beta^{-1}}dB_t \quad (\text{勾配ランジュバン動力学})$$

$$\text{定常分布: } \pi \propto \exp(-\beta L(X))$$

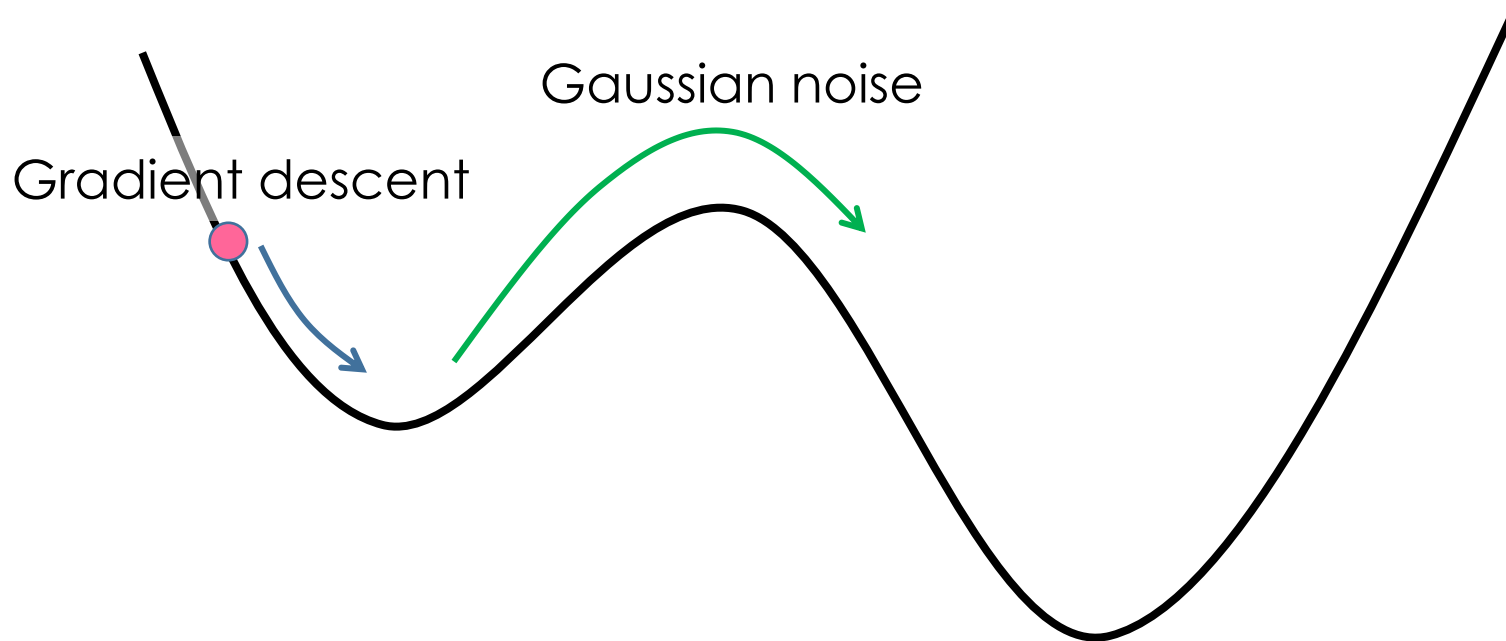


時間離散化

[Gelfand and Mitter (1991); Borkar and Mitter (1999); Welling and Teh (2011)]

GLD: $X_{t+1} = X_t - \eta \nabla L(X_t) + \sqrt{2\eta\beta^{-1}}\xi_t$ (Euler-Maruyama scheme)
 $\xi_t \sim N(0, I)$

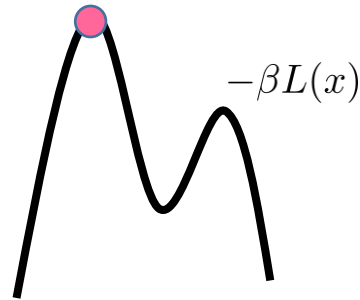
SGLD: $X_{t+1} = X_t - \eta \frac{1}{|I_B|} \sum_{i \in I_B} \nabla \ell_i(X_t) + \sqrt{2\eta\beta^{-1}}\xi_t$
確率的勾配



$$dX_t = -\nabla L(X_t)dt + \sqrt{2\beta^{-1}}dB_t$$

適当な条件のもとGLDの定常分布は以下で与えられる:

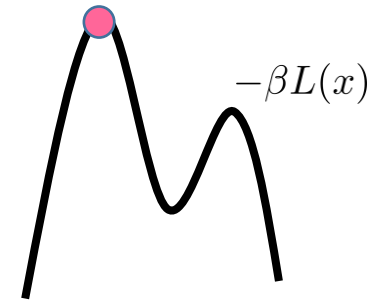
$$\pi_\infty(dx) \propto \exp(-\beta L(x))dx$$



- β が大きければ, 定常分布は最適解周りに集中する.
- 定常分布からのサンプリングにも使える (というかこっちが本来の目的).

定常分布

$$\pi_\infty(dx) \propto \exp(-\beta L(x)) dx$$



対数ソボレフ不等式 (π_∞ の性質):

- 例:
- 二次関数+有界関数
 - Weak Morse型関数

任意の(π_∞ に対して絶対連続な)確率分布 $d\nu = f d\pi_\infty$ に対し,

$$\int f \log(f) d\pi_\infty \leq 2c_{LS} \int \frac{\|\nabla f\|^2}{f} d\pi_\infty$$

$$(D(\nu||\pi_\infty) \leq 2c_{LS} I(\nu||\pi_\infty))$$

KL-div $D(\mu||\nu) = \int \log\left(\frac{d\mu}{d\nu}\right) d\mu$, Fisher-div $I(\mu||\nu) = \int \left\| \nabla \log \frac{d\mu}{d\nu} \right\|^2 d\mu$

➡ 幾何的エルゴード性

ρ_t : X_t の周辺分布

$$D(\rho_t||\pi_\infty) \leq \exp(-2t/c_{LS}) D(\rho_0||\pi_\infty)$$

定常分布へKL-divergenceの意味で指数オーダーの収束

$$L(x) = \frac{1}{n} \sum_{i=1}^n \ell_i(x) + \lambda_1 \|x\|^2$$

Bounded perturbation lemma:

$$|\ell_i(x)| \leq B \ (\forall x) \quad \longrightarrow \quad c_{\text{LS}} \leq \frac{1}{2\lambda_1\beta} \exp(4\beta B)$$

[R. Holley and D. Stroock. Logarithmic sobolev inequalities and stochastic Ising models. Journal of statistical physics, 46(5-6):1159–1194, 1987.]

