^{九州大学集中講義} 深層学習および機械学習の数理

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

2020年9月2日~4日

機械学習と人工知能の歴史



ネオコグニトロン



LeNet



- <u>誤差逆伝搬法</u>でパラメータを更新
- 手書き文字認識データセット(MNIST)で99%の精度を達成

Deep Learning 深層学習

ImageNet



ImageNet Challenge

IM GENET

- 1,000 object classes (categories).
- Images:
 - 1.2 M train
 - 100k test. 0



ImageNet: 1,000カテゴリ,約120万枚の訓練画像データ ILSVRC (ImageNet Large Scale Visual Recognition Competition) [J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In CVPR09, 2009.]

ImageNetデータにおける識別精度の変遷



ImageNet: 21841クラス, 14,197,122枚の訓練画像データ そのうち1000クラスでコンペティション



様々なタスクで高い精度

AlphaGo/Zero



deep neural networks and tree search, Nature, 529, 484-489, 2016]

画像認識



自動翻訳

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



[Wu et al.: Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. arXiv:1609.08144]

画像の生成



[Glow: Generative Flow with Invertible 1x1 Convolutions. Kingma and Dhariwal, 2018]





[Zhu, Park, Isola, and Efros: Unpaired image-to-image translation using cycle-consistent adversarial networks. ICCV2017.]

諸分野への波及

Tumor ...

Macrophag





[タオル畳み、サラダ盛り付け 「指動く」ロボット初公開, ITMedia:http://www.itmedia.co.jp/news/articles/1711/30/news089 .html]



[Niepert, Ahmed&Kutzkov: Learning Convolutional Neural Networks for Graphs, 2016]

[Gilmer et al.: Neural Message Passing for Quantum Chemistry, 2017] [Faber et al.:Machine learning prediction errors better than DFT accuracy, 2017.]



- 人を超える精度 (FROC73.3% -> 87.3%) - 悪性腫瘍の場所も特定

> [Detecting Cancer Metastases on Gigapixel Pathology Images: Liu et al., arXiv:1703.02442, 2017]





- 「〇〇法が良い」という様々な仮説の氾濫.
- <u>世界的課題</u>
- 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)

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.



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



深層学習の理論概観



基礎

応用

教師あり学習



モデル:関数の集合(例:深層NNの表せる関数の集合)



 $h_1(u) = [h_{11}(u_1), h_{12}(u_2), \dots, h_{1d}(u_d)]^T$

活性化関数は通常要素ごとにかかる. Poolingのように要素ごとでない非線形変換もある.

- ☆ReLU (Rectified Linear Unit) :
- シグモイド関数:







●第ℓ層

$$\phi_{\ell+1}(x) = \eta(W^{(\ell)}\phi_{\ell}(x) + b^{(\ell)})$$
$$W^{(\ell)} \in \mathbb{R}^{m_{\ell+1} \times m_{\ell}} \quad b^{(\ell)} \in \mathbb{R}^{m_{\ell+1}}$$



活性化関数の例

☆ReLU (Rectified Linear Unit)

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





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



※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 $L(\hat{\theta}) - \inf_{f} \mathbb{E}[\ell(Y, f(X))]$ $L(\theta)$ $\widehat{\widehat{L}}(\theta)$ 汎化ギャップ 余剰誤差 θ^*

学習機の複雑さと学習能力

• No free lunch theorem

「あらゆる問題で性能の良い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」



^{機械学習への教訓} 「必要以上に複雑なモデルを当てはめると失敗する」

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

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]



 どこを贔屓するか? ▶モデルの外に真があればそもそも学習方法の優劣を論じ ることは難しい(どれもドングリの背比べ) ▶仮にモデルに真が入っていた場合にどこが贔屓されるか P:真の分布のモデル $P^* \in \mathcal{P}$:真の分布 $\hat{\theta}, \tilde{\theta}:$ 推定量 $D^n = (x_i, y_i)_{i=1}^n$:訓練データ ■ ミニマックス最適性 $\sup_{P^* \in \mathcal{P}} \mathbb{E}_{D^n \sim P^*} [L(\hat{\theta})] = \inf_{\tilde{\theta}: \text{Estimator}} \sup_{P^* \in \mathcal{P}} \mathbb{E}_{D^n \sim P^*} [L(\tilde{\theta})]$ ■ 許容性 つぎのようなÕが存在しない: $\mathbb{E}_{D^n \sim P^*}[L(\tilde{\theta})] \le \mathbb{E}_{D^n \sim P^*}[L(\hat{\theta})] \quad (\forall P^* \in \mathcal{P})$ $\mathbb{E}_{D^n \sim P^*}[L(\hat{\theta})] < \mathbb{E}_{D^n \sim P^*}[L(\hat{\theta})]$ $(\exists P^* \in \mathcal{P})$ ■ 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 $\hat{\theta}$. (→ ベイズ推定量)

損失関数最小化

経験損失(訓練誤差)

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

$$\ell(y,y') = (y - y')^{2} \qquad 二乗損失 (回帰) \quad (y,y' \in \mathbb{R})$$

$$\ell(y,y') = -\sum_{k=1}^{K} y_{k} \log(y'_{k}) \quad \text{Cross-entropy損失 (多値判別)} \\ (y_{k} \in \{0,1\}, y'_{k} \in [0,1], \text{ともに和が1})$$

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

- 基本的には確率的勾配降下法 (SGD) で最適化を実行
- AdaGrad, Adam, Natural gradientといった方法で高速化
 微分はどうやって求める? → 誤差逆伝搬法

• 勾配降下法



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

誤差逆伝搬法

合成関数

$$f(x;W) = f_{3,W_3}(f_{2,W_2}(f_{1,W_1}(x)))$$

= $f_{3,W_3} \circ f_{2,W_2} \circ f_{1,W_1}(x)$

合成関数の微分

$$\frac{\partial f}{\partial W_1}(x) = \frac{\partial f_{3,W_3}}{\partial f_{2,W_2}} \frac{\partial f_{2,W_2}}{\partial f_{1,W_1}} \frac{\partial f_{1,W_1}}{\partial W_1}(x)$$



連鎖律を用いて微分を伝搬

$$\frac{\partial f}{\partial W_3}(x) = \frac{\partial f_{3,W_3}}{\partial W_3} \left(f_{2,W_2} \circ f_{3,W_3}(x) \right)$$
$$\frac{\partial f}{\partial W_2}(x) = \frac{\partial f_{3,W_3}}{\partial f_{2,W_2}} \frac{\partial f_{2,W_2}}{\partial W_2} \left(f_{3,W_3}(x) \right)$$
$$\frac{\partial f}{\partial W_1}(x) = \frac{\partial f_{3,W_3}}{\partial f_{2,W_2}} \frac{\partial f_{2,W_2}}{\partial f_{1,W_1}} \frac{\partial f_{1,W_1}}{\partial W_1}(x)$$

パラメータによる微分と入力による微分は違うが、情報をシェアできる. $f_{1,W}(x) = h(Wx) \text{ O 場合} \qquad u = Wx$ $\frac{\partial f_{1,W}}{\partial W_{ij}}(x) = \frac{\partial h}{\partial u_i}(u)x_j \qquad \frac{\partial f_{1,W_1}}{\partial x_j}(x) = \sum_i \frac{\partial h}{\partial u_i}(u)W_{ij}$

確率的勾配降下法 (SGD)

(Stochastic Gradient Descent)

沢山データがあるときに強力



大きな問題を分割して個別に処理

普通の勾配降下法:

$$W^{t} = W^{t-1} - \alpha \nabla L(W)$$

= $W^{t-1} - \alpha \frac{1}{n} \sum_{i=1}^{n} \nabla \ell(z_{i}, W)$
全データの計算

確率的勾配降下法 (SGD)

(Stochastic Gradient Descent)

沢山データがあるときに強力



大きな問題を分割して個別に処理

普通の勾配降下法:

確率的勾配降下法:

毎回の更新でデータを一つ(または少量)しか見ない







理論的課題

表現能力 (第一章) どれだけ難しい問題まで学習でき るようになるか?





要点のまとめ



▶重要な関数クラス(Barron, Holder, Sobolev, Besov) はほ ぼ最適な効率性で近似できる.

▶ 適応的な関数近似によりカーネル法を優越する.

·汎化能力

 ▶重要な関数クラスの推定精度はミニマックス最適レート を達成できる→特徴抽出機能によりカーネル法を優越.
 ▶データサイズがパラメータ数より小さくても過学習しない→実質的統計的自由度が小さい.

▶陰的正則化等により、自由度が小さく抑えられる→汎化する.

•最適化能力

 >横幅を十分広くとれば大域的最適解が勾配法で求まる.
 >初期パラメータのスケーリングによってNeural Tangent KernelとMean fieldの二つの状況に大きく分けられる.
 >ノイズ付加により平滑化・局所解からの脱出を実現.

第1章 深層学習の表現能力

Kolmogorovの加法定理

任意の連続関数は横幅固定の4層の"ニューラルネット" で表現できる.

定理 (Kolmogorov's superposition theorem)

- 定数 $\lambda_j > 0 \ (j = 1, ..., d), \sum_{j=1}^d \lambda_j \le 1,$
- *I* = [0,1]から*I*への狭義単調増大連続関数 *ϕ*_q (*q* = 1,...,2*d*+1)

が存在して,任意の連続関数 $f \in C([0,1]^d)$ が次のように表現できる:

$$f(x_1,\ldots,x_d) = \sum_{q=1}^{2d+1} g(\lambda_1\phi_q(x_1) + \cdots + \lambda_d\phi_q(x_d))$$

なお, *g* ∈ *C*([0,1])は*f*にのみ依存した関数.

- 一変数関数の合成で多変数関数が作れてしまう。
- 任意の連続関数は4層ニューラルネットの最終層だけを学習すればよいことになる。しかし、gの滑らかさはfおよび入力の次元dに強く依存し、最適な学習精度は達成できない。
- ヒルベルトの23の問題の13番目を一般化した定理



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	ĥ
1987	Hecht-Nielsen	対象毎に構成	$C(R^d)$	
1988	Gallant & White	Cos	$L_2(K)$	
	Irie & Miyake	integrable	$L_2(\mathbb{R}^d)$	
1989	Carroll & Dickinson	Continuous sigmoidal	$L_2(K)$	
	Cybenko	Continuous sigmoidal	$\mathcal{C}(K)$	
	Funahashi	Monotone & bounded	C(K)	
1993	Mhaskar + Micchelli	Polynomial growth	C(K)	
2015	Sonoda + Murata	Unbounded, admissible	$L_{1}(R^{d})/L_{2}(R^{d})$	

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












連続関数の近似

• Cybenkoの理論

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



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

とできる.

証明の直感的概略

・シグモイド型の関数に対し,

$$h\left(a(\alpha^{\top}x+\beta)+\theta\right) \stackrel{a\to\infty}{\longrightarrow} \begin{cases} 1 & (\alpha^{\top}x+\beta>0)\\ h(\theta) & (\alpha^{\top}x+\beta=0)\\ 0 & (\alpha^{\top}x+\beta<0) \end{cases}$$

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、 $\cos(\alpha^{\mathsf{T}}x + \beta)$ や $\sin(\alpha^{\mathsf{T}}x + \beta)$ 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.





(Sonoda & Murata, 2015)

- Ridgelet変換による解析(Fourier変換の親戚)
- 3層NNはridgelet変換で双対空間(中間層)に行って から戻ってくる(出力層)イメージ

積分表現の概略(Ridgelet変換)

ある $\psi: \mathbb{R} \rightarrow \mathbb{R}$ が存在して、以下の「許容条件」が成り立つとする:

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

定理

 $f, \hat{f} \in L^1(\mathbb{R}^d)$ の時,許容条件を満たす η, ψ に対し以下がほとんどいたるところの $x \in \mathbb{R}^d$ に対して成り立つ:

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

なお,連続点においては等式が常に成り立つ.

つまり, $f(x) = \int \frac{T(a,b)}{K_{\psi,\eta} \|a\|} \eta(a^{\top}x - b) \mathrm{d}a \mathrm{d}b$ と書ける.



40

Ridgelet変換 = Radon変換 +Wavelet変換

カーネル法との関係



- 例: リッジ回帰 (Tikhonov正則化)

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\arg\min} \frac{1}{n} \|X\beta - Y\|_2^2 + \lambda_n \|\beta\|_2^2$$

変数変換:

- 正則化項のため、 $\hat{\beta} \in \text{Ker}(X)^{\perp}$. つまり、 $\hat{\beta} \in \text{Im}(X^{\top})$.
- ある $\hat{\alpha} \in \mathbb{R}^n$ が存在して、 $\hat{\beta} = X^{\top} \hat{\alpha}$ と書ける.

$$\hat{\alpha} \leftarrow \arg\min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \| X X^\top \alpha - Y \|_2^2 + \lambda_n \alpha^\top (X X^\top) \alpha.$$

新しい入力xに対しては $y = x^{\top} X^{\top} \hat{\alpha}$ で予測.

