

JDLA E資格認定プログラム

全人類がわかるE資格コース

MLP基礎（バッチ処理編）



AVILEN

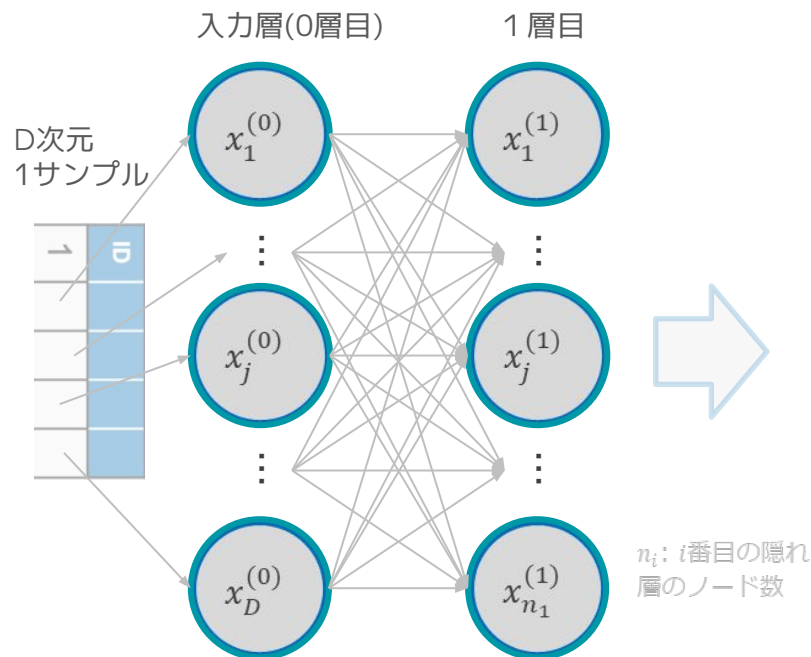
- 1) バッチ処理（順伝播）
- 2) バッチ処理（逆伝播）
- 3) 最適化手法

【Chapter03】MLP基礎（バッチ処理編）

バッチ処理（順伝播）

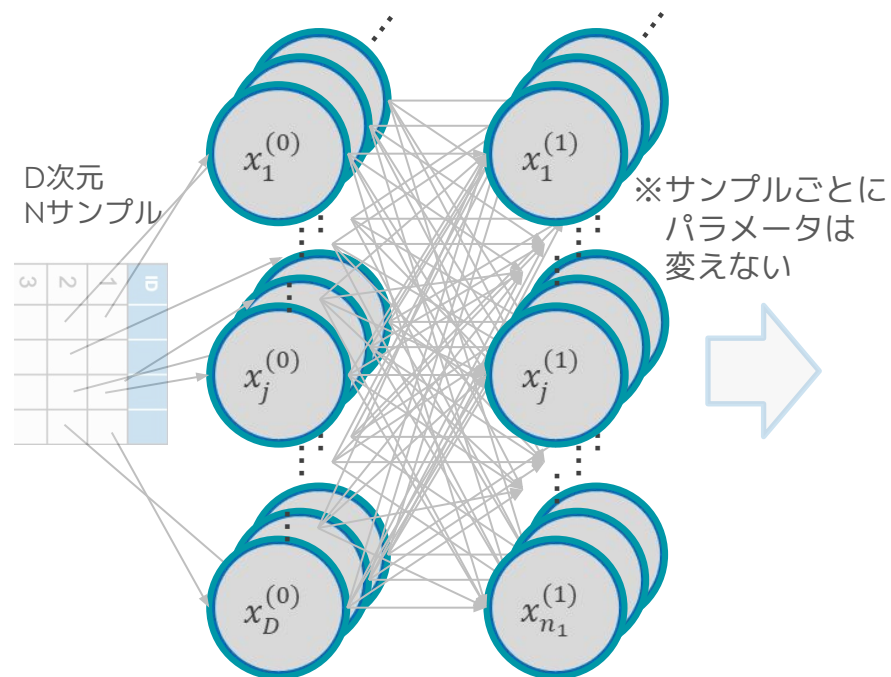
オンライン学習

1サンプルずつ学習する



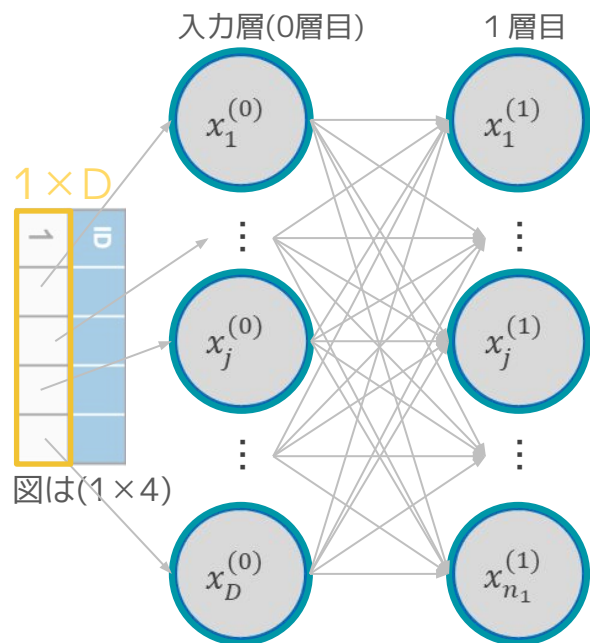
バッチ処理 (バッチ学習)

Nサンプルまとめて学習する



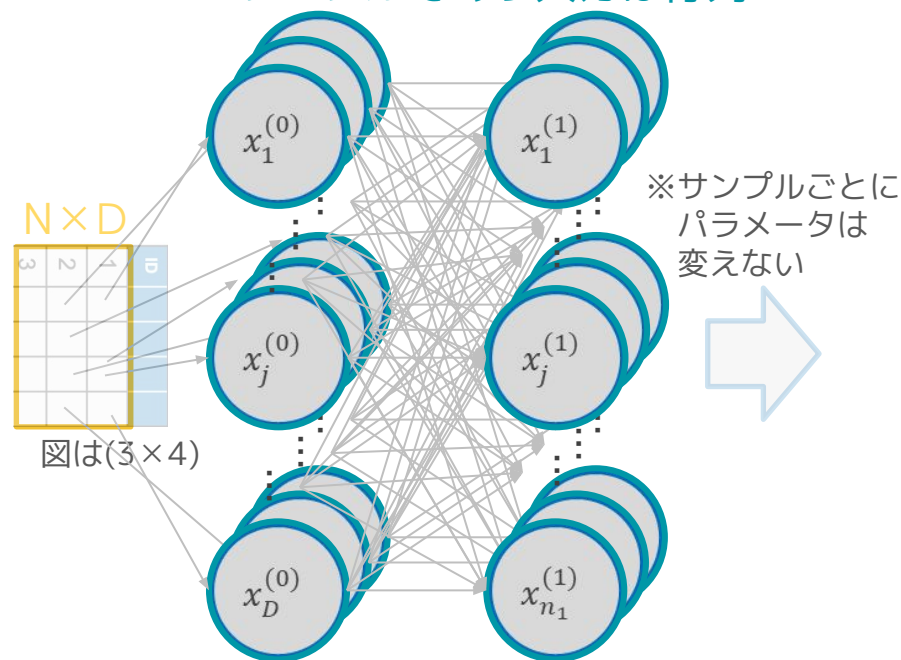
オンライン学習

1サンプルなので入力はベクトル



バッチ処理 (バッチ学習)

Nサンプルであり入力は行列



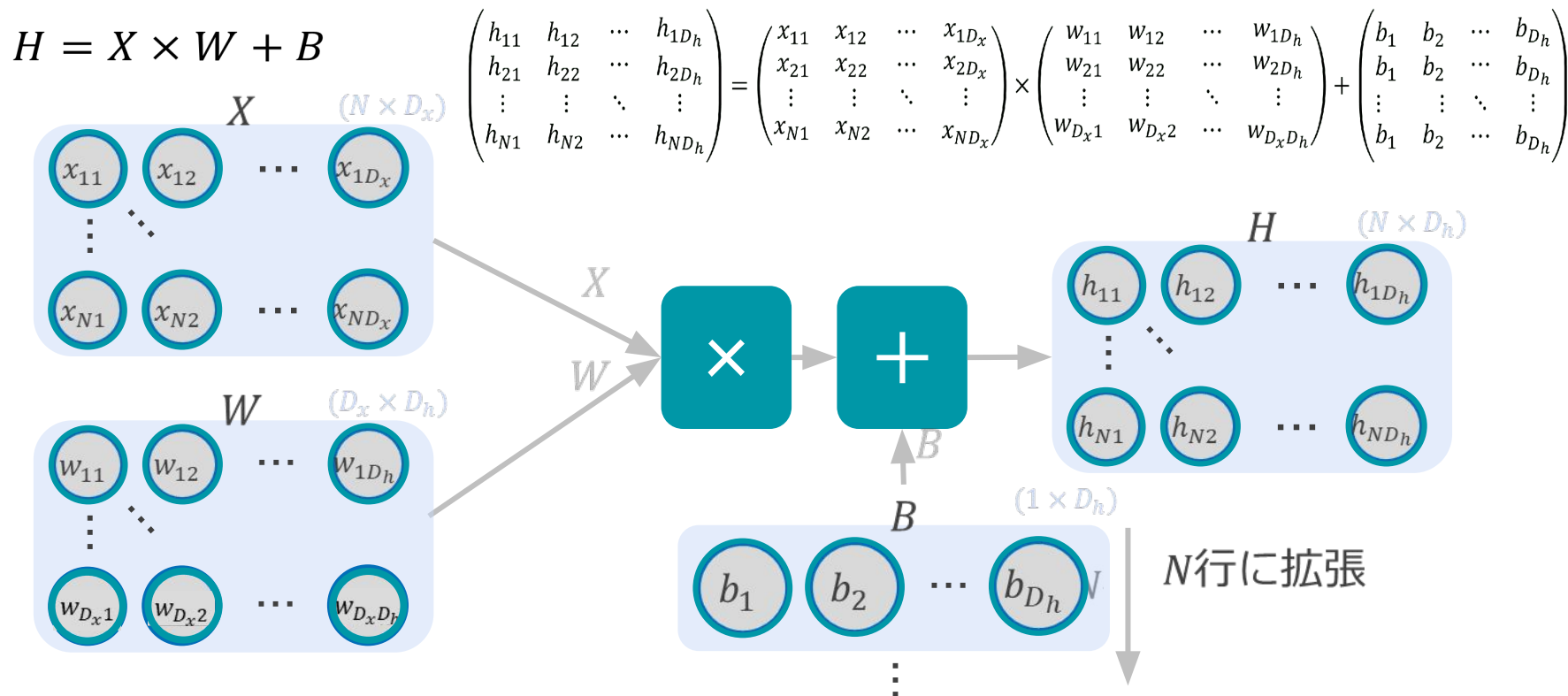
バッチサイズとは、パラメータを一度更新するための損失計算に必要なデータ数

【言い換え】ひとまとまりのグループのことをバッチサイズという

これまでの資料でのNはバッチサイズと呼ばれる。
すなわち・・・

	オンライン学習	バッチ処理（バッチ学習）
バッチサイズN	1	全データ数
損失関数 （二乗和誤差関数）	$\sum_n^1 (t_n - y_n)^2$	$\sum_n^N (t_n - y_n)^2$

バッチ版Affine変換の順伝播を記す



$$\begin{pmatrix} h_{11} & h_{12} & \cdots & h_{1D_h} \\ h_{21} & h_{22} & \cdots & h_{2D_h} \\ \vdots & \vdots & \ddots & \vdots \\ h_{N1} & h_{N2} & \cdots & h_{ND_h} \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1D_x} \\ x_{21} & x_{22} & \cdots & x_{2D_x} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{ND_x} \end{pmatrix} \times \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1D_h} \\ w_{21} & w_{22} & \cdots & w_{2D_h} \\ \vdots & \vdots & \ddots & \vdots \\ w_{D_x1} & w_{D_x2} & \cdots & w_{D_xD_h} \end{pmatrix} + \begin{pmatrix} b_1 & b_2 & \cdots & b_{D_h} \\ b_1 & b_2 & \cdots & b_{D_h} \\ \vdots & \vdots & \ddots & \vdots \\ b_1 & b_2 & \cdots & b_{D_h} \end{pmatrix}$$

