# 深層神経回路の数理

## ランダム行列・無限次元近似・自由確率論

Hayase Tomohiro
早瀬 友裕

# Table of Contents

## S1

### Overview

Roles of random matrices in deep learning

## S2

### Jacobian

The first application of free probability to deep learning

## S3

### Fisher Information

Learning dynamics needs a fruitful random matrix theory

## S4

### Conclusion

Current works, future works, and the other topics

2

深層神経回路の数理

ランダム行列・無限次元近似・自由確率論

## Deep Neural Networks

Rossenblatt '57~62    Multilayer perceptron

$(x, y)$ : given data.



$w_1$    $w_2$

$z = f_{w_1, w_2}(x)$

$L(f_{w_1, w_2}) = |y - f_{w_1, w_2}(x)|^2$

parameter.

Stochastic gradient decsents    Amari, Tsypkin '66~67

Error Back Propagation            '76

§1 Overview

## A standard setting

### 1. Multilayer Perceptron is a parametric family $(f_\theta)_{\theta \in \Theta}$

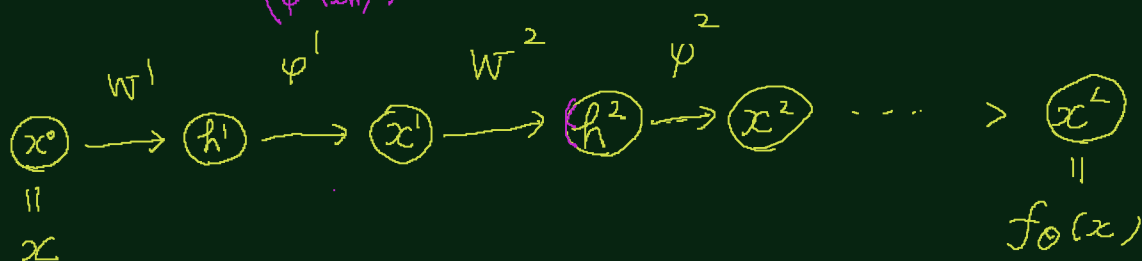$$w/ \quad f_\theta : \mathbb{R}^M \to \mathbb{R}^M \quad (M \in \mathbb{N})$$

$$\hookrightarrow \varphi^L \circ g_{\theta^L}^L \circ \cdots \circ \varphi^1 \circ g_{\theta^1}^1$$

$$g_{\theta^\ell}^\ell(x) = \underbrace{W^\ell}_{M \times M} x + \underbrace{b^\ell}_{\in \mathbb{R}^M} \quad (\theta^\ell = (W^\ell, b^\ell))$$

activation $\rightarrow$ $\varphi^\ell \in C(\mathbb{R})$ : differentiable, except for finite number of points

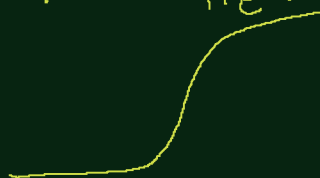$$\varphi^\ell(x) = \begin{pmatrix} \varphi^\ell(x_1) \\ \vdots \\ \varphi^\ell(x_n) \end{pmatrix}$$

$$x^0 \xrightarrow{W^1} h^1 \xrightarrow{\varphi^1} x^1 \xrightarrow{W^2} h^2 \xrightarrow{\varphi^2} x^2 \cdots > x^L$$

$$x^0 \;\|\; x \qquad\qquad x^L \;\|\; f_\theta(x)$$

### Examples of Activation

#### ReLU '15

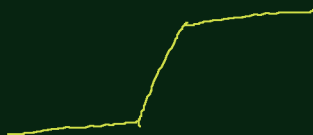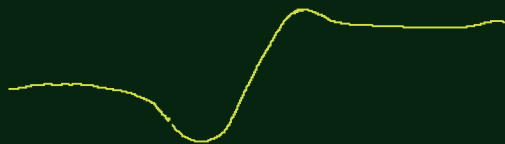$$\varphi(x) = \max(0, x)$$

#### Sigmoid

$$\varphi(x) = \frac{1}{1 + e^{-x}}$$

#### Hard tanh

#### SiLU (Sigmoid Linear Unit) '17
(or Swish)

$$\varphi(x) = \frac{x}{1 + e^{-x}}$$

testing dataset : training dataset
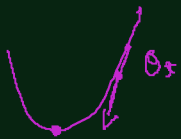に含まれてはいないが
あることが望しい dataset。

2. A training dataset is a set of pairs $(x_n, y_n)_{n=1}^N$

3. Gradient Descents

(MSE)