• リッジ回帰の変数変換版

$$\hat{\alpha} \leftarrow \operatorname*{arg\,min}_{\alpha \in \mathbb{R}^n} \frac{1}{n} \| X X^{\top} \alpha - Y \|_2^2 + \lambda_n \alpha^{\top} (X X^{\top}) \alpha.$$

※ $(XX^{\mathsf{T}})_{i,j} = x_i^{\mathsf{T}} x_j \, \mathrm{d} x_i \, \mathrm{d} x_j$ の内積.

・<u>カーネル法のアイディア</u>

xの間の内積を他の関数で置き換える:

 $x_i^{\top} x_j \rightarrow k(x_i, x_j)$

この $k: \mathbb{R}^{p} \times \mathbb{R}^{p} \to \mathbb{R}$ をカーネル関数と呼ぶ.

<u>カーネル関数の満たすべき条件</u>

- 対称性: k(x, x') = k(x', x)
- 正値性: $\sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j k(x_i, x_j) \ge 0, \forall (\{x_i\}_{i=1}^{m}, \{\alpha_i\}_{i=1}^{m}, m)$

この条件さえ満たしていればなんでも良い

カーネルリッジ回帰

カーネルリッジ回帰:
$$K = (k(x_i, x_j))_{i,j=1}^n$$
 として,

$$\hat{\alpha} \leftarrow \arg\min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \| \mathbf{K} \alpha - \mathbf{Y} \|_2^2 + \lambda_n \alpha^\top \mathbf{K} \alpha.$$

新しい入力 x に対しては,

$$y = \sum_{i=1}^{n} k(x, x_i) \hat{\alpha}_i$$

で予測.





[https://scikit-learn.org/stable/auto_examples/plot_kernel_ridge_regression.html]

カーネル関数の例

• ガウシアンカーネル

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

● 多項式カーネル

$$k(x,x') = \left(1 + x^{\top}x'\right)^{p}$$

- $\chi^2 \neg \neg \neg \ddot{\chi}$ $k(x, x') = \exp\left(-\gamma^2 \sum_{j=1}^d \frac{(x_j - x'_j)^2}{(x_j + x'_j)}\right)$
- Matérn-kernel

$$k(x,x') = \int_{\mathbb{R}^d} e^{\mathrm{i}\lambda^\top (x-x')} \frac{1}{(1+\|\lambda\|^2)^{\alpha+d/2}} \mathrm{d}\lambda$$

グラフカーネル,時系列カーネル,…
 (ユークリッド空間でなくてもカーネルが定義できればカーネル法が使える)



カーネル関数 \Leftrightarrow 再生核ヒルベルト空間 (RKHS) k(x,x') \mathcal{H}_k



集合 Ω 上の再生核ヒルベルト空間(Reproducing kernel Hilbert space, RKHS) \mathcal{H} とは、 Ω 上の関数からなるヒルベルト空間であって、任意の $x \in \Omega$ に対し $\phi_x \in \mathcal{H}$ が 存在し、

$$f(x) = \langle \phi_x, f \rangle_{\mathcal{H}} \quad (f \in \mathcal{H})$$

を満たすものをいう.

- $k(x,y) \coloneqq \langle \phi_x, \phi_y \rangle_{\mathcal{H}}$ は正定値対称カーネル関数
- 逆に正定値対称カーネルが与えられたら対応するRKHSが一章に存在

定理 (Moore-Aronszajnの定理)

k(x,y): 正定値対称カーネル (given)

 Ω 上の関数からなるヒルベルト空間 \mathcal{H}_k で以下の条件を満たすものが一意に存在:

 $f(x) = \left\langle \sum_{i=1}^{n} \alpha_i k(x_i, \cdot), k(x, \cdot) \right\rangle_{\mathcal{A}_i} = \sum_{i=1}^{n} \alpha_i k(x_i, x)$

1.
$$k(x,\cdot) \in \mathcal{H}_k$$

2. $f = \sum_{k=1}^{n} \alpha_k k(x_k) t \geq 7 \pm 1$

- 2. $f = \sum_{i=1}^{n} \alpha_i k(x_i, \cdot)$ なる有限和は \mathcal{H}_k 内で稠密
- 3. 再生成:

 $f(x) = \langle k(x, \cdot), f \rangle_{\mathcal{H}_k} \ (\forall x \in \Omega, \forall f \in \mathcal{H}_k).$

再生核ヒルベルト空間のイメージ

47

- 高次元(無限次元)特徴空間にφで写像して推論を行う.
- 再生核ヒルベルト空間では線形な処理が元の空間では非線形 処理になる。



 $k(\mathbf{x},\mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}}$

多項式カーネル



http://www.eric-kim.net/eric-kim-net/posts/1/kernel_trick.html

カーネルリッジ回帰の再定式化

再生性: *f* ∈ *H_k* に対し

$$f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}_k}.$$

カーネルリッジ回帰

$$\hat{f} \leftarrow \min_{f \in \mathcal{H}_k} \quad \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 + C \|f\|_{\mathcal{H}_k}^2$$

• 表現定理

$$\exists \alpha_i \in \mathbb{R} \quad \text{s.t.} \quad \hat{f}(x) = \sum_{i=1}^n \alpha_i k(x_i, x),$$
$$\Rightarrow \quad \|\hat{f}\|_{\mathcal{H}_k} = \sqrt{\sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j)} = \sqrt{\alpha^\top K \alpha}$$

さきほどのカーネルリッジ回帰の定式化と一致.



カーネル関数の分解とRKHSの表現

50

• カーネル関数は対象行列の対角化に対応する分解を許す.

ある非負測度
$$\nu$$
に対して、 $\int_{\Omega} k(x,x) d\nu(x) < \infty$ が成り立つなら、
高々加算個の正実数列(μ_i)_{i \in I}および $L_2(\nu)$ 内の正規直交基底(e_i)_{i \in I}が存在して、
 $k(x,x') = \sum_{i \in I} \mu_i e_i(x) e_i(x')$ ($\forall x,x'$)
と分解できる (各点収束).

このような分解は他にもいろいろとバージョンがある, e.g., Mercer展開.
 (詳細は[Steinwart&Scovel: Constructive Approximation, 35(3):363—417, 2012])

• さらに $(\sqrt{\mu_i} e_i)_{i \in I}$ はRKHS内の正規直交基底になる. $\mathcal{H}_k = \left\{ f(x) = \sum_{i \in I} \alpha_i \sqrt{\mu_i} e_i(x) \mid \sum_{i \in I} \alpha_i^2 < \infty \right\}$ $\|f\|_{\mathcal{H}_k} = \sqrt{\sum_{i \in I} \alpha_i^2}, \quad \langle f, f' \rangle_{\mathcal{H}_k} = \sum_{i \in I} \alpha_i \alpha_i'$

カーネル回帰の再定式化



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.







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



再生核ヒルベルト空間の理論



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

$$\eta(\hat{F}_{\ell-1}(x))$$

$$\hat{w}^{\top}\eta(\hat{F}_{\ell-1}(x)) \in \mathcal{H}_{\ell}$$

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

\hat{k}_{ℓ} に対応した再生核ヒルベルト空間

$$\mathcal{H}_{\ell} = \{ f(x) = w^{\top} \eta(\hat{F}_{\ell-1}(x)) \mid w \in \mathbb{R}^{m_{\ell}} \}, \\ \|f\|_{\mathcal{H}_{\ell}} = \|w\|. \quad (\texttt{s} \ \texttt{b} \ \texttt{I} \ \texttt{c} \ \texttt{i} \ \texttt{c} \ \texttt{b} \ \texttt{c} \ \texttt{i} \ \texttt{f} \ \texttt{b} \ \texttt{f} \ \texttt{f$$

See Montavon et al. (2011); Bach (2015)

る)

- 深層学習は各層でカーネル関数を逐次的に構築 する学習方法であるとも言える。
- データから適応的にカーネル関数を学習する点が通常のカーネル法と異なる。

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?





→「表現力」を比べてみる.

万能近似理論より二層ニューラルネットでも中間層のユ ニット数を無限に増やせば任意の関数を任意の精度で近 似できる.

歴史的には後にSVMの理論に繋がってゆく. (例:Gaussian kernelの万能性)

Q:ではなぜ深い方が良いのか? A:深さに対して**指数的に表現力が増大**するから.

表現力と層の数

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<u>指数的</u>に上がる.



折り紙のイメージ

Montufar, Guido F., et al. "On the number of linear regions of deep neural networks." 2014.



他にも同様の結論を出している論文多数

- 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





対称性の高い関数は,特に層を深く することで得をする.









区分線形関数の表現

- 任意の区分線形関数($\mathbb{R}^d \to \mathbb{R}$)は深さ[$\log_2(d+1)$]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





有理関数の近似

有理関数をReLU-DNNで近似

 $p: [0,1]^d \to [-1,+1], q: [0,1]^d \to [2^{-k},+1]$: <u>r次</u>多項式 p/qをReLU-DNNで近似したい

あるReLU-DNNf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:



• 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-DNNf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:



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

有理関数の近似

有理関数をReLU-DNNで近似

$$p: [0,1]^d \to [-1,+1], q: [0,1]^d \to [2^{-k},+1]$$
 : r次多項式
 p/q をReLU-DNNで近似したい

あるReLU-DNNf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:



統計的推定理論による比較

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

$$y_i = f^{o}(x_i) + \xi_i \quad (i = 1, ..., n)$$

 $\xi_i \sim N(0, \sigma^2)$ は観測誤差



平均二乗誤差:

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

※実はこれは<u>二乗損失の平均余剰誤差</u>になっている. $\mathbb{E}[L(\hat{\theta}) - \inf_{f} \mathbb{E}[\ell(Y, f(X))]]$

なぜ深層学習が良いのか?

• 真の関数f°の形状によって深層が有利になる





(複雑な関数形状に適応的にフィットすることができる)



- ・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

[Suzuki, ICLR2019]



- 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.



非線形回帰問題

非線形回帰モデル
$$y_i = f^{o}(x_i) + \xi_i$$
 $(i = 1, ..., n)$ ただし、 $\xi_i \sim N(0, \sigma^2)$ かつ $x_i \sim P_X(X)$ (i.i.d.).

$$f^{o}$$
をデータ $(x_i, y_i)_{i=1}^n$ から推定したい



※以下の理論は判別問題でも展開可能

バイアスとバリアンスのトレードオフ 69



推定精度 = バイアス (モデル誤差) + バリアンス (分散)



- ・Barronクラス
- Hölderクラス
- Sobolevクラス
- Besovクラス

真の関数が各クラスに含まれているときに 近似誤差はどれくらいになるか調べたい.

Hölder, Sobolev, Besov空間

 $\Omega = [0, 1]^d \subset \mathbb{R}^d$ • Hölder space ($\mathcal{C}^\beta(\Omega)$)

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

• Sobolev space $(W_p^k(\Omega))$

$$\|f\|_{W_p^k} = \left(\sum_{|\alpha| \le k} \|D^{\alpha}f\|_{L^p(\Omega)}^p\right)^{\frac{1}{p}}$$

• Besov space $(B_{p,q}^{s}(\Omega))$ $(0 < p, q \le \infty, 0 < s \le m)$ $\varpi_{m}(f,t)_{p} := \sup_{\|h\| \le t} \left\| \sum_{j=0}^{m} (-1)^{m-j} {m \choose j} f(\cdot + jh) \right\|_{L^{p}(\Omega)}$, $\|f\|_{B_{p,q}^{s}(\Omega)} = \|f\|_{L^{p}(\Omega)} + \left(\int_{0}^{\infty} [t^{-s} \omega_{m}(f,t)_{p}]^{q} \frac{\mathrm{d}t}{t} \right)^{1/q}$. $\|f\|_{B_{p,q}^{s}(\Omega)} = \|f\|_{L^{p}(\Omega)} + \left(\int_{0}^{\infty} [t^{-s} \omega_{m}(f,t)_{p}]^{q} \frac{\mathrm{d}t}{t} \right)^{1/q}$.

