深層学習の数理: 統計力学的アプローチ

産業技術総合研究所 人工知能研究センター 機械学習研究チーム研究員

唐木田 亮

Deep learning and Physics 2020 June. 18

目次

・イントロダクション

- ・統計力学的アプローチの紹介
- 深層学習とランダムネス
- 平均場理論
- ランダム行列理論
- ・さらなる発展: Fisher情報行列とNeural Tangent Kernel
 - Fisher情報行列と学習率

[Karakida, Akaho & Amari, Universal Statistics of Fisher Information in Deep Neural Networks: Mean Field Approach, AISTATS 2019]

- Neural Tangent Kernelとの対応
- Batch Normalizationの解析

[Karakida, Akaho & Amari, *The Normalization Method for Alleviating Pathological Sharpness in Wide Neural Networks*, NeurIPS 2019]





Deep Network Network (DNN):

Fully-connected, CNN, ResNet, Transformer ...

使いやすい訓練法: SGD, Adam, K-FAC ... Dropout, Batch Normalization ...

+ Adversarial attacks & defences, 解釈性, contrastive learning ...

深層生成モデル: VAE, GAN, Flow-based … 既存の機械学習手法との融合: 強化学習, カーネル法…

DNNを基盤とした機械学習の集合





Deep Network Network (DNN):

Fully-connected, CNN, ResNet, Transformer ···

自然な疑問: 深層ニューラルネットワークはなぜ/どのような設定で 性能が高いのか

"手持ちのデータで性能が出ない"

原因は, データ?正規化?層数?アルゴリズム?ステップ数?…

数理の必要性、近年発展が著しい、大きく分けると,

表現能力 (Expressivity/Representation power), 訓練性 (Trainability), 汎化能力 (Generalization)

表現能力:次元の呪いと階層性

- 万有近似性 (Universal Approximation) [1990年代~]
 3層(shallow)ニューラルネットは入力を任意の精度で近似可能
- Barronの定理 (1993)

(ある特定のクラスの) 関数f(x) に対して, $\mathrm{E}_x[(f(x)-f_M(x))^2] \leq C_f/M$ を満たす W_1, W_2 が存在

一方, 基底(
$$W_1$$
)を学習しない場合,
 $E_x[(f(x) - f_M(x))^2] \ge C'_f/M^{2/D}$

 $f_M(x) = W_2$ Sigmoid(W₁x) 入力x 入力次元 D

階層性が近似精度における次元の呪いを克服.3層で十分?

表現能力

Deep vs. Shallow

非線形性や入力の種類(位相)に応じた様々な理論

- The number of monomials [Delalleau, Bengio, NIPS2011]
- The number of *linear regions* [Montufar+, NIPS2014]
- Betti numbers [Bianchini, Scarselli, IEEE (2014)] etc...
- Barronの定理の拡張 [Lee+ COLT2017]

階層型ニューラルネットが表現できる関数の 複雑さは, 幅に対してベキ的に増加, 層数に対して**指数関数的に増加**

非線形変換





有限の計算資源(メモリ)では,深層モデルの方が効率的

訓練性の問題



- 一般には, global minimaにたどり着く保証はない
- DNNは勾配が消失 (or 発散)しやすい

$$D_l W_l^{\top} D_{l+1} W_{l+1}^{\top} \cdots W_L^{\top} (y-f)$$

訓練性とloss Landscape

[RK, Okada, Amari, Neural Networks (2016)]

統計モデル Hidden h $p(\mathbf{h}, \mathbf{v}) = \exp(\mathbf{h}^T W \mathbf{v} - ||\mathbf{h}||^2 / 2 - ||\mathbf{v}||^2 / 2) / Z$ И 負の対数尤度 Visible v 他の固定点 = saddle points M唯一のlocal minima = global minima 同様のランドスケープはテンソル分解,線形ネットなど でも知られる、わかりやすい描像、 8

訓練性とloss landscape

- ・いくつかの描像
- DNNの幅が十分に大きければ, すべての 固定点がglobal minima

Extremely wide nets [Nguyen&Hein, ICML '17&'18] [Gori&Tesi IEEE 1992]

<u>M(DNNの幅) ≧ T (訓練サンプル数)</u>



Sigmoid, softplus (極限としてReLU), CNN 幅広く成立

- モデルだけでなくアルゴリズムも貢献している SGDがLocal minimaを抜け出す条件 [Kleinberg+ICML '18]

$$\theta_{t+1} = \theta_t - \frac{\partial E}{\partial \theta}(\theta, x_t)$$
 Mini-batch由来のノイズが効く

