



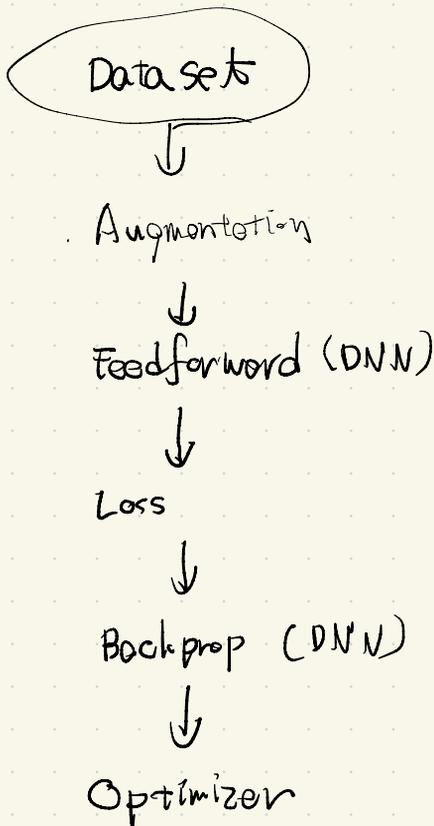
プログラミングと深層学習

3/23, 24

Day 1

深層学習練習

- ✓ framework of deep learning
- ✓ examples.



Framework

103

$$\mathcal{D} \subset X \times Y \subset \mathbb{R}^n \times \mathbb{R}^{n_0}$$

finite

μ : prob. dist. on $X \times Y$

$$L: \mathbb{R}^{n_0} \times \mathbb{R}^{n_0} \rightarrow \mathbb{R}$$

Find $f: \mathbb{R}^n \rightarrow \mathbb{R}^{n_0}$

s.t. minimize $\int_{(x,y) \in X \times Y} L(f(x), y) \mu(dx dy)$

μ
unknown

with information \mathcal{D} .

★ $\mathcal{D} \subset X \times Y \subset \mu \in \mathcal{P} \subset \text{approx}$

★ $f = f_\theta \quad (\theta \in \Theta \subset \mathbb{R}^p)$
parametrize

★ $\mathcal{P} \subset \text{approx}$

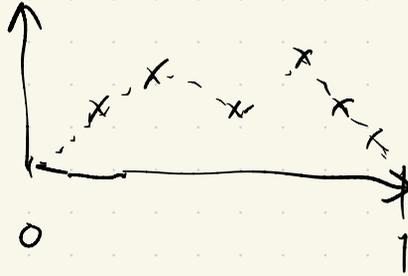
★ "minimize" $\in \text{approx}$

Dataset D

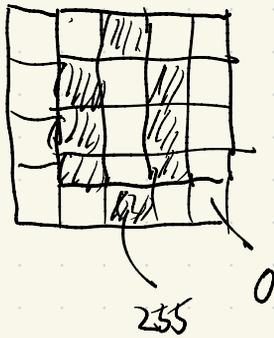
CV

1-dim

$$X, Y = [0, 1] \subset \mathbb{R}$$



Computer Vision



Pixel

RGB

$$H \times W \times C$$

↑
3

→ vectorize

HWC-dim
vector

output: $0 \rightsquigarrow \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = e_1 \in \mathbb{R}^{10}$

$\vee 0 \sim 9 \rightarrow e_1 \sim e_{10}$

NLP

MNIST: ~~手写字~~

Natural Language Processing

(NLP)

I am eating pizza.