計算フローはオンライン版(バッチサイズ=1の時)とほとんど変わらない

```

1  N = 10 # バッチサイズ
2  Dx = 3 # 入力のカラムは3次元
3  Dh = 2 # outputのカラムは2次元
4  linear = nn.Linear(in_features=Dx, out_features=Dh)
5  input = torch.randn(N, Dx)
6  output = linear(input)
7  print(output)

```

```

tensor([[ 0.1459,  0.7930],
        [ 0.7778,  0.4139],
        [ 0.4035,  0.6504],
        [ 0.0781, -0.0228],
        [ 0.0316,  0.0502],
        [-0.0242,  0.2527],
        [-0.3394,  0.5411],
        [-0.1472,  0.3472],
        [-0.6998,  0.1616],
        [ 0.1409,  0.6075]], grad_fn=<AddmmBackward>)

```

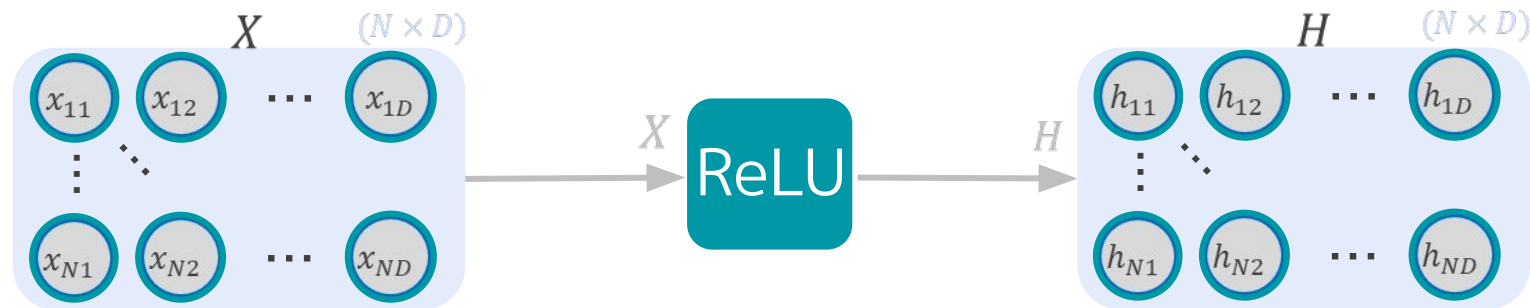
inputの次元(N, D_x)=(10, 3)
outputの次元(N, D_h)=(10, 2)

ReLU関数のバッチ版順伝播を記す

$$h_{ij} = \text{ReLU}(x_{ij}) = \begin{cases} x_{ij} & (x_{ij} > 0) \\ 0 & (x_{ij} \leq 0) \end{cases}$$

Point

1要素ずつReLUに通す



$$h_{ij} = \text{ReLU}(x_{ij}) = \begin{cases} x_{ij} & (x_{ij} > 0) \\ 0 & (x_{ij} \leq 0) \end{cases}$$

```
1  import torch
2  import torch.nn as nn
3
4  N = 10 # バッチサイズ
5  D = 3 # カラムが3次元のデータ
6  relu = nn.ReLU()
7  input = torch.randn(N, D)
8  output = relu(input)
9  print(output)
```

```
tensor([[1.6895, 0.7285, 0.0000],
        [0.9630, 0.0000, 0.9106],
        [0.1444, 0.0000, 0.1002],
        [0.2353, 0.0000, 0.0000],
        [0.0000, 0.7483, 0.8810],
        [0.0000, 0.0000, 0.8269],
        [0.0000, 0.0000, 0.4797],
        [0.2976, 0.0607, 0.2066],
        [0.9198, 0.0000, 0.5978],
        [0.6547, 0.0000, 0.3788]])
```

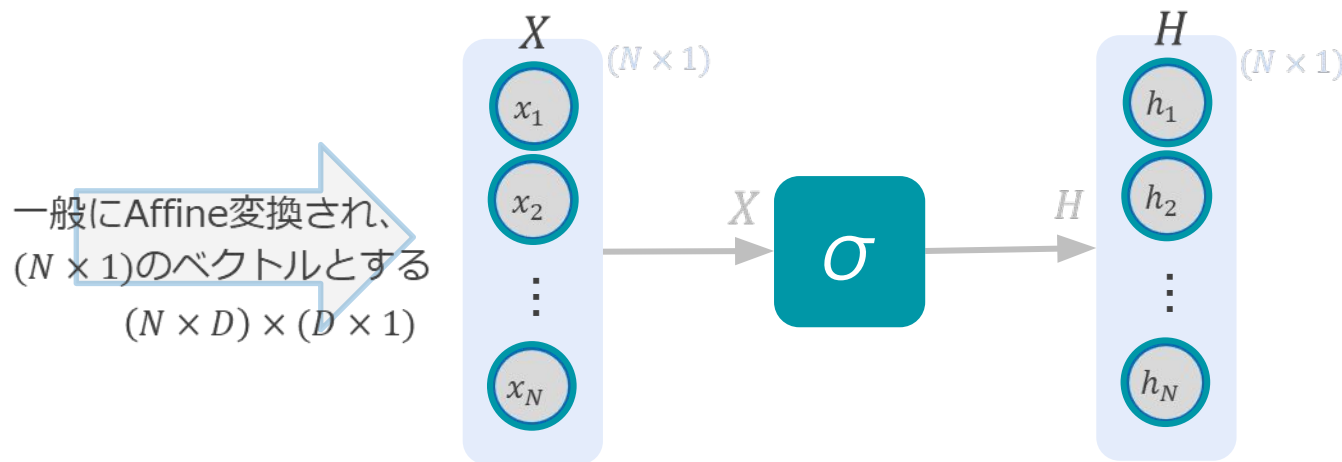
inputの次元(N, D_x)=(10, 3)
outputの次元(N, D_h)=(10, 3)

シグモイド関数のバッチ版順伝播を記す

$$h_i = \sigma(x_i) = \frac{1}{1 + e^{-x_i}}$$

Point

1要素ずつシグモイドに通す



$$h_i = \sigma(x_i) = \frac{1}{1 + e^{-x_i}}$$

```

1  N = 10 # バッチサイズ
2  D = 1 # カラムが1次元のデータ
3  sigmoid = nn.Sigmoid()
4  input = torch.randn(N, D)
5  output = sigmoid(input)
6  print(output)

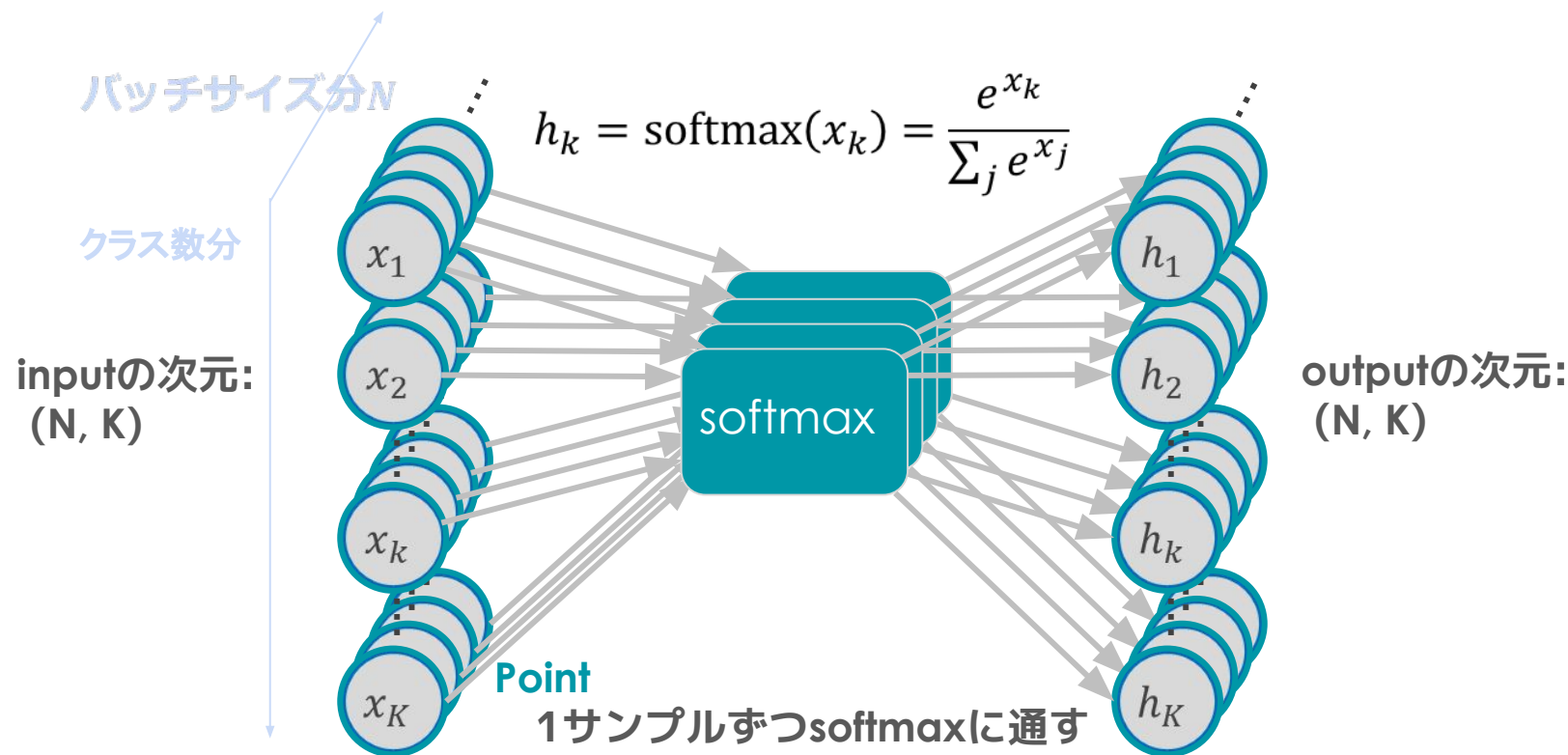
```

```

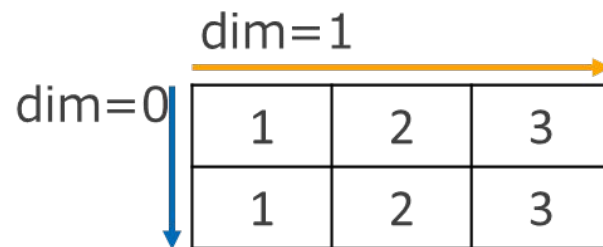
tensor([[0.2597],
        [0.0730],
        [0.2210],
        [0.1687],
        [0.5634],
        [0.6582],
        [0.8052],
        [0.1817],
        [0.4726],
        [0.6259]])

```

inputの次元(N, D_x)=(10, 1)
 outputの次元(N, D_h)=(10, 1)



$$h_k = \text{softmax}(x_k) = \frac{e^{x_k}}{\sum_j e^{x_j}}$$



```

1  sm0 = nn.Softmax(dim=0) # 列方向に足し算
2  sm1 = nn.Softmax(dim=1) # 行方向に足し算
3  input = torch.tensor([[1, 2, 3], [1, 2, 3]], dtype=torch.float64)
4  output0 = sm0(input)
5  output1 = sm1(input)

```

```

6
7  print(input)      tensor([[1., 2., 3.],
8                    [1., 2., 3.]], dtype=torch.float64)
9  print(output0)    tensor([[0.5000, 0.5000, 0.5000],
10                        [0.5000, 0.5000, 0.5000]], dtype=torch.float64)
11 print(output1)    tensor([[0.0900, 0.2447, 0.6652],
12                        [0.0900, 0.2447, 0.6652]], dtype=torch.float64)

```

inputの次元:(2, 3)

outputの次元: (2, 3)

通常、データは以下のような形のため、dim=1に設定する

データ

説明変数

ID				
1				
2				
3				
4				
5				
...				