汎化の問題

表現能力が高い ≠ 学習させやすい



DNNの数理の最前線. 汎化性能を測るさまざまな指標の提案. 網羅的な実験検証も行われつつある:

"Fantastic generalization measures and ..." [Jiang+ ICLR 2020]

40以上の汎化指標

目次

・イントロダクション

- ・統計力学的アプローチの紹介
- 深層学習とランダムネス
- 平均場理論
- ランダム行列理論
- ・さらなる発展: Fisher情報行列とNeural Tangent Kernel
 - Fisher情報行列と学習率
- Neural Tangent Kernelとの対応
- Batch Normalizationの解析

深層学習の多様性

- ネットワーク構造
 素子数, 層数, skip connection, 畳み込み …
- 活性化関数 Sigmoid, (Leaky-)ReLU, ELU …





• 勾配法

constant rate, step-wise, Adam, K-FAC …

$$\theta_{t+1} = \theta_t - \eta_t \frac{\partial E}{\partial \theta}$$

個別のモデルでの個別の数理 (表現能力,訓練性,汎化性能,…)

多様な深層学習に統一的・普遍的な視点は与えられないか?

深層学習とランダムネス

DNN最適化の初期値は
 ランダム行列

M



広く使われている初期値の例: 一様乱数 [Glorot&Bengio AISTATS 2010] $W \sim U \Big[- \frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}, \frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}} \Big]$

ある層の素子数Mに対し、パラメータ分散はO(1/M)

 $\sum_{i} W_{ij} x_{j} \sim O(1)$ 次の層の**出力**が素子数に依存しない. ほどよいスケールに規格化される



平均ゼロ, 共分散: $\mathbb{E}(L(\mathbf{w}_1)L(\mathbf{w}_2)) = f\left(\|\mathbf{w}_1 - \mathbf{w}_2\|^2\right)$ f: 原点まわりで滑らかな関数

2次元グリッド @wikipedia

深層ネットの誤差関数ランドスケープとの定性的な類似

- ・固定点は鞍点が大半で, local minimaはほとんどglobal minima. 誤差が大きいほど負の固有値の割合が増加. [Dauphin+ NIPS '15]
- ・学習中における結合パラメータのノルムの増加はランダムガウス場上の ランダムウォークと同じ($\|\mathbf{w}_t - \mathbf{w}_0\| \sim \log t$) [Hoffer+ NIPS '17]

ランダム深層ニューラルネットの平均場理論

[Amari (1970-)][Poole+ NIPS '16]

$$u_i^l = \sum_j W_{ij}^l h_j^{l-1} + b_i^l, \quad h_i^l = \phi(u_i^l)$$

・ ランダムパラメータ e.g. ガウス分布 $W_{ij}^l \sim \mathcal{N}(0, \sigma_w^2/M)$ $b_i^l \sim \mathcal{N}(0, \sigma_b^2)$



• ニューラルネットの典型的な挙動を表す秩序変数の導出

第*l*層の平均活動度 $q^{l} = \sum_{i}^{M} (u_{i}^{l})^{2}/M$ $M \gg 1$ $q^{l} = \sigma_{w}^{2} \int Dz \phi^{2} \left(\sqrt{q^{l-1}}z \right) + \sigma_{b}^{2}$ **ガウス積分 (ランダム変数の和)**

Mean Field Theory (統計神経力学ともいう)

ランダム深層ニューラルネットの平均場理論



基本的には, 深層学習の "平均場"理論 = 大数の法則と中心極限定理



様々なarchitecture (shallow&deep, sigmoid, ReLU, ResNet, CNN ...) に対して共通して行える計算

ランダム深層ニューラルネットの平均場理論



基本的には, 深層学習の "平均場"理論 = 大数の法則と中心極限定理

参考: Recurrent Neural Network (RNN)の平均場 "近似" [Rozonoer (1969)] [Amari (1970-)][Sompolinsky+ PRL(1988)]

秩序変数の更新則が深層ネットとRNNで類似している ので, アナロジーで"平均場"理論と呼ばれている. RNNの(真の意味の)平均場は, セルフコンシストに解く 必要がある. また, 解が厳密な計算と一致するとは限らない. "Amari solution" [Crisanti, Sommers, Sompolinsky (2008)]



個々の素子が独立で同じ統計性 を持つと仮定

[Amari (1970-)][Poole+ NIPS '16]

異なる入力に対して、ニューロンの活動は独立とは限らない(結合を 共有しているため).入力間の相関に応じて活動も相関.