- 直感:
 - ・非整数回の微分も定義したい.
 - ・ 整数回微分を"つなげる"→実補間
 - qはそのつなげ方を制御
 - *s*で関数の滑らかさを制御
 - pで滑らかさの空間的一様性を制御
空間の間の関係

• For $m \in \mathbb{N}$,

$$B_{p,1}^{m} \hookrightarrow W_{p}^{m} \hookrightarrow B_{p,\infty}^{m},$$
$$B_{2,2}^{m} = W_{2}^{m}.$$

• For $0 < s < \infty$ and $s \notin \mathbb{N}$,

$$\mathcal{C}^s = B^s_{\infty,\infty}.$$



• 連続関数の領域: s > d/p

$$B^s_{p,q} \hookrightarrow C^0$$

• L^r -可積分な領域: $s > d(1/p - 1/r)_+$

 $B^s_{p,q} \hookrightarrow L^r$



• \emptyset : $B_{1,1}^1([0,1]) \subset \{\text{bounded total variation}\} \subset B_{1,\infty}^1([0,1])$

不連続な領域 d/p > s



滑らかさが非一様的な場合:
 *p*が小さい状況



これらの性質にも関わらず深層学習は良い学習ができるか?

スパース性との関係







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

$$\mathcal{F}(L, W, S, B)$$
 $\begin{bmatrix} \bullet & 縦幅: L \\ \bullet & 横幅: W \\ \bullet & 枝の数: S \\ \bullet & A パラメータの上限: B \end{bmatrix}$ の深層NNモデルの集合

・活性化関数はReLUを仮定



関数近似能力

• $0 < p,q,r \le \infty \ge 0 < s < \infty$ が以下を満たすとする:

 $s > d(1/p - 1/r)_+$ (L^r-可積分性)

• $m \epsilon s < \min\{m, m - 1 + 1/p\}$ 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:
$$L = O(\log(N)), \qquad \qquad W = O(N), \\S = O(N\log(N)), \qquad B = O(N^{(d/p-s)_+}),$$
すると,深層NNは以下の誤差でBesov空間の元を近似できる:
$$\frac{\sup_{f^{\circ} \in U(B_{p,q}^{s}([0,1]^{d}))} \inf_{\tilde{f} \in \mathcal{F}(L,W,S,B)} \|f^{\circ} - \tilde{f}\|_{L^{r}([0,1]^{d})} \lesssim N^{-s/d}.$$

Pinkus (1999), Mhaskar (1996): $p = r \rightarrow 01 \le p$, 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 ($s \le k + 1 + (d - 1)/2$).

B-spline

$$\mathcal{N}(x) = egin{cases} 1 & (x \in [0,1]), \ 0 & (ext{otherwise}). \end{cases}$$

次数mのcardinal B-spline:

$$\mathcal{N}_m(x) = (\underbrace{\mathcal{N} * \mathcal{N} * \cdots * \mathcal{N}}_{m+1 \text{ times}})(x)$$

→ 区分多項式



$$\mathcal{N}_{k,j}^{(d)}(x_1,\ldots,x_d) = \prod_{i=1}^d \mathcal{N}_m(2^k x_i - j_i)$$

Cardinal B-splineによる展開 (DeVore & Popov, 1988)

• Atomic decomposition:

 $f \in B^s_{p,q}$ の必要十分条件:

$$f = \sum_{k \in \mathbb{N} + j \in J(k)} \alpha_{k,j} \mathcal{N}_{k,j}^{(d)}$$

と分解できて
$$(ただし, J(k) = \{j \in \mathbb{Z}^d \mid -m < j_i < 2^{k_i+1} + m\})$$

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

• ノルムの同値性: $\|f\|_{B^s_{p,q}} \simeq N(f)$



(see also Bolcskei, Grohs, Kutyniok, Petersen: Optimal Approximation with Sparsely Connected Deep Neural Networks. 2018)



• 真の関数f°が次のように展開できるとする:

$$f^{o} = \sum_{k=1}^{\infty} \sum_{j \in J(k)} \alpha_{k,j} \phi_{k,j}(x)$$

- 各基底関数 $\phi_{k,i}$ をReLU-NNでよく近似できるなら、f°も良く近似できる.
- Cardinal B-splineはReLU-NNでよく近似できる: log(1/ε)層で近似誤差εを達成する.
- B-splineに関する定理を深層学習に持ち込める.
 → Besov空間に限らない理論を展開可能

[Bölcskei et al.: arXiv:1705.01714]



$s > d(1/p - 1/r)_+$ なる仮定のもとで

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

- *p* = *q* = ∞の時, Yarotsky (2016)の結果に帰着 (Hölder空間)
- <u>適応的</u>な<u>非線形</u>近似が必要 (Dung, 2011)



関連研究

Chui et al. (1994) and Bölcskei et al. (2017) dealt with a "smooth" activation with lim_{x→∞} η(x)/x^k → 1, lim_{x→-∞} η(x)/x^k = 0 with k ≥ 2 under 1 ≤ p. Mhaskar and Micchelli (1992) studied s = k + 1. Mhaskar (1993) studied k ≥ 2 and s = k + 1, Mhaskar (1996) considered the Sobolev space W^m_p with a "bump" activation function (excluding ReLU).

バイアスとバリアンス分解



- これまで示したこと:バイアス(近似誤差)
- これから示すこと:経験誤差最小化のバリアンス

$$\hat{f} = \operatorname*{argmin}_{f \in \mathcal{F}(L,W,S,B)} \sum_{i=1}^{n} (y_i - f(x_i))^2.$$

$$E[\|f^{\circ} - \hat{f}\|_{L^{2}(P_{X})}^{2}]$$

$$\lesssim \underbrace{\frac{S[L\log(BW) + \log(Ln)]}{n}}_{\text{Variance}} + \underbrace{\inf_{f \in \mathcal{F}(L,W,S,B)} \|f - f^{\circ}\|_{L^{2}(P_{X})}^{2}}_{\text{Bias}}$$

深さ 横幅
$$スパース性$$

 $(\#starset network = starset n$

⇒ バイアスとバリアンスのトレードオフをバランスすればよい.



•最小二乗解(訓練誤差最小化)

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

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

定理(推定精度)

$$\begin{split} \|f^{\mathbf{o}}\|_{B^{s}_{p,q}} &\leq 1, \|f^{\mathbf{o}}\|_{\infty} \leq 1 \, \text{m o } 0 < p, q \leq \infty, \, s > d(1/p - 1/2)_{+} \, \text{o } \varepsilon \neq s, \\ N &\asymp n^{\frac{d}{2s+d}} \, \varepsilon \neq \delta \, z \, \varepsilon \, \varepsilon, \end{split}$$

$$||f^{\circ} - \hat{f}||^{2}_{L^{2}(P_{X})} \leq n^{-\frac{2s}{2s+d}} \log(n)^{2}.$$

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

線形推定量との比較

- 線形推定法 (Donoho & Johnstone, 1994) $\hat{f}(x) = \sum_{i=1}^{n} y_i \varphi_i(x_1, \dots, x_n; x) + b \geq 書ける <u>任意の推定量</u>$ (<u>カーネルリッジ回帰</u>, Sieve法, Nadaraya-Watson推定量...) $\hat{f}(x) = K_{x,X}(K_{X,X} + \lambda I)^{-1}Y \quad \hat{f} = \underset{f \in \mathcal{H}_k}{\operatorname{argmin}} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{H}_k}^2$ $n^{-\frac{2s - 2(1/p - 1/2) + 2(1/p - 1/2) + 2}{2s + 1 - 2(1/p - 1/2) + 2}} \quad (d=10 \text{ B})$
- 深層学習

p < 2で差が出る



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

间凤



スパース推定との違い:

スパース:

あらかじめ用意した多数の基底の中から<u>重要な</u> 基底を選択

Deep:
 直接,基底を構築する





David Donoho: ガウス賞 (2018) スパース推定, wavelet-shrinkage, 圧縮センシング, ...

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

例:ジャンプが3か所の区分定数関数

深層:1/n, カーネル: $1/\sqrt{n}$



→ さらには「**スパース推定**」という観点からも説明できる.

[Satoshi Hayakawa and Taiji Suzuki: On the minimax optimality and superiority of deep neural network learning over sparse parameter spaces. arXiv:1905.09195.]

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<u>任意の推定量</u> 例: カーネルリッジ回帰 $\hat{f}(x) = K_{x,X}(K_{X,X} + \lambda \mathbf{I})^{-1}Y$ ("浅い"学習法とみなす)



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]







例 (1): 低次元データ構造



例 (2): 縮小ランク回帰

縮小ランク回帰

$$Y_i = UVX_i + \xi_i \quad (i = 1, \ldots, n)$$

ただし, $Y_i \in \mathbb{R}^M, X_i \in \mathbb{R}^N$, $U \in \mathbb{R}^{M \times r}, V \in \mathbb{R}^{r \times N}$ $(r \ll M, N)$.

- Linear estimator $\hat{f}(x) = \sum_{i=1}^{n} Y_i \varphi(X_1, \dots, X_n, x)$,
- Deep learning $\hat{f}(x) = \hat{U}\hat{V}x$.



低ランク行列の凸包はフルランク行列

Barronクラスの汎化誤差

Sobolevクラスの近似理論

$$\Pi_{N} := \left\{ \sum_{j=1}^{N} \alpha_{j} \eta(a_{j}^{\top} x + b_{j}) \mid \alpha_{j} \in \mathbb{R}, a_{j} \in \mathbb{R}^{d}, b_{j} \in \mathbb{R} \right\} \quad \begin{bmatrix} \mathbf{p} \parallel \mathbb{R} & \mathbf{p} \parallel \mathbb{R} \\ \square \mathbb{R} = \neg \mathbf{p} \mid \mathbf{p} \parallel \mathbf{p} \parallel$$

$$\eta$$
がある開区間で無限回微分可能であり、その開区間のある点bにおいて
 $\frac{\partial^k \eta}{\partial x^k}(b) \neq 0 \ (\forall k \in \mathbb{Z}, k \ge 0)$
とする. すると、 $\forall f \in W_p^s([0,1]^d)$ に対してある $g \in \Pi_N$ が存在して、
 $\|f - g\|_p \leq N^{-\frac{s}{d}} \|f\|_{W_p^s}$
(ノード数Nの中間層を用いた近似誤差)

[Mhaskar: Neural networks for optimal approximation of smooth and analytic functions. Neural Computation, 8(1):164–177, 1996]

- この近似誤差はN個の基底を用いた近似法の中で最適なオーダーを達成.
- シグモイド関数は条件を満たす。ReLUは満たさない。
- 滑らかな関数はより近似しやすい.

中田



- 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。



「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



「過去問は解けるけれども本番の試験は解けない」 という状況を回避したい

従来の学習理論









従来の学習理論





「**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.



[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.



深層学習の汎化誤差

 $\psi(y_i, f(x_i))$:損失関数 (1-Lipschitz continuous w.r.t. f)

$\hat{\Psi}(f) := \frac{1}{n} \sum_{i=1}^{n} \psi(y_i, f(x_i))$	$\Psi(f) := \mathbb{E}[\psi(Y, f(X))]$
経験誤差(訓練誤差)	期待誤差(汎化誤差)

• 汎化ギャップ

任意の $\hat{f} \in \mathcal{F}$ に対して成り立つ一様バウンドが欲しい:

$$\Psi(\hat{f}) \le \hat{\Psi}(\hat{f}) + \underbrace{\bar{R}_n(\psi \circ \mathcal{F})}_{}$$

Complexity

E.g., Rademacher 複雑度: $P(\epsilon_i = -1) = P(\epsilon_i = 1)$ (i.i.d.)

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

Uniform bound



Uniform bound



"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



深層学習の汎化誤差バウンド(抜粋)

<u>ノルム型バウンド</u>

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)	$\frac{\prod_{\ell=1}^{L} R_{\ell,2}}{\sqrt{n}} \left(\frac{R_{\ell,2 \to 1}^{2/3}}{R_{\ell,2}^{2/3}} \right)^{3/2}$	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)	$\sqrt{\frac{L^2 \max_{1 \le i \le n} \widehat{f}(x_i) ^2 \sum_{\ell=1}^{L} \frac{1}{\mu_{\ell}^2 \mu_{\ell}^2}}{n}}}{n}$	Compression
Baykal et al. (2018)	$\sqrt{\frac{L^2 \max_{1 \le i \le n} \widehat{f}(x_i) ^2 \sum_{\ell=2}^{L} (\hat{\Delta}^{\ell} \rightarrow)^2 \sum_{i=1}^{W} S_i^{\ell}}{n}}$	Compression
Suzuki et al. (2018)	$\sum_{\ell=2}^{L} \sqrt{\lambda_{\ell}} + \sqrt{\frac{\sum_{\ell=1}^{L} m_{\ell+1}^{\sharp} m_{\ell}^{\sharp}}{n}}$	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)



 ・ パラメータ数 ∑_{l=1}^L m_l m_{l+1} がそのままバウンドに現れてしまう.