強凸な場合: $L(x)$ が μ -強凸 $\Rightarrow c_{\text{LS}} \leq 1/(\mu\beta)$

[Bakry and Émery, 1985]

• **散逸的 (dissipative):**

$$\exists m > 0, b \geq 0, \langle x, \nabla L(x) \rangle \geq m\|x\|^2 - b$$

• **平滑性:**

$$\exists M, \|\nabla L(x) - \nabla L(y)\| \leq M\|x - y\|$$



$$\longrightarrow \quad c_{\text{LS}} \leq \frac{2m^2 + 8M^2}{m^2 M \beta} + \left(\frac{6M(d + \beta) + 2}{m} \right) e^{O(\beta+d)}$$

[Raginsky, Rakhlin and Telgarsky, 2017]

過程: L は M -平滑: $\exists M, \|\nabla L(x) - \nabla L(y)\| \leq M\|x - y\|$

定理

[Vempala and Wibisono, 2019]

ν_k : Marginal distribution of X_k (discrete time dynamics)

$$D(\nu_k || \pi_\infty) \lesssim \exp(-k\eta/c_{LS})D(\nu_0 || \pi_\infty) + 8c_{LS}dM^2\eta$$

定理 (informal)

散逸性と平滑性の条件のもと (and other technical condition),

$$E[L(X_k)] - L(X^*) \lesssim \exp(-ck\eta/c_{LS}) + c_{c_{LS},\beta,d}\eta + \frac{d \log(\beta + 1)}{\beta}$$

幾何的エルゴード性

時間離散化の誤差

$E_{\pi_\infty}[L(X)] - L(X^*)$

定常分布が最適解まわりにどれだけ集中しているか

where $c, c_{c_{LS},\beta,d} > 0$ are constants.

[Raginsky, Rakhlin and Telgarsky, 2017; Xu, Chen, Zou, and Gu, 2018; Erdogdu, Mackey and Shamir, 2018]

- 逆温度パラメータ β が十分大きければ、目的関数が非凸でも最適解の近くに到達できる。
- ただし、一般には対数ソボレフ不等式は β に指数的に依存することに注意。
(そうでない場合もある: 強凸目的関数, Weak Morse関数)

- **Finite dimensional Langevin dynamics:**

- [Convergence in low \(convex case\)](#): Dalalyan and Tsybakov, 2012; Dalalyan, 2016; Durmus and Moulines, 2015, ..
- [Non-convex Optimization](#): Raginsky et al., 2017; Xu et al., 2018; Erdogdu, Mackey and Shamir, 2018
- [Log-Sobolev inequality](#): Vempala and Wibisono, 2019.

- **Infinite dimensional Langevin dynamics:**

- Continuous time:
 - [Existence & Uniqueness of invariant measure](#): Da Prato and Zabczyk, 1992; Maslowski, 1989; Sowers, 1992.
 - [Geometric ergodicity](#): Jacquot and Royer, 1995; Shardlow, 1999; Hairer, 2002, Its explicit rate: Goldys and Maslowski, 2006.
- Discrete time:
 - [Weak approximation rate of discretized scheme](#): Hausenblas, 2003; Debussche, 2011; Bréhier, 2014; Bréhier and Kopec 2016.

Other topics (MCMC in Hilbert space):

- [preconditioned Crank–Nicolson \(pCN\)](#): Hairer et al., 2014; Eberle, 2014; Vollmer, 2015; Rudolf and Sprungk, 2018.
- [Metropolis-Adjusted Langevin Algorithm \(MALA\)](#): Durmus and Moulines, 2015; Beskos et al., 2017.

GLDによるNNの最適化

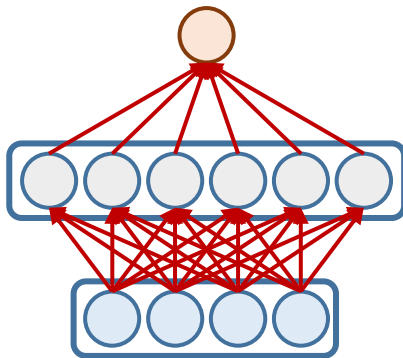
観測モデル:

$$y_i = f^\circ(x_i) + \xi_i \quad (i = 1, \dots, n)$$

where $x_i \sim \text{Unif}(\mathbb{S}^{d-1})$, $\xi_i \sim (\text{mean } 0, \text{variance } \sigma^2, \text{bounded})$

教師生徒設定 with **ReLU活性化関数**:

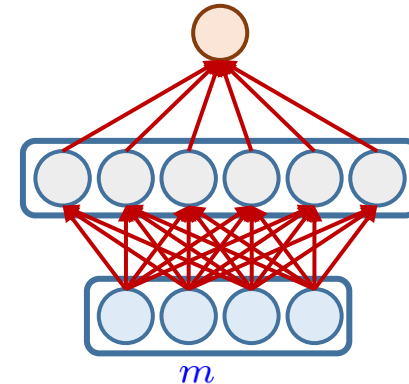
Teacher



$$f^\circ(x) = \sum_{j=1}^m a_j^\circ \sigma(\langle w_j^\circ, x \rangle)$$

$\sigma : \text{ReLU}$

Student (overparameterization)



$$f_\Theta(x) = \sum_{j=1}^m a_j \sigma(\langle w_j, x \rangle)$$

- 教師と生徒で同じサイズとする. (おそらく緩和可能)
- 生徒は教師を多項式時間で推定できるか？

L2-正則化あり経験損失関数:

$$\hat{\mathcal{R}}_\lambda(\Theta) = \frac{1}{n} \sum_{i=1}^n \left(y_i - \sum_{j=1}^M a_j \sigma(\langle w_j, x \rangle) \right)^2 + \lambda \sum_{j=1}^M (a_j^2 + \|w_j\|^2)$$

L2-正則化

$$\sum_{j=1}^M |a_j| \|w_j\| \leq \frac{1}{2} \sum_{j=1}^M (a_j^2 + \|w_j\|^2) \quad \leftarrow \text{ReLUではL2-正則化はスパース正則化 (L1-正則化) でもある。}$$

二段階最適化

(1) 探索フェーズ: GLD

$$\Theta^{(k+1)} = \Theta^{(k)} - \eta^{(1)} \nabla \hat{\mathcal{R}}_\lambda(\Theta^{(k)}) + \sqrt{\frac{2\eta^{(1)}}{\beta}} \zeta^{(k)}$$

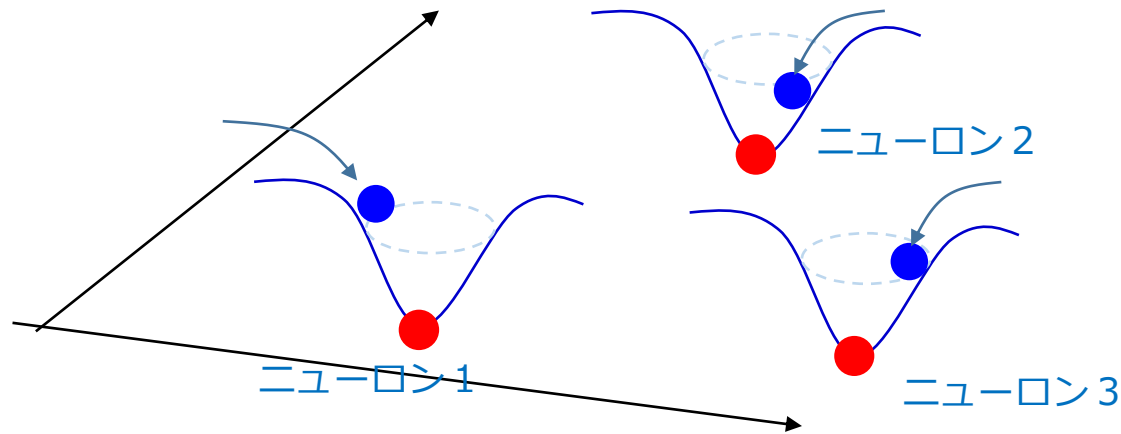
ガウスノイズ

(2) 収束フェーズ: ノイズ無し・正則化なしの勾配法

$$\Theta^{(k+1)} = \Theta^{(k)} - \eta^{(1)} \nabla \hat{\mathcal{R}}_0(\Theta^{(k)})$$