信号相関の伝達と秩序-カオス相転移



同一の発火パターンに引き込まれ, 入力信号が区別できない

異なる信号間の差を **拡大**する処理 ₁₉

Backpropagationの平均場理論

[Schoenholz+ ICLR '17]

勾配:
$$\frac{\partial f}{\partial W_{ij}^l} = \delta_i^l \phi(z_j^{l-1})$$
 逆誤差伝播: $\delta_i^l = \phi'(z_i^l) \sum_j \delta_j^{l+1} W_{ji}^{l+1}$
勾配ノルム $\tilde{q}^l := \sum_i^M (\delta_i^l)^2$ $\tilde{q}^l = \tilde{q}^{l+1} \sigma_w^2 \int \mathcal{D}z [\phi'(\sqrt{q_l}z)]^2$

- 指数関数的な**勾配の消失・発散**が起こる <u>転移点はfeedforwardと同じ</u> $\chi = \sigma_w^2 \int \mathcal{D}z \left[\phi' \left(\sqrt{q^*}z \right) \right]^2$
- 実は, 導出は非自明

仮定: Gradient independence W_{ji}^{l+1} をfeedforwardと独立なアンサンブル におきかえる

厳密導出: 拘束条件(u_l = W_lh_{l-1})つきでのランダムガウス平均 [Yang 2019] [Arora+ 2019]

Backpropagationの平均場理論

[Schoenholz+ ICLR '17]



ランダム行列理論とDynamical Isometry

 $u_l = W_l h_{l-1} + b_l, \quad h_l = \phi(u_l)$ Input-output Jacobian: $\frac{\partial h_L}{\partial h_0} = \prod_{l=1}^L D_l W_l \quad D_l = \text{diag}(\phi'(u_l))$

平均場理論: 勾配の消失/発散を防ぐには χ = 1**が必要** 特異値の二乗<u>平均</u> $\sim \chi^L$

Dynamical Isometry:

特異値の平均だけでなく,分布の形状を層数と独立にしたい

- 自由確率論によるJacobianのスペクトル解析 [Pennington+ NIPS '17, AISTATS '18] [Pastur, arXiv2001.06188]
- ガウス初期値ではなく, 直交初期値が必要
- ReLUには補正(e.g. shift)が必要



さまざまなアーキテクチャでの検証

Convolutional Neural Network [Xiao+ ICML 2018]



そのほか ResNet [Yang+ NIPS '17], (gated) RNNs [Chen+ ICML '18] ...

おまけ - 誰が平均場理論に再注目したか? -

Career opportunities [edit]

[https://en.wikipedia.org/wiki/Google_Brain]

Google Brain Residency Program [edit] (2016-2017期)

Google Brain Residency Program^[29] is targeted at people who are eager to devote own passion to machine learning and artificial intelligence. This is an opportunity to get hands-on experience in Google team and have chance to keep in touch with professional researchers with Google Brain team. The program lasted 12 months.

Within the program were groups of new graduates from top universities with degree of BAs or Ph.Ds in computer science, physics, mathematics, and neuroscience, or others who come from years of industry experience. They were picked to work with researchers in Google Brain Team at the forefront of machine learning.

The breadth of topics in this program allowed team members to flexibly combine their professional knowledge with the application of algorithms, natural language understanding, robotics, neuroscience, genetics and more. Just several months, these new future researchers have already made a great impact in the research field.

Some of the recently published technical papers resulting from the residency program are listed below:

Unrolled Generative Adversarial Networks 🔉

Conditional Image Synthesis with Auxiliary Classifier GANs 🔊

Regularizing Neural Networks by Penalizing Their Output Distribution₽

Mean Field Neural Networks&

Learning to Remember_₽

Towards Generating Higher Resolution Images with Generative Adversarial Networks₽

Multi-Task Convolutional Music Models

Audio DeepDream: Optimizing Raw Audio With Convolutional Networks

理論を検証できる数値実験が必要.系統的かつ大規模に.