 ・ パラメータ数≫データサイズの状況を説明できていない.



"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

 $\mathsf{NN} \in \mathcal{F} \mathcal{V}: \ f(x) = (W^{(L)} \eta(\cdot)) \circ (W^{(L-1)} \eta(\cdot)) \circ \cdots \circ (W^{(1)} x)$

• Bartlett et al. (2017): 正規化マージンバウンド

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

$$R_{\ell,2} := \|W^{(\ell)}\| = \sigma_1(W^{(\ell)}) \qquad \qquad R_{\ell,2\to1} := \sum_j \|W_{j,:}^{(\ell)}\|$$

• Golowich et al. (2018)

$$\prod_{\ell=1}^{L} R_{\ell,F} \min\left\{\frac{1}{n^{1/4}}, \sqrt{\frac{L}{n}}\right\}$$
$$R_{\ell,F} := \|W^{(\ell)}\|_{\mathrm{F}} : \operatorname{Frobenius} \mathcal{I} \mathrel{\text{IV L}}_{(R_{\ell,2} \downarrow \Downarrow \downarrow \bigstar \lor)}$$

○ 横幅に依存しない. → Overparametrizationの状況を説明!
 ※ 縦幅に指数的に依存する場合がある.(バウンドによる)

圧縮型バウンド



[Arora et al., 2018; Zhou et al., 2019; Baykal et al., 2019; Suzuki et al., 2018]

注:これらのバウンドはfの汎化誤差は与えていない. "圧縮していない"ネットワークfの汎化誤差も与えられる (xページ).

<u>|非圧縮ネット</u>の圧縮型バウンド

仮定: \hat{f} がより小さなネットワーク $f^{\#}$ に圧縮できるとする. ($\hat{f} \in \hat{F}, f^{\#} \in F^{\#}; \hat{F}$ は学習済みネットの集合, $F^{\#}$ は圧縮したネットの集合) [Suzuki, Abe, Nishimura: ICLR2020]

非圧縮ネットの圧縮型バウンド:

$$\Psi(\widehat{f}) \leq \widehat{\Psi}(\widehat{f}) + \frac{1}{\sqrt{n}}\widehat{r} + \sqrt{\frac{\sum_{\ell=1}^{L} m_{\ell+1}^{\sharp} m_{\ell}^{\sharp}}{n}}\log(n)$$

 $\|\widehat{f} - f^{\sharp}\|_{n}^{2} \leq \widehat{r}^{2}$ (a.s.) <u>:圧縮可能性は訓練データのみから判断.</u> (学習結果が圧縮可能なように学習するのが好ましい)



$$\Psi(\mathbf{f}^{\sharp}) \lesssim \hat{\Psi}(\hat{f}) + \hat{\mathbf{r}} + \sqrt{\frac{\sum_{\ell=1}^{L} m_{\ell}^{\sharp} m_{\ell+1}^{\sharp}}{n}} \log(n)$$
既存のバウンド)

<u>|非圧縮ネット</u>の圧縮型バウンド

仮定: \hat{f} がより小さなネットワーク $f^{\#}$ に圧縮できるとする. ($\hat{f} \in \hat{F}, f^{\#} \in F^{\#}; \hat{F}$ は学習済みネットの集合, $F^{\#}$ は圧縮したネットの集合) [Suzuki, Abe, Nishimura: ICLR2020]



より正確なステートメント

- 仮定: \hat{f} がより小さなネットワーク $f^{\#}$ に圧縮できるとする. ($\hat{f} \in \hat{F}, f^{\#} \in F^{\#}; \hat{F}$ は学習済みネットの集合, $F^{\#}$ は圧縮したネットの集合)
- $\|\widehat{f} f^{\sharp}\|_{n}^{2} \le \widehat{r}^{2}$ (a.s.)
- $\dot{R}_r(\mathcal{G}') := \bar{R}_n(\{f \in \mathcal{G}' \mid \|f\|_{L_2(P_X)} \le r\}) : \underline{\beta \beta \beta Rademacher複雑度}$
- $r_* := \inf\{r > 1/n \mid \dot{R}_r(\psi(\widehat{\mathcal{F}}) \psi(\mathcal{F}^{\sharp})) \le r^2\}$:局所Rad.の不動点 $\psi(\widehat{\mathcal{F}}) - \psi(\mathcal{F}^{\sharp}) := \{\psi(\widehat{f}) - \psi(f^{\sharp}) \mid \widehat{f} \in \widehat{\mathcal{F}}, f^{\sharp}\mathcal{F}^{\sharp}\}$

Theorem(非圧縮ネットの圧縮型バウンド)

いつどれくらい圧縮できるか?

- 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.



[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]

[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.

	元サイズ	圧縮可能 サイズ
Layer	Original	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	$2,\!359,\!296$	$3,\!886$
	*	/\





Near low rank covariance

Distribution of eigenvalues of the covariance matrix in an internal layer



9-th layer in VGG-13 trained on CIFAR-10



Compression error



If the eigenvalue drops rapidly, then the network can be compressed into much smaller one.

Then
$$\exists f^{\#}$$
 s.t. $\|f^{\sharp} - \hat{f}\|_{n}^{2} \leq C \left(\sum_{\ell=2}^{L} \sqrt{\lambda_{\ell}}\right)^{2} \|f\|_{n}^{2} = \frac{1}{n} \sum_{i=1}^{n} f(x_{i})^{2}$

Singular value distribution

Singular value distribution of trained weight matrix



[Martin&Mahoney: Traditional and Heavy-Tailed Self Regularization in Neural Network Models. arXiv:1901.08276]

Near low rank weight and covariance¹²⁰

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

• Near low rank weight matrix:

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

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

Both of weight and covariance are near low rank

 $(\sigma_j(\cdot): j$ -th largest eigenvalue)

+ Other boundedness condition.



Comparison with existing work

Comparison of intrinsic dimensionality between our degree of freedom and that in Arora et al. (2018). They are computed on VGG-19 network trained on CIFAR-10.

Layer	Original	Arora et al. (2018)	Our bound
1	1,728	$1,\!645$	567
4	$147,\!456$	$644,\!654$	$75,\!240$
6	$589,\!824$	$3,\!457,\!882$	$394,\!200$
9	$1,\!179,\!648$	$36,\!920$	$10,\!395$
12	$2,\!359,\!296$	22,735	756
15	$2,\!359,\!296$	$26,\!584$	882
		larger	smaller

Implicit number of parameters in each internal layer:

 $\hat{N}_{\ell}(\lambda_{\ell})\hat{N}_{\ell+1}(\lambda_{\ell+1})k^2,$

where k is the size of filter. We used $\lambda_{\ell} = 10^{-2} \text{Tr}[\widehat{\Sigma}_{(\ell)}]$ (sufficiently small).

[S. Arora, R. Ge, B. Neyshabur, and Y. Zhang. Stronger generalization bounds for deep nets via a compression approach. ICML2018.]

BigNAS

[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







理論の副産物

自由度の理論解析により、ネットワークのどこ に着目すればどれだけ圧縮できるかがわかる.



深層ニューラルネットワークの圧縮技術への応用



ニューラルネットワークの圧縮



提案手法:

従来手法より良い精度

94%の圧縮 (精度変わらず)

約半分に圧縮しても精度落ちず

74.04%

Spec-ResB w/ ft

86.67%

91.77%

20.69M

20.69M

5.25G

124

転移学習のネットワーク構造決定

ある閾値以上の固有値をカウント (e.g., 10⁻³).
 → 縮小したネットワークのサイズとして使う.

[Shinya, Simo-Serra, and Suzuki: Understanding the Effects of Pre-training for Object Detectors via Eigenspectrum. ICCV2019, Neural Architects Workshop]

• その後、スクラッチから学習(*S*)もしくはImageNet事前学習モデルをファイン チューニングする(*J*).



Backbone	Normalization -	Classification		COCO ($2 \times$ schedule)					
Dackoone		MACs	#params	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L
ResNet-50 [35]	SyncBN	3.8 G	—	34.5	55.2	37.7	20.4	36.7	44.5
ResNet-50*	GN	4.09 G	25.5 M	35.5	55.6	38.5	21.3	37.5	45.3
ResiaxNet $S3-50$ (MACs)	GN	4.06 G	18.6 M	35.4	55.4	38.6	21.5	37.3	45.2
ResiaxNet $\mathcal{I}1$ -50 (MACs)	GN	4.05 G	21.7 M	35.5	55.5	38.6	21.4	37.3	46.0
ResiaxNet $\mathcal{I}3$ -50 (MACs)	GN	4.07 G	22.0 M	35.4	55.6	38.4	21.3	37.8	45.5
ResiaxNet \mathcal{I} 3-50 (params)*	GN	4.92 G	24.7 M	35.8	55.9	38.9	21.8	38.0	45.6
DetNet-59 [35]	SyncBN	4.8+ G		36.3	56.5	39.3	22.0	38.4	46.9
DetNet-59 [†]	GN	5.00+ G	18.3 + M	36.2	56.0	39.3	22.1	38.3	46.0
DetiaxNet I2-59 (MACs)	GN	4.94+ G	17.4+ M	36.2	56.0	39.3	22.5	38.1	46.0

Overparameterizeされた ネットワークの統計学



"新しい"バイアス-バリアンスのトレードオフ



- モデルがある複雑度(サンプルサイズ)を超えた後,第二の降下が起きる.
- モデルサイズがデータより多いと推定量の<u>バリアンスがむしろ減る</u>.

※設定によるので注意が必要.

線形回帰における二重降下

- Hastie et al.: Surprises in High-Dimensional Ridgeless Least Squares Interpolation, arXiv:1903.08560.
- 線形回帰を考察

$$y_{i} = x_{i}^{\top}\beta + \epsilon_{i}$$

$$\mathbb{E}[\epsilon_{i}] = 0, \operatorname{Var}(\epsilon_{i}) = \sigma^{2}, \operatorname{Cov}(x_{i}) = \Sigma \qquad \|\beta\| = r$$
• 最小ノルム解: $\hat{\beta} = (X^{\top}X)^{+}X^{\top}Y$

$$(X = [x_{1}, \dots, x_{n}]^{\top}, Y = [y_{1}, \dots, y_{n}]^{\top})$$
• 期待予測誤差: $\mathbb{E}_{\epsilon}[\|\hat{\beta} - \beta\|_{\Sigma}^{2}|X]$

定理

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

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



直感:

 次元(d)>サンプルサイズ(n)だとデータの張る部分空間は全体の一部 →実質的自由度がdより低く、バリアンス小

注意:

- 次元が大きくなると真の関数も変化している設定.
- 単なる線形回帰なので第一層も学習する深層学習とは異なる.

2層NNの二重降下理論

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



[Ba, Erdogdu, Suzuki, Wu, Zhang: Generalization of Two-layer Neural Networks: An Asymptotic Viewpoint. ICLR2020]

二層NNの学習.今度は a_m を固定して w_m を学習. w_m の初期値が大きい状況 (NTK)→二重降下現れる w_m の初期値が小さい状況 (平均場)→二重降下が弱い

小さな初期値から始めると真の関数を表す最小限の表現を獲得する 大きな初期値から始めるとカーネル法と似た状況で過学習しやすい



Implicit regularization (陰的正則化)¹³¹

- ニューラルネットワークの学習では様々な「陽的正則化」を用いる:
 バッチノーマリゼーション, Dropout, Weight decay, ...
- 実は深層学習の構造が自動的に生み出す「陰的正則化」も強く効いているという説。

例:線形ネットワーク $f_W := W^{(L)} W^{(L-1)} \dots W^{(1)} x$

(L2正則化学習)
$$\widehat{W} = \underset{W}{\operatorname{argmin}} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \lambda \sum_{\ell=1}^{L} \|W^{(\ell)}\|_{\mathrm{F}}^2.$$

任意の局所最適解は<u>低ランク</u>になる:

$$\widehat{W}^{(L-1)} = \cdots = \widehat{W}^{(1)} = \widehat{u}\widehat{v}^{ op}$$
 (rank 1)

<u>モデルの複雑さが大幅に削減されている.</u>

(見た目のパラメータ数) $Lm^2 \rightarrow 2m$ (実質的パラメータ数)

※非線形活性化関数がある場合は完全には解明されていない(論文は多数)



 小さな初期値から勾配法を始めるとノルム最小 化点に収束しやすい→陰的正則化



[Gunasekar et al.: Implicit Regularization in Matrix Factorization, NIPS2017] [Soudry et al.: The implicit bias of gradient descent on separable data. JMLR2018] [Gunasekar et al.: Implicit Bias of Gradient Descent on Linear Convolutional Networks, NIPS2018] [Moroshko et al.: Implicit Bias in Deep Linear Classification: Initialization Scale vs Training Accuracy, arXiv:2007.06738]



ResNetのODE解釈

ResNetの各層は特徴の最適化の一反復,常微分方程式の離散化とみなせる.