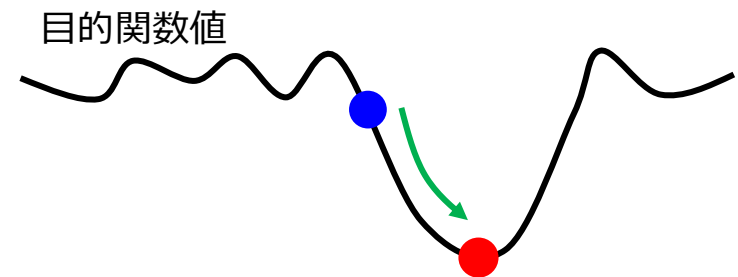
探索フェーズ:

最適解近傍への収束



収束フェーズ:

最適解まわりでは大体強凸
→ 普通の勾配法で線形収束



二段階の収束フェーズ：実験的にも観測されている (Hidden progress)

収束解析

[Akiyama and Suzuki: Excess Risk of Two-Layer ReLU Neural Networks in Teacher-Student Settings and its Superiority to Kernel Methods. arXiv:2205.14818]

- $\hat{\mathcal{R}}_\lambda(\Theta) = \frac{1}{n} \sum_{i=1}^n \left(y_i - \sum_{j=1}^M a_j \sigma(\langle w_j, x \rangle) \right)^2 + \lambda \sum_{j=1}^M (a_j^2 + \|w_j\|^2)$
- $\mathcal{R}(\Theta) = \mathbb{E}[(f^\circ(X) - f_\Theta(X))^2]$
- $W^\circ = (w_1^\circ, \dots, w_m^\circ)$ の m 番目の特異値は σ_{\min} で下から抑えられる。

定理 (informal)

(m, d, σ_{\min}) 固定の元で) 適当な定数時間 K_1 が存在して,

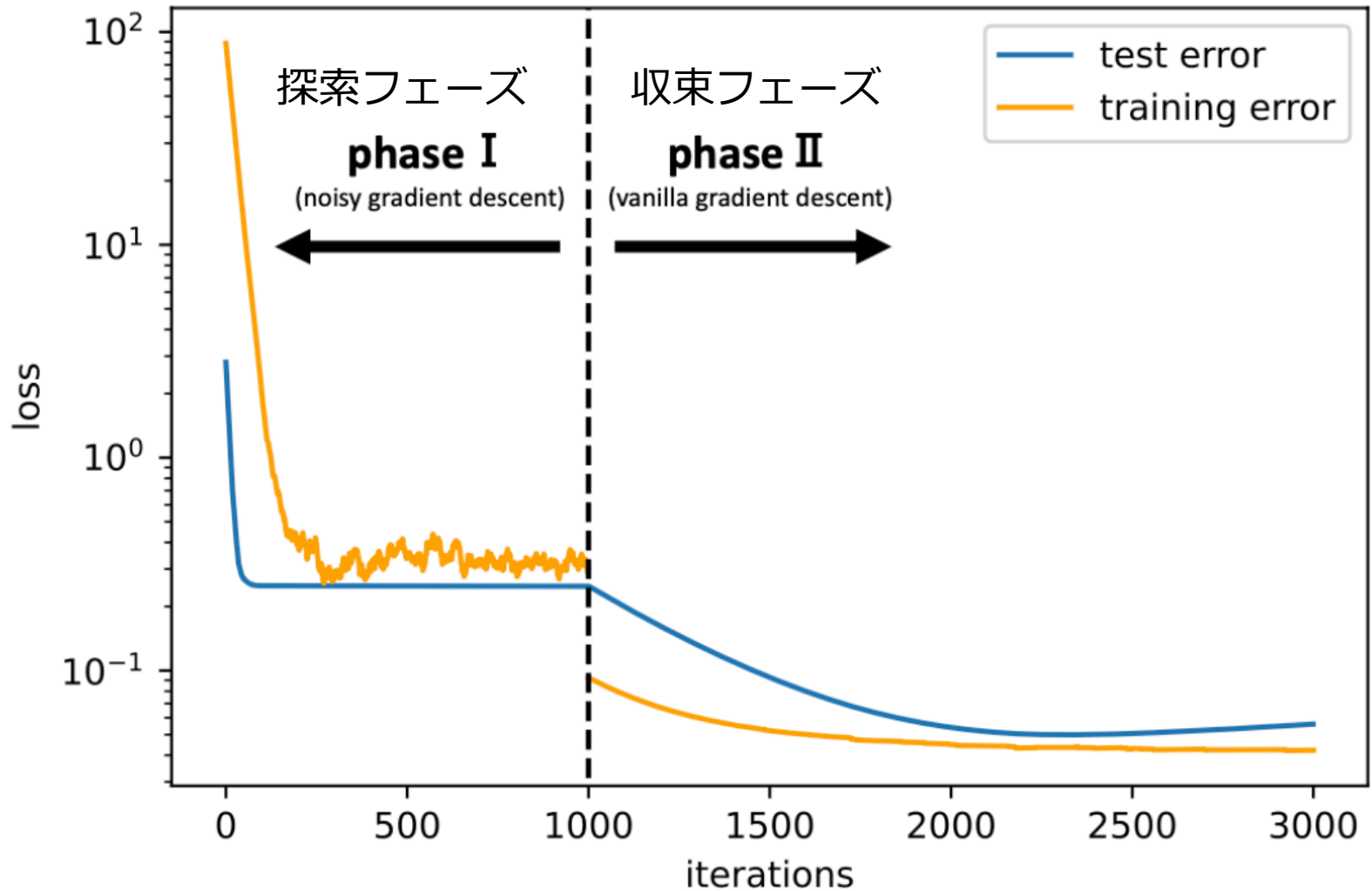
Phase 1 (大域的探索フェーズ): しばらく探索すると期待損失をある閾値以下まで減らせる (真に近くなる)

$$R(\Theta^{(K_1)}) \leq \epsilon_0$$

Phase 2 (収束フェーズ): 大域的最適解へ線形収束

$$\hat{R}_0(\Theta^{(k)}) - \hat{R}_0(\Theta^*) \leq c_1 \exp(-c_2(k - K_1)).$$

数値実験：二段階の学習ダイナミクス



予測誤差の比較

- 前ページのアルゴリズムで最適化して得られた解 Θ^k

$$\|f_{\Theta^k} - f^\circ\|_{L^2(P_X)}^2 \lesssim \frac{\sigma_{\min}^{-4} m^5 \log(n)}{n}$$

- 線形推定量の予測誤差の下限:

$$R_{\text{lin}} \gtrsim n^{-\frac{d+2}{2d+2}}$$

For $d = 2$: $n^{-2/3}$

For large d : $n^{-1/2}$

線形推定量は次元の呪いを受ける
→ 特徴量も学習することで改善

Related work

Non-overparameterized setting:

- [Li & Yuan, 2017] showed global convergence under $M = m$ and a special network structure (ResNet like structure).
- [Zhong et al., 2017] showed local convergence under $M = m$, i.e., they showed convergence when the initial solution is close to the true parameter.

