# 深層学習の汎化誤差のための近似性能と複雑性解析

2019/11/22 IBIS企画セッション

今泉允聡

(統計数理研究所 / 理化学研究所 / JSTさきがけ)

# 深層学習と問題設定

多層二ユーラルネット(DNN)モデル

• 多くの層を持つ(変換を繰り返す)

ス x2 x3 x4 **全**上層 ℓ層目の変換

$$h_{\ell} = \sigma(\Theta_{\ell}^{\mathsf{T}} h_{\ell-1})$$
 $\Theta_{\ell}$ : パラメタ (行列)
 $\sigma$ : 活性化関数 sigmoid, ReLUなど

データ上の経験誤差を最小化し、汎化誤差を評価

出

### 経験誤差(訓練誤差)

$$\mathcal{L}(\Theta) = n^{-1} \sum_{i=1}^{n} \ell(y_i, f_{\Theta}(x_i))$$
  
  $f_{\Theta}$ : DNN,  $\Theta = (\Theta_1, ..., \Theta_L)$ 

Θを 学習



汎化誤差(精度の尺度)  $E[\mathcal{L}(\widehat{\Theta})]$ 

2

# 汎化誤差の分解

誤差を三つの要素で説明

$$E[\mathcal{L}(\widehat{\Theta})] \leq \inf_{\Theta'} \mathcal{L}(\Theta') + \left| E[\mathcal{L}(\widehat{\Theta})] - \mathcal{L}(\widehat{\Theta}) \right| + \left| \mathcal{L}(\widehat{\Theta}) - \inf_{\Theta'} \mathcal{L}(\Theta') \right|$$

汎化誤差

近似誤差

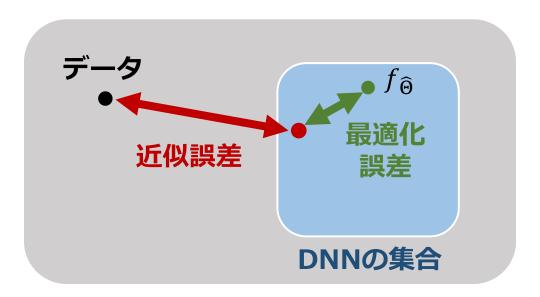
DNNの表現力

複雑性誤差

DNNの集合の大きさ

最適化誤差

学習がうまくいくか

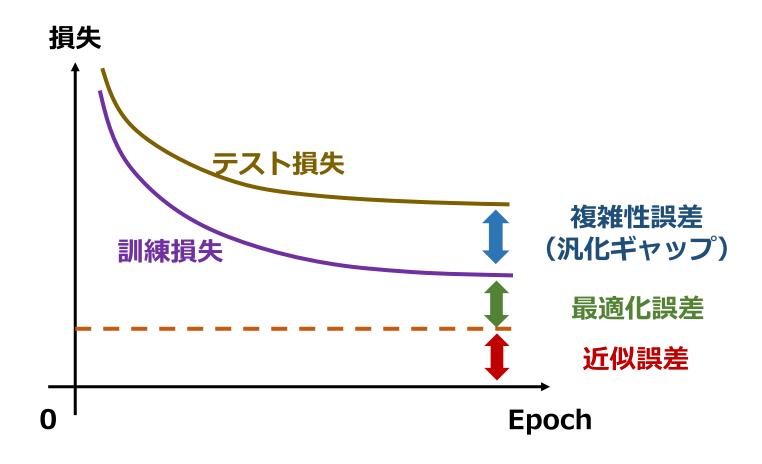


### 本講演のトピック

- ・近似誤差
- ・複雑性誤差

# 汎化誤差の分解

・実際の学習の軌跡との対応



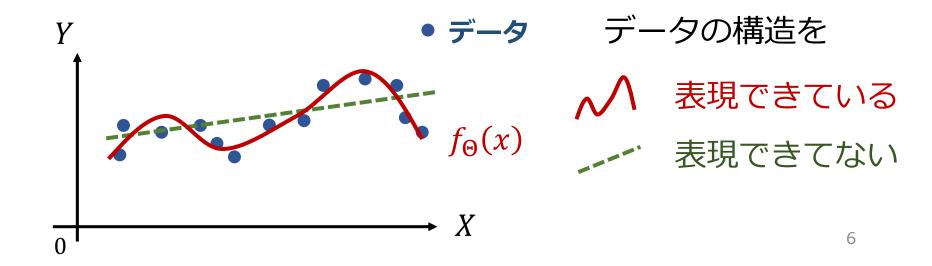
# 近似誤差

DNNと既存法の差を示す関数近似理論

# 近似誤差とは

• DNN  $f_{\Theta}$  が表現できるデータの構造で決まる

回帰の例(損失 $\ell$ は二乗損失)  $f_{\Theta}$ : DNN データが  $Y = f^*(X) + \varepsilon$  から生成されている時:  $\inf_{\Theta'} \mathcal{L}(\Theta') \leq \inf_{\Theta} \|f^* - f_{\Theta}\|_{\infty}^2 + \text{Noise Terms}$ 