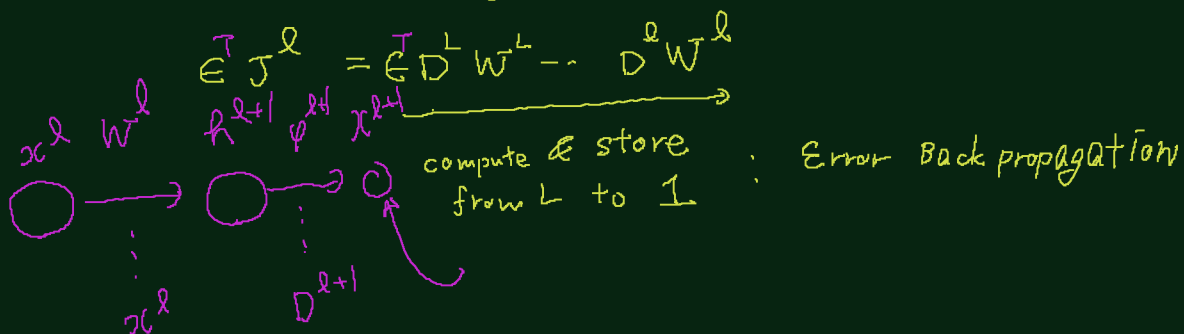$$L(f_\theta) = \frac{1}{2N} \sum_{n=1}^N \| f_\theta(x_n) - y_n \|_2^2$$

$t = 0, 1, 2, \cdots$
$$\theta_{t+1} = \theta_t - \eta_t \frac{d}{d\theta} L(f_\theta) \Big|_{\theta = \theta_t} \qquad (\eta_t > 0)$$

.✕. In practice, $(x_n, y_n)$ are picked randomly from the training dataset.

4. Error Back-Propagation

$$\frac{\partial L}{\partial W^\ell} = \underbrace{(f_\theta(x) - y)^T}_{\epsilon^T} \underbrace{\frac{\partial f_\theta}{\partial x^{\ell+1}}}_{J^{\ell \times 1}} D^{\ell+1} x^\ell$$

$$\underbrace{\epsilon^T J^\ell = \epsilon^T D^L W^L \cdots D^\ell W^\ell}_{}$$

$x^\ell$ $W^\ell$ $h^{\ell+1}$ $\phi^{\ell+1}$ $x^{\ell+1}$

compute & store : Error Back propagation
from L to 1

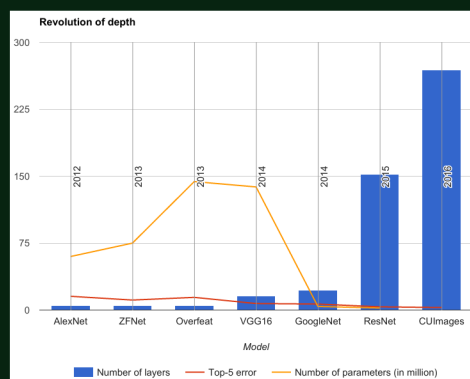$x^\ell$ $D^{\ell+1}$
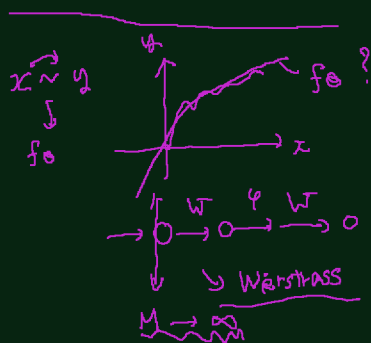
# Deep

- ILSVRC 2012
  ( ImageNet Large Scale Visual Recognition Challenge)

- AlexNet ( Hinton e.t.c )
  ImageNet Classification with Deep Convolutional Neural Networks
  ( NIPS 2012 )

- Deep である必要
  - 表現能力の指数的向上
  - 性能は上昇する(?)

$x \tilde{\sim} y$       $f_\theta$ ?

$f_\theta$

$\to \bigcirc \xrightarrow{W} \bigcirc \xrightarrow{\varphi} \xrightarrow{W} \bigcirc$

→ Weierstrass



Revolution of depth

## S2. Dynamical Isometry

Fix $L$ ($=$ the number of Layers).

$$J = D^L W^L \cdots D^1 W^1$$

$D^\ell$ : activation ⓐ Jacobian.