Overparameterized setting:

- [Li, Ma & Zhang, 2020] showed global convergence of GD for an overparameterized setting $M > m$.

$$f^\circ(x) = \sum_{j=1}^m |\langle w_j, x \rangle| \quad \mathcal{L}(\hat{f}) - \mathcal{L}(f^\circ) = O(d^{-(1+Q)})$$

where Q is a small constant.

True network

- Tensor decomposition technique is used.
 - The true network has a special structure.
 - The convergence is not exactly shown (it converges as $d \rightarrow \infty$).
- [Chizat, 2019]: Convergence to sparse solution with sparse reg.
 - BLASSO [De Castro & Gamboa, 2012]

2-homogeneous activation + NDSC condition

ReLU

Guarantee ← [Akiyama&Suzuki, 2021]

2層NNの学習・粒子勾配法 (参考資料：最終日に説明予定)

- 前のスライドでは横幅が小さいNNの最適化を考えた.
- 実際は過剰パラメータ化した横幅が広いNNを使うことが多い (表現力を上げるため).
- 先のGLDの収束レートには次元が現れるので, そのままでは使えない.
(実際, 次元に指数関数的に依存する)

→ 平均場ランジュバン動力学

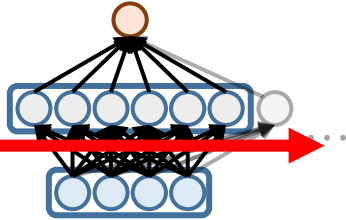
2層NNのGLDによる最適化

$$f_{\Theta}(x) = \frac{1}{M} \sum_{j=1}^M h_{\theta_j}(x) \quad \text{例: } h_{\theta}(x) = r\sigma(w^{\top}x) \text{ for } \theta = (r, w)$$

$$d\theta_{j,t} = -\nabla_{\theta_j} L(f_{\Theta_t}) dt + \sqrt{2/\beta} dB_t$$

ニューロンが沢山あると普通のGLDの理論が適用できない。
 しかし、平均場ランジュバン動力学の理論により理論保証ができる。
 (逆にニューロン数無限大の極限を考えると理論保証可能になる)

多粒子化 (平均場) :

$$f_{\Theta}(x) = \frac{1}{M} \sum_{j=1}^M h_{\theta_j}(x) \xrightarrow{M \rightarrow \infty} f_{\mu}(x) = \int h_{\theta}(x) d\mu(\theta)$$


定理 (Hu, Ren, Šiška, and Szpruch, 2021; Mei, Montanari, and Ngyue, 2018)

$M \rightarrow \infty, t \rightarrow \infty$ の極限で粒子 θ_j の分布 μ_t は以下の分布に収束:

$$\mu_{\infty} = \arg \min_{\mu \in \mathcal{P}} L(f_{\mu}) + \frac{1}{\beta} \text{Ent}(\mu)$$

エントロピー
 $(\text{Ent}(\mu) = \int \log(\mu) d\mu)$

重要 : 分布 μ に対しては凸関数 ! (if 損失が凸)

MF-LDの収束

目的関数 $F(\mu) = \frac{1}{n} \sum_{i=1}^n \ell_i(f_\mu) + \lambda_1 \mathbb{E}_\mu[\|x\|^2]$ (損失関数+正則化)

$$\frac{\delta F(\mu)}{\delta \mu}(x) = \frac{1}{n} \sum_{i=1}^n \ell'_i(f_\mu) h_x(z_i) + \lambda_1 \|x\|^2 \quad f_\mu(x) = \int h_\theta(x) d\mu(\theta)$$

平均場ランジュバン動力学:

$$dX_t = -\nabla \frac{\delta F(\mu_t)}{\delta \mu}(X_t) dt + \sqrt{2\lambda_2} dB_t \quad \mu_t = \text{Law}(X_t)$$

(わかりにくいが単純に各ニューロンを勾配法で動かして微小ノイズを加えていることに対応)

近接点更新解: (c.f., Mirror descent, exponentiated gradient)

$$p_\mu(x) \propto \exp\left(-\frac{1}{\lambda_2} \frac{\delta F(\mu)}{\delta \mu}(x)\right) \quad p_\mu = \arg \min_{\nu \in \mathcal{P}} (\nu - \mu) \frac{\delta F(\mu)}{\delta \mu} + \lambda_2 \text{Ent}(\nu)$$

損失を線形化して得られる解

定理 (Entropy sandwich) [Nitanda, Wu, Suzuki (AISTATS2022)][Chizat (2022)]

p_{μ_t} が一様に対数ソボレフ不等式 (定数 α)を満たすとすると,

$$\mathcal{L}(\mu_t) - \mathcal{L}(\mu^*) \leq \exp(-2\alpha\lambda_2 t)(\mathcal{L}(\mu_0) - \mathcal{L}(\mu^*)). \quad (\text{線形収束!})$$

ただし, $\mathcal{L}(\mu) = F(\mu) + \lambda_2 \text{Ent}(\mu)$

応用例

平均場二層ニューラルネットワーク

$$f(x) = \frac{1}{M} \sum_{j=1}^M r_j \sigma(w_j^\top x) \xrightarrow{M \rightarrow \infty} f_\mu(x) = \int r \sigma(w^\top x) d\mu(r, w)$$

$$f_\mu(z) = \int h_\theta(z) d\mu(\theta) \quad \text{Example } h_\theta(x) = r \sigma(w^\top x) \text{ for } \theta = (r, w)$$

目的関数 $F(\mu) = \frac{1}{n} \sum_{i=1}^n \ell_i(f_\mu) + \lambda_1 \mathbb{E}_\mu[\|x\|^2]$

$\xrightarrow{\mu \text{ で微分}}$

$$\frac{\delta F(\mu)}{\delta \mu}(x) = \frac{1}{n} \sum_{i=1}^n \ell'_i(f_\mu) h_x(z_i) + \lambda_1 \|x\|^2$$

MMD最小化によるノンパラメトリック密度推定

$$F(\mu) = \text{MMD}^2(g * \mu, \hat{\mu}_n) + \lambda_1 \mathbb{E}_\mu[\|x\|^2]$$