# 普遍近似定理

### よく知られている結果

**普遍近似定理** (Cybenko (1989)など)

層が2つのニューラルネットワーク(NN)は、十分な数のパラメタがあれば、連続関数を任意の精度で近似できる。

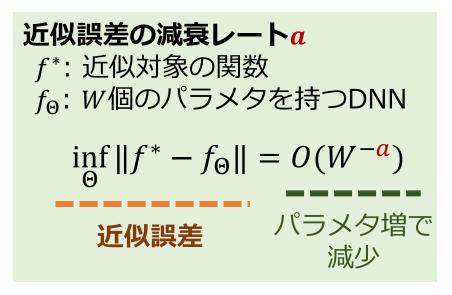


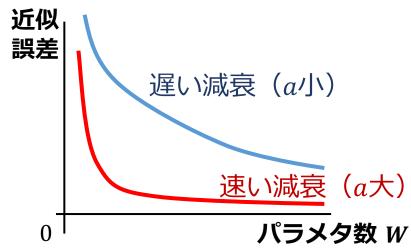
- 層2つで成立
  - → 深層でなくて良い
- 普遍近似は多くの他手法でも成立
  - → NNでなくて良い
- → なんでDNNを使うの?

# より詳細に近似を調べるには

### 近似誤差の減衰レート

・パラメタ(エッジ)が増えるときの誤差減少スピード





- レートを出すには、f\*が滑らかである必要
  - **→** *f*\*が<mark>微分可能</mark>である状況を調べる

# 滑らかな関数に対する近似レート

 $f^*$ : 近似対象(入力d次元、 $\beta$ 回微分可能)

 $f_{\Theta}$ : DNN (L層, パラメタW個、活性化関数 $\sigma$ )

### 

DNNは L=2 のもとで以下を達成:

$$\inf_{\Theta} \|f^* - f_{\Theta}\| = O(W^{-\beta/d})$$



### 活性化関数がReLUの場合 (Yarotsky (2017)など)

L層のDNNは以下を達成:

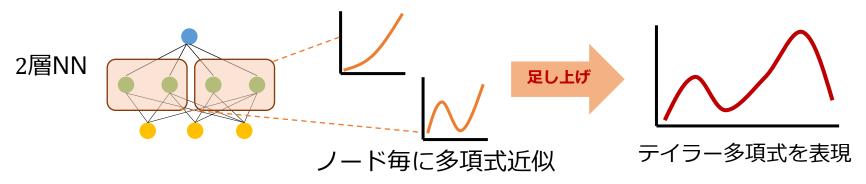
$$\inf_{\Theta} \|f^* - f_{\Theta}\| = O(W^{-\beta/d} + 2^{-L})$$

ReLUの尖りから 来る影響

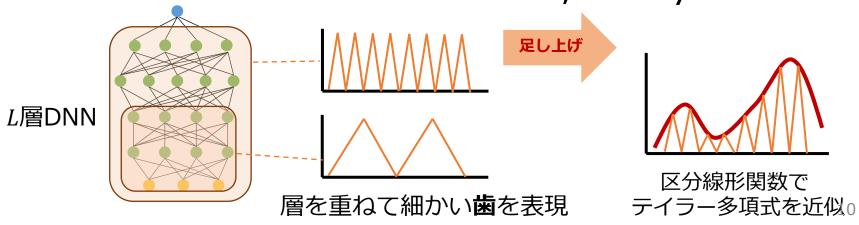
誤差レート $\beta/d$ は、 $f^*$ の滑らかさで増加、入力次元で減少。

# 活性化関数のによる違い

• σが滑らかな場合(sigmoid, softplus)



• σが滑らかでない場合(ReLU, LeakyReLU)



# 良いレートなの?

誤差レートβ/dは理論的に最適

近似誤差の最適性 (DeVore+ (1989)など) 近似誤差レート $O(W^{-\beta/d})$ は理論上の最適値。

• しかし、他手法も同じように最適

他手法の近似レート (Newman+ (1964)など) フーリ工基底、多項式基底などによる近似は レート $O(W^{-\beta/d})$ を達成する。 すごいぞ! やはりDNNは 最適なんだ!