目次

・イントロダクション

- ・統計力学的アプローチの紹介
- 深層学習とランダムネス
- 平均場理論
- ランダム行列理論
- ・さらなる発展: Fisher情報行列とNeural Tangent Kernel
 - Fisher情報行列と学習率

[Karakida, Akaho & Amari, Universal Statistics of Fisher Information in Deep Neural Networks: Mean Field Approach, AISTATS 2019]

- Neural Tangent Kernelとの対応
- Batch Normalizationの解析

[Karakida, Akaho & Amari, *The Normalization Method for Alleviating Pathological Sharpness in Wide Neural Networks*, NeurIPS 2019]

パラメータ空間の重要性

最急勾配における学習率の設計



- 学習率の決め方はヒューリスティックが多い e.g. Training lossの減少がε以下になったら,学習率を1/a倍にする



二乗回帰では,

$$F = \mathbf{E}[\nabla_{\theta} f(x;\theta) \nabla_{\theta} f(x;\theta)^{\top}]$$
$$\mathbf{E}[\cdot]: 入力サンプル平均 x$$



・ 訓練誤差ゼロの大域解まわりのヘシアンに対応

 $Loss(\theta) = \mathbb{E}[||y - f||^{2}]$ Hessian: $\nabla_{\theta} \nabla_{\theta} Loss(\theta) = F - \sum_{k=1}^{C} \mathbb{E}[(y_{k} - f_{k}) \nabla_{\theta} \nabla_{\theta} f_{k}]$



• 情報幾何の基本量. パラメータ空間は曲がっている $D_{\mathrm{KL}}(p(x,y;\theta)||p(x,y;\theta+d\theta)) \sim d\theta^{\top}Fd\theta$

$$u_i^l = \sum_j W_{ij}^l h_j^{l-1} + b_i^l, \quad h_i^l = \phi(u_i^l) \quad (l = 1, 2, ..., L)$$

・ パラメータはrandom Gaussian $W_{ij}^l \sim \mathcal{N}(0, \sigma_w^2/M) \quad b_i^l \sim \mathcal{N}(0, \sigma_b^2)$

入力/中間層: 素子数 M (≫ 1) 出力層: 素子数 C (= 0(1))



- 入力もrandom Gaussian $x_i(t) \sim \mathcal{N}(0,1)$ (t = 1, ..., T)訓練サンプル数
- non-centered network (e.g. ReLU, Tanh with bias terms, ...) バイアスが非ゼロ($\sigma_b \neq 0$)あるいは活性化関数のガウス平均が非ゼロ ($\int Dz\phi(z) \neq 0$)

Fisher 情報行列の巨視的理解

・FIMの計算方法 = chain rule (backprop.と同様)

$$F = \mathbf{E}[\nabla_{\theta} h^{L}(x)^{\top} \nabla_{\theta} h^{L}(x)]$$

 $\tilde{q}^l := \sum_i (\delta^l_i(t))^2$

$$\frac{\partial h_k^L}{\partial W_{ij}^l} = \delta_i^l \phi(u_j^{l-1}), \quad \delta_i^l = \phi'(u_i^l) \sum_j \delta_j^{l+1} W_{ji}^{l+1} \quad (l = 1, ..., L-1)$$

(i) Feedforwardの秩序変数 [Amari 1970-, Poole+ NIPS '16]

$$\hat{q}^{l} := \sum_{i} (h_{i}^{l}(t))^{2} / M$$
 $\hat{q}_{st}^{l} := \sum_{i} h_{i}^{l}(s) h_{i}^{l}(t) / M$

(ii) Backpropagationの秩序変数 [Schoenholz+ ICLR '17]

$$\tilde{q}_{st}^l := \sum_i \delta_i^l(s) \delta_i^l(t)$$

Fisher情報行列の固有値

$$M$$
 (width) が十分に大きいとき, 漸近的に
 $\lambda_{max} \sim (L-1) \left(\frac{T-1}{T} \kappa_2 + \frac{\kappa_1}{T} \right) M$ L: 層数
 T : 訓練サンプル数

$$\kappa_1 := \sum_{l=1}^{L} \tilde{q}^l \hat{q}^{l-1} / (L-1) \qquad \kappa_2 := \sum_{l=1}^{L} \tilde{q}^l_{st} \hat{q}^{l-1}_{st} / (L-1)$$

・非常に大きい孤立した最大固有値 **O**(**M**)

$$\mathbf{E}[\nabla_{\theta} f \nabla_{\theta} f^{\top}] = \mathbf{Cov}(\nabla_{\theta} f, \nabla_{\theta} f) + \mathbf{E}[\nabla_{\theta} f] \mathbf{E}[\nabla_{\theta} f]^{\top}$$
$$\longrightarrow \lambda_{max} \sim ||\mathbf{E}[\nabla_{\theta} f]||^{2}$$