➤ $\text{MMD}^2(\nu_1, \nu_2) := \int k(x, x') d\nu_1(x) d\nu_1(x') - 2 \int k(x, x') d\nu_1(x) d\nu_2(x') + \int k(x, x') d\nu_2(x) d\nu_2(x')$

k : 正定値カーネル

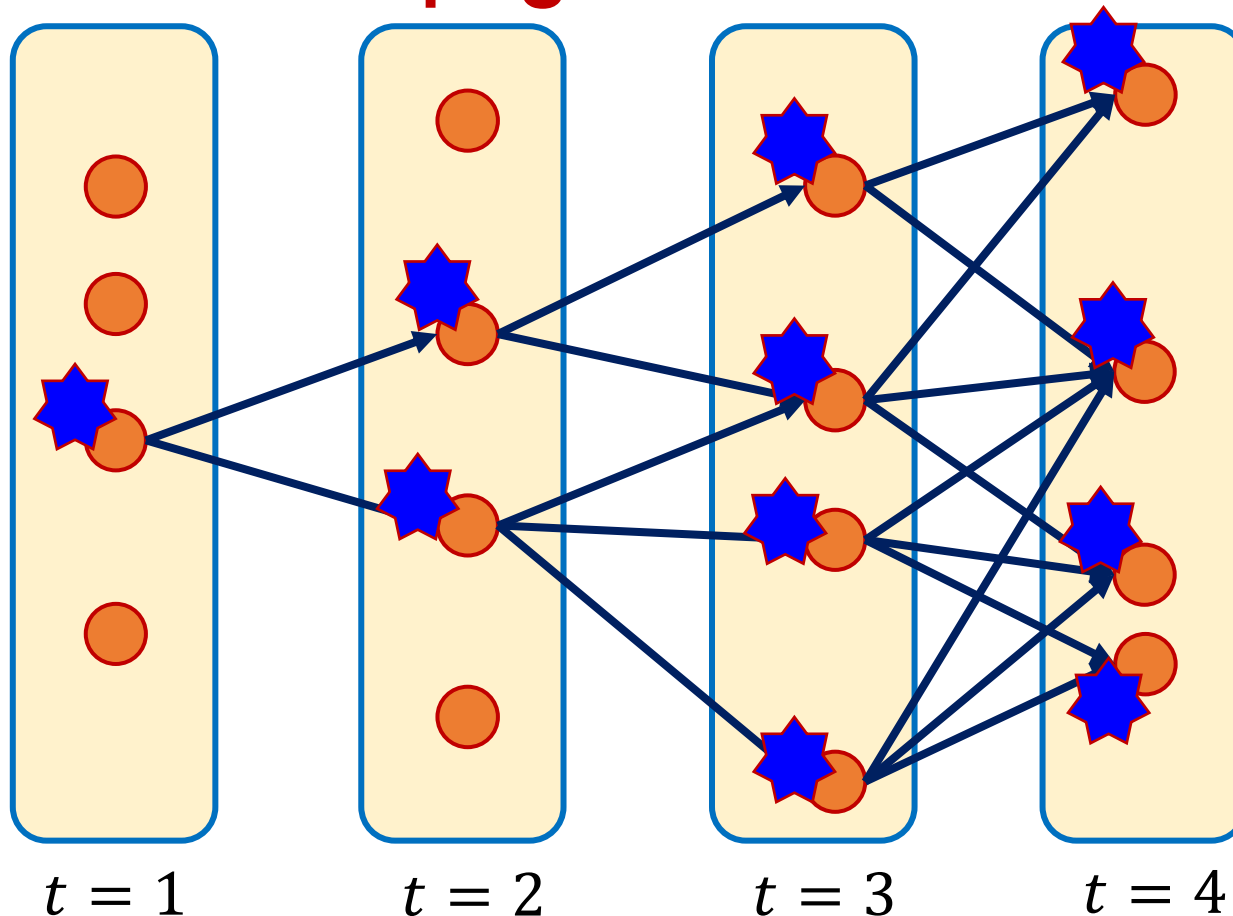
➤ $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$: 経験分布 (訓練データ) ➤ $g(x) = \frac{1}{\sqrt{(2\pi\sigma^2)^d}} \exp\left(-\frac{\|x\|^2}{2\sigma^2}\right)$

ベイズ事後分布の変分推論

難しさ: McKean-Vlasov過程

- 粒子間相互作用のある確率微分方程式はMcKean-Vlasov過程として知られている。 (McKean, Kac,..., 60年代)
- 離散時間・有限粒子での収束を示す際にはPropagation of chaosの評価が難しい。 (粒子を増やすことでそれぞれがあたかも独立に振る舞う現象)

Propagation of chaos



一つの粒子の微小な変化が他の粒子に伝播して増幅される可能性がある。

平均場NNの線形収束
連続時間・無限粒子

[Nitanda, Wu, Suzuki
(AISTATS2022)]
[Chizat (2022)]

時間・空間離散化：「二重ループの手法」

- PDA法 [Nitanda, Wu, Suzuki: NeurIPS2021]
- P-SDCA法 [Oko, Suzuki, Wu, Nitanda: ICLR2022]
- 無限次元拡張 [Nishikawa, Suzuki, Nitanda: NeurIPS2022]

難しい：Propagation of chaos
(McKean, Kac, ..., 60年代より)

空間離散化・連続時間：
Uniform-in-time propagation of chaos

- Super対数Sobolev不等式
[Suzuki, Nitanda, Wu (ICLR2023)]
- Leave-one-out型評価
[Chen, Ren, Wang (arXiv2022)]

時間・空間離散化・確率的勾配：
「一重ループの手法」

Suzuki, Wu, Nitanda
(arXiv:2306.07221)

粒子双対平均化法

(Particle Dual Averaging; PDA)

[Nitanda, Wu, Suzuki: NeurIPS2021]

$$\min_{q:\text{prob.density}} \frac{1}{n} \sum_{i=1}^n \ell(\mathbb{E}_q[h_\theta(x_i)], y_i) + \lambda_1 \mathbb{E}_q[\|\theta\|^2] + \lambda_2 \mathbb{E}_q[\log(q)]$$

近似

↓
 q に関する線形汎関数で近似 (勾配を用いる)
 $\mathbb{E}_{\theta \sim q}[\bar{g}^{(t)}(\theta)]$ (線形近似; $\bar{g}^{(t)}$ は基本的に勾配)

$\bar{g}^{(t)}$ の決定に双対平均化法のルールを用いる

$$\min_{q:\text{prob.density}} \mathbb{E}_{\theta \sim q}[\bar{g}^{(t)}(\theta)] + \lambda_2 \mathbb{E}_q[\log(q)]$$

↳ 解: $q^{(t+1)}(\theta) \propto \exp(-\bar{g}^{(t)}(\theta)/\lambda_2)$ 具体形が得られる.