他のも最適だから結局同じ?

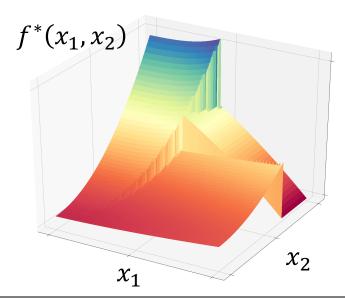


f\*が滑らかなら、DNNと他手法の理論的性能は同等。11

# DNNが重要になる状況1

汎化性能も優越するが、 ここでは近似性能のみ記述

• f\*が滑らかでない場合、DNNが他に優越



### 区分上でのみ滑らかな関数

$$f^* = \Sigma_m f_m \otimes 1_{R_m}$$

 $f_m$ : 滑らかな関数,  $1_{R_m}$ : 区分上の指示関数

### 近似レートの差別化

(Imaizumi & Fukumizu (2019))

DNNのレート:

 $O(\max\{W^{-\beta/d}, W^{-\alpha/2(d-1)}\})$ 

他手法(カーネル等)のレート:

 $O(\max\{W^{-\beta/d}, W^{-\alpha/4(d-1)}\})$ ( $\alpha$ は区分の境界線の滑らかさ)

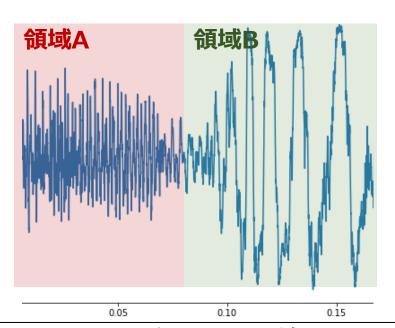
区分の境界が複雑な形

→ DNNが速いレートを達成

### DNNが重要になる状況2

汎化性能も優越するが、 ここでは近似性能のみ記述

• f\*が不均一な滑らかさを持つ場合、DNNが優越



### Besov空間の関数

$$f^* = \sum_{j} c_{j} \phi_{j} + \sum_{j,k} c_{j,k} \psi_{j,k}$$

$$\|c_{\cdot}\|_{p} + \left(\sum_{k} 2^{qk(\beta+1/2-1/p)} \|c_{\cdot,k}\|_{p}^{q}\right)^{1/q} < \infty$$

### 近似レートの差別化

(Suzuki (2019))

DNNのレート:

$$O(W^{-\beta/d})$$

他手法(カーネル等)のレート:

$$O(W^{-(\beta-(1/p-1/2)_+)/d})$$

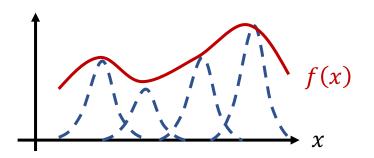
(pは不均一さの程度)