$W^1, \cdots, W^L$ : initial state

Error back propagation $\|\vec{\epsilon}^T J\|_2$ can

explode / vanish as $L \to \infty$ , $\|$

~ Uniform.
~ Gaussian
~ Haar orthogonal
(uniformly picking up from $O(M)$)

"Expoding/Vanishing Gradient Problem"

Pennington, Schoenholz, Ganguli [NIPS '17, AISTATS '18]

$\sqrt{JJ^T}$'s eigenvalues

⌐ Ensuring $J$'s singular values $\sim O(1)$ as $L \to \infty$ (Dynamical Isometry) is essential for avoiding the exploding/vanishing of gradients ⌐

⌐ Gaussian Initialization ✗

Orthogonal Initialization + normalization of $\rho \Rightarrow$ D.I.

$$S_{\lim_{M \to \infty} \mu_{JJ^T}}(z) \xrightarrow{L \to \infty} \exp\left(-\frac{z\sigma_0^2}{1+z}\right) \quad \text{(Voiculescu)}.$$

↳ Free Infinite Multiplicative Infinite Divisible

free prob.



limit spectral distribution

### Sketch

Assume that $(W_1, W_1^*), \cdots, (W_L, W_L^*), (D_1, \cdots, D_L)$ are asymptotic free as $M \to \infty$.

free probability

Let $(u_1, u_1^*), \cdots (u_L, u_L^*), (d_1, \cdots, d_L)$ be free family in a $C^*$-prob. sp. $(\mathcal{A}, \tau)$.

$$\tilde{j} = d_L u_L \cdots d_1 u_1$$

$$\tilde{j}\tilde{j}^* = d_L u_L (\underbrace{d_{L-1} u_{L-1} \cdots d_1^2 \cdots u_{L-1} d_{L-1}}_{\tilde{j}_{L-1}\tilde{j}_{L-1}^*}) u_L^* d_L$$

$$S_{\tilde{j}\tilde{j}^*}(z) = S_{d_L^2}(z) S_{\tilde{j}_{L-1}\tilde{j}_{L-1}^*}(z),$$

$$= S_{d_1^2}(z) \cdots S_{d_L^2}(z)$$

↳ unitary だから消えた

§3. Fisher Information Matrix & Neural Tangent Kernel

$$I(\theta) = \frac{1}{N}\sum_{n=1}^{N}\frac{\partial f_\theta(x_n)}{\partial \theta}^T\frac{\partial f_\theta(x_n)}{\partial \theta}$$

$: M^2L \times M^2L$ 行列

(Empirical)
Fisher Information Matrix

$b^\ell = 0$ とした行列

$D_{KL}(P_\theta \| P_{\theta+d\theta})$
$\simeq (d\theta)^T I(\theta)\, d\theta$

( Information Geometry (Amari))

Neural Tangent Kernel

$$\textcircled{H} = \frac{\partial f_\theta(x)}{\partial \theta}\frac{\partial f_\theta(x)}{\partial \theta}^T$$

$P_\theta = \exp\left(-L(f_\theta)\right)/Z.$

$: M \times M$ 行列
$C \times C \quad (C \ll M) :$ 名もの次元

# NTK describes learning dynamics

Learning dynamics of parameters is given by:
$$\frac{d\theta_t}{dt} = \eta(\nabla_\theta f_{\theta_t})^T(y - f_{\theta_t})$$

Learning dynamics of DNN is given by:
$$\frac{df_{\theta_t}}{dt} = \eta\Theta_t(y - f_{\theta_t})$$
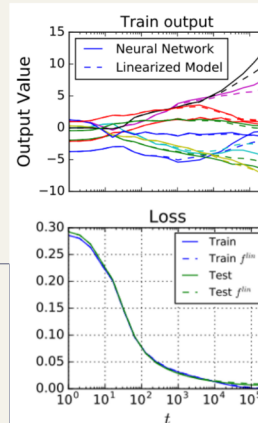$$\Theta_t = \nabla_\theta f_{\theta_t}(\nabla_\theta f_{\theta_t})^T$$

where

**Informal [Jacot+NeurIPS2018, Lee+NeurIPS2019]:** Under the wide limit M \to \infty, the learning of the DNN is approximated by
$$\frac{df_{\theta_t}}{dt} = \eta\Theta(y - f_{\theta_t})$$
where
$$\Theta = \nabla_\theta f_\theta(\nabla_\theta f_\theta)^T$$
( **Neural Tangent Kernel)**

Train output — Neural Network, Linearized Model; Output Value

Loss — Train, Train $f^{lin}$, Test, Test $f^{lin}$

14

$\Theta :$ 初期値

Gaussian

$\textcircled{?}$ orthogonal

$C = O(1)$ as $M \to \infty$

$\frac{}{M!}$

$D.I\textcircled{?}$

Spectrum of $\frac{\partial f_\theta}{\partial \theta}$ $\left(= b^L \frac{\partial x^L}{\partial \theta}\right)$  per sample.

arXiv: 2006, 07814

T.H. & Ryo Karakida [AISTATS '21]