→ この分布からは以下の勾配ランジュバン動力学を用いてサンプリング可能:

$$d\theta_t = -\nabla(\bar{g}^{(t)}(\theta)/\lambda_2)dt + \sqrt{2}d\xi_t.$$

時間離散化 → $\theta_k = \theta_{k-1} - \eta \nabla \bar{g}^{(t)}(\theta)/\lambda_2 + \sqrt{2\eta} \xi_{k-1}$

計算量解析:

1. 外側ループ: $\mathcal{L}(\hat{q}^{(t)}) - \mathcal{L}(q^*) \leq O(1/t)$
 2. 内側ループ: $T_t = \tilde{O}(t^2 \exp(8/\lambda_2)/(\lambda_1 \lambda_2))$ (GLDによる)
- ⇒ 合計: $O(\epsilon^{-3})$ の勾配アップデートで十分.

➤ 初の多項式オーダー最適化手法

粒子確率的双対座標上昇法

(Particle Stochastic Dual Coordinate Ascent; P-SDCA)

[Oko, Suzuki, Wu, Nitanda: ICLR2022]

主問題

$$\min_p P(p) = \frac{1}{n} \sum_{i=1}^n \ell_i \left(\int p(\theta) h_i(\theta) \right) + \lambda_1 \int \|\theta\|^2 p(\theta) d\theta + \lambda_2 \int p(\theta) \log(p(\theta)) d\theta$$

|| by Fenchelの双対定理

双対問題

$$\ell_i^*(g) := \sup_{u \in \mathbb{R}} \{ug - \ell_i(u)\}$$

$$- \min_{g \in \mathbb{R}^n} D(g) = \frac{1}{n} \sum_{i=1}^n \ell_i^*(g_i) + \lambda_2 \log \left(\int q[g](\theta) d\theta \right)$$

ただし $q[g](\theta) := \exp \left\{ -\frac{1}{\lambda_2} \left(\frac{1}{n} \sum_{i=1}^n h_i(\theta) g_i + \lambda_1 \|\theta\|^2 \right) \right\}$

- 双対変数の座標をランダムに選択し, その座標に関して最適化.
 → 確率的双対座標上昇法

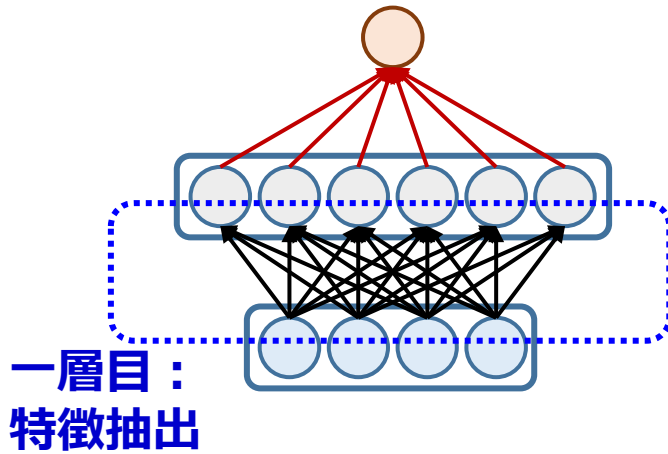
計算量解析:

双対ギャップ ϵ_P を達成するのに必要な外側ループ数:

$$t_{\text{end}} = 2 \left(n + \frac{1}{\lambda_2 \gamma} \right) \log \left(\frac{nC}{\epsilon_P} \right)$$

- 指数オーダーでの収束を達成
- サンプルサイズ n への依存を緩和

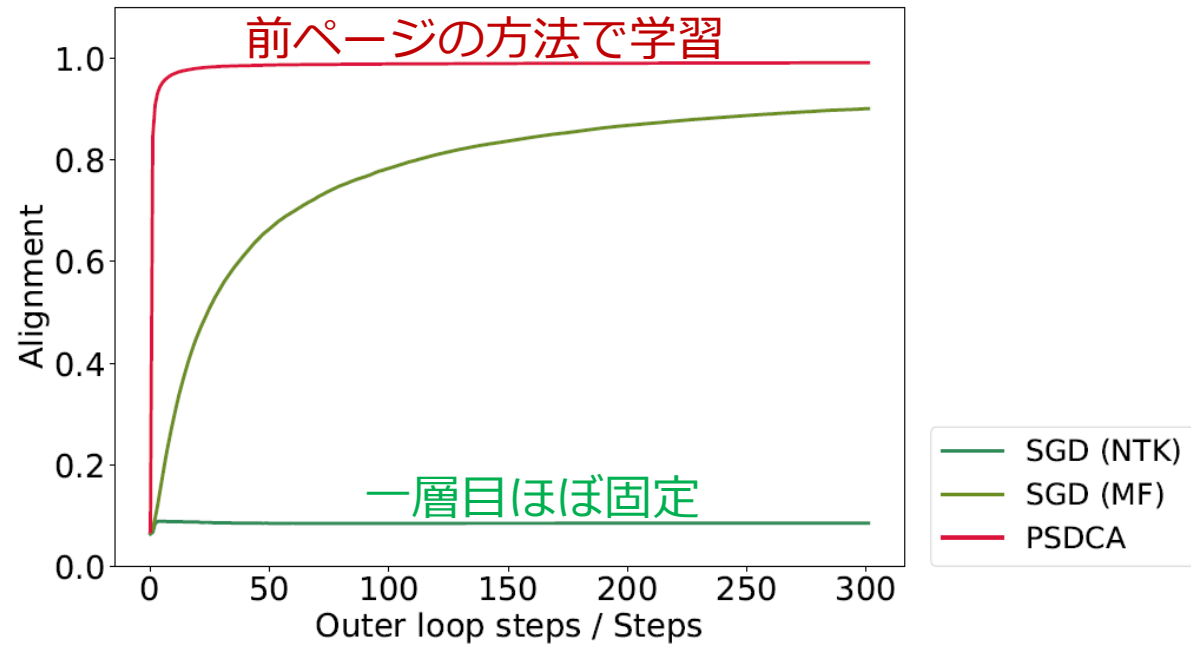
Kernel alignment



カーネルalignment:

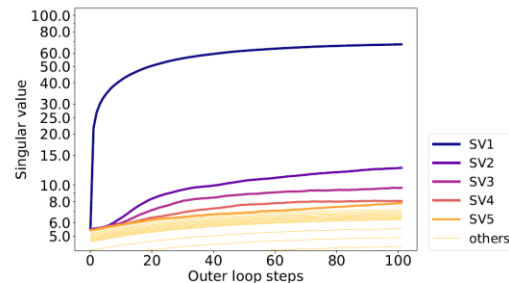
$$A(k_W) := \frac{\langle K_W, yy^\top \rangle_F}{\sqrt{\langle K_W, K_W \rangle_F \langle yy^\top, yy^\top \rangle_F}}$$