tokenize

$\rightarrow [0, 10, 201, 1018]$

vector

$\rightarrow x \in X.$

$\gamma?$

mark

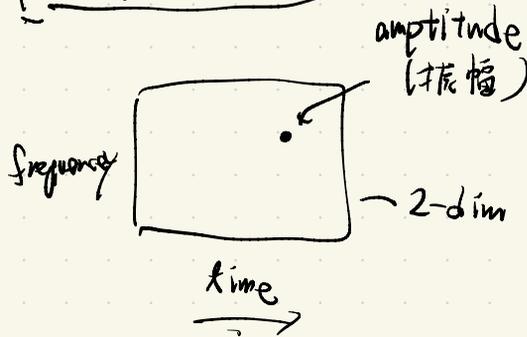
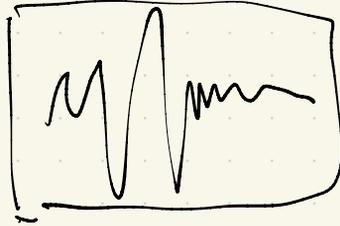


x : I am  pizza

y : I am eating pizza.

Audio

Audio

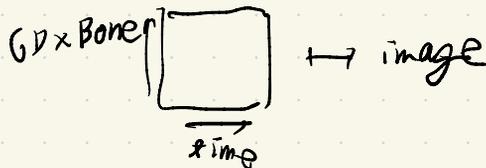
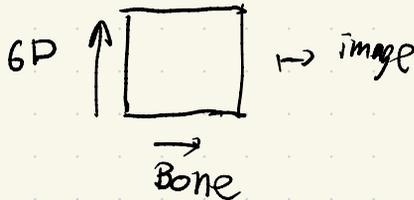


画像として扱う
→

3D Bone Data

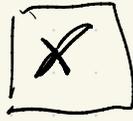


3D(+3D) x Bone

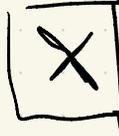


Augmentation

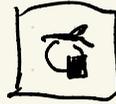
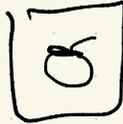
Image



flip



Random
Flip



Random
Erasing

data に依存して増やす。

Model

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^{n_0}$$

$$f = f_{\theta}$$

$$f_{\theta}(x) = Wx + b$$

$$\theta = (W, b)$$

$$\begin{cases} W \in \mathbb{R}^{n_0 \times n} \\ b \in \mathbb{R}^{n_0} \end{cases}$$

Multilayer-Perceptron

Aug 2

(MLP)

$$x_0 = x$$

$$\left\{ \begin{array}{l} h_l = W_l x_{l-1} + b_l \\ x_l = \varphi(h_l) \\ l = 1, \dots, L \end{array} \right.$$

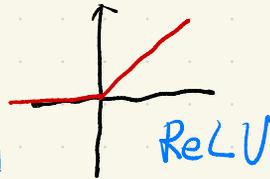
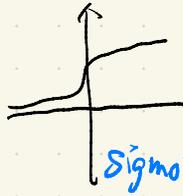
$L \in \mathbb{N}$

φ : activation ftn.

$$\varphi: \mathbb{R} \rightarrow \mathbb{R}$$

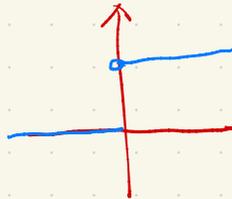
↳ conti.

↳ diff' except for finite pts.



$$f_{\Theta}(x) = h_L$$

$$\Theta = (W_1, b_1, \dots, W_L, b_L)$$



* L : Loss fn

LF

↳ L2-fn

$$\hookrightarrow L(f(x), y) = \|f(x) - y\|^2$$

↳ Softmax Cross Entropy --

↳ L1 - - -

Backpropagation

L_0

$$\left[\begin{array}{l} \text{minimize } \sum \|f_\theta(x) - y\|^2 \\ \theta \in \Theta \mid (x, y) \in \mathcal{D} \end{array} \right] + \lambda \|\theta\|^2$$