softmax関数の設定

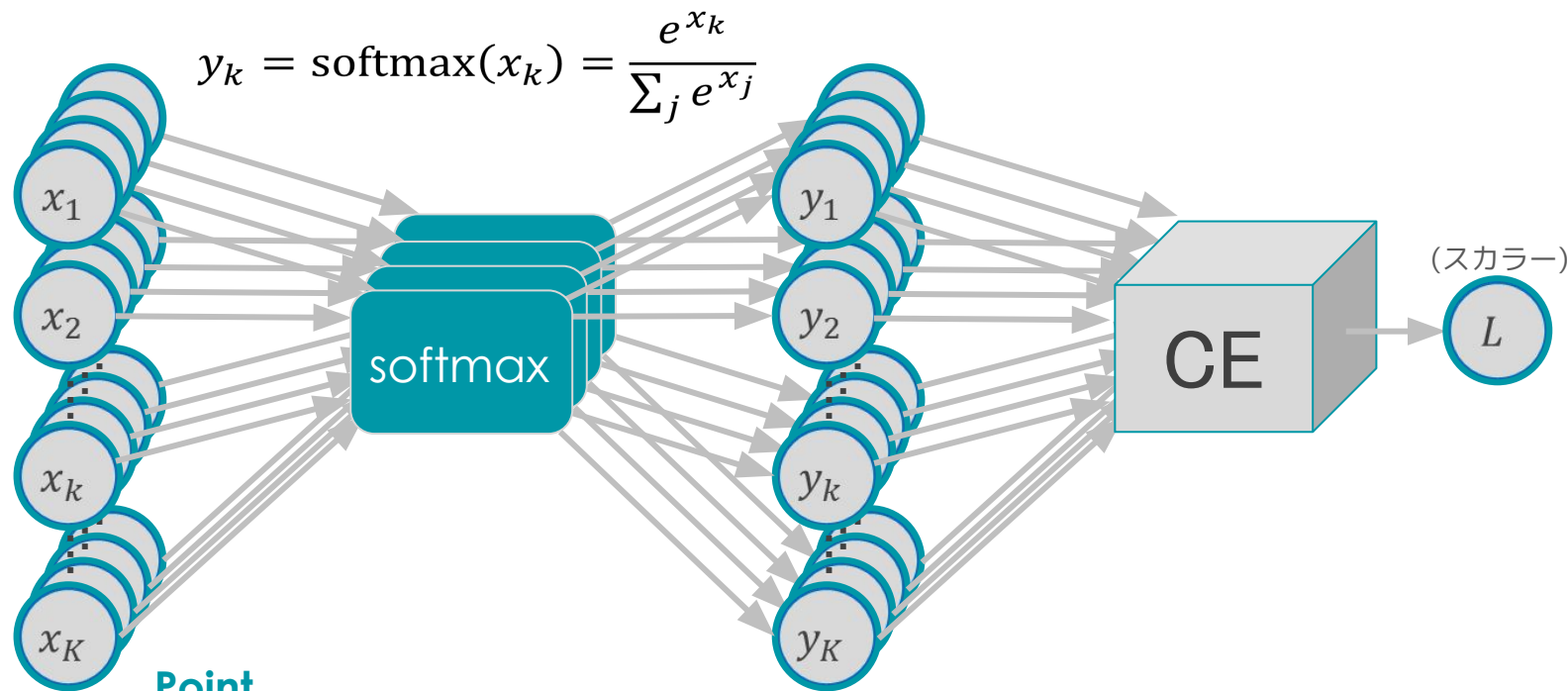
dim=1

dim=0

1	2	3
1	2	3

データサンプルの方向(dim=0)ではなく、
変数が並ぶ方向(dim=1)に設定する。

その結果、1サンプルずつsoftmaxに通すことができる



1サンプルずつSoftmax+CEに通し、
CEでバッチサイズ分の損失平均を算出する

【Chapter03】MLP基礎（バッチ処理編）

バッチ処理（逆伝播）

スカラー量 L の行列 $W(2 \times 3)$ における偏微分行列 $\frac{\partial L}{\partial W}$ を考える

$$W = \begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{pmatrix} \text{としたとき} \quad \frac{\partial L}{\partial W} = \begin{pmatrix} \frac{\partial L}{\partial w_{11}} & \frac{\partial L}{\partial w_{12}} & \frac{\partial L}{\partial w_{13}} \\ \frac{\partial L}{\partial w_{21}} & \frac{\partial L}{\partial w_{22}} & \frac{\partial L}{\partial w_{23}} \end{pmatrix}$$

(2×3) の同じ形式

一般に、スカラー量 L の行列 $\mathbf{W}(D \times H)$ における偏微分行列 $\frac{\partial L}{\partial \mathbf{W}}$ は

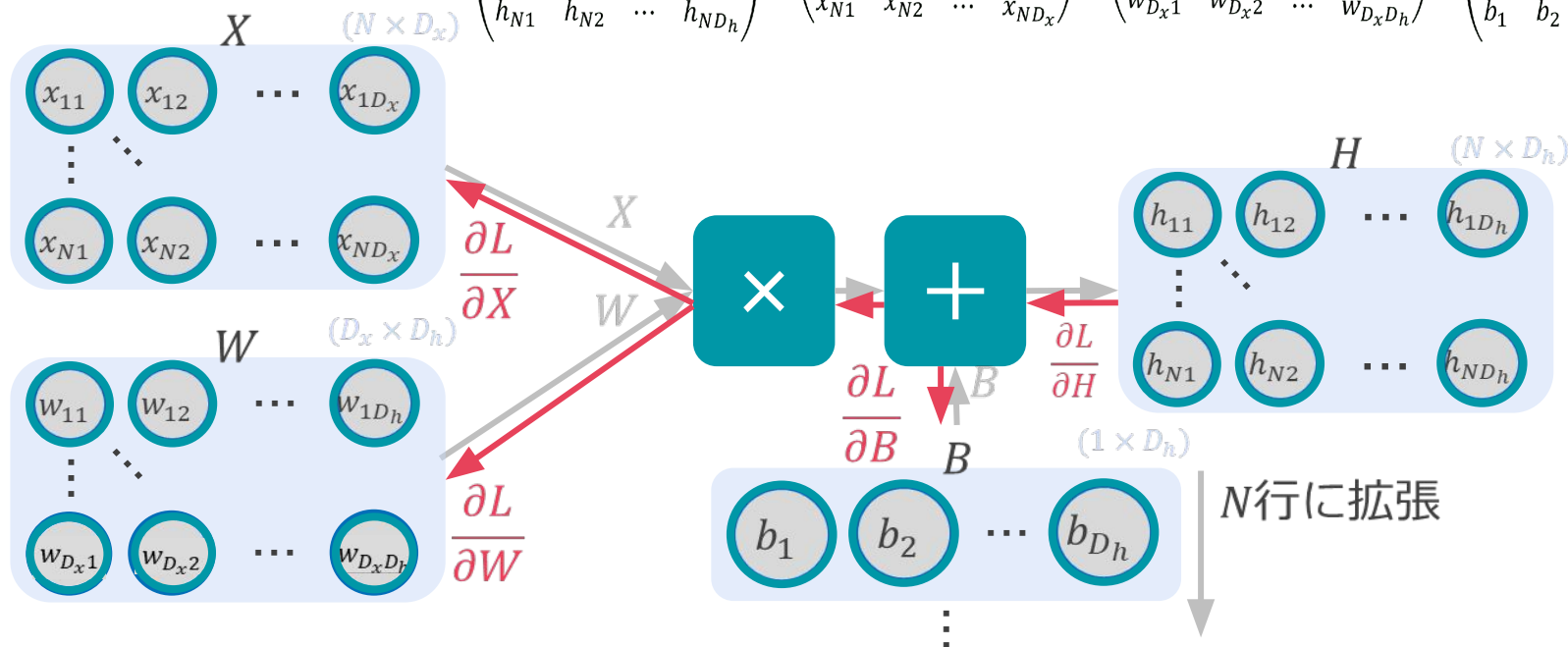
$$\mathbf{W} = \begin{pmatrix} w_{11} & \cdots & w_{1H} \\ \vdots & \ddots & \vdots \\ w_{D1} & \cdots & w_{DH} \end{pmatrix} \text{としたとき} \quad \frac{\partial L}{\partial \mathbf{W}} = \begin{pmatrix} \frac{\partial L}{\partial w_{11}} & \cdots & \frac{\partial L}{\partial w_{1H}} \\ \vdots & \ddots & \vdots \\ \frac{\partial L}{\partial w_{D1}} & \cdots & \frac{\partial L}{\partial w_{DH}} \end{pmatrix}$$

$(D \times H)$ の同じ形式

バッチ版Affine変換の逆伝播を記す

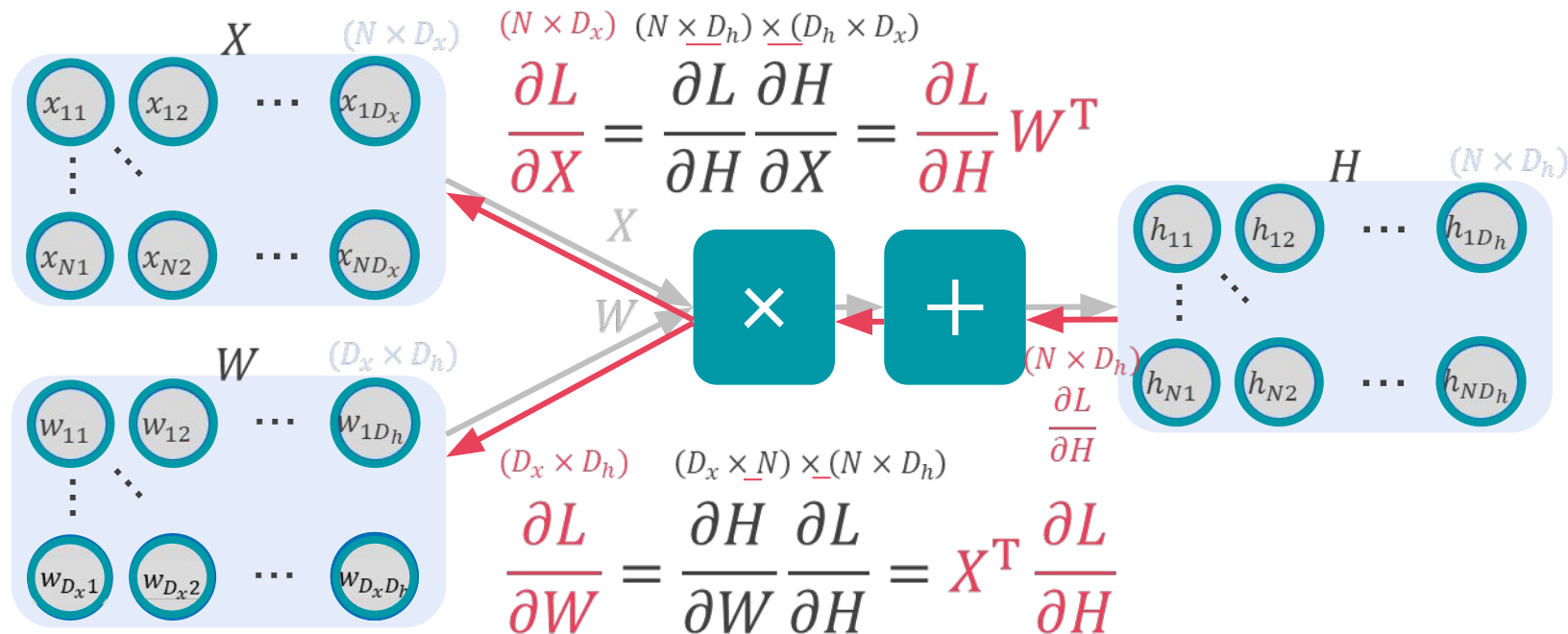
$$H = X \times W + B$$

$$\begin{pmatrix} h_{11} & h_{12} & \cdots & h_{1D_h} \\ h_{21} & h_{22} & \cdots & h_{2D_h} \\ \vdots & \vdots & \ddots & \vdots \\ h_{N1} & h_{N2} & \cdots & h_{ND_h} \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1D_x} \\ x_{21} & x_{22} & \cdots & x_{2D_x} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{ND_x} \end{pmatrix} \times \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1D_h} \\ w_{21} & w_{22} & \cdots & w_{2D_h} \\ \vdots & \vdots & \ddots & \vdots \\ w_{D_x1} & w_{D_x2} & \cdots & w_{D_xD_h} \end{pmatrix} + \begin{pmatrix} b_1 & b_2 & \cdots & b_{D_h} \\ b_1 & b_2 & \cdots & b_{D_h} \\ \vdots & \vdots & \ddots & \vdots \\ b_1 & b_2 & \cdots & b_{D_h} \end{pmatrix}$$