一層目で抽出された特徴量が教師信号(y)とどれだけ相関しているか?
→ 高いほど特徴量が真の関数の成分を多く含んでいる。

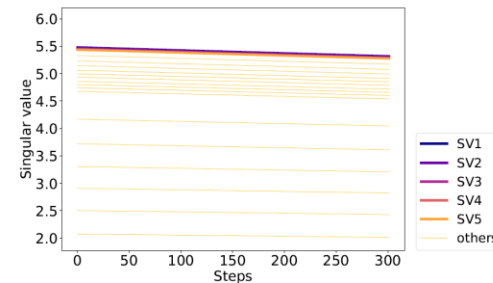


一層目も学習することで真の関数により適合した特徴量が学習できている。

固有値の分布:



(a) P-SDCA (mean field regime)



(b) SGD (NTK regime)

平均場NNの線形収束
連続時間・無限粒子

[Nitanda, Wu, Suzuki
(AISTATS2022)]
[Chizat (2022)]

時間・空間離散化：「二重ループの手法」

- PDA法 [Nitanda, Wu, Suzuki: NeurIPS2021]
- P-SDCA法 [Oko, Suzuki, Wu, Nitanda: ICLR2022]
- 無限次元拡張 [Nishikawa, Suzuki, Nitanda: NeurIPS2022]

難しい：Propagation of chaos
(McKean, Kac, ..., 60年代より)

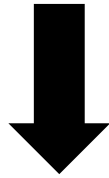
空間離散化・連続時間：
Uniform-in-time propagation of chaos

- Super対数Sobolev不等式
[Suzuki, Nitanda, Wu (ICLR2023)]
- Leave-one-out型評価
[Chen, Ren, Wang (arXiv2022)]

時間・空間離散化・確率的勾配：
「一重ループの手法」

Suzuki, Wu, Nitanda
(arXiv:2306.07221)

$$dX_t = -\nabla \frac{\delta F(\mu_t)}{\delta \mu}(X_t) dt + \sqrt{2\lambda_2} dB_t$$



(時間離散化)

$$X_{k+1}^{(i)} = X_k^{(i)} - \eta v_k^i + \sqrt{2\eta\lambda_2} \xi_k^{(i)}$$

ただし $\mathbb{E}[v_k^i] = \nabla \frac{\delta F(\hat{\mu}_k)}{\delta \mu}(X_k^i)$ かつ $\hat{\mu}_k = \frac{1}{N} \sum_{i=1}^N \delta_{X_k^{(i)}}$

(確率的勾配)

(空間離散化)

- 時間離散化: $X_t \rightarrow X_k^{(i)}$
- 空間離散化: N 粒子で近似 ($\hat{\mu}_k$) [もっとも難しい]
- 確率的勾配: 勾配計算を軽量化 ($\nabla \frac{\delta F(\mu)}{\delta \mu}(x) = \mathbb{E}[v_k(x; \mu)]$)

収束解析

$$p_\mu(x) \propto \exp\left(-\frac{1}{\lambda_2} \frac{\delta F(\mu)}{\delta \mu}(x)\right) : \text{proximal Gibbs measure}$$

定理 (1ステップ更新の減少)

p_μ は対数Sobolev不等式を定数 α で満たすとする。
損失関数の凸性と平滑性の仮定のもと、

$$\begin{aligned} & \mathcal{L}^{(N)}(\hat{\mu}_{k+1}) - \mathcal{L}(\mu^*) \\ & \leq \exp(-\lambda_2 \eta_k \alpha) \left(\mathcal{L}^{(N)}(\hat{\mu}_k) - \mathcal{L}(\mu^*) \right) \\ & \quad + C \left(\underbrace{\eta_k^3 + \lambda_2 \eta_k^2}_{\substack{\text{時間} \\ \text{離散化}}} + \underbrace{\frac{\eta_k}{N}}_{\substack{\text{空間} \\ \text{離散化}}} + \underbrace{\eta_k^{\frac{3}{2}} \lambda_2^{\frac{1}{2}} \sigma_k \tilde{\sigma}_k}_{\substack{\text{確率的} \\ \text{勾配}}} \right) \left[\begin{array}{l} \sigma_k^2 = \max_i \mathbb{E} [\|v_k^i - \mathbb{E}[v_k^i]\|^2] \\ \tilde{\sigma}_k^2 = \max_i \mathbb{E} \left[\left\| \nabla v_k^{i\top}(\mathcal{X}) - \nabla \nabla^\top \frac{\delta F(\mu \mathcal{X})}{\delta \mu}(X^i) \right\|_{\text{op}}^2 \right] \end{array} \right] \end{aligned}$$

既存研究では粒子数は時間に対して指数関数的に依存

[Mei et al., 2018; Javanmard et al., 2019; De Bortoli et al., 2020]

Assumption:

1. $F: \mathcal{P} \rightarrow \mathbb{R}$ is convex and has a form of $F(\mu) = L(\mu) + \lambda_1 \mathbb{E}_\mu[\|x\|^2]$.
2. (smoothness) $\left\| \nabla \frac{\delta L(\mu)}{\delta \mu}(x) - \nabla \frac{\delta L(\nu)}{\delta \mu}(y) \right\| \leq C(W_2(\mu, \nu) + \|x - y\|)$ and
(boundedness) $\left\| \nabla \frac{\delta L(\mu)}{\delta \mu}(x) \right\| \leq R$.

SGD-MFLD:

$$F(\mu) = \frac{1}{n} \sum_{j=1}^n f_j(\mu) \quad (\text{有限和}),$$

$$v_k^i = \frac{1}{B} \sum_{j \in I_k} \frac{\delta f_j(\hat{\mu}_k)}{\delta \mu}(X_k^i) \quad (\text{確率的勾配})$$

(Mini-batch size = B)

$$\mathcal{L}^{(N)}(\hat{\mu}_k) - \mathcal{L}(\mu^*) \lesssim \exp(-\lambda_2 \eta k \alpha) + \frac{1}{\alpha \lambda_2} \left(\underbrace{\eta^2 + \lambda_2 \eta}_{\text{時間 離散化}} + \underbrace{\frac{1}{N}}_{\text{空間 離散化}} + \underbrace{\frac{(n-B)\sqrt{\eta \lambda_2}}{B(n-1)}}_{\text{確率的 勾配}} \right)$$

更新回数のバウンド:

By setting $\eta = O\left(\epsilon \alpha \wedge (\lambda_2 \epsilon \alpha)^2 \frac{B^2(n-1)^2}{(n-B)^2 \lambda_2} \wedge \sqrt{\lambda_2 \epsilon \alpha}\right)$,
the iteration complexity becomes

$$k = O\left(\frac{1}{\epsilon \alpha} + \left(\frac{1}{\lambda_2 \epsilon \alpha}\right)^2 \frac{\lambda_2 (n-B)^2}{B^2 (n-1)^2} + \sqrt{\frac{1}{\lambda_2 \alpha \epsilon}}\right) \frac{1}{\lambda_2 \alpha} \log(\epsilon^{-1})$$

to achieve $\epsilon + O(1/(\lambda_2 \alpha N))$ accuracy.

➤ $B = n \wedge \sqrt{1/(\lambda_2 \alpha \epsilon)}$ is the optimal mini-batch size. $\rightarrow k = O(\log(\epsilon^{-1})/\epsilon)$.