 $h_{j+1} = h_j + F_j(h_j)$ $\bigvee_{\substack{dh_t \\ dt}} = \underbrace{\nabla L(h_t)}_{(=F_t)} \circ h_t$

[E, 2017][Sonoda & Murata, 2017][Li & Shi, 2017]



 F_i

 h_i

ResNetと常微分方程式をつなげることで 常微分方程式の数値解法をネットワーク 構造の決定に持ち込める.

→ PolyNet, FractalNet, RevNet, Linear-Multistep-ResNet, ...

[Lu et al.: Beyond Finite Layer Neural Networks: Bridging Deep Architectures and Numerical Differential Equations, ICML2018]

ODE-Net

 $\mathbf{h}_{t+1} = \mathbf{h}_t + f_{\theta_t}(\mathbf{h}_t)$ ResNet



- 層を連続化することですべての層が一つのネットワークに集約される。
- ODEにすることで汎用ODEソルバーを適用できる.
 - ▶ 入力点に応じて離散化の粒度を変えられる. (層の概念がもはやない)



ODE Network

Chen et al. Neural ordinary differential equations. NeurIPS2018 Best Paper.]

この方向性は現在「流行っている」 しかし, ODE的観点が汎化性能等の意 味で本当に意味があり有用であるかは 慎重な考察が必要である。

ResNetの平均場理論と最適制御

- ResNetは入力から出力へ"最短"で繋いでいる.
 - → 最適制御理論による特徴付け (HJB方程式).



[E, Han, Li: A Mean-Field Optimal Control Formulation of Deep Learning. Research in the Mathematical Sciences, 6-10, 2019]
[Benning, Celledoni, Ehrhardt, Owren, Schönlieb: Deep learning as optimal control problems: Models and numerical methods. Journal of Computational Dynamics, 2019, 6(2) : 171-198.]

• 最適制御を用いたResNetの平均場における最適化理論研究もある.

[Lu, Ma, Lu, Lu, and Ying: A mean-field analysis of deep resnet and beyond: Towards provable optimization via overparameterization from depth. ICML2020.]

ResNet型勾配ブースティング法

[Nitanda&Suzuki: Functional Gradient Boosting based on Residual Network Perception. ICML2018] [Nitanda&Suzuki: Gradient Layer: Enhancing the Convergence of Adversarial Training for Generative Models. AISTATS2018]

ResNetと勾配ブースティング法の類似性

Residual Network

スキップコネクションを持つ巨大な深層ニューラルネットワーク. 画像認識タスク等でSOTA.

勾配ブースティング (XGBoost, LightGBM) 予測器についての関数勾配によるブースティング法 (アンサンブル学習法).
 データマイニング系のコンペティションで最も有用とされるモデル.

ResFGB=両者の関連性に着目したブースティング法

[Veit, Wilber, &Belongie 2016, Littwin & Wolf 2016, Weinan 2017, Haber, Ruthotto, & Holtham 2017, Jastrzebski, Arpit, Ballas, Verma, Che, & Bengio 2017, Chang, Meng, Haber, Tung, and Begert 2017. etc]

1層加える=勾配法1反復

$$x$$
 y $h_{j+1} = h_j + F_j(h_j)$

$$h_{j+1} = h_j + F_j(h_j)$$

$$\sum_{\substack{l \\ dh_t \\ dt}} = \sum_{\substack{r \\ (=F_t)}} h_j \circ h_t$$

層を重ねるごとに目的関数を減少 (関数空間での無限次元勾配法) ResFGB: ResNet型勾配ブースティング法¹³⁹ [Nitanda&Suzuki: Functional Gradient Boosting based on Residual Network Perception. ICML2018] • ResNetの特性を備えた勾配ブースティング法

(識別問題の経験損失) $\min_{\phi,W} \mathcal{L}_n(\phi, W) = \widehat{\mathbb{E}}_{x,y}[l(W^{\top}\phi(x), y)].$

特徴写像 ϕ につい $\mathcal{L}_n(\phi) = \min_W \mathcal{L}_n(\phi, W)$ の関数勾配を用いた最適化.

関数勾配 $abla_{\phi} \mathcal{L}_{n}(\phi)$: 訓練データの特徴 $\phi(x)$ の線形分類 可能性を向上させるための方向の群れ.

未知データに適用するには平滑化が必要.
→ カーネル
$$k(x,x') = \iota(\phi(x))^{\mathsf{T}}\iota(\phi(x'))$$
で畳み込む.
 ι はNNで十分な降下方向を与えるよう学習.
 $T_k \nabla_{\phi} \mathcal{L}_n(\phi) = \widehat{\mathbb{E}}_{x,y} \left[\nabla_{\phi} \mathcal{L}_n(\phi)(x)\iota(\phi(x))^{\mathsf{T}} \right] \iota(\phi(\cdot))$

残差ブロック:平滑化された関数勾配. 作り方から直交しない限り常に損失関数の降下方向.



ResFGBの概要

勾配ブースティングの一反復=ResNetの層追加. 関数空間での最適化の観点からマージン分布の最小化を示せる. → 新しい勾配ブースティングの理論に基づいた汎化保証付き深層ResNetの学習.

提案手法ResFGBの数値実験

中~大規模データでの多値識別問題.以下の手法と比較.

Random Feature + SVM, Random Forest, Gradient Boosting (LightGBM)

	Method	LETTER	USPS	ijcnn1	MNIST	COVTYPE	SUSY
提案手法	RESFGB (LOGISTIC)	0.975 (0.0016)	0.954 (0.0006)	0.987 (0.0011)	0.985 (0.0007)	0.968 (0.0017)	0.804 (0.0000)
	ResFGB (smooth hinge)	0.975 (0.0012)	0.950 (0.0022)	0.988 (0.0018)	0.987 (0.0010)	0.965 (0.0058)	0.804 (0.0004)
	SUPPORT VECTOR MACHINE	0.959 (0.0062)	0.948 (0.0023)	0.977 (0.0015)	0.969 (0.0041)	0.824 (0.0059)	0.754 (0.0534)
	RANDOM FOREST	0.964 (0.0012)	0.939 (0.0018)	0.980 (0.0005)	0.972 (0.0005)	0.948 (0.0005)	0.802 (0.0004)
	GRADIENT BOOSTING	0.964 (0.0011)	0.938 (0.0039)	0.982 (0.0010)	0.981 (0.0004)	0.972 (0.0005)	0.804 (0.0005)

SOTAとされるLightGBM以上の精度を確認.

幾つかのデータでは数反復で収束.通常の勾配ブースティングより**効率的な最適化**. ↑関数勾配の平滑化の仕方の差異による.



第3章 深層学習の最適化

深層学習の"学習"



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?



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



誤差逆伝搬法

合成関数

$$f(x;W) = f_{3,W_3}(f_{2,W_2}(f_{1,W_1}(x)))$$

= $f_{3,W_3} \circ f_{2,W_2} \circ f_{1,W_1}(x)$

合成関数の微分

$$\frac{\partial f}{\partial W_1}(x) = \frac{\partial f_{3,W_3}}{\partial f_{2,W_2}} \frac{\partial f_{2,W_2}}{\partial f_{1,W_1}} \frac{\partial f_{1,W_1}}{\partial W_1}(x)$$

深層学習で主に使われる確率的最適化法14

SGD

● モーメンタム SGD (Nesterov の加速法と類似, 一番よく利用されている)

$$g_t = \nabla L(w_t)$$

 $\Delta w_t = \theta \Delta w_{t-1} - (1 - \theta) \eta g_t$
 $w_{t+1} = w_t + \Delta w_t$

• Nesterovの加速法 (凸関数における加速法と同様,パラメータの設定は異なる)

$$g_t = \nabla L(w_t + \theta \Delta w_{t-1}), \Delta w_t = \theta \Delta w_{t-1} - (1 - \theta)\eta g_t, w_{t+1} = w_t + \Delta w_t$$

- AdaGrad (Duchi et al., 2011)
- Adam (Kingma and Ba, 2014): AdaGrad と加速法を組み合わせたような方法. AdaGrad と違い、勾配のノルムを減衰させて次に伝える. モーメンタム SGD と並んでよく使われている.
- RMSprop (Hinton et al.): AdaGrad において勾配のノルムを減衰させて和 を取る方法.




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.

大域的最適性



 <u>線形深層NN</u>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<u>線形NN</u>のみ.

→ 臨界点が大域的最適解であることの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.

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

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 $(M \ge n)$ 二層 $NNf_{(a,W)}(x) = \sum_{m=1}^M a_m \eta(w_m^T x)$ 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 Landscapes. JMLR, 20:1-34, 2019.]







二層目の重みを固定する設定

(Tian, 2017; Brutzkus and Globerson, 2017; Li and Yuan, 2017; Soltanolkotabi, 2017; Soltanolkotabi et al., 2017; Shalev-Shwartz et al., 2017; Brutzkus et al., 2018)

148

$$y = \sum_{j=1}^{k} v_j \eta (w_j^\intercal x + b_j)$$

- Li and Yuan (2017): ReLU,入力はガウス分布を仮定
 > SGDは多項式時間で<u>大域的最適解</u>に収束
 > <u>学習のダイナミクスは2段階</u>
 → 最適解の近傍へ近づく段階 + 近傍での凸最適化的段階
- Soltanolkotabi (2017): ReLU,入力はガウス分布を仮定

 <u>過完備(横幅>サンプルサイズ)なら勾配法で最適解に線形収束</u> (Soltanolkotabi et al. (2017)は二乗活性化関数でより強い帰結)
- Brutzkus et al. (2018): ReLU
 - ▶ 線形分離可能なデータなら過完備ネットワークで動かしたSGDは 大域的最適解に有限回で収束し、過学習しない。

(線形パーセプロトロンの理論にかなり依存)

Li and Yuan (2017): Convergence Analysis of Two-layer Neural Networks with ReLU Activation.

Soltanolkotabi (2017): Learning ReLUs via Gradient Descent.

Brutzkus, Globerson, Malach and Shalev-Shwartz (2018): SGD learns over parameterized networks that provably generalized on linearly separable data.

ーバーパラメトライゼーション オ

横幅が広いと局所最適解が大域的最適解になる.



自由度が上がるため,初期値から最適解 (完全フィット)へ到達しやすい.



149

- 二種類の解析手法
 - Neural Tangent Kernel
 - ➤ Mean-field analysis (平均場解析)

二つのスケーリング

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

• Neural Tangent Kernelのregime (lazy learning)

 $\succ a_j = \mathbf{O}(1/\sqrt{M})$

[Jacot+ 2018][Du+ 2019][Arora+ 2019]

150

・平均場解析のregime
 ▶ a_i = 0(1/M)

[Nitanda & Suzuki (2017), Chizat & Bach (2018), Mei, Montanari, & Nguyen (2018)]

 $%NTKの1/\sqrt{M}$ 自体はそこまで本質ではない、1/Mより大きいことが重要.

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

NTK

 $f_W(x) \simeq (W - W^{(0)})^\top \nabla_W f_{W^{(0)}}(x)$

初期値のスケールが大きいので,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.