連鎖律を用いて、バッチ版Affine変換の逆伝播の詳細を記す

$$H = X \times W + B$$

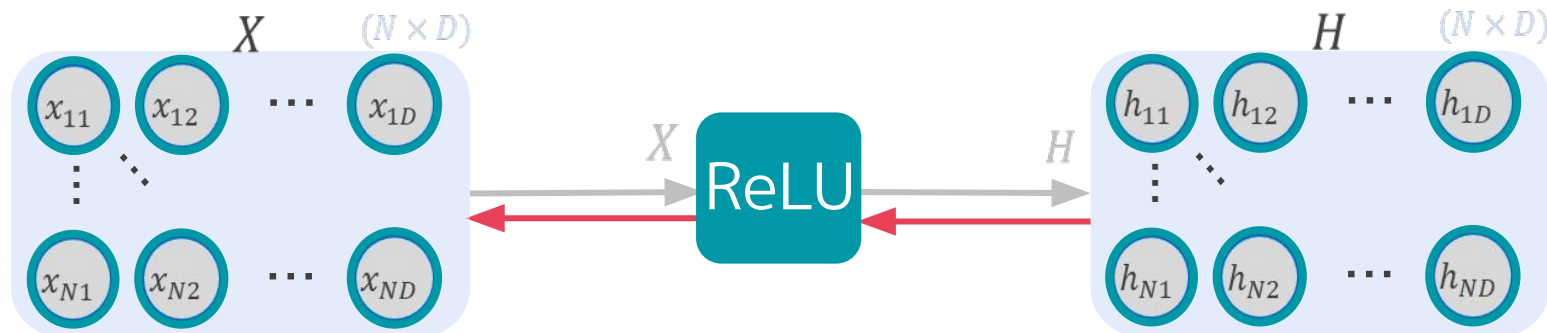


ReLU関数のバッチ版逆伝播を記す

$$h_{ij} = \text{ReLU}(x_{ij}) = \begin{cases} x_{ij} & (x_{ij} > 0) \\ 0 & (x_{ij} \leq 0) \end{cases}$$

Point

順伝播では1要素ずつReLUに通したので、逆伝播でも「要素ごとの偏微分」を伝播させる

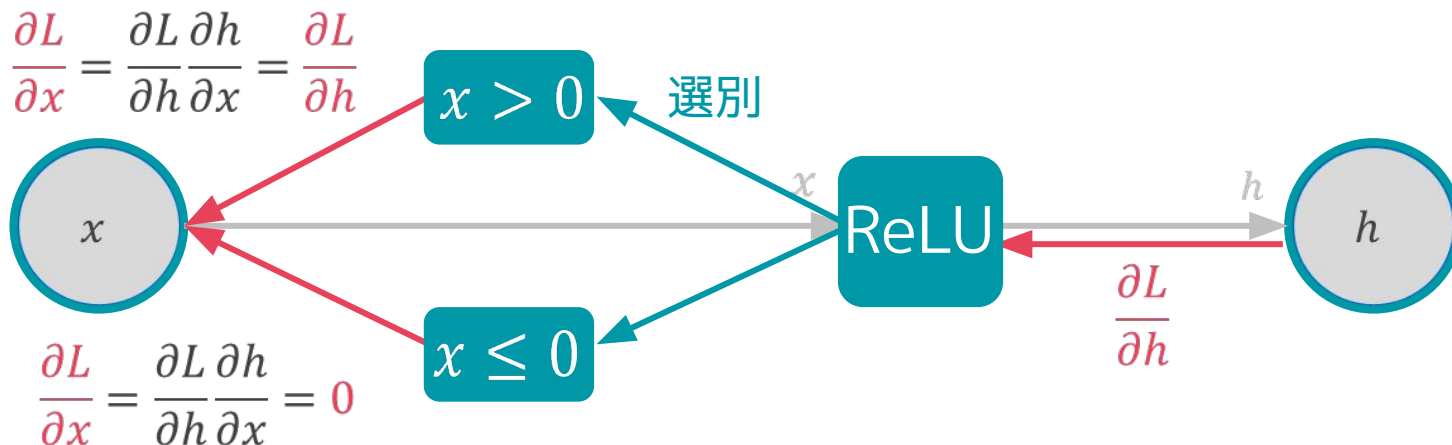


バッチ版ReLU関数の微分

$$h_{ij} = \text{ReLU}(x_{ij}) = \begin{cases} x_{ij} & (x_{ij} > 0) \\ 0 & (x_{ij} \leq 0) \end{cases}$$

要素ごとの偏微分 $\frac{\partial h_{ij}}{\partial x_{ij}}$ を求める
これはオンライン版と同じである

オンライン版の再掲



ReLU関数のバッチ版逆伝播の詳細を記す

$$h_{ij} = \text{ReLU}(x_{ij}) = \begin{cases} x_{ij} & (x_{ij} > 0) \\ 0 & (x_{ij} \leq 0) \end{cases}$$

Point

要素ごとの偏微分を逆伝播させる

$$\frac{\partial L}{\partial x_{ij}} = \frac{\partial L}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial x_{ij}} = \frac{\partial L}{\partial h_{ij}}$$



$$\frac{\partial L}{\partial x_{ij}} = \frac{\partial L}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial x_{ij}} = 0$$