分散縮小勾配法

SVRG-MFLD:

$$F(\mu) = \frac{1}{n} \sum_{j=1}^n f_j(\mu) \quad (\text{有限和}),$$

$$v_k^i = \frac{1}{B} \sum_{j \in I_k} \nabla \frac{\delta f_j(\hat{\mu}_k)}{\delta \mu}(X_k^{(i)}) - \frac{1}{B} \sum_{j \in I_k} \nabla \frac{\delta f_j(\dot{\mu})}{\delta \mu}(\dot{X}^{(i)}) + \nabla \frac{\delta F(\dot{\mu})}{\delta \mu}(\dot{X}^{(i)})$$

(分散縮小勾配)
(\dot{X} は m 回(に一回更新)

$$\begin{aligned} & \mathcal{L}^{(N)}(\hat{\mu}_k) - \mathcal{L}(\mu^*) \\ & \lesssim \exp(-\lambda_2 \eta k \alpha) \\ & \quad + \frac{1}{\lambda_2 \alpha} \left(\underbrace{\eta^2}_{\text{時間}} + \underbrace{\lambda_2 \eta}_{\text{空間}} + \frac{1}{N} + \frac{n-B}{B(n-1)} \lambda_2^{1/2} \eta \sqrt{m(\eta + \lambda_2)} \right) \end{aligned}$$

確率的
勾配の誤差

線形GLDの既存解析
[Kinoshita, Suzuki:
NeurIPS2022] の非線
形への拡張/改善

更新回数: $\eta = \epsilon \alpha \wedge \sqrt{\lambda_2 \alpha \epsilon},$

$$k = \frac{1}{\lambda_2 \alpha \eta} \log(1/\epsilon) = O\left(\frac{1}{\epsilon \alpha} + \sqrt{\frac{1}{\lambda_2 \alpha \epsilon}}\right) \frac{1}{\lambda_2 \alpha} \log(\epsilon^{-1}) \quad \text{ただし } B = \sqrt{m} = n^{1/3}.$$

総勾配計算回数: $Bk + \frac{nk}{m} \lesssim n^{1/3} \left(\frac{1}{\alpha \epsilon} + \sqrt{\frac{1}{\lambda_2 \alpha \epsilon}}\right) \frac{1}{\lambda_2 \alpha} \log(\epsilon^{-1}).$ \sqrt{n} in Kinoshita&Suzuki (2022)

統計的性質

- ℓ_i : ロジスティック損失
- $h_z(x) = \bar{R} \cdot [\tanh(\langle x_1, z \rangle + x_2) + x_3]/2$

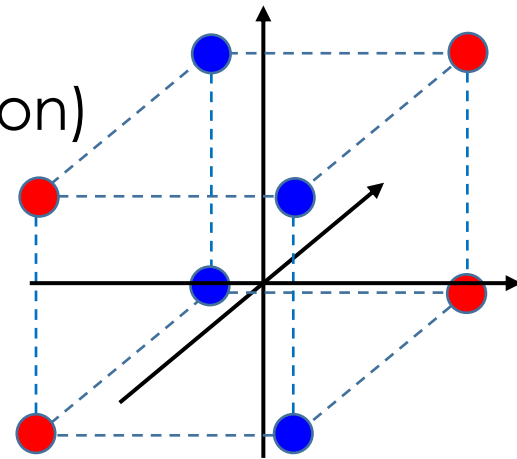
• k -スパースパリティ問題

- $X \sim \text{Unif}(\{-1, 1\}^d)$ (up to freedom of rotation)
- $Y = X_{i_1} X_{i_2} \dots X_{i_k}$ for $i_j \in [d]$ with $i_j \neq i_l$.

Q: この問題設定でカーネル法を上回る？

A: Yes.

[Suzuki, Wu, Oko, Nitanda: Feature learning via mean-field Langevin dynamics: classifying sparse parities and beyond. 2023]



Authors	regime/method	width	class error	number of iterations
Ji and Telgarsky (2019)	NTK/SGD	d^8	d^2/n	d^2/ϵ
Telgarsky (2023)	NTK/SGD	d^2	d^2/n	d^2/ϵ
Barak et al. (2022)*	Two phase SGD	$O(1)$	$d^{(k+1)/2}/\sqrt{n}$	d/ϵ^2
Telgarsky (2023)	mean-field/GF	d^d	d/n	∞
Wei et al. (2019)	mean-field/WF	∞	d/n	∞
Ours*	mean-field/MFLD	$e^{O(d)}$	$\exp(-O(\sqrt{n}/d))$	$e^{O(d)}$
Ours*	mean-field/MFLD	$e^{O(d)}$	d/n	$e^{O(d)}$

統計的性質

- ℓ_i : ロジスティック損失
- $h_z(x) = \bar{R} \cdot [\tanh(\langle x_1, z \rangle + x_2) + x_3]/2$

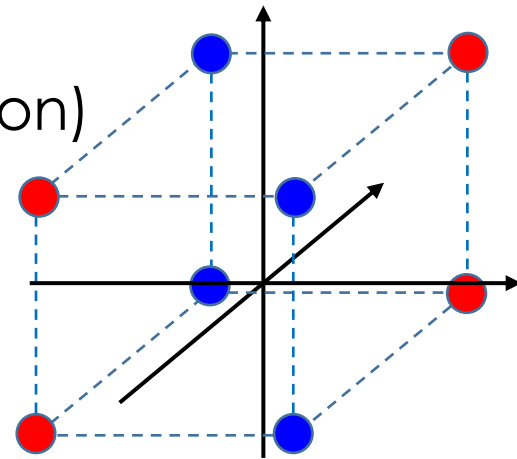
• k -スパースパリティ問題

- $X \sim \text{Unif}(\{-1, 1\}^d)$ (up to freedom of rotation)
- $Y = X_{i_1} X_{i_2} \dots X_{i_k}$ for $i_j \in [d]$ with $i_j \neq i_l$.

Q: この問題設定でカーネル法を上回る？

A: Yes.

[Suzuki, Wu, Oko, Nitanda: Feature learning via mean-field Langevin dynamics: classifying sparse parities and beyond. 2023]



特徴学習によって次元への依存性が改善されている。

Authors	regime/method	width	class e	
Ji and Telgarsky (2019)	NTK/SGD	d^8	d^2/ϵ	
Telgarsky (2023)	NTK/SGD	d^2	d^2/n	d^2/ϵ
Barak et al. (2022)*	Two phase SGD	$O(1)$	$d^{(k+1)/2}/\sqrt{n}$	d/ϵ^2
Telgarsky (2023)	mean-field/GF	d^d	d/n	∞
Wei et al. (2019)	mean-field/WF	∞	d/n	∞
Ours*	mean-field/MFLD	$e^{O(d)}$	$\exp(-O(\sqrt{n}/d))$	$e^{O(d)}$
Ours*	mean-field/MFLD	$e^{O(d)}$	d/n	$e^{O(d)}$