NTKと平均場の違い

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

 η : ReLUとする. $a_j = O(1), w_j = O(1/\sqrt{M})$ または $w_j = O(1/M)$ とスケール変換

•各 w_j がO(1/M)だけ動けば、全体としてO(1)の変化(データにフィットできる). •横幅は十分大きく取る: $M \gg n$ (overparameterization)



Neural Tangent Kernel

連続時間ダイナミクスを考える.

$$Model: f_{W}(x) = \sum_{j=1}^{M} a_{j} \eta(w_{j}^{\top}x) \cdot a_{j} \text{ tble} \\ \cdot w_{j} \text{ を学習}$$

$$[Jacot, Gabriel&Hongler, NeurIPS2018]$$

$$\frac{dw_{j}}{dt} = -\nabla_{w_{j}} \hat{L}(f_{W}) \quad (Gradient descent, GD)$$

$$= -\frac{1}{n} \sum_{i=1}^{n} \ell'_{i}(f_{W}(x_{i}))a_{j}\nabla_{w_{j}}\eta(w_{j}^{\top}x_{i}) \quad \overline{\nabla_{w_{j}}\eta(w_{j}^{\top}x_{i})} = x_{i}\eta'(w_{j}^{\top}x_{i})$$

$$(Jacot, Gabriel&Hongler, NeurIPS2018]$$

$$= -\frac{1}{n} \sum_{i=1}^{n} \ell'_{i}(f_{W}(x_{i}))a_{j}\nabla_{w_{j}}\eta(w_{j}^{\top}x_{i}) \quad \overline{\nabla_{w_{j}}\eta(w_{j}^{\top}x_{i})} = x_{i}\eta'(w_{j}^{\top}x_{i})$$

$$(Jacot, Gabriel&Hongler, NeurIPS2018]$$

$$= -\frac{1}{n} \sum_{i=1}^{n} \ell'_{i}(f_{W}(x_{i}))a_{j}\nabla_{w_{j}}\eta(w_{j}^{\top}x_{i}) \quad \overline{\nabla_{w_{j}}\eta(w_{j}^{\top}x_{i})} = x_{i}\eta'(w_{j}^{\top}x_{i})$$

$$(Jacot, Gabriel&Hongler, NeurIPS2018]$$

$$(Jac$$

目的関数の減少速度

$$\frac{\mathrm{d}\hat{L}(f_W)}{\mathrm{d}t} = \frac{1}{n} \sum_{i=1}^n \frac{\mathrm{d}f_W(x_i)}{\mathrm{d}t} \ell'_i(f_W(x_i))$$

$$= -\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \ell'_i(f_W(x_i)) k_W(x_i, x_j) \ell'_j(f_W(x_j))$$

$$= -\frac{1}{n^2} \|\nabla_f \hat{L}(f_W)\|^2_{K_W}$$

$$\leq -\lambda_{\min} \frac{1}{n^2} \|\nabla_f \hat{L}(f_W)\|^2 \quad (\lambda_{\min} \colon \mathcal{T} \ni \Delta \mathcal{T} \mathfrak{I}) \mathfrak{O}_{\frac{1}{2}} \mathfrak{L}(f_W)$$

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

• <u>ランダム初期化</u>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 $(\exp(-\lambda_{\min}t))$

ランダム初期値とNTKの正定値性

$$K_{\infty,i,j} = \mathbb{E}_{w \sim N(0,I)} [x_i^{\top} x_j \eta'(w^{\top} x_i) \eta'(w^{\top} x_j)]$$

(横幅無限大のNTK)

$$||x_i|| = 1, ||x_i - x_j|| \ge \phi \implies K_{\infty} \succeq C\phi n^{-2}$$

Hoeffdingの不等式より
$$P\left(|K_{W^{(0)},i,j} - K_{\infty,i,j}| \le \sqrt{\frac{\log(2/\delta')}{2M}}\right) \ge 1 - \delta'$$

補題

一様バウンドを取って
$$P\left(\|K_{W^{(0)}}-K_{\infty}\|_{\mathrm{F}}^{2} \leq n^{2} rac{\log(2n^{2}/\delta)}{2M}
ight) \geq 1-\delta$$

十分横幅Mが広ければ、ランダム初期化した $K_{W^{(0)}}$ の正定値性が保証される.

Optimization in NTK regime

以下のように初期化する:

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

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 $\sqrt{y^{\mathsf{T}}(K_{W^{(0)}})^{-1}y/n}$ で抑えられる.

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

- <u>訓練誤差0</u>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)

 $f_W(x) = \sum a_j \eta(w_j^\top x)$

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/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

158

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

Frechet 微分 in L₂(P_n): ∇_f L(f)
∇_f L(f) = (ℓ'_i(f(x_i)))ⁿ_{i=1}
L(f + h) = L(f) + ⟨∇_f L(f), h⟩_{L₂(P_n)} + o(||h||²_{L₂(P_n)})
平滑化積分作用素:

$$T_k f(x) := \int k(x, x') f(x') \mathrm{d}P_n(x')$$

 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).

NTK regimeでのSGD

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

(We train both of first and second layers)



目的関数:

$$Y = f^*(X) + \epsilon$$
 (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 $j \in \{1, ..., M\}$: $a_j^{(t+1)} = a_j^{(t)} - \alpha_t [\nabla_a \ell(y_t, f_{a^{(t)}, W^{(t)}}(x_t)) + \lambda(a^{(t)} - a^{(0)})]$ $W_j^{(t+1)} = W_j^{(t)} - \alpha_t [\nabla_W \ell(y_t, f_{a^{(t)}, W^{(t)}}(x_t)) + \lambda(W^{(t)} - W^{(0)})]$ end for Return $\bar{a}^{(T)} = \frac{1}{T} \sum_{t=0}^{T-1} a^{(t)}, \ \bar{W}^{(T)} = \frac{1}{T} \sum_{t=0}^{T-1} W^{(t)}.$

NTKにおける余剰誤差の速い収束

160

[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.





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:

$$\mu_{j} = O(j^{-\beta})$$
.
 $j = \lambda_{\infty}(x, x') = \sum_{j=1}^{\infty} \mu_{j}\phi_{j}(x)\phi_{j}(x')$

$$\frac{1}{\sqrt{2}}$$

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).





NTKの固有値固有関数分解 (ϕ_m)[∞]_{m=1}:固有関数. L₂(P_X)内の 正規直交基底.

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)]

足される.

高周波成分

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

Beyond kernel

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)

平均場解析

ニューラルネットワークの最適化をパラメータの分布最適化としてみなす。

$$f(x) = \frac{1}{M} \sum_{j=1}^{M} a_j \eta(w_j^{\top} x) \xrightarrow{M \to \infty} \int a \eta(w^{\top} x) \rho(a, w) dadw$$



$$f$$
の最適化 $\Leftrightarrow \rho$ の最適化

$$\frac{\partial \rho_t}{\partial t} = -\nabla \cdot (v_t \rho_t)$$
 連続方程式

Wasserstein勾配流

[Atsushi Nitanda and Taiji Suzuki: Stochastic Particle Gradient Descent for Infinite Ensembles. arXiv:1712.05438.]

$$v_t(a,w) = -\frac{1}{n} \sum_{i=1}^n \nabla_{(a,w)}(a\eta(w^{\top}x_i))\ell'(y_i, f_{\rho_t}(x_i))$$
 (各粒子は勾配降下方向へ移動)

MMDとの関係

MMD: Maximum Mean Discrepancy

[Gretton et al., NIPS2006]

$f_{\nu^*}(x)$:真の関数

$$\mathbb{E}_{X}[(f_{\nu}(X) - f_{\nu^{*}}(X))^{2}] = \mathbb{E}_{X}\left[\int \int a\eta(w^{\top}X)a'\eta(w'^{\top}X)(d\nu(a,w)d\nu(a',w') - 2d\nu(a,w)d\nu^{*}(a',w') + d\nu^{*}(a,w)d\nu^{*}(a',w'))\right] \\ = \int \int \mathbb{E}_{X}\left[a\eta(w^{\top}X)a'\eta(w'^{\top}X)\right](d\nu(a,w)d\nu(a',w') - 2d\nu(a,w)d\nu^{*}(a',w') + d\nu^{*}(a,w)d\nu^{*}(a',w'))\right]$$

k((a,w),(a',w')) = $\langle \phi_k(a,w), \phi_k(a',w') \rangle_{\mathcal{H}_k}$ ↑ 正定値カーネルになっている.

 $= \langle \mathbb{E}_{\nu}[\phi_{k}(a,w)], \mathbb{E}_{\nu}[\phi_{k}(a,w)] \rangle_{\mathcal{H}_{k}} - 2 \langle \mathbb{E}_{\nu}[\phi_{k}(a,w)], \mathbb{E}_{\nu^{*}}[\phi_{k}(a,w)] \rangle_{\mathcal{H}_{k}} \\ + \langle \mathbb{E}_{\nu^{*}}[\phi_{k}(a,w)], \mathbb{E}_{\nu^{*}}[\phi_{k}(a,w)] \rangle_{\mathcal{H}_{k}}$

 $= \left\| \mathbb{E}_{\nu}[\phi_k] - \mathbb{E}_{\nu^*}[\phi_k] \right\|_{\mathcal{H}_k}^2 : \mathsf{MMD}$

L2距離最小化⇔ MMD最小化

[Arbel et al. arXiv:1906.04370][Sonoda, arXiv:1902.00648]





M→∞の極限で,最適解への収束が成り立つ場合がある. [Nitanda&Suzuki, 2017][Chizat&Bach, 2018][Chizat, 2019]

> ノイズありのダイナミクス: McKean-Vlasov過程 [Mei, Montanari&Nguyen, 2018]

Wasserstein距離について

μ, ν :距離空間(\mathcal{X}, c)上の確率測度($\exists x \exists x \exists r b land 2 \exists$)

$$W_p(\mu,\nu) = \left(\inf_{\pi \in \Pi(\mu,\nu)} \int_{\mathcal{X} \times \mathcal{X}} c(x,y)^p \mathrm{d}\pi(x,y)\right)^{1/p}$$

Π(μ,ν): 周辺分布がμ,νであるX × X上の同時分布の集合
周辺分布を固定した同時分布の中で最小化

$$(\mathcal{X} = \mathbb{R}^d: c(x, y) = ||x - y||)$$

- 分布のサポートがずれていてもwell-defined
- 底空間の距離が反映されている ※KL-divergenceは距離が反映されない.

(双対表現: Kantorovich双対)

$$\inf_{\pi \in \Pi(\mu,\nu)} \int c(x,y)^p d\pi(x,y) = \sup_{\psi,\phi} \left\{ \int \psi d\mu + \int \phi d\nu \mid \psi(x) + \phi(y) \le c(x,y)^p \right\}$$

「輸送距離」とも言われる

W2距離と粒子勾配降下法の関係

168

$$\begin{split} W_{2} 距離による近接点アルゴリズムを考える: f_{\nu}(x) = \int h(w,x) d\nu(w) \\ L(\nu) + \frac{W_{2}^{2}(\mu,\nu)}{2\delta} \\ &= \mathbb{E}_{X} \left[\ell \left(\int h(w,X) d\nu(w) \right) \right] + \frac{W_{2}^{2}(\mu,\nu)}{2\delta} \\ &\simeq \left\langle \frac{\ell' \left(\int h(w,X) d\mu(w) \right)}{E^{2}(\mu,\nu)} \right\rangle \int h(v,X) d(\nu-\mu)(v) \right\rangle_{L^{2}(P)} + \inf_{\pi \in \Pi(\mu,\nu)} \int \frac{||w-v||^{2}}{2\delta} d\pi(w,v) \\ &= \inf_{\pi \in \Pi(\mu,\nu)} \int \left(\left\langle \Delta_{\mu}, h(v,\cdot) \right\rangle_{L^{2}(P)} + \frac{||w-v||^{2}}{2\delta} \right) d\pi(w,v) + \text{const.} \end{split}$$

$$v = \underset{v'}{\operatorname{arg\,min}} \left\{ \langle \Delta_{\mu}, h(v', \cdot) \rangle_{L^{2}(P)} + \frac{\|w - v'\|^{2}}{2\delta} \right\}$$

$$\simeq w - \delta \nabla_{w} \langle h(w, \cdot), \ell'(f_{\mu}) \rangle$$

: 最急降下法

- 各粒子ごとにみると単純な最急降下法.
- 粒子勾配降下法は W_2 距離を近接項とした近接点アルゴリズムの一次近似 $\rightarrow \delta \rightarrow 0$ の極限 (連続時間): <u>Wasserstein gradient flow</u>

連続の方程式と勾配流

「連続の方程式」
$$\frac{\partial \rho_t}{\partial t} = -\nabla \cdot (v_t \rho_t)$$
 の意味

$$\frac{\mathrm{d}}{\mathrm{d}t} \int f(w) \mathrm{d}\rho_t(w) = \int (\nabla f(w))^\top v_t(w) \mathrm{d}\rho_t(w)$$
$$(\forall f: \exists \geq n^2 \not = h, \ \mathcal{C}^{\infty} - \mathcal{W})$$

 今, ρ_t は写像 $T_t: \mathbb{R}^d \to \mathbb{R}^d$ による ρ_0 の<u>押し出し</u>であるとする: $\rho_t = T_{t\#}\rho_0$. つまり, $w \sim \rho_0$ に対する $T_t(w)$ の分布が ρ_t であるとする.
 写像 T_t を生成するベクトル場を $\frac{dT_t}{dt}(w) = v_t(T_t(w))$ とする.

$$\frac{\mathrm{d}}{\mathrm{d}t} \int f(w) \mathrm{d}\rho_t(w) = \frac{\mathrm{d}}{\mathrm{d}t} \int f(T_t(w)) \mathrm{d}\rho_0(w)
= \int \nabla f(T_t(w))^\top \frac{\mathrm{d}T_t(w)}{\mathrm{d}t} \mathrm{d}\rho_0(w)
= \int \nabla f(T_t(w))^\top v_t(T_t(w)) \mathrm{d}\rho_0(w)
= \int \nabla f(w)^\top v_t(w) \mathrm{d}\rho_t(w). \quad (連続の方程式)$$

 $w_t = T_t(w)$ に対し、 $v_t(w_t) = -\nabla_w \langle h(w, \cdot), \ell'(f_{\rho_t}) \rangle|_{w=w_t}$ としたのが前ページの更新式.