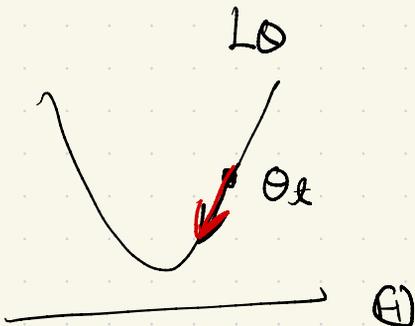
weight decay

Gradient Descent

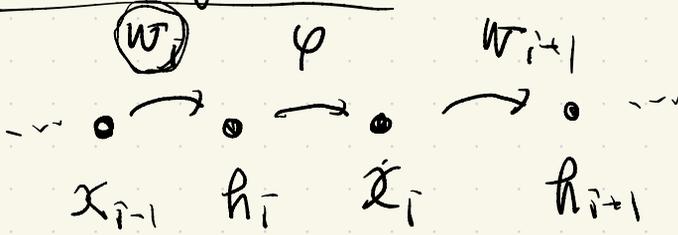
$$\theta_{t+1} = \theta_t - \eta_t \cdot \nabla_{\theta} L_0 \theta_t$$

$$\eta_t > 0$$

: learning rate

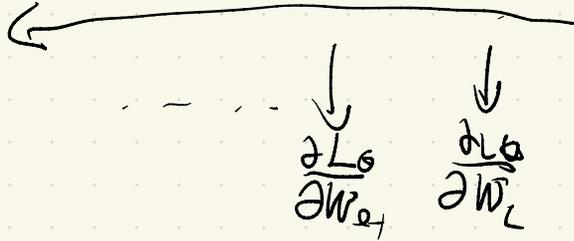


Back propagation



$$\frac{\partial L_0}{\partial w_i} = \frac{\partial L_0}{\partial h_i} \cdot \frac{\partial h_i}{\partial w_i}$$

The terms $\frac{\partial L_0}{\partial w_i}$, $\frac{\partial L_0}{\partial h_i}$, and $\frac{\partial h_i}{\partial w_i}$ are circled in the original image. A long arrow points from the right side of the equation back towards the left, indicating the direction of backpropagation.