不均一さがより強い

→ DNNが速いレートを達成

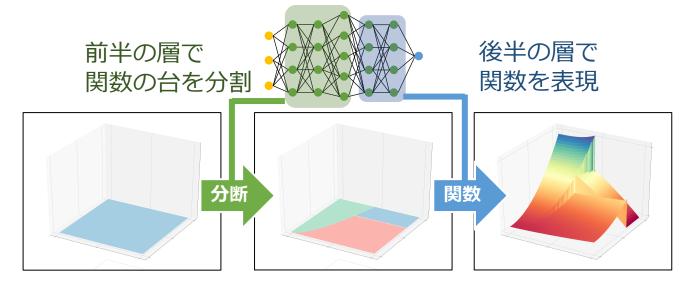
# DNNは局所構造を表現できる

これまで:均一な滑らかさ・構造を表現



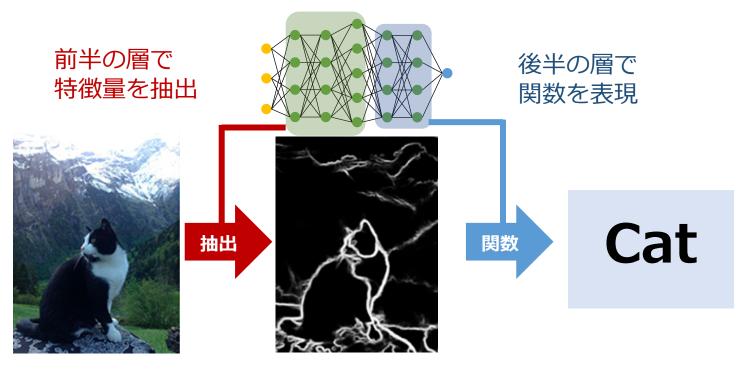
カーネル法 スプライン法 フーリエ法 など

DNN:局所的に滑らかさを変えられる



# 更なるDNNの役割:特徴量抽出

### 特徴量変換+関数表現を考える

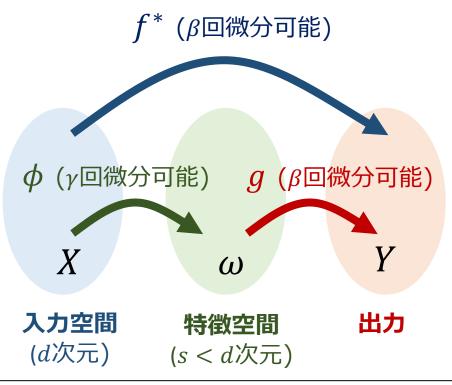


DNNは複数回の四則演算をする

→ 既存の変換(フーリエ変換など)が近似的に可能

# DNNが重要になる状況3

•特徴空間への写像がある場合、DNNのレート改善

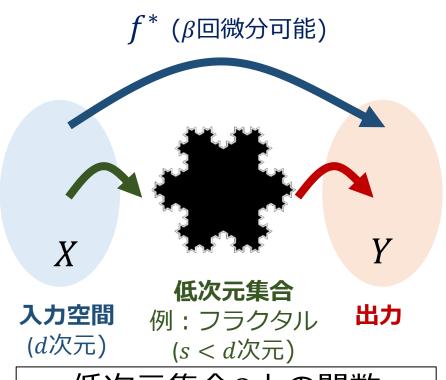


特徴写像 $\phi$ (未知)と関数gの合成  $f^* = g \circ \phi$ 

特徴量が低次元・ *φ*がシンプル **→** DNNのレート改善

### DNNが重要になる状況3+

• 特徴空間の具体例: 低次元集合



低次元集合 $\Omega$ 上の関数  $f^*(X)$ ,  $Supp(X) = \Omega$ 

### 低次元特徴量がある時の 近似レート

(Nakada & Imaizumi (2019))

DNNのレート:

 $\tilde{O}(W^{-\beta/s})$ 

一般的なレート:

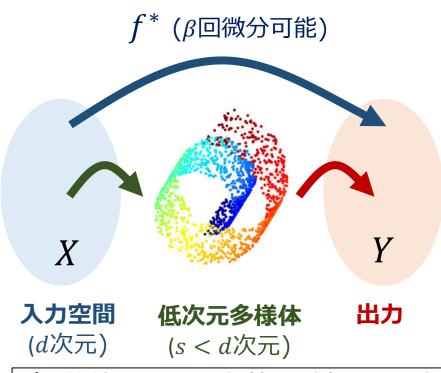
 $O(W^{-\beta/d})$ 

(s: 低次元集合の次元)

→ DNNのレート改善

### DNNが重要になる状況3+

• 特徴空間の具体例:低次元多様体



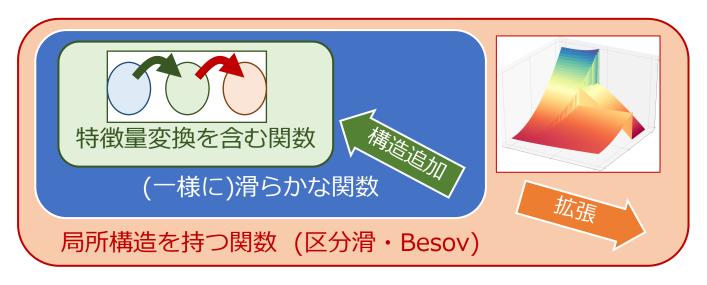
多様体上への変換を持つ関数  $f^* = g \circ \phi$ 

# 多様体特徴量がある時の近似レート (Schmidt-Hieber (2019b)) DNNのレート: $\tilde{O}(W^{-\beta/s})$ 一般的なレート: $O(W^{-\beta/d})$

→ DNNのレート改善

# 近似誤差のまとめ

• DNNが優位・改善する状況の発見



- ・知見1:局所構造を持つ関数ではDNNが優位
  - 既存法(カーネル法など)に対する優位性も示される
- ・知見2:特徴量抽出が有効ならDNNは改善
  - 狭い関数クラスだが、特徴量の構造がDNNに合う

# 近似誤差の未解決点

まだ未解決な点も多い

結局、層は すごく多いのが いいの?



### 超深層(100層など)の意義は不明

- ・上記の結果は3~5層で成立
- ・ReLUの場合もO(log n)層くらい

議論: Yarotsky(2018)など

どんな特徴量が 必要なの?