- $\rho_t = T_{t\#}\rho_0$
- $\frac{\mathrm{d}T_t}{\mathrm{d}t}(w) = v_t(T_t(w))$
- ある ϕ_t を用いて $v_t = \nabla \phi_t$ と書けるとする. この時,以下が成り立つ:

$$\lim_{\delta \to 0} \frac{W_2(\rho_{t+\delta}, (\mathrm{id} + \delta v_t)_{\#} \rho_t)}{\delta} = 0$$

詳細は以下を参照:

Ambrosio, Gigli, and Savaré. *Gradient Flows in Metric Spaces and in the Space of Probability Measures*. Lectures in Mathematics. ETH Zürich. Birkhäuser Basel, 2008.







Brenierの定理

同条件のもと

ρ_{0}, ρ_{1} が確率密度関数を持つ時、以下が成り立つ: $W_{2}^{2}(\rho_{0}, \rho_{1}) = \inf_{T:T_{\#}\rho_{0}=\rho_{1}} \mathbb{E}_{X \sim \rho_{0}}[\|X - T(X)\|^{2}]$

- Infを達成する写像*T**が存在する.
- しかも、ある凸関数 ψ が存在して $T^*(x) \in \partial \psi(x)$ と書ける.
- このT*を<u>最適輸送写像</u>という.

Benamou-Brenier formula (連続の方程式とW2距離の関係):

$$W_2^2(\rho_0, \rho_1) = \inf_{\{v_t\}_t} \int_0^1 \|v_t\|_{L_2(\rho_t)}^2 \mathrm{d}t$$

ただし、 $\inf lap_0 n \delta \rho_1 \land \bar{\rho}_1 \land \bar{\rho}$

• $\rho_t = T_{t\#}\rho_0$ • $\frac{\mathrm{d}T_t}{\mathrm{d}t}(w) = v_t(T_t(w))$



- Wasserstein勾配流は W_2 距離を用いた近接点ア ルゴリズムで特徴づけられることが分かった.
- 目的関数がW2距離に関する凸性(displacement convexity)が成り立つなら、大域的最適解への 収束が示せる (エントロピーなど):

$W_2(\rho_t, \rho^*) \le e^{-\lambda t} W_2(\rho_0, \rho^*).$

- しかし、NNの最適化では凸性は成り立たない。
 そのため、大域収束を示すことが難しい。
- •もっとも、局所的には凸性が成り立ちうる.

例:スパースな最適解 [Chizat, 2019]





局所最適性条件

定理 (Nitanda&Suzuki, 2017)

ある解 $\hat{\mu}$ がコンパクトな台の確率密度関数を持つとする. この時,ある μ^* s.t. $L(\mu^*) < L(\hat{\mu})$ が存在して

- supp(*µ*^{*}) ⊆ supp(*µ*̂)かつ*µ*^{*}は確率密度を持つ, or
- μ^* は密度を持たず $supp(\mu^*)$ は $supp(\hat{\mu})$ に内部に含まれる, が満たされるとき,降下方向が存在して粒子降下法によって目的関 数値を減らすことができる.



平均場解析と陰的正則化

二値判別をexp-損失を用いて解く (ラベルノイズなしとする):

 $\min_{\rho} \sum_{i=1}^{n} \exp\left(-y_i f_{\rho}(x)\right) \qquad \text{ただし} \qquad f_{\rho}(x) = \int \eta(w^{\top} x) d\rho(w)$ 符号付測度の中で最適化

平均場解析の設定で最適化する. 初期値が小さいので判別に必要なニューロンだけが「生えてくる」.



0.0

-0.5

0.5

→ スパースな解:陰的正則化

[Chizat&Bach:Implicit Bias of Gradient Descent for Wide Two-layer Neural Networks Trained with the Logistic Loss. COLT2020.]

最適化の結果として「単純な」 解が求まってしまう.

判別平面はL1-正則化解マージン最大化元に収束する: $\max_{\rho:\|\rho\|_{\mathcal{F}_1}} \min_{i \in \{1,...,n\}} y_i f_{\rho}(x_i) \qquad \|\rho\|_{\mathcal{F}_1} = |\rho|(\mathbb{R}^d)$

勾配法と陰的正則化(再掲)

 小さな初期値から勾配法を始めるとノルム最小 化点に収束しやすい→陰的正則化



[Gunasekar et al.: Implicit Regularization in Matrix Factorization, NIPS2017] [Soudry et al.: The implicit bias of gradient descent on separable data. JMLR2018] [Gunasekar et al.: Implicit Bias of Gradient Descent on Linear Convolutional Networks, NIPS2018] [Moroshko et al.: Implicit Bias in Deep Linear Classification: Initialization Scale vs Training Accuracy, arXiv:2007.06738]

各regimeにおける陰的正則化

各regimeにおける陰的正則化の種類

Regime	対応する正則化
NTK, カーネル法 with early stopping	L2-正則化
平均場理論	L1-正則化

- ニューラルネットワークの学習では様々な「陽的正則化」を用いる:
 バッチノーマリゼーション, Dropout, Weight decay, MixUp, ...
- 一方で、深層学習の構造が自動的に生み出す「陰的正則化」も強く効いていると考えられる。
 - → オーバーパラメタライズしても過学習しない.

ノイズあり勾配法と大域的最適性

Sharp minima vs flat minima



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 ⇒ 解にノイズを乗せている ⇒ 目的関数の平滑化

 $\begin{aligned} x_t &= x_{t-1} - \eta (\nabla L(x_{t-1}) + \xi_t) & (y_t = x_t + \eta \xi_t) \\ \Rightarrow y_t &= y_{t-1} - \eta \xi_{t-1} - \eta \nabla L(y_{t-1} - \eta \xi_{t-1}) \\ \Rightarrow \mathbb{E}_{\xi_{t-1}}[y_t] &= y_{t-1} - \eta \nabla \mathbb{E}_{\xi_{t-1}}[L(y_{t-1} - \eta \xi_{t-1})] \end{aligned}$

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

関連研究: Graduated optimization

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.

$$\sigma$$
-nice性の導入. 多項式オーダーでの収束
 $\hat{L}_{\delta}(x) = \mathrm{E}_{u \sim U(\mathrm{B}(\mathrm{R}^d))}[L(x + \delta u)]$

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.



180
GLD/SGLD

• Stochastic Gradient Langevin Dynamics (SGLD)

$$\begin{split} \min_{x \in \mathbb{R}^d} L(x) &= \min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell_i(x) \qquad (非凸) \\ \\ \mathbf{d}X_t &= -\nabla L(X_t) \mathbf{d}t + \sqrt{2\beta^{-1}} \mathbf{d}B_t \qquad (勾配Langevin動力学) \\ \\ & \hat{\mathbf{c}} \ensuremath{\hat{\pi}} \ensuremath{\hat{\pi}} : \ \pi \propto \exp(-\beta L(X)) \\ \\ \\ \mathbf{d}X_t &= X_t - \eta \nabla L(X_t) + \sqrt{2\eta\beta^{-1}} \xi_t \qquad (Euler-Maruyama fill) \\ \\ & \mathbf{GLD}: \quad X_{t+1} = X_t - \eta \nabla L(X_t) + \sqrt{2\eta\beta^{-1}} \xi_t \qquad (Euler-Maruyama fill) \\ \\ & \mathbf{f}_t \sim N(0, I) \\ \\ \\ & \mathbf{SGLD}: \quad X_{t+1} = X_t - \eta \frac{1}{|I_B|} \sum_{i \in I_B} \nabla \ell_i(X_t) + \sqrt{2\eta\beta^{-1}} \xi_t \\ \hline \\ & \overline{\text{d}x} \ensuremath{\hat{\pi}} \ensur$$





収束定理 (有限次元)

- f_i : 有界, Lipschitz連続, 滑らかな勾配 $\|\ell_i\|_{\infty} \leq A, \|\nabla \ell_i\|_{\infty} \leq B, \|\nabla \ell_i(x) - \nabla \ell_i(y)\| \leq M \|x - y\|$
- <u>散逸条件</u>:

$$\langle \nabla L, w \rangle \ge m \|w\|^2 - b \quad (\forall w \in \mathbb{R}^d)$$

(+ その他細かい条件)

Thm [Raginsky, Rakhlin and Telgarsky, COLT2017]

$$E[L(X_k)] - L(X^*) \le \tilde{O}\left((\beta + d)k\eta^{5/4} + \frac{\beta + d}{\sqrt{\lambda^*}} \exp\left(-\tilde{\Omega}\left(\frac{\lambda^* k\eta}{\beta(d + \beta)}\right)\right) + \frac{d\log(\beta + 1)}{\beta}\right)$$

- λ_* はスペクトルギャップと言われる量. $\rightarrow \underline{\lambda_* \alpha \gamma} - \underline{\lambda_* \alpha \gamma}$
- 逆温度パラメータが十分大きくて、更新を十分な回数回せば最適解付近に近づける。
- Xu et al. (NeurlPS2018) は収束レートを改善しているが, 証明にいくつかの 間違いあり.





対数Sobolev不等式

 $\pi_{\infty}(dx) \propto \exp(-\beta L(x)) dx$:連続時間ダイナミクスの定常分布

〔2. 平滑性

対数Sobolev不等式

$$d\nu = f \, d\pi_{\infty} \quad \text{(probability)}$$
$$\int f \log(f) d\pi_{\infty} \leq 2c_{\text{LS}} \int \frac{\|\nabla f\|^2}{f} d\pi_{\infty} \quad (D(\nu || \pi_{\infty}) \leq 2c_{\text{LS}} I(\nu || \pi_{\infty}))$$

Geometric ergodicity $\rho_t: X_t$ の周辺分布 $D(\rho_t || \pi_\infty) \le \exp(-2t/c_{\text{LS}})D(\rho_0 || \pi_\infty)$

定常分布へ線形収束

[Bakry, Gentil, and Ledoux: Analysis and Geometry of Markov Diffusion Operators. Springer, 2014. Th. 5.2.1]

連続の方程式再び

• 勾配ランジュバン動力学に対応する連続の方程式

$$\partial_t \rho_t = \nabla \cdot \left(\rho_t \nabla \log(\rho_t / \pi_\infty) \right)$$

 勾配ランジュバン動力学は相対エントロピー (KL-ダイ バージェンス)をWasserstein勾配流で最適化しているこ とに対応:

$$D(\rho || \pi_{\infty}) = \int \log \left(\frac{\mathrm{d}\rho}{\mathrm{d}\pi_{\infty}} \right) \mathrm{d}\rho$$

Remark:

通常の (相対でない) エントロピーは W_2 -距離に関して凸 (displacement convexity). つまり、 W_2 -距離に関する測地線上で凸関数になる.



[Muzellec, Sato, Massias, Suzuki, arXiv:2003.00306][Suzuki, arXiv:2007.05824]



E.g., Bayesian optimization on infinite dimensional space

[Zimmermann and Toussaint. Bayesian functional optimization. AAAI, 2018] [Vellanki, Rana, Gupta, de Celis Leal, Sutti, Height, and Venkatesh: Bayesian functional optimisation with shape prior. AAAI, 2019]



例: NNの学習

ldea: 分布の学習 → **輸送写像の学習**



2層NNの学習: 直接表現

$$L(W) = \frac{1}{n_{\rm tr}} \sum_{i=1}^{n_{\rm tr}} \ell_i(f_W(x_i)) + \frac{\lambda_0}{2} \|W\|_{\rm F}^2$$

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

$$\begin{cases} \bullet \ a_j \leq j^{-\gamma} \text{ for } \gamma > 1/2 \\ \bullet \ \eta \text{ is a smooth activation, e.g., sigmoid.} \end{cases}$$



NTKと違い,
$$a_j$$
はデータサイズ
にも横幅にも依存させずスケー
ルを固定できる.
(TNKは $a_j = 1/\sqrt{M}$ とする)

無限次元ランジュバン動力学
$$x = \sum_{j=1}^{\infty} x_j f_j \in \mathcal{H}$$

$$\min_{x \in \mathcal{H}} L(x) \implies \min_{x \in \mathcal{H}} \left\{ L(x) + \frac{\lambda}{2} \|x\|_{\mathcal{H}_K}^2 \right\} \qquad \mathcal{H}_K : \mathsf{RKHS with kernel } K.$$

$$\mathcal{H}_K \hookrightarrow \mathcal{H}$$

$$\mathrm{d}X_t = -\nabla\left(L(X_t) + \frac{\lambda}{2} \|X_t\|_{\mathcal{H}_K}^2\right) \mathrm{d}t + \sqrt{\frac{2}{\beta}} \mathrm{d}\xi_t$$

ノルム: For
$$x=\sum_{j=1}^\infty x_j f_j\in \mathcal{H}$$
 , we let $\|x\|_{\mathcal{H}_K}^2=\sum_{j=1}^\infty \mu_j^{-1}x_j^2$ where $\mu_j\sim j^{-2}$.