$x_{ij} > 0$

選別

$x_{ij} \leq 0$

x_{ij}

ReLU

h_{ij}

h_{ij}

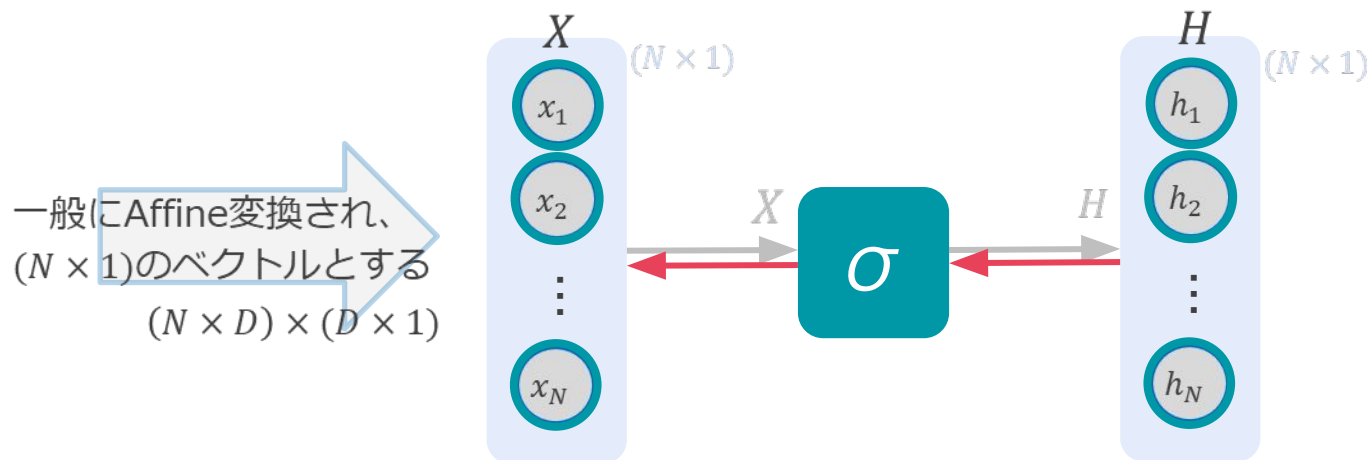
$$\frac{\partial L}{\partial h_{ij}}$$

シグモイド関数のバッチ版逆伝播を記す

$$h_i = \sigma(x_i) = \frac{1}{1 + e^{-x_i}}$$

Point

ReLUと同様に、順伝播では1要素ずつ σ に通したので、逆伝播でも「要素ごとの偏微分」を伝播させる



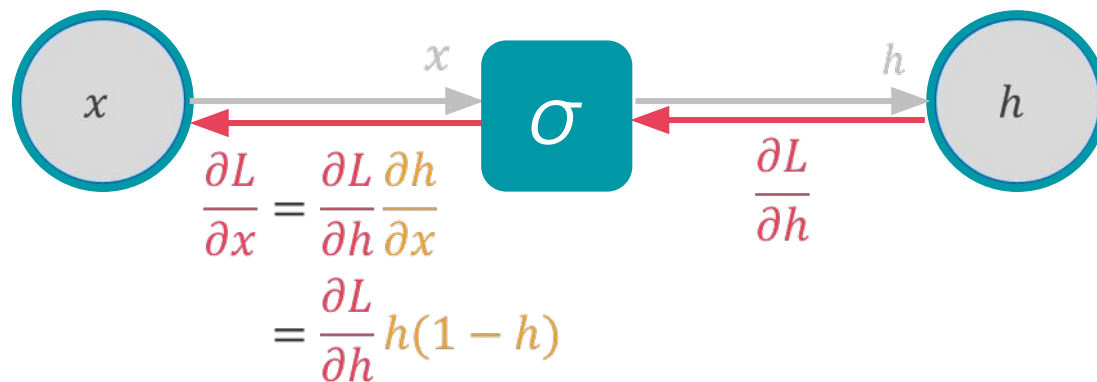
バッチ版シグモイド関数の微分

$$h_i = \sigma(x_i) = \frac{1}{1 + e^{-x_i}}$$

要素ごとの偏微分 $\frac{\partial h_{ij}}{\partial x_{ij}}$ を求める

これはオンライン版と同じである

オンライン版の再掲

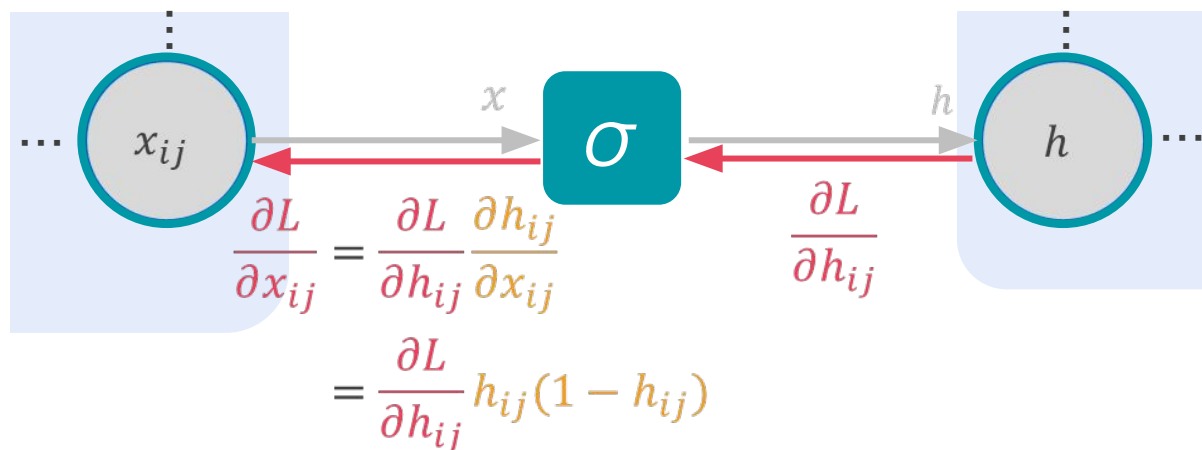


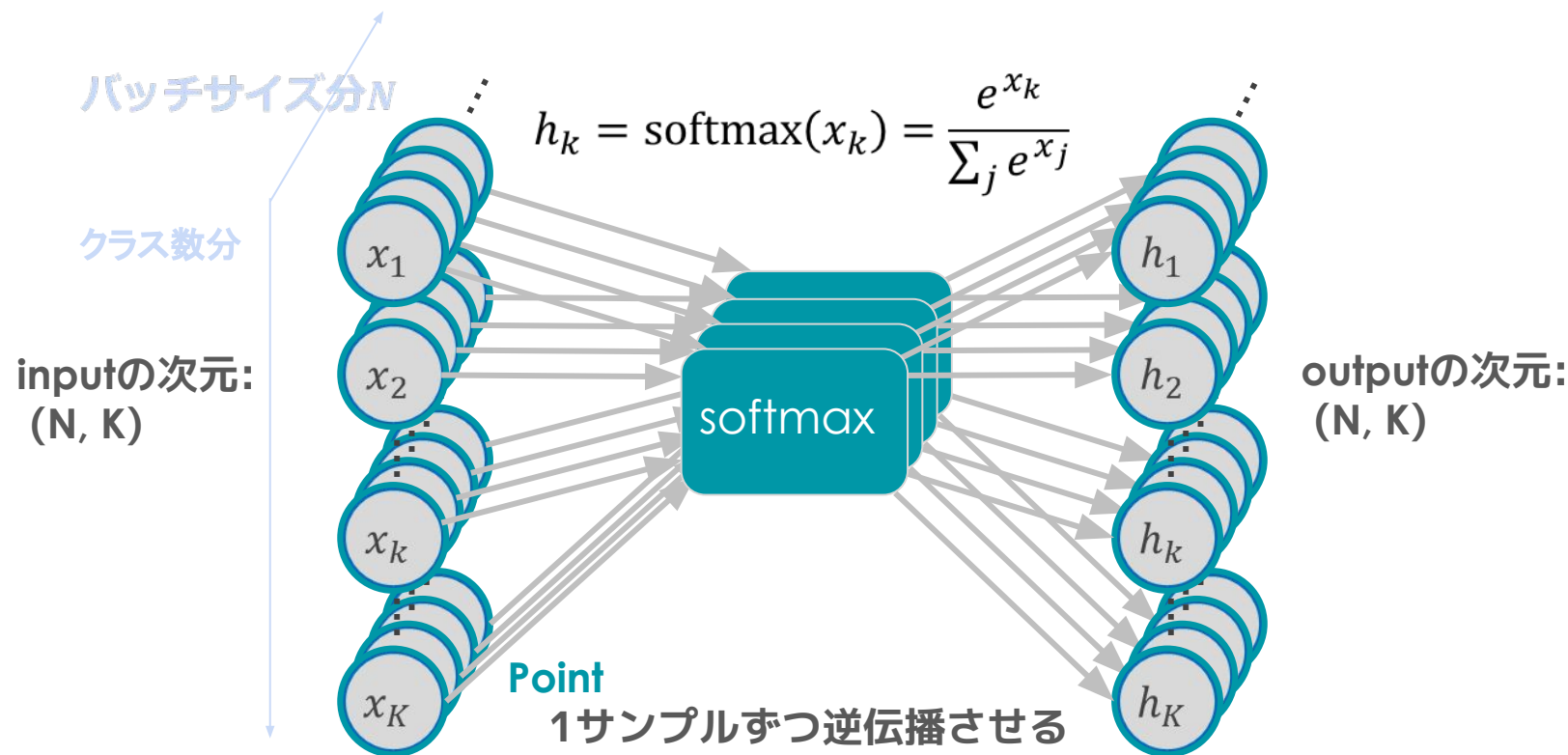
シグモイド関数のバッチ版逆伝播の詳細を記す

$$h_i = \sigma(x_i) = \frac{1}{1 + e^{-x_i}}$$

Point

要素ごとの偏微分を逆伝播させる

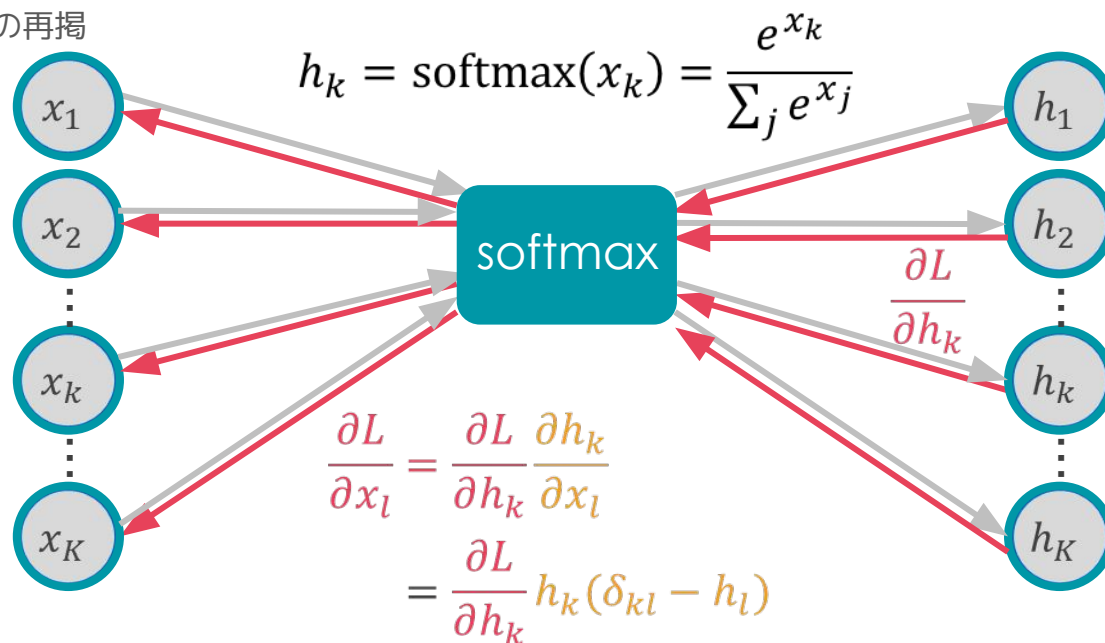




バッチ版softmax関数の微分

1サンプルごとに独立の計算なのでオンライン版×Nと考える

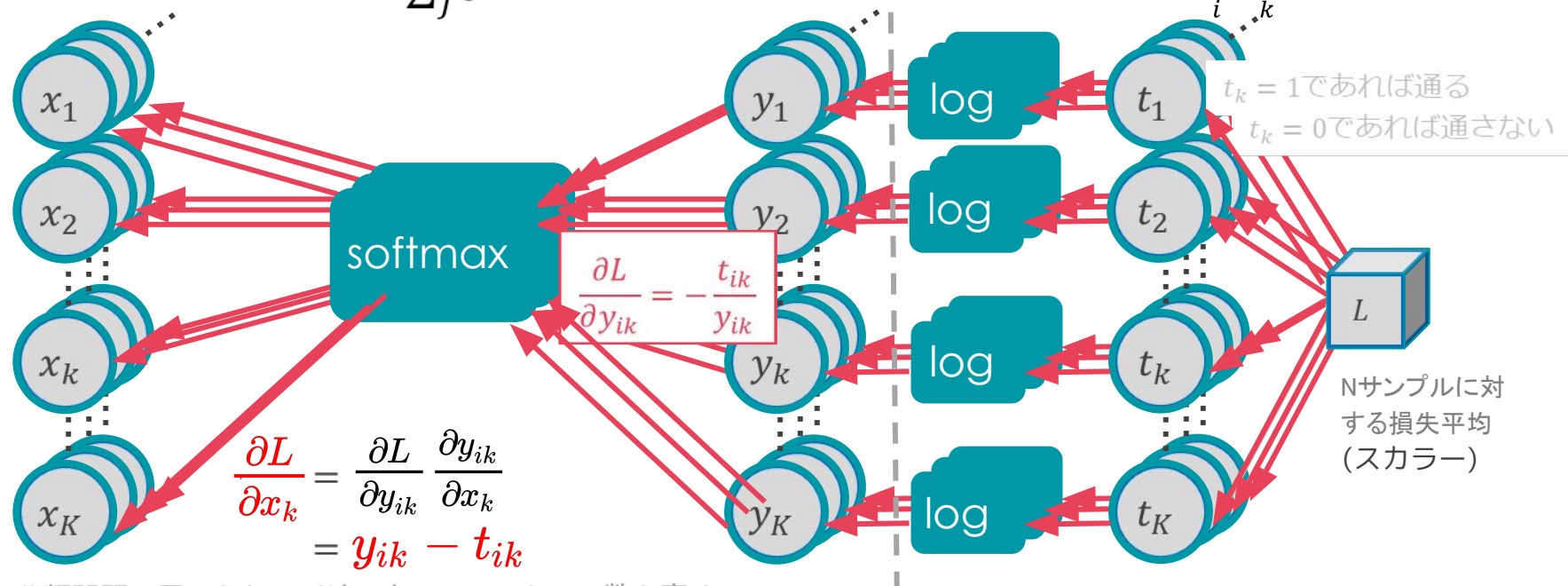
オンライン版の再掲



バッチ版Softmax With CrossEntropyの逆伝播の結果のみ記す

$$h_k = \text{softmax}(x_k) = \frac{e^{x_k}}{\sum_j e^{x_j}}$$

$$L = CE(t, y) = -\frac{1}{N} \sum_i \sum_k t_{ik} \log y_{ik}$$



分類問題に用いられる手法である。Kはクラス数を表す

タスクによって出力ユニットが異なり、適切な損失関数も異なる

タスク		出力ユニット	活性化関数	損失関数	損失関数の式
回帰	線形回帰	線形ユニット	恒等関数	二乗和誤差関数	$\sum_i (t_i - y_i)^2$
分類	二値分類	シグモイドユニット	sigmoid関数	バイナリークロスエントロピー関数	$-\sum_i (t_i \log y_i + (1 - t_i) \log(1 - y_i))$
	多値分類	ソフトマックスユニット	softmax関数	クロスエントロピー誤差関数	$-\sum_i \sum_k t_{ik} \log y_{ik}$

t_i : 目的変数

y_i : NN出力値

k : 所属クラスを表す添え字 (one-hot encoding後)

i : データサンプルの添え字 (何番目のサンプルかを表す)

バッチ処理ver.

【補足】線形ユニットに対する損失関数

バッチ処理でも、ガウス分布の対数尤度の最大化と二乗和誤差の最小化は等しい

前提

線形ユニットはガウス分布 $\mathcal{N}(t; y, I)$ を出力するためのユニットである

$$\begin{aligned} L(\theta) &= -\log N(t; y, I) \\ &= -\log e^{-\frac{\sum_i (t_i - y_i)^2}{2\sigma^2}} \\ &= \sum_i (t_i - y_i)^2 + \text{const.} \end{aligned}$$

二乗和誤差関数

t : 目的変数

y : NN出力値

$y_i = b + W^T h$ (i 番目のデータの予測値)

損失関数 $L(\theta)$ が最小となるような W を求めることを「損失関数の最適化」という

【補足】シグモイドユニットに対する損失関数

前提

シグモイドユニットはベルヌーイ分布 $\text{Bern}(t; y)$ を出力するためのユニットである

$$L(\theta) = -\log \text{Bern}(t; y)$$

t : 目的変数

y : NN出力値

$$= -\sum_i (t_i \log y_i + (1 - t_i) \log(1 - y_i)) = \sigma(b + W^T h)$$

(i番目のデータの予測値)

バイナリークロスエントロピー誤差関数

【補足】 one-hot encodingされている場合、
正解の予測確率の対数をとって和を求めるだけ

【補足】ソフトマックスユニットに対する損失関数

バッチ処理でも、
マルチヌーイ分布の対数尤度の最大化とクロスエントロピー誤差関数の最小化は等しい
前提

ソフトマックスユニットはマルチヌーイ分布 $\text{Muln}(t; y)$ を出力するためのユニットである

$$L(\theta) = -\log \text{Muln}(t; y)$$

$$= - \sum_i \sum_k t_{ik} \log y_{ik}$$

クロスエントロピー誤差関数

【補足】 one-hot encodingされている場合、
正解の予測確率の対数をとって和を求めるだけ

t : 目的変数

t_{ik} : i 番目のデータのクラス k
のone-hot label

y : NN出力値

$y_{ik} = \text{softmax}^{(k)}(b + W^T h)$
(i 番目のデータのクラス k
に対する予測値)

【Chapter03】MLP基礎（バッチ処理編）

最適化手法

バッチ処理では、すべての訓練データを一度にNNに流した

しかし【問題点】がある

データ数が膨大であれば計算量が膨大となりパラメータ更新が進みづらい

そこで「ミニバッチ処理」

訓練データをいくつかのグループに分けて
NNに投入する

ミニバッチ処理の【利点】

分割すると細かい更新が多く収束が速い

イメージ

ID	col1	col2	col3	col4
1				
2				
...				
100				
101				
...				

1回目の学習

2回目の学習

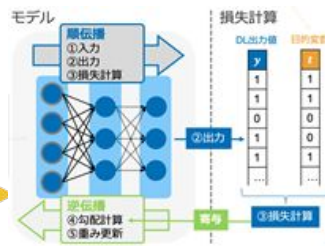
バッチサイズとは、パラメータを一度更新するための損失計算に必要なデータ数

例)バッチサイズ=100のとき

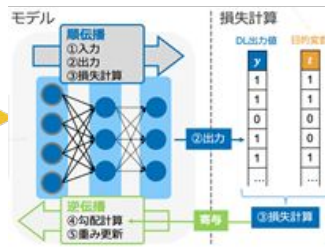
ID	col1	col2	col3	col4
1				
2				
...				
100				
101				
...				