### 使える特徴量概念は抽象的

- ・解析できる特徴は次元などに限定
- ・具体的な特徴付けは今後の課題

議論: Alemi+(2016)など

ラフな差別化は出来つつあるが、現象を説明するにはまだ壁





# 複雑性誤差

なぜ巨大モデルでも汎化する?

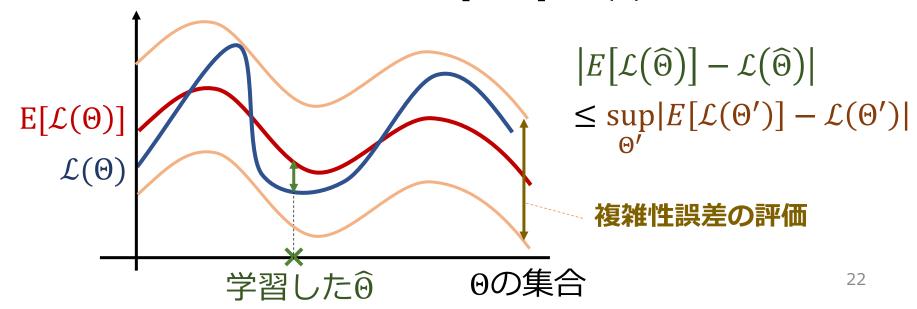
### そもそも複雑性誤差とは?

• 汎化誤差(期待値)と訓練誤差(経験平均)との差

$$|E[\mathcal{L}(\widehat{\Theta})] - \mathcal{L}(\widehat{\Theta})|$$

評価方法:一様収束誤差

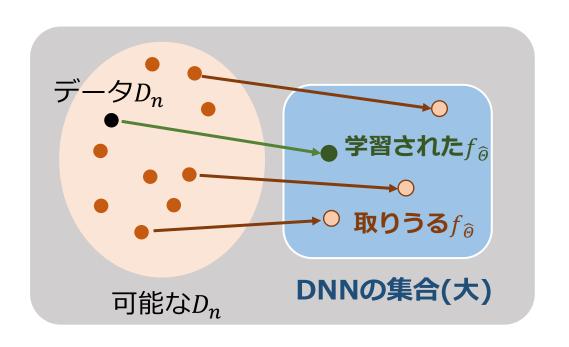
可能な全てのΘ上でのE[L(Θ)]とL(Θ)の差



# なぜ全てのΘを考える?

アルコリスムの ランダムを考慮することもある

- 汎化誤差E[L(Q)](期待値)を考えるとは?
  - 可能な $D_n$ すべての場合の平均値を考えること



汎化誤差(期待値)を考える



取りうる $D_n$ をすべて考慮

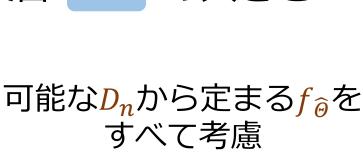


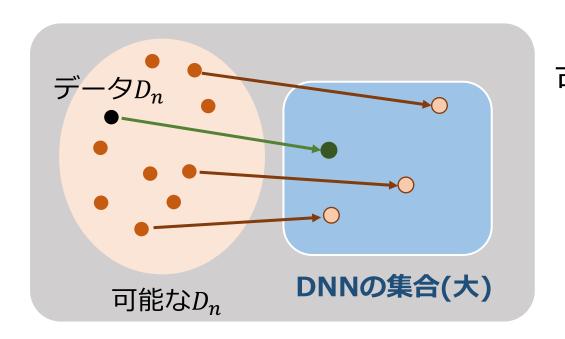
可能な*D<sub>n</sub>*から定まる*f<sub>ô</sub>を* **すべて**考慮

# 既存理論の考え方

### 既存理論

• 複雑性誤差 = 可能な $f_{
ho}$ の集合 の大きさ





可能な**f**<sub>0</sub>の候補集合が 大きいほど 複雑性誤差が増加

# 複雑性評価の数学的方法

レートは改善可だが 大きさへの依存は不可

複雑性評価 (e.g. Anthony & Bartlett (1999))

$$\sup_{\Theta} |E[\mathcal{L}(\Theta)] - \mathcal{L}(\Theta)| = O\left(\frac{1}{\sqrt{n}} \int_{0}^{\infty} \sqrt{\log N_{\delta}} \, d\delta\right)$$