Cylindrical Brownian motion: $\xi_t = \sum_{j=1}^{\infty} \xi_{j,t} f_j$

時間離散化:

$$\begin{pmatrix}
X_{n+1} = S_{\eta} \left(X_n - \eta \nabla L(X_n) + \sqrt{2\frac{\eta}{\beta}} \xi_n \right) & \left(S_{\eta} := (I + \eta \lambda A)^{-1} \right) \\
(準陰的Eulerスキーム) & A = \operatorname{diag}(\mu_1^{-1}, \mu_2^{-1}, \ldots) \\
\xi_n = \sum_{j=1}^{\infty} \gamma_{n,j} f_j \text{ where } \gamma_{n,j} \sim N(0, 1) \text{ (i.i.d.).}
\end{cases}$$

定常分布

$$dX_t = -\nabla \left(L(X_t) + \frac{\lambda}{2} \|X_t\|_{\mathcal{H}_K}^2 \right) dt + \sqrt{\frac{2}{\beta}} d\xi_t$$

$$\frac{\mathrm{d}\pi_{\infty}}{\mathrm{d}\mu_{*}}(x) \propto \exp\left(-\beta L(x)\right)$$

 $\mu_* = N(0, C)$ (Hilbert空間上のガウス過程) where $C = (\beta \lambda)^{-1} \operatorname{diag}(\mu_0, \mu_1, \dots).$

$$\pi_{\infty}(x) \propto \exp\left(-\beta L(x) - \frac{1}{2}x^{\top}C^{-1}x\right)$$
 と解釈しても良い.

(無限次元)勾配ランジュバン動力学の定常分布は
 <u>ガウス過程事前分布を用いたベイズ事後分布</u>に対応する.
 → 過学習を防ぎ汎化する [Suzuki, arXiv:2007.05824]

無限次元の設定

ヒルベルト空間

$$\mathcal{H} = \left\{ \sum_{k=0}^{\infty} \alpha_k f_k \mid \sum_{k=0}^{\infty} \alpha_k^2 < \infty \right\}$$

 $\langle x, y \rangle = \sum_{k=0}^{\infty} \alpha_k \beta_k$ for $x = \sum_k \alpha_k f_k, \ y = \sum_k \beta_k f_k.$

RKHS構造

$$\mathcal{H}_K = \left\{ \sum_{k=0}^{\infty} \alpha_k f_k \mid \sum_{k=0}^{\infty} \alpha_k^2 / \mu_k < \infty \right\}$$

 $\langle x, y \rangle_{\mathcal{H}_K} = \sum_{k=0}^{\infty} \alpha_k \beta_k / \mu_k \quad \text{for } x = \sum_k \alpha_k f_k, \ y = \sum_k \beta_k f_k.$

仮定(固有値の減少)
$$\mu_k \simeq k^{-2}$$

(あまり本質的ではない. $\mu_k \sim k^{-p} (p > 1)$ としても良い.)

$$\min_{x \in \mathcal{H}} L(x) = \min_{x \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \ell_i(x) + \left(\frac{\lambda_0}{2} \|x\|^2\right)$$

$$\mathcal{H}$$

 \mathcal{H}_{K}

Assumption (1)

- It either holds:
 - (Strict Dissipativity) $\lambda > M\mu_0$, or (強):強凸
 - (Bounded gradients) $\|\nabla L(\cdot)\| \leq B$, for B > 0. (33)



Assumption (2)





 π_{∞} :定常分布

Thm (informal) [Muzellec, Sato, Massias, Suzuki, 2020]

上記の条件のもと,次が成り立つ:

$$L(X_n) - \int L(x) d\pi_{\infty}(x) \lesssim \exp\left(-\Lambda_{\eta}^* n\eta\right) + \frac{c_{\beta}}{\Lambda_0^*} \eta^{1/2-\kappa}$$
(geometric ergodicity + time discretization)

ただし $\kappa > 0$ は任意の正の実数, $c_{\beta} = \sqrt{\beta}$ (有界な勾配), $c_{\beta} = 1$ (強散逸条件).

Remark: $\int L(x) d\pi_{\infty}(x) \simeq L(\tilde{x}) \quad \text{for} \quad \tilde{x} := \operatorname*{arg\,min}_{x \in \mathcal{H}} \left\{ L(x) + \frac{\lambda}{2} \|x\|_{\mathcal{H}_{K}}^{2} \right\}$

証明は以下の論文のテクニックを援用: Brehier 2014; Brehier&Kopec 2016; Mattingly et al., 2002; Goldys&Maslowski, 2006.

誤差の解析(2)



 Λ_{η}^* : スペクトルギャップ, β に対して指数的依存がある.

証明は以下の論文のテクニックを援用: Brehier 2014; Brehier&Kopec 2016; Mattingly et al., 2002; Goldys&Maslowski, 2006.

• 深層学習の最適化への応用と汎化誤差解析: Suzuki, arXiv:2007.05824.

ノイズのコントロール

- 大域的最適解を得るためには $\beta \rightarrow \infty$ が必要.
- スペクトルギャップはβに指数的に依存.
- 大域的最適解まわりで局所的に凸になっていて、 離れた場所より目的関数値が真に小さければ途 中で勾配法に切り替えても良い。
- 例えば2層NNでは訓練誤差の形状が局所的に 強凸になることがある [Li and Yuan, 2017][Chizat, 2019]
 (各ニューロンが適度にばらけている場合はそうなる)





 $|\mathbb{E}[\phi(X_n)] - \phi(x^*)| \leq ?$ for a smooth function ϕ .



-項のバウンド 弟

 $\mathbb{E}[\phi(X_n) - \phi(x^*)] = \left[\mathbb{E}[\phi(X_n) - \phi(X^{\mu_\eta})]\right] + \mathbb{E}[\phi(X^{\mu_\eta}) - \phi(X^{\pi_\infty})] + \mathbb{E}[\phi(X^{\pi_\infty}) - \phi(x^*)]$

補題 (離散時間ダイナミクスのGeometric ergodicity)

ある定常分布 μ_{η} がだた一つ存在して (極限分布), geometric ergodicity (定常分布への線形収束) が成り立つ:

 $\mathbb{E}[\phi(X_n) - \phi(X^{\mu_\eta})] \le C(1 + ||x_0||) \exp\left(-\Lambda_{\eta}^* n\eta\right)$

ただし, "スペクトルギャップ" Λ^*_η は以下のように与えられる,

(i) (Strict dissipative) $\Lambda_{\eta}^{*} = \frac{\frac{\lambda}{\mu_{0}} - M}{1 + \eta \frac{\lambda}{\mu_{0}}}$ (ii) (Bounded gradient) $\Lambda_{\eta}^{*} = C \min\left(\frac{\lambda}{2\mu_{0}}, \frac{1}{2}\right) \delta$ for $\delta = \exp(-O(\beta))$

 $X^{\mu_{\eta}}$: r.v. obeying μ_{η} $X_0 = x_0$ (constant)

- 有限次元の場合と違い, 強平滑条件がないとおそらく成り立たない.
- Coupling argument: Lyapunov条件, majorization条件より (Mattingly et al. (2002)とGoldys&Maslowski (2006)のテクニックを合わせる)

Geometric ergodicity

• Coupling argument



第二項のバウンド

 $\mathbb{E}[\phi(X_n) - \phi(x^*)] = \mathbb{E}[\phi(X_n) - \phi(X^{\mu_\eta})] + \mathbb{E}[\phi(X^{\mu_\eta}) - \phi(X^{\pi_\infty})] + \mathbb{E}[\phi(X^{\pi_\infty}) - \phi(x^*)]$

 $X^{\mu_{\eta}}:$ <u>離散時間</u>ダイナミクスの定常分布 $X^{\pi}:$ <u>連続時間</u>ダイナミクスの定常分布 (存在と一意性は保証されている)

補題(連続・離散時間ダイナミクスの定常分布の違い)

任意の $0 < \kappa < 1/2$ に対し、ある定数Cが存在して、

$$\mathbb{E}[\phi(X^{\mu_{\eta}}) - \phi(X^{\pi})]| \le C \|\phi\|_{0,2} \frac{c_{\beta}}{\Lambda_0^*} \eta^{1/2-\kappa}$$

- $\|\phi\|_{0,2} = \max\{\|\phi\|_{\infty}, \|D\phi\|_{\infty}, \|D^2\phi\|_{\infty}\}$
- $c_{\beta} = \sqrt{\beta}$ for bounded gradient condition, and $\beta = 1$ otherwise
- Malliavin解析
- ステップサイズηを0に近づけると、離散時間ダイナミク スが連続時間ダイナミクスに近づく.
- βはΛ₀に影響している.

有限次元バージョンとの関係

$$\mathcal{H}_{K} = \left\{ \sum_{k=0}^{\infty} \alpha_{k} f_{k} \mid \sum_{k=0}^{\infty} \alpha_{k}^{2} / \mu_{k} < \infty \right\}$$
• $\mu_{k} \simeq 1/k^{2}$ (我々の状況)

$$|\mathbb{E}[\phi(X_{n}) - \phi(X^{\pi})]| \leq C \left[\exp\left(-\Lambda_{\eta}^{*} n\eta\right) + \frac{c_{\beta}}{\Lambda_{0}^{*}} \frac{\eta^{1/2-\kappa}}{\dots} \right] \quad \text{(optimal)}$$
• $\mu_{k} \simeq 1/k^{p}$ (予想) see [Andersson,Kruse&Larsson, 2016] for finite time horizon.
 p が大きくなるほど関数クラスは"単純"になる.

$$|\mathbb{E}[\phi(X_{n}) - \phi(X^{\pi})]| \leq C \left[\exp\left(-\Lambda_{\eta}^{*} n\eta\right) + \frac{c_{\beta}}{\Lambda_{0}^{*}} \frac{\eta^{p-1} - \kappa}{\dots} \right]$$
有限次元の解析は $p \to \infty$ に対応 (定数を無視すれば):

 $\left|\mathbb{E}[\phi(X_n) - \phi(X^{\pi})]\right| \le C \left[\exp\left(-\Lambda_{\eta}^* n\eta\right) + \frac{c_{\beta}}{\Lambda_0^* \dots}\right]$

 $\simeq 1/k^{\nu}$

203

[Xu et al. (2018)]

判別問題における速い収束

Assumption

- ・強低ノイズ条件:
 - $|P(Y = 1|X) 1/2| \ge \delta$ (a.s.)

• $\operatorname{supp}(P_X) \subset [0,1]^d$ and P_X has density p such that $p(x) \ge c_0 \quad (\forall x \in \operatorname{supp}(P_X)).$

●活性化関数はなめらか:

 $\sigma \in \mathcal{C}^m(\mathbb{R}) \quad \text{for } 2m > d$

•真の関数はモデルに入っているとする: $f^* = f_{W^*}$.

十分大き $\alpha n \ge \beta \le n$ に対し,

$$\mathbb{E}[P_{\pi_k}(\{W_k \in \mathcal{H} \mid P_X[\operatorname{sign}(f_{W_k}(X)) = \operatorname{sign}(f^*(X))] \neq 0\})]$$

$$\lesssim \exp(-c\beta\delta^{2m/(2m-d)}) + \frac{\Xi_k}{\delta^{2m/(2m-d)}}$$

ベイズ最適な判別機が高い確率で求まる. (Butz most stock)



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



- 深層学習の理論
 - ▶ 表現能力
 ▶ 汎化能力
 ▶ 最適化能力
- 表現力
 - 万能近似性
 - 層を深くすることで指数的に表現力増大
- 汎化能力
 - 適応的な学習手法を実現→真の関数に非凸性・スパース性があれば、カーネル法のような特徴写像を固定する方法に優越
 - 陰的正則化等でoverparametrizedな状況でも汎化
- 最適化理論
 - Overparametrizeされていれば大域的最適解を得る
 - Neural Tangent Kernelと平均場
 - スケールの仕方によって凸らしさが変わる→最適化の難しさと陰的正則化の種類が変わる(L2 vs L1)

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」



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