1回目の学習

2回目の学習



- 1.入力・出力
- 2.損失計算
- 3.勾配更新
- 4.重み更新



- 1.入力・出力
- 2.損失計算
- 3.勾配更新
- 4.重み更新

つまり・・・

オンライン学習

バッチサイズ=1

バッチ学習

バッチサイズ
= 全データ数

バッチサイズとは、パラメータを一度更新するための損失計算に必要なデータ数

	オンライン学習	ミニバッチ処理	バッチ処理（バッチ学習）
バッチサイズN	1	中間の値 (例: 64/128/256)	全データ数
損失関数 (二乗和誤差関数)	$\sum_n^1 (t_n - y_n)^2$	$\sum_n^{128} (t_n - y_n)^2$	$\sum_n^N (t_n - y_n)^2$

バッチサイズのトレードオフ

一般に、バッチサイズが大きければ勾配推定はより正確になるが、計算コストは高まる

一方で、バッチサイズが小さければ勾配推定の誤差を高める代わりに、更新頻度が早く収束が早くなると言われている

trade-offまとめ	勾配推定	収束まで
バッチサイズ大	より正確	遅い
バッチサイズ小	誤差を高める	早い

よく用いられる有効なバッチサイズ

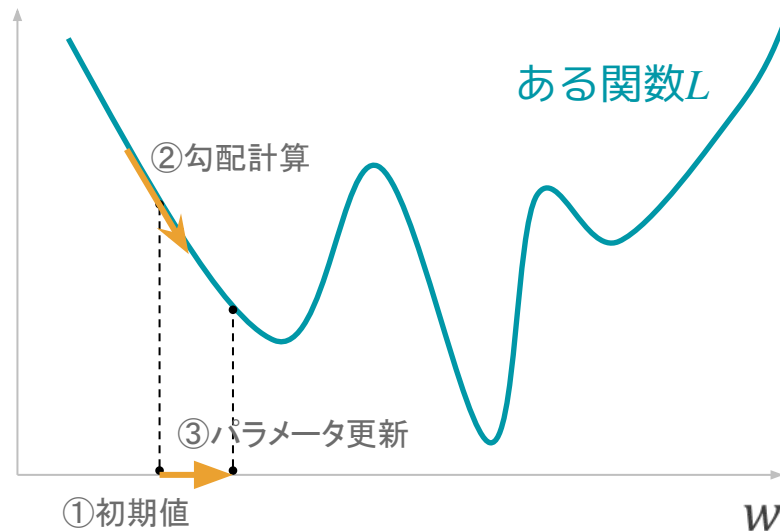
計算リソースが限られている場合	16 / 32 / 64 / 128 / 256
計算リソースに余裕がある場合	512 / 1024 またはそれ以上

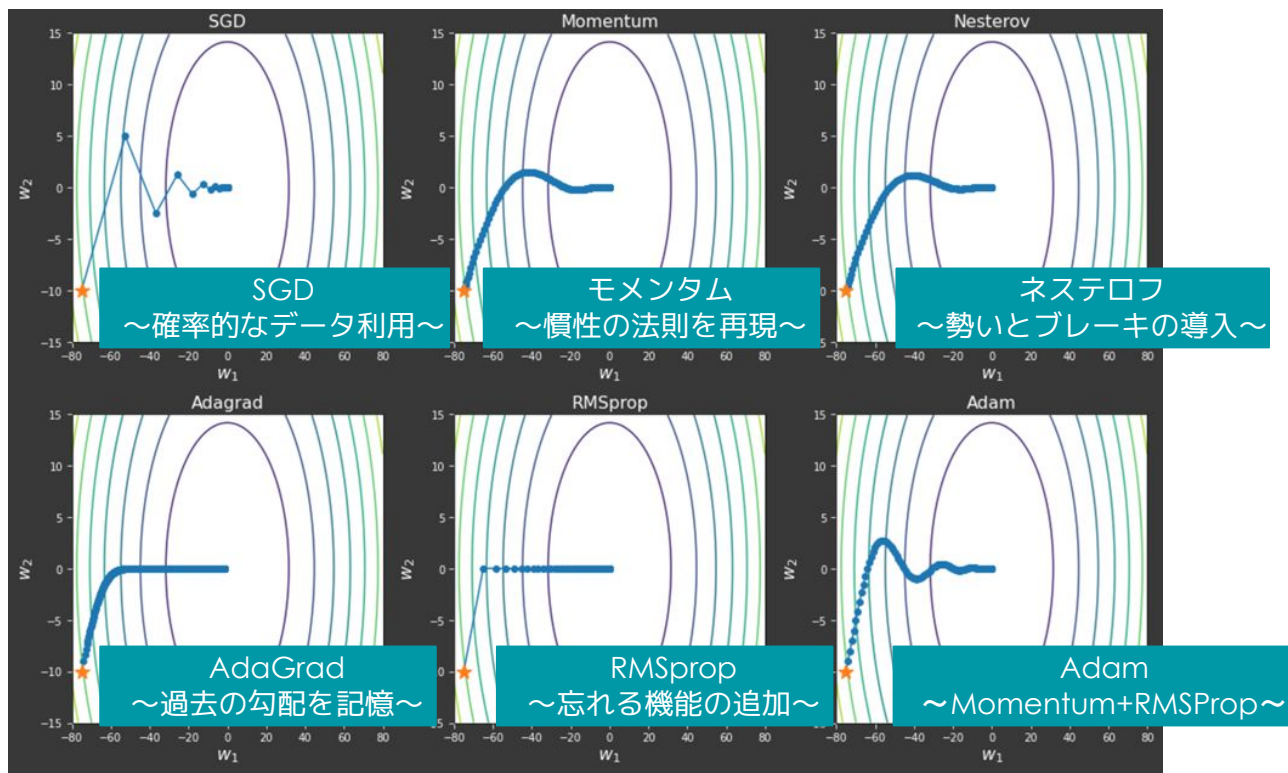
最適化とは、勾配の情報をヒントにパラメータをアップデートすること

最適化の基本的流れ

- ①パラメータ w の初期値を決める
- ②その初期値における勾配 $\frac{\partial L}{\partial w}$ を計算
- ③勾配の大きさやその方向を頼りにパラメータ を更新する

③の方法が各最適化手法によって異なる

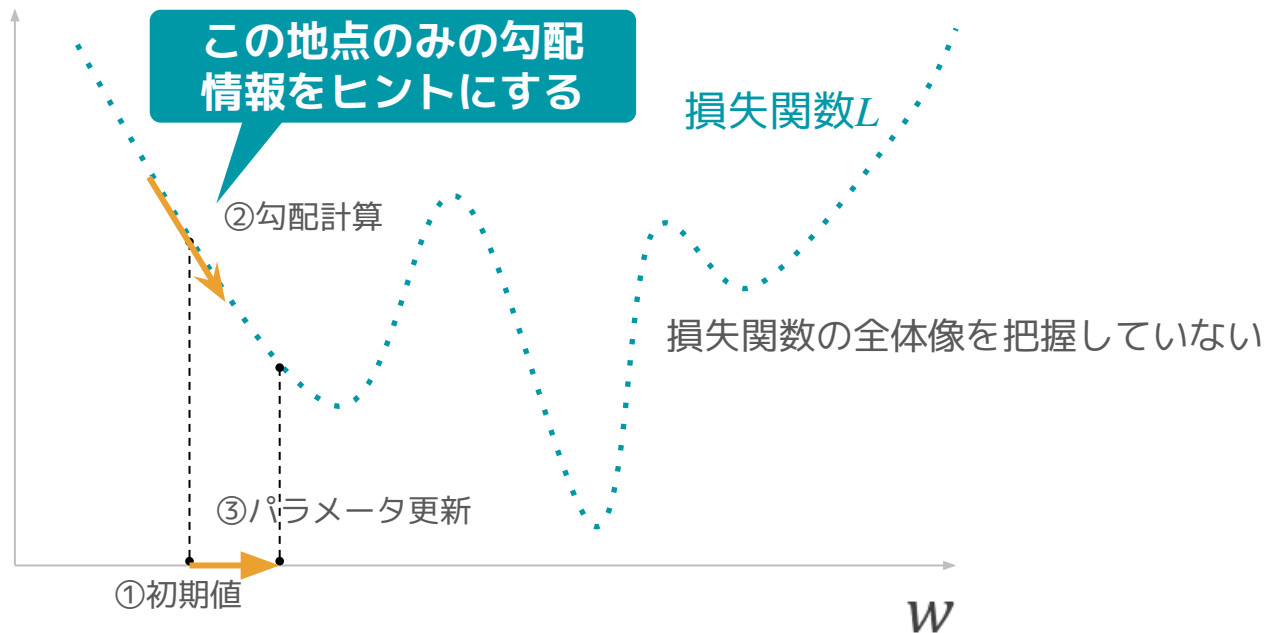




※対象の損失関数によって最適化手法の相性は異なる。図の損失関数の定義式は $L = w_1^2 + 5w_2^2$

最急降下法は、その地点での勾配方向に更新

イメージ図



最急降下法は、その地点での勾配方向に更新

数式

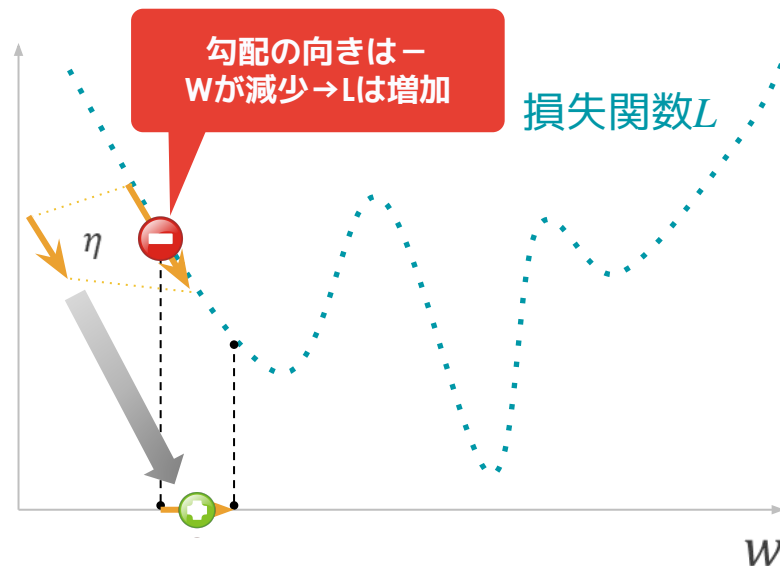
$$w_{t+1} = w_t - \eta \left(\frac{\partial L}{\partial w} \right)_{w=w_t}$$

ハイパーパラメータ(学習率)