### 可能なf<sub>e</sub>の集合の**大きさ**

### 導出の流れ

一様誤差

$$\sup_{\Theta} |E[\mathcal{L}(\Theta) - \mathcal{L}(\Theta)]|$$



Rademacher複雑性

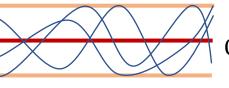
$$n^{-1/2} \mathbb{E} \left[ \sup_{\Theta} \Sigma_{i=1}^n \sigma_i \ell(y_i, f_{\Theta}(x_i)) \right]$$



Dudley積分

$$n^{-1/2} \int_0^\infty \sqrt{\log N_\delta} \, d\delta$$

→集合の大きさを評価



 $N_{\delta}$ :  $\{f_{\Theta}\}$ の最小 $\delta$ 被覆数

 $\sigma_i$ : Rademacher変数

×: 離散点(被覆球の中心)

# 複雑性はパラメタ数が主

DNNの複雑性評価 (e.g. Anthony & Bartlett (1999))

$$O\left(\frac{1}{\sqrt{n}} \int_{0}^{\infty} \sqrt{\log N_{\delta}} \, d\delta\right) = O\left(\frac{\sqrt{W \log L}}{\sqrt{n}}\right)$$

$$\rightarrow \mathcal{N} \ni \mathcal{A} \not \Rightarrow \mathcal{W} \quad \mathcal{N} \hat{\Xi} \Leftrightarrow \mathbb{Z}$$

• この理論はDNNの実性能を説明できない



大量のパラメタは 複雑性誤差を上げる



高精度DNNは 膨大なパラメタ数

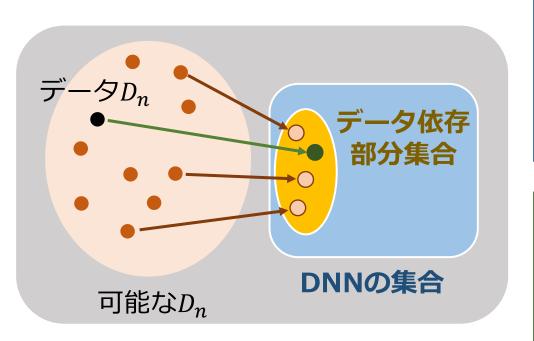
Alex Net  $\rightarrow$  6 千万 VGG Net  $\rightarrow$  1 億

統計・学習理論の(大)原則

# 新しい複雑性評価の着想

**着想**:可能な  $f_{\Theta}$ すべてを考える必要は無いのでは?

• データ依存で実現しうる $f_{\Theta}$ の部分集合がありそう







# 新しい着想を支持する実験

- データ依存集合の重要性を示す**実験**(Zhang+ (2017))
  - ・全然違うデータ $D_n, D'_n$ でも、 DNNの近似誤差・複雑性誤差は両方とも**小さい**

既存理論:複雑性誤差= の大きさ  $D_n$   $D_n'$   $D_n'$   $D_n'$   $D_n'$   $D_n'$   $D_n'$   $D_n'$  大きな近似誤差

近似誤差: どちらも小

複雑性誤差:大

近似誤差: どちらかは大

複雑性誤差:小

28

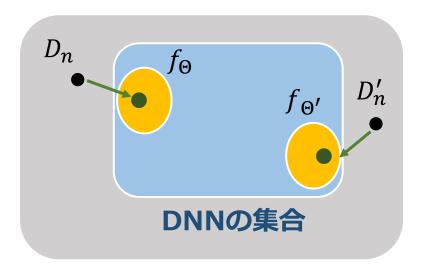
# データ依存集合は実験を説明

- データ依存集合の重要性を示す**実験**(Zhang+ (2017))
  - 全然違うデータ $D_n, D'_n$ でも、 DNNの近似・複雑性誤差は両方とも小さい

新しい着想:複雑性誤差=データ依存集合



の大きさ



近似誤差:どちらも小

複雑性誤差:常に小



実験・実現象と一致!

# データ依存集合を考える

- データ依存集合の重要性を示す実験(Zhang+ (2017))
  - 全然違うデータ $D_n, D'_n$ でも、 DNNの近似・複雑性誤差は両方とも小さい

新しい着想:複雑性誤差=データ依存集合



の大きさ



何によって



決まっているの?