$H_L := \frac{1}{M} \frac{\partial x^L}{\partial \theta} \frac{\partial x^{L*}}{\partial \theta}$ .... NTK ..... A dual of Fisher Information Matrix per sample

Then limit spectral distribution of $H_L$

concentrates on the points $\underline{g L}$ as $M \to \infty$
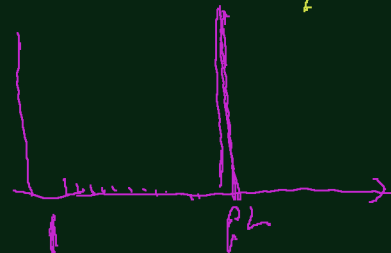
w/ $g = \lim_{M \to \infty} \frac{1}{M} \|x^0\|_2^2$

under some assumptions to achieve dynamical isometry.

Sketch

$\frac{1}{M}\|x^\ell\|_2^2$

$H_{\ell+1} = \hat{g}_\ell + W_{\ell+1} D_\ell H_\ell D_\ell W_{\ell+1}^*$

asymptotic freeness

$\Rightarrow \mu_{\ell+1} = (g + \cdot) * (\nu_\ell \boxtimes \mu_\ell)$

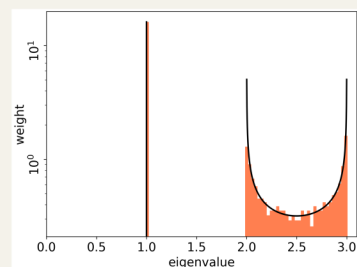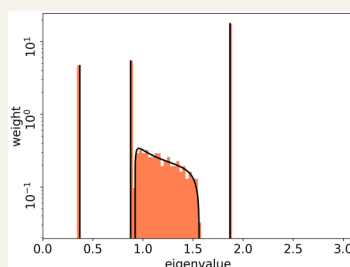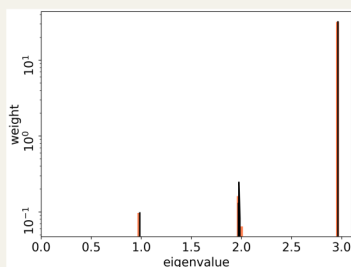$H_{\ell+1}$    $D_\ell^2$

0  $\partial^{-1}$

Point : If activation function is ReLU-like ftn,

· $\mu_\ell$ has an atom as a maximal point spectrum.

· $\nu \boxtimes \mu (\{a b\}) = \nu(\{a\}) \mu(\{b\}) + \mu(\{b\}) - 1$.



# Limit spectral distributions : L = 3

$D \longrightarrow p\,\hat{g}$   $p\cdot\hat{g}$ : free

0 or 1
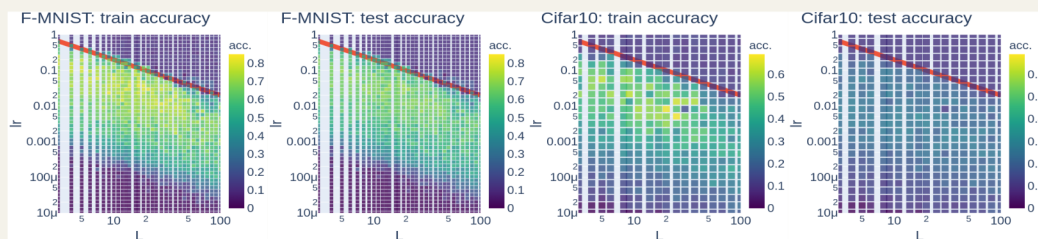
16

## Training Dynamics

- $\eta > \dfrac{1}{\lambda_{max}(H_L)}$  $\Rightarrow$ The training dynamics may not converge

# Training under D. Isometry

Red line (the boarder line of the exploding gradients) :

$$\eta = 2/L$$

This line is expected by our theory !

# Conclusion

- Random Matrices appear in the theory of deep neural networks ( e.g. dynamical isometry, Fisher Information, NTK )

- Since Jacobian ( input or parameter ) are (noncommutative) polynomial of Random Matrices.

Free Probability Theory provides tools for handling them!

To use this, we have to prove asymptotic freeness { a.s. expected. strong (operator norm)

Gaussian { Hanin-Nica '19
Yang '19, '20
Pastur '20

Orthogonal → H. )    (under preparation)
'67

gradient independence

面白いよ { $W^1, \cdots, W^L$ ⟩  independent
それらい     $D^1, \cdots, D^L$ ⟨  と思って
とする.

freeness → $tr( Q(W^1, \cdots, W^L, D^1, \cdots, D^L))$ : ok (?)