η : どれだけ更新させるか

特徴

求めた勾配方向とは**逆向き**にその
大きさだけパラメータを更新する

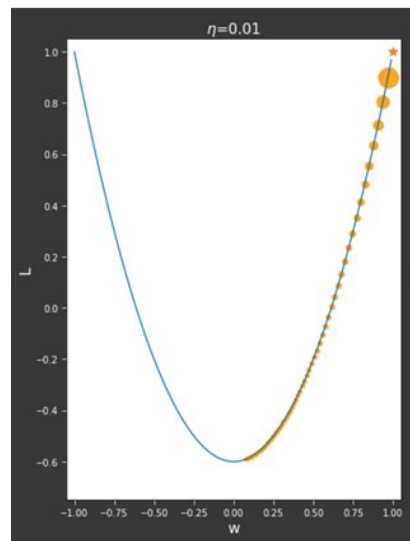


最急降下法は、全データを用いた損失関数を最適化するのでバッチ学習である

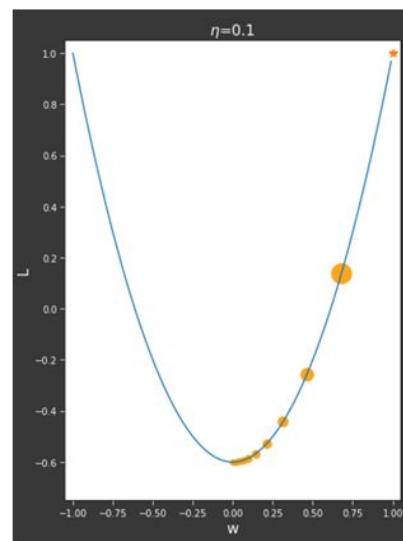
収束までの経路図



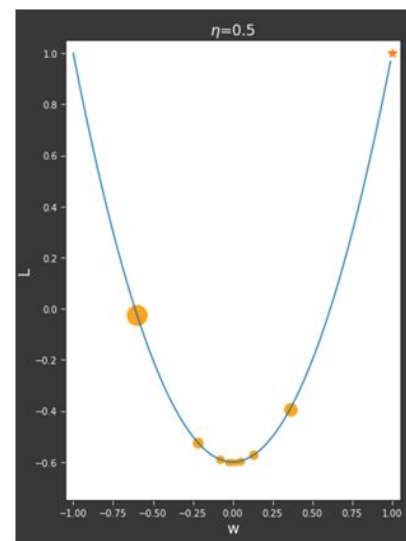
学習率 η が小さすぎ



学習率 η がちょうどいい



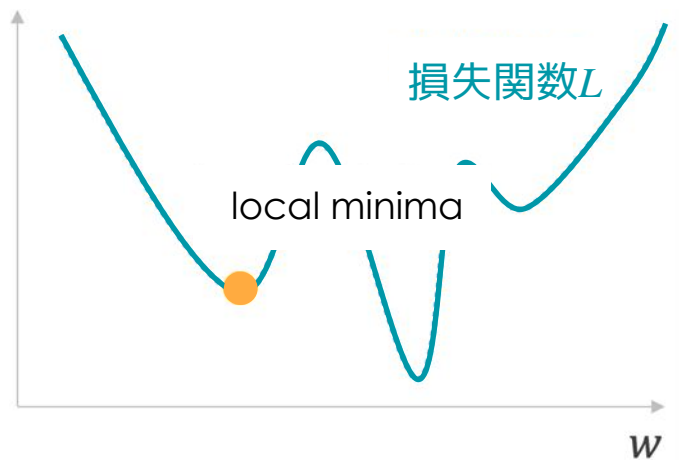
学習率 η が大きすぎ



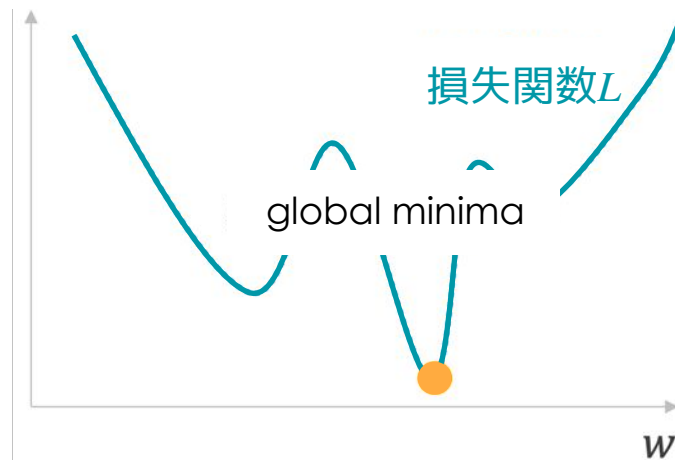
動画で詳細に理解したい方：<https://ai-trend.jp/basic-study/neural-network/optimizer/>

最急降下法では「局所的最適解(local minima)」にはまってしまい、
本当に低い谷へ落ち着くことが難しい

(\Leftrightarrow 全体の最適解を大域的最適解(global minima)と呼んだりもする)



ここから抜け出せなくなる



こちらが望ましい

SGDは、確率的な山の一地点での勾配方向に更新

数式

$$w_{t+1} = w_t - \eta \left(\frac{\partial L(X)}{\partial w} \right)_{w=w_t}$$

L はデータ X に依存

ハイパーパラメータ(学習率)

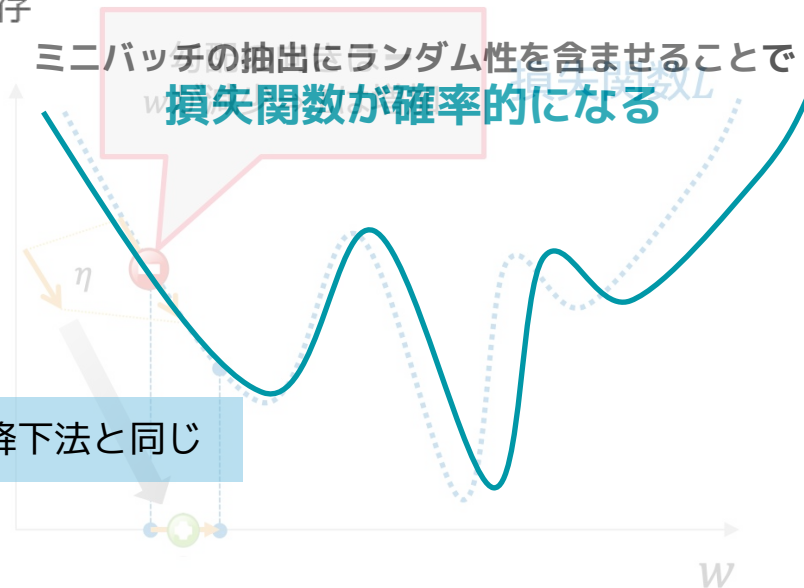
η : どれだけ更新させるか

特徴

求めた勾配方向とは逆向きにその
大きさだけパラメータを更新する

手法は最急降下法と同じ

ミニバッチの抽出にランダム性を含ませることで
損失関数が確率的になる

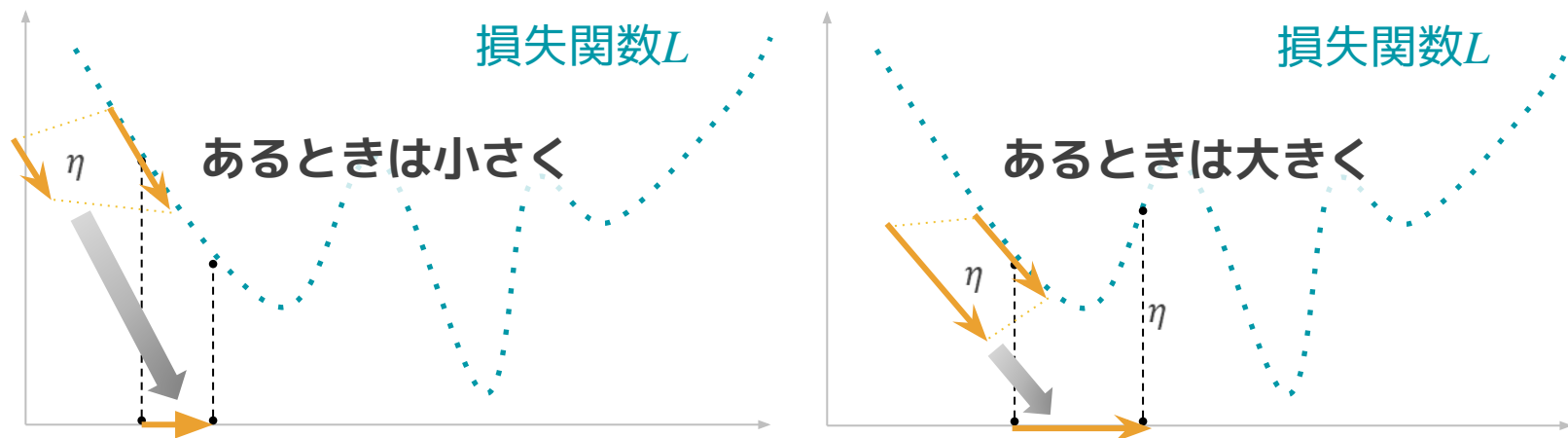


SGDは、グループ化されたサンプルを用いた損失関数を最適化するのでミニバッチ処理である

SGDは、データを変えて損失関数自体を確率的に与えた

【背景】SGDが登場するまでの解決策

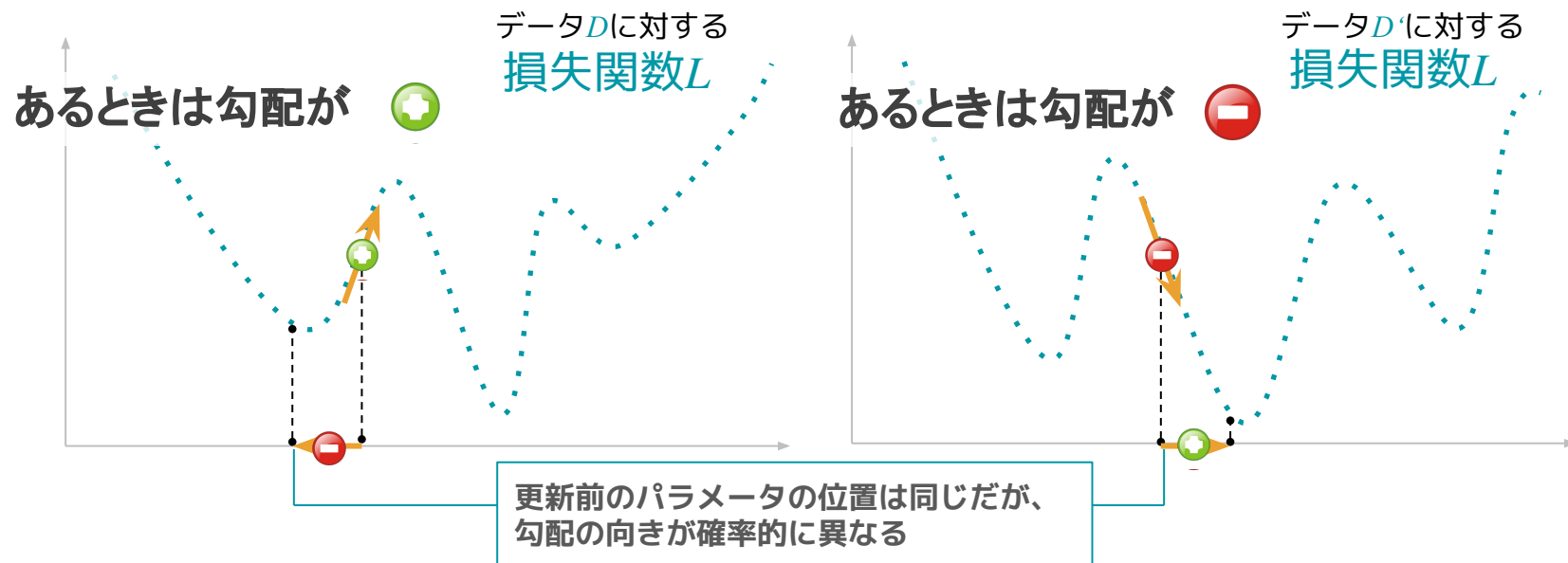
例) 学習率を確率的に変える



最適化手法を変えようとしていた
しかし、損失関数が複雑すぎてうまくいかない

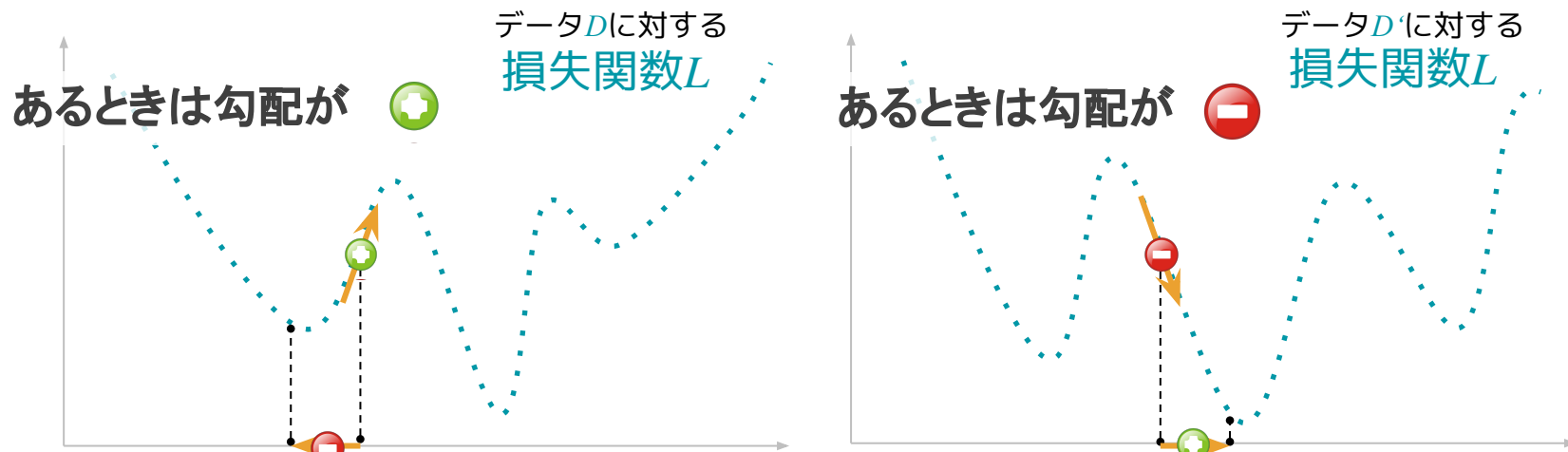
SGDは、データを変えて損失関数自体を確率的に与えた

SGDは学習率ではなく、データを変えた



SGDは、最適化の途中で振動して収束が遅くなる

確率的に変化する損失関数を最適化することにも弱点がある。
それは、進む方向がランダムになり振動を助長し収束が遅くなること



【解決策】振動を抑制するように、前回までの勾配によって「勢い」を変化させる。
それがモメンタム

関数平面上で“ボール”が転がるようにパラメータ更新

数式

$$v_{t+1} = \alpha v_t - \eta \left(\frac{\partial L(X)}{\partial w} \right)_{w=w_t}$$