Summary

✓ Dataset

✓ Augmentation

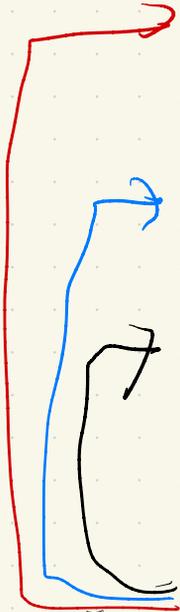
✓ NN

✓ Forward

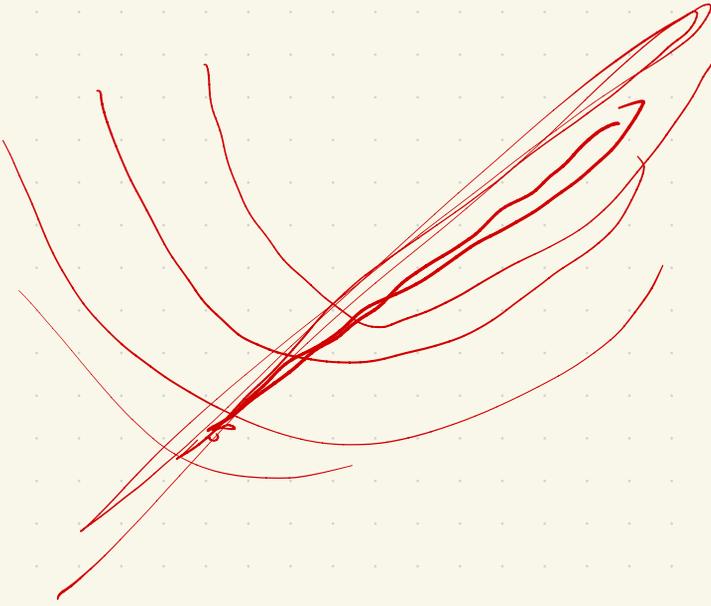
✓ Loss

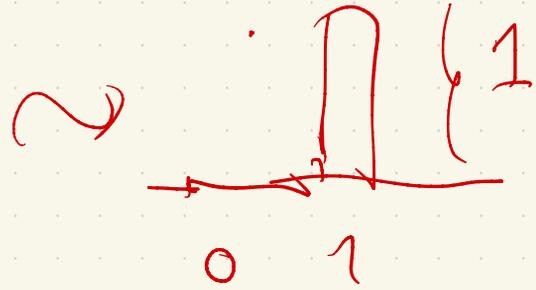
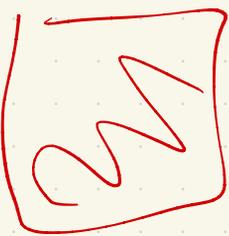
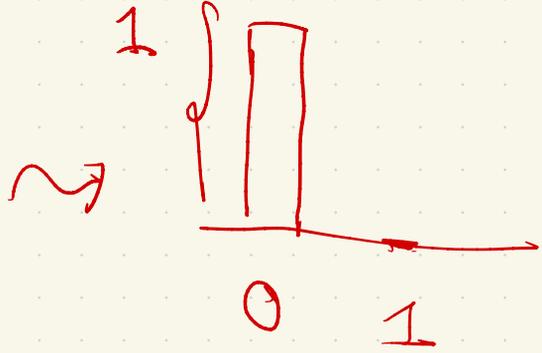
optimizer

✓ Back propagation

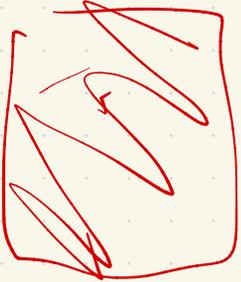


$$y = ABx$$

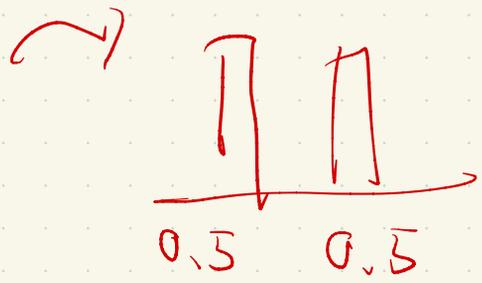




Mixup



$\alpha = 0.5$



Day 2 Signal Propagation

(1) Back Propagation

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} L_0$$

$$L_0 = \sum_{x,y} L(f_{\theta}(x), y) + \text{[scribble]}$$

$$\theta_t = (w_1^{(t)}, \dots, w_L^{(t)}) \\ b_1^{(t)}, \dots, b_L^{(t)}$$

Initialization が重要.

$$W_{\ell, \vec{i}j} \sim N(0, \frac{1}{n}) \quad \text{i.i.d.}$$

$$(\vec{i}j = 1, \dots, n)$$

$$(\ell = 1, \dots, L)$$

$$b_{\ell, \vec{i}j} \sim N(0, \frac{1}{n}) \quad \text{i.i.d.}$$

$$(\vec{i}j = 1, \dots, n)$$

$$(\ell = 1, \dots, L)$$

$$\Rightarrow \mathbb{E}[\|h_{\ell}\|^2] = \mathbb{E}[\|x_{\ell-1}\|^2] + 1$$



$x_i = \varphi(h_i)$ 2倍ぐらい増す

ランダム行列

: 行列値 (R-valued)

- 確率変数

$W_{ij} \sim N(0, \sigma^2) \quad (i, j = 1, \dots, n)$

i.i.d

Ginibre

Random

Matrix.