・クラス数Cだけ λ_{max} に縮退 Eigenvectors $E[\nabla_{\theta} f_k]$ (k = 1, ..., C)

Fisher情報行列の固有値

$$M$$
 (width) が十分に大きいとき, 漸近的に $\lambda_{max} \sim (L-1) \left(rac{T-1}{T} \kappa_2 + rac{\kappa_1}{T}
ight) M$ L: 層数 T : 訓練サンプル数

- ・同様にして,固有値の平均値は **0(1/M)**
- クラス数だけ大きな固有値の 実験的報告
 [Sagun+ 2017] [Papyan, ICML '19]





Cross-entropy [Papyan, 2019]



• 最急勾配 (バッチ学習)
$$\theta \leftarrow \theta - \underline{\eta} \frac{\partial E(\theta)}{\partial \theta}$$

学習率

最急勾配が大域解 $\theta^* = \{W^{l*}, b^{l*}\}$ s.t. $E(\theta^*) = 0$ の近傍で収束する必要条件

グレー:

 $\eta < 2/\lambda_{max}$ [LeCun, Kanter & Solla, PRL 1991]など

(右図) 1 epoch後の訓練誤差 L=4 ReLU on MNIST, SGD

- 発散/収束する学習率に明確な境界 理論線 $(= 2/\lambda_{max})$





Neural Tangent Kernel (NTK) 理論

[Jacot+ NeurIPS '18]

$$\frac{d\theta_t}{dt} = \eta \nabla_{\theta} f_t^{\top} (y - f_t) \quad \nabla_{\theta} f_t : CN \times P(\mathcal{N} \ni \mathcal{X} - \mathcal{P} \otimes \mathcal{Y})$$
行列

学習率のオーダーを1/Mでとる (or NTK parameterizationを使う)

• 対応する関数勾配をみる

$$\frac{df_t}{dt} = \eta \nabla_{\theta} f_t \nabla_{\theta} f_t^{\top} (y - f_t)$$
=: Θ_t NTK (*CN*×*CN*行列)

(Informal)

10

隠れ層幅が無限大のとき、
$$f_t$$
のダイナミクスは
$$\frac{df_t^{\text{lin}}}{dt} = \eta \Theta_0 (y - f_t^{\text{lin}})$$
のダイナミクスと一致



MNIST; ReLU (L=2), M=2048, SGD+momentum [Lee+ NeurIPS '19]

Neural Tangent Kernel (NTK) 理論

[Jacot+ NeurIPS '18] [Lee+ NeurIPS '19]

- ・ 線形モデルの訓練と等価 $f_t^{\text{lin}}(x) = f_0(x) + \nabla_{\theta} f_0(x)^{\top} (\theta_t \theta_0)$ 微小変化するパラメータがたくさんあるので 出力は $\mathcal{O}(1)$ で変化 $\eta \nabla_{\theta} f \sim M^{-1} \cdot M^{-1/2}$ $P \cdot M^{-2} \sim 1$
 - 大域収束 $f_t = f_0 + (I \exp(-\Theta_0 t))(y f_0)$
 - 未知データx'に対しても可解. Gaussian Processと等価.
 特に, 訓練されたモデルはKernel ridge-less regression

 $\langle f_{\infty}(x') \rangle_{\text{ini.}} = \Theta(x', x) \Theta(x, x)^{-1} y$

※ NTKの最小固有値は正と仮定 (たとえば,入力の正規化とnon-poly.のactivationで成立) ※ 非漸近論 (十分大きいが有限のM) による収束証明も多数

 $M\gtrsim T^3~$ [Huang & Yau ICML 2020]

そのほかの収束証明の概要: [Zou & Gu, "An improved analysis of ...", NeurIPS '19]

NTKの構造

• カーネルの具体的な計算は秩序変数のrecurrence eq.を使っている

$$\Theta_{ab} = \sum_{l=1}^{L} \tilde{q}_l(a, b) \hat{q}_{l-1}(a, b)$$



FIMと固有値を共有

- "Duality"
$$F/M = \begin{bmatrix} \nabla_{\theta} f & NTK = \end{bmatrix}$$

 $\nabla_{\theta} f^{\top}$
パラメータ空間 関数空間

NTK regimeの "外側"

- 初期値から大きく離れる状況では NTK理論の成立が非自明 (e.g. 有限幅の効果, ノイズのような外乱…)
 - 外側こそが積極的に"特徴抽出"