$$w_{t+1} = w_t + v_{t+1}$$

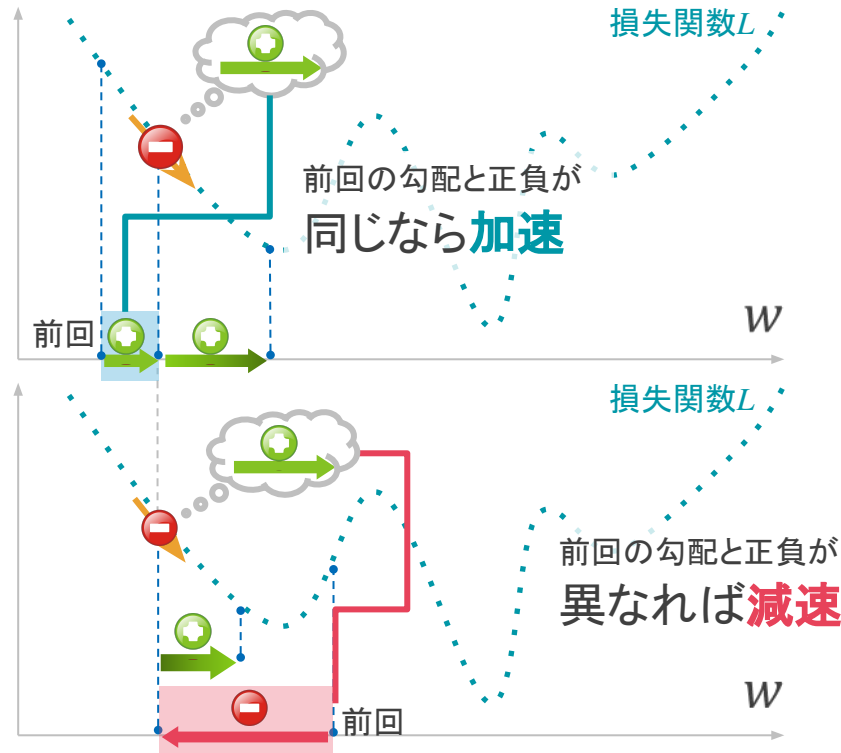
ハイパーパラメータ

η : 今回の損失の勾配の割合(学習率)

α : 前回までの勾配の割合

特徴

SGDに「慣性の法則」を導入
移動平均をとることで振動幅を抑制する



過去の勾配の二乗和を全て記憶し学習率を調整する

数式

$$h_t = h_{t-1} + \left(\frac{\partial L}{\partial w} \right)^2_{w=w_t}$$

$$\eta_t = \frac{\eta_0}{\sqrt{h_t}}$$

$$w_{t+1} = w_t - \eta_t \odot \left(\frac{\partial L}{\partial w} \right)_{w=w_t}$$

行列

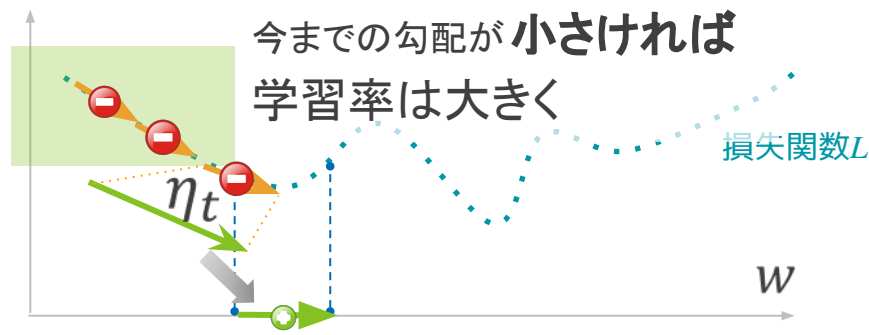
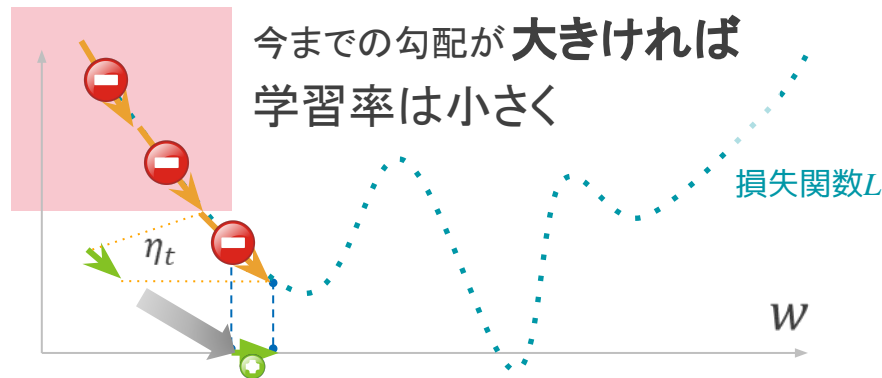
要素ごとの積であり、
パラメータごとに異なる
学習率を設定する

ハイパーパラメータ(学習率)

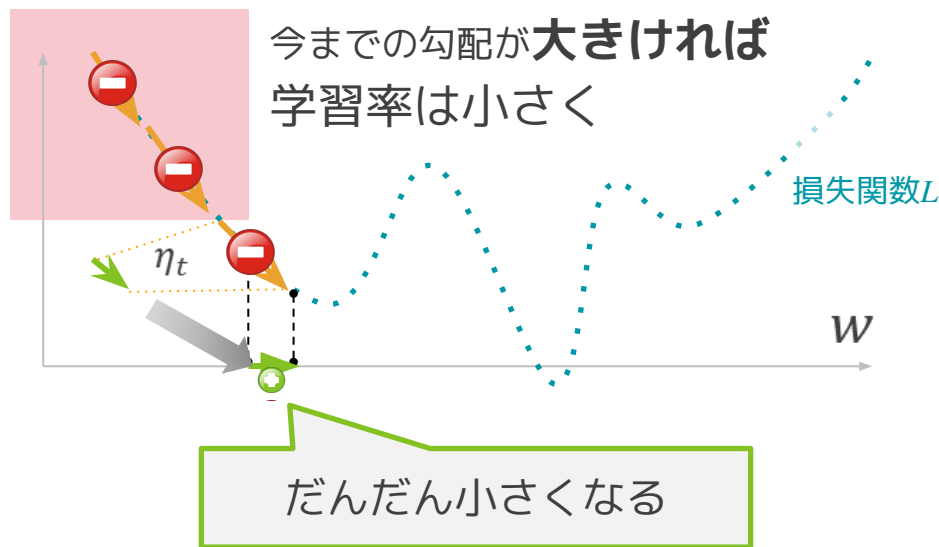
η_0 : 学習率の初期値

特徴

- ✓ **パラメータ**ごとに学習率を更新
- ✓ 前回までに大きく更新されたパラメータの学習率は、より小さく調整される。その逆も然り。(右図参照)



学習が進むにつれて、学習率が小さくなり、やがて0になる
ゆえに学習が進まなくなる可能性がある



古い情報は忘れて、新しい勾配情報が反映されるように記憶し学習率を調整する

数式

$$h_t = \alpha h_{t-1} + (1 - \alpha) \left(\frac{\partial L}{\partial w} \right)^2_{w=w_t}$$
$$\eta_t = \frac{\eta_0}{\sqrt{h_t}}$$
$$w_{t+1} = w_t - \eta_t \left(\frac{\partial L}{\partial w} \right)_{w=w_t}$$

ハイパーパラメータ(学習率)

η_0 : 学習率の初期値

α : 新古の勾配情報の割合

特徴

過去の勾配情報と新しい勾配情報の割合を調整して記憶する。



勾配の2乗を記憶した「ボール」が関数平面上を転がる

数式

Algorithm 1: *Adam*, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation. g_t^2 indicates the elementwise square $g_t \odot g_t$. Good default settings for the tested machine learning problems are $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. All operations on vectors are element-wise. With β_1^t and β_2^t we denote β_1 and β_2 to the power t .

Require: α : Stepsize

Require: $\beta_1, \beta_2 \in [0, 1]$: Exponential decay rates for the moment estimates

Require: $f(\theta)$: Stochastic objective function with parameters θ

Require: θ_0 : Initial parameter vector

$m_0 \leftarrow 0$ (Initialize 1st moment vector)

$v_0 \leftarrow 0$ (Initialize 2nd moment vector)

$t \leftarrow 0$ (Initialize timestep)

while θ_t not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep t)

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)

$\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)

$\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)

$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$ (Update parameters)

end while

return θ_t (Resulting parameters)

初期段階の
不安定性の解消
(後述)

特徴

Momentum + RMSProp

とりあえず「Adam」といった風潮がある

引用: <https://arxiv.org/pdf/1412.6980.pdf>



忘れる機能の追加



忘れる機能の追加



イテレーション数 t でバイアス補正することで、初期段階の不安定性を解消

改善前の【課題】初期段階では $v_1 = 0$ のため、変動が大きくなり不安定になる

$$\begin{aligned} g_t &\leftarrow \nabla_{\theta} f_t(\theta_{t-1}) && \text{(イテレーション数 } t \text{ における勾配)} \\ m_t &\leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t && \text{(新旧の一次モーメントの加重平均)} \\ v_t &\leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2 && \text{(新旧の二次モーメントの加重平均)} \\ \theta_t &\leftarrow \theta_{t-1} - \alpha \cdot m_t / (\sqrt{v_t} + \epsilon) && \text{(パラメータ更新)} \end{aligned}$$

【改善施策】イテレーション数 t でバイアス修正

$$\begin{aligned} g_t &\leftarrow \nabla_{\theta} f_t(\theta_{t-1}) && \text{(イテレーション数 } t \text{ における勾配)} \\ m_t &\leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t && \text{(新旧の一次モーメントの加重平均)} \\ v_t &\leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2 && \text{(新旧の二次モーメントの加重平均)} \\ \widehat{m}_t &\leftarrow m_t / (1 - \beta_1^t) && \text{(イテレーション数 } t \text{ でバイアス補正された一次モーメント)} \\ \widehat{v}_t &\leftarrow v_t / (1 - \beta_2^t) && \text{(イテレーション数 } t \text{ でバイアス補正された二次モーメント)} \\ \theta_t &\leftarrow \theta_{t-1} - \alpha \cdot \widehat{m}_t / (\sqrt{\widehat{v}_t} + \epsilon) && \text{(パラメータ更新)} \end{aligned}$$

```
1  # 最適化したい関数の定義
2  def func(x):
3      return x**2
4
5  # パラメータ初期値
6  x = torch.tensor(-2.0, requires_grad=True)
7  params = [x]
8  # 最適化手法を定義
9  # lr: 学習率
10 opt = optim.SGD(params, lr=0.01) # SGD
11 opt = optim.SGD(params, lr=0.001, momentum=0.9) # Momentum
12 opt = optim.SGD(params, lr=0.001, momentum=0.9, nesterov=True) # Nesterov
13 opt = optim.Adagrad(params, lr=1.0) # AdaGrad
14 opt = optim.RMSprop(params, lr=1.0) # RMSProp
15 opt = optim.Adam(params, lr=1.0) # Adam
16
17 #関数の出力に対して勾配を求め、optimizerで更新を繰り返す
18 for _ in range(80):
19     opt.zero_grad() # 勾配を0にセット(PyTorchのデフォルトは勾配が加算されていく)
20     outputs = func(x)
21     outputs.backward() # 勾配計算
22     opt.step() # パラメータ更新
```


- 1) バッチ処理（順伝播）
- 2) バッチ処理（逆伝播）
- 3) 最適化手法



AVILEN