### 試みの一部

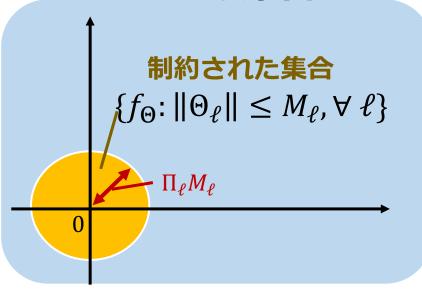
- 1. ノルム制約(暗黙的正則化)
- 2. 学習アルゴリズム
- 3. 各種正則化手法

他多数

# 試み1. ノルム制約下の複雑性

- パラメタΘが十分に小さい場合
  - 各層 $\ell=1,...,L$ の $\|\Theta_{\ell}\|$ が $M_{\ell}$ で抑えられると仮定

### DNNの集合



### ノルム制約下での複雑性誤差

(Neyshabur+ (2015a)) (Golowich+ (2019))

$$O\left(\frac{B\sqrt{L}\prod_{\ell=1}^{L}M_{\ell}}{\sqrt{n}}\right)$$

 $B = \max_{i} ||x_i||$ : データの大きさ

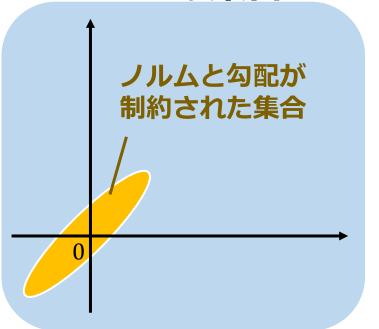
他、Bertlett+ (2017)など

- パラメタ数Wには(陽には)依存しない。
- M<sub>e</sub>が小さいほど、複雑性誤差は小さくなる。

# 試み1. ノルム制約下の複雑性

• 他ノルムや、各層の勾配を抑えることも

### DNNの集合



### ノルム・ヤコビアン制約下の誤差

(Wei & Ma (2019))

$$\tilde{O}\left(\frac{\left(\Sigma_{\ell=1}^{L}(M_{\ell}'K_{\ell})^{2/3} + (M_{\ell}''K_{\ell}')^{2/3}\right)^{3/2}}{\sqrt{n}}\right)$$

 $M'_{\ell}, M''_{\ell}$ :  $\Theta_{\ell}$ の別ノルムでの上限

 $K_{\ell}, K'_{\ell}$ :  $\ell$ 層のヤコビアンの上限

# 試み1. ノルム制約下の複雑性

・ノルム制約の実現性に課題

パラメタのノルムが常に小さければ 複雑性誤差も小さい!





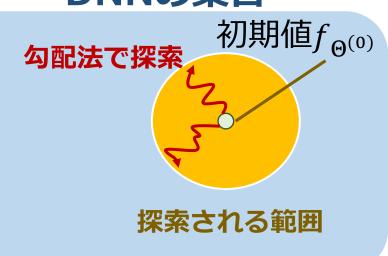
じゃあどういう時にノルムが小さいの? それはDNN特有なの?

⇒ まだ明確な特徴づけは明らかでない

数学的発見による仮説はある (線形モデルの収束など: Neyshabur+ (2014), Arora+(2019))

- ・学習アルゴリズム(勾配法)の性能を評価
  - 探索される範囲=実現するf<sub>e</sub>の集合

### DNNの集合



確率的勾配法(SGD)にも 拡張可能

### 勾配法による最適化

初期値を<u>設定</u> Θ<sup>(0)</sup>

パラメタ更新 t = 1, ..., T  $\Theta^{(t+1)} \leftarrow \Theta^{(t)} - \eta_t \nabla_{\Theta} \mathcal{L}(\Theta^{(t)})$   $\eta_t: ステップサイズ$ 

• 探索性能を下げる  $\rightarrow$  実現する $f_{\Theta}$ も減る

### DNNの集合

ステップサイズ・更新回数を減



探索範囲は縮小

### 早く学習をやめれば誤差減

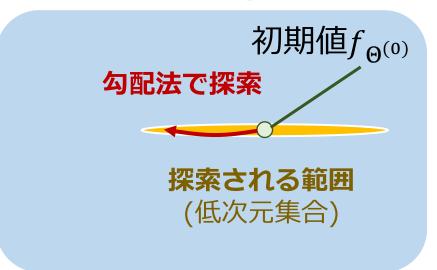
(Hardt+ (2016)) (Kuzborskij+ (2017))  $\eta_t = O(1/t)$ とする  $O\left(\frac{(W/L)^L T^q}{n}\right)$ 

 $T \ge 1$ : パラメタの更新回数  $q \in (0,1)$ : 減衰率

- 探索範囲を狭めれば、複雑性誤差は小さくなる。
- パラメタ数Wの影響は受ける。

・実際の探索の低次元性を考慮

### DNNの集合



### 低次元集合上の探索で誤差減

(Imaizumi & Sriperumbudur (2019))  $\eta_t = O(1/t)$ とする

$$\tilde{O}\left(\frac{\mathbf{d}_{\mathbf{W}}Ld\log T^{q'}}{n}\right)$$

 $d_W < W$ : 探索される空間の次元  $q' \in (0,1)$ : 減衰率

- 低次元構造で、パラメタ数Wの影響を緩和。
- ただし活性化関数に制約が入る。

・ノルム制約の実現性に課題

学習を控えめにすれば 複雑性誤差は小さい!