$W \sim$ uniformly sample

from $O(n)$

Haar Orthogonal

Random Matrix

W : Gram matrix

$$W = U D V^T$$

$$U, V \in O(n)$$

Haar
Orthogonal

$$D = \text{diag} \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}$$

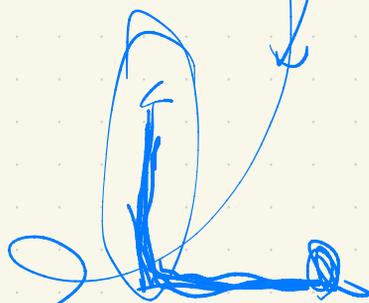
$\lambda_1 \geq \dots \geq \lambda_n$: singular value of W

($W^T W$ の固有値の square root)

Back propagation

$$J = \frac{\partial h_L}{\partial x_0} = W_L D_{L-1} \underbrace{W_{L-1}}_{\dots W_1} D_{L-2}$$

J の singular value の性質が良くなる
→ 学習が速くなる



Explosion / Vanishing gradient

→ J が orthogonal に近くなる

→ J の特異値が I 周辺

に集中しているような

状況をつかいませんか?

↳ Dynamical Isotometry
(Saxe, Ponnagunton)

$$J = W_L D_{L-1} W_{L-1} \dots W_1$$

$$D_e = \psi'(h_e)$$

⇐

↑

Random
Matrix.

⇐

Random
vector.

RM の 種 の singular value ϵ controle
できるか?

→ free probability theory

prob. theory

random variable X

prob. distribution μ_X

expectation
/ variance

Independence

$$\mu_{X,Y} = \mu_X \otimes \mu_Y$$

Noncommutative
Probability element of
 $A \in A$ algebra

$n \mapsto \text{tr}(a^n)$

$\text{tr}(a^n)$

free product

free probability theory

RM eigenvale $X = W^T W$

~~singular~~ value, of

$$\left\{ n \mapsto \text{tr}(X^n) \right\}$$

\leadsto eigenvalue

固有空間 $U \sim O(n)$
uniform

\Leftrightarrow non-commutative
with other RM

i.e. $U, V \sim$ uniform on $O(n)$

U, V : non-commutative

\leadsto "asymptotically" free

3.7.3 a prob.

$$X \perp\!\!\!\perp Y \rightarrow \mu_{X+Y} = \mu_X \boxplus \mu_Y$$

$$\mathbb{E}[e^{it(X+Y)}]$$

$$= \mathbb{E}[e^{itX}] \mathbb{E}[e^{itY}]$$

W_1, W_2 : asyn. free

$$\mu_{W_1 + W_2} \approx \mu_{W_1} \boxplus \mu_{W_2}$$

free additive
convolution

$$\mu_{W_1 W_2} \approx \mu_{W_1} \boxtimes \mu_{W_2}$$

$$J = \sum_{L=1}^L D_L W_L \dots D_1 W_1$$

$$\Rightarrow \mu_{JJ} \approx \mu_{W_1^T W_1} \otimes \mu_{D_1} \otimes \dots \otimes \mu_{W_{L-1}^T W_{L-1}} \otimes \mu_{D_{L-1}}$$

asym. freeness

B. Collins

Comm.

Math.

Phys.

trivial
or



Marčenko-Pastur

AISTATS

Pennington,

arXiv 2018
2019

Emergence university

in Neural Network

$\partial \theta f$

$$\theta_{t+1} = \theta_t - \boxed{\nabla_{\theta} L_{\theta}}$$

$$(\nabla_{\theta} L_{\theta})^T (\nabla_{\theta} L_{\theta}) ?$$

$$\textcircled{H} := \sum_m \frac{\partial f(x(m))}{\partial \theta} \frac{\partial f(x(m))}{\partial \theta}^T$$

$$\mathcal{D} = \{ (x_m, y_m) \mid m=1, \dots, N \}$$

Jacobi
2018

Neural Tangent Kernel

$$\frac{df_{\theta}}{dt} = \textcircled{H} (Y - f_{\theta}(X))$$

arXiv

Conjugate Kernel

2022 F. Weng.

$$\begin{cases} Y = (y^{(1)} \dots y^{(N)}) \\ X = (x^{(1)} \dots x^{(N)}) \end{cases}$$