"Active" regime
 ×
 t = 0
 Lazy (NTK) regime

"Large learning rate phase"

[Lewkowycz+ arXiv2003.02218]





Batch Normalization (BN)

[loffe&Szegedy ICML '15]

$$\begin{aligned} u_i^l(t) \leftarrow \frac{u_i^l(t) - \mu_i^l}{\sigma_i^l} \gamma_i^l + \beta_i^l \\ \mu_i^l &:= \mathbf{E}[u_i^l(t)] \quad \sigma_i^l := \sqrt{\mathbf{E}[u_i^l(t)]^2 - (\mu_i^l)^2} \end{aligned}$$

経験的には - 大きい学習率での高速な収束 - 汎化しやすい



BNのメカニズム [Santurkar+ NeurIPS '18]

- ・通説 (Internal covariate shiftの抑制)への反証
- ・Loss landscapeの急激な変化を抑えている可能性



Batch norm in the last layer

簡単のため, 最終層のmean subtractionのみを考える $\bar{f}_k(t) := (u_k^L(t) - \mu_k(\theta))\gamma_k + \beta_k \quad k = 1, \dots, C$ (クラス数)

仮定(I) 幅Mと訓練サンプル数Tが十分大きく, $\rho := M/T$ (固定) (II) gradient independence

$$\rho\alpha(\kappa_{1} - \kappa_{2}) + c_{1} \leq \lambda_{max} \leq \sqrt{(C\alpha^{2}\rho(\kappa_{1} - \kappa_{2})^{2} + c_{2})M}$$

 $\kappa_{1}, \kappa_{2}: 秩序変数から計算, \quad \alpha = L - 1$
 $c_{1}, c_{2}: 非負値定数$

・ λ_{max} のオーダーが $\Theta(M)$ から高々 $\Theta(\sqrt{M})$ へ減少

導出の方針: $E[\nabla_{\theta} f \nabla_{\theta} f^{\top}] = Cov(\nabla_{\theta} f, \nabla_{\theta} f) + E[\nabla_{\theta} f]E[\nabla_{\theta} f]^{\top}$

$$\bar{f}_k(t) = f_k(t) - \mathbf{E}[f_k] \qquad \mathbf{E}[\nabla_\theta \bar{f}_k] = 0$$

Batch norm in the last layer

簡単のため, 最終層のmean subtractionのみを考える

$$\rho\alpha(\kappa_1 - \kappa_2) + c_1 \le \lambda_{max} \le \sqrt{(C\alpha^2\rho(\kappa_1 - \kappa_2)^2 + c_2)M}$$

 κ_1, κ_2 :秩序変数から計算, $\alpha = L - 1$
 c_1, c_2 :非負値定数

- λ_{max}のオーダーが Θ(M) から高々 Θ(√M) へ減少
 理論の上限は保守的. 経験的には Θ(1)
- ・ 中間層のみにBN Deep ReLU net では, $\lambda_{max} = \Theta(M)$
- Layer normalization $\lambda_{max} = \Theta(M)$

$$\mathbb{E}[\nabla_{\theta} f_k(t)] \neq \frac{1}{C} \sum_{k=1}^{C} [\nabla_{\theta} f_k(t)]$$



学習率 η と λ_{max}



最終層BNだけで,幅に依存しない大きな学習率が許容された

Dynamical isometry成立下でのFIM

[Hayase & Karakida arXiv:2006.07814]

Input-output Jacobian:

 $\frac{\partial h_L}{\partial h_0} = \prod_{l=1}^{L} D_l W_l$ の統計性が層数に非依存でも, FIMは依存する



オンライン学習 (mini-batch size =1)

 *θ*_{t+1} = *θ*_t - η∇_θ[M⁻¹L(f_{θt}(x(t) - y(t))]
 [右図] 学習初期(500 step)の挙動

Hard Tanh on Fashion-MNIST



まとめ

ランダムネスに基づいた,大自由度(幅無限大)極限 における深層ネットの数理を紹介した

- DNNの挙動を統一的にとらえられる, 普遍性のある枠組み
- さらに, 個別の現象に定量的な示唆を与える
 - ・勾配の発散/消失を防ぐ初期値
 - ・学習率の設定
 - ・NTK regimeの発見

今後の課題

- 各種モデル・アルゴリズムのNTK regimeでの理解
- 固有値解析 (2/λ_{max}境界, NTK regime の成立条件 λ_{min} > 0)
 の精緻化
- NTK regimeの外側への展開