それは良いパラメタが 学習できないんじゃないの?

→ 近似・最適化誤差との関係は明らかでない。

損失関数の形状の研究が必要。

実験的反論もある: Hoffer+ (2018)

# 複雑性誤差の理論まとめ

問題:DNNは膨大なパラメタを持つが、複雑性誤差が小さい

理論化の方向:複雑性誤差=データ依存のモデル部分集合

• その部分集合の正体は明らかでない



### 正体の候補

- ・ノルム制約集合
- ・学習アルゴリズムの影響
- ・各種正則化手法 (バッチ正規化など), 他多数



多くの仮説・実験・数学的発見

・暗黙的正則化 / バイアス / 損失関数の形状, etc…

# まとめと展望

# まとめ・展望

### 近似誤差

• 進捗しているが更なる発展の余地

### 複雑性誤差

- ・実現象との大きな矛盾
- モデル部分集合のスキーム
  - 未解明な点は多い
  - 他のスキームも提案

実験・数学的発見から理論の拡張へ

### 実運用

数学的発見



Lottery Ticket, Implicit bias, Double Descent, Flat minima, etc.

DNNの 理論

ご静聴ありがとうございました。

### Reference

- Cybenko, G. (1989). Approximations by superpositions of a sigmoidal function. Mathematics of Control, Signals and Systems, 2, 183-192.
- Mhaskar, H. N. (1996). Neural networks for optimal approximation of smooth and analytic functions. Neural computation, 8(1), 164-177.
- Yarotsky, D. (2017). Error bounds for approximations with deep ReLU networks. Neural Networks, 94, 103-114.
- DeVore, R. A., Howard, R., & Micchelli, C. (1989). Optimal nonlinear approximation. Manuscripta mathematica, 63(4), 469-478.
- Newman, D. J. and Shapiro, H. S., (1964) Jackson's theorem in higher dimensions, On Approximation Theory (Proceedings of Conference in Oberwolfach, 1963), pp. 208–219.
- Imaizumi, M., & Fukumizu, K. (2018). Deep neural networks learn non-smooth functions effectively. AI & Statistics.
- Suzuki, T. (2018). Adaptivity of deep ReLU network for learning in Besov and mixed smooth Besov spaces: optimal rate and curse of dimensionality.
- Petersen, P., & Voigtlaender, F. (2018). Optimal approximation of piecewise smooth functions using deep ReLU neural networks. Neural Networks, 108, 296-330.
- Schmidt-Hieber, J. (2019a). Nonparametric regression using deep neural networks with ReLU activation function. Annals of Statistics, to appear.

- Nakada, R., & Imaizumi, M. (2019). Adaptive approximation and estimation of deep neural network to intrinsic dimensionality. arXiv preprint arXiv:1907.02177.
- Schmidt-Hieber, J. (2019b). Deep ReLU network approximation of functions on a manifold. arXiv preprint arXiv:1908.00695.
- Yarotsky, D. (2018). Optimal approximation of continuous functions by very deep relu networks. Conference on Learning Theory.
- Alemi, A. A., Fischer, I., Dillon, J. V., & Murphy, K. (2016). Deep variational information bottleneck. International Conference on Learning Representations.
- Anthony, M., & Bartlett, P. L. (1999). Neural network learning: Theoretical foundations. cambridge university press.
- Zhang, C., Bengio, S., Hardt, M., Recht, B., & Vinyals, O. (2016). Understanding deep learning requires rethinking generalization. International Conference on Learning Representations.
- Neyshabur, B., Tomioka, R., & Srebro, N. (2014). In search of the real inductive bias: On the role of implicit regularization in deep learning. arXiv preprint arXiv:1412.6614.
- Arora, S., Cohen, N., Hu, W., & Luo, Y. (2019). Implicit Regularization in Deep Matrix Factorization. Neural Information Processing Systems.
- Hardt, M., Recht, B., & Singer, Y. (2016). Train faster, generalize better: Stability
  of stochastic gradient descent. In International Conference on Machine
  Learning.
- Kuzborskij, I., & Lampert, C. H. (2017). Data-dependent stability of stochastic gradient descent. arXiv preprint arXiv:1703.01678.