

# Neural Network の基礎

本谷 秀堅

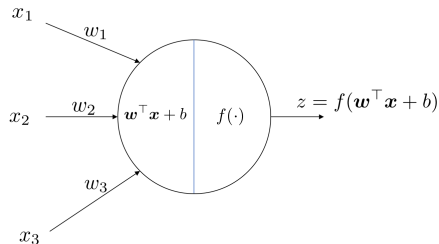
名古屋工業大学

## ニューラルネットワークによる関数の表現と学習

## ユニットの入出力

- 入力： $\mathbf{x} = (x_1, x_2, \dots)^\top$
- 出力： $y = f(u)$ 
  - まず重みとの内積＋定数を計算：
$$u = w_1x_1 + w_2x_2 + \dots + b = \mathbf{w}^\top \mathbf{x} + b$$
  - 次に活性化関数： $f(u) = f(\mathbf{w}^\top \mathbf{x} + b)$

# ニューラルネットワークの構成単位



## ユニットの入出力

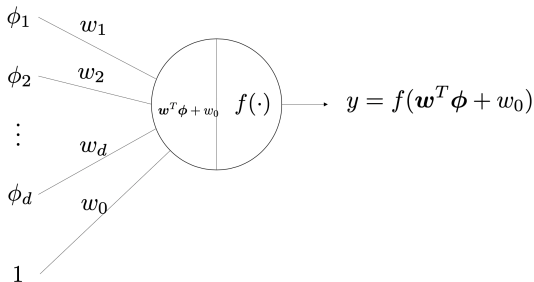
- 入力： $\mathbf{x} = (x_1, x_2, \dots)^\top$
- 出力： $y = f(u)$ 
  - まず重みとの内積＋定数を計算： $u = w_1x_1 + w_2x_2 + \dots + b = \mathbf{w}^\top \mathbf{x} + b$
  - 次に活性化関数： $f(u) = f(\mathbf{w}^\top \mathbf{x} + b)$

# これまでとの比較

$$\phi(\mathbf{x}) = [\phi_1, \phi_2, \dots, \phi_d]^\top$$

活性化関数,  $f(\cdot)$

$$y = f(\mathbf{w}^\top \phi + w_0)$$



# これまでの復習

## 2クラス識別（ロジスティック回帰）

入力,  $\mathbf{x}$ , と出力,  $z \in \mathbb{R}$  の関係：

$$z(\mathbf{x}) = p(y = 1|\mathbf{x}) = \sigma(a) = \frac{1}{1 + \exp(-(\mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}) + w_0))}$$

学習データ  $\mathcal{D} = \{(\mathbf{x}_i, y_i) | i = 1, 2, \dots, N, y_i \in \{0, 1\}\}$  が与えられている  
最尤法による係数  $\mathbf{w}$  の求める：

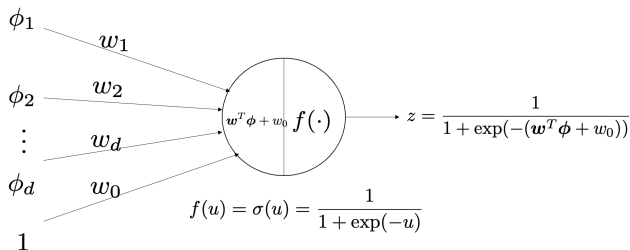
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left[ - \sum_{i=1}^N y_i \log p(y_i | \mathbf{x}_i) \right]$$

$$= \arg \min_{\mathbf{w}} \left[ - \sum_{i=1}^N \{y_i \log z(\mathbf{x}) + (1 - y_i) \log(1 - z(\mathbf{x}))\} \right]$$

# ニューラルネットワークによる2クラス識

入力,  $\mathbf{x}$ , と出力,  $z \in (0, 1)$  の関係

$$z(\mathbf{x}) = p(y = 1|\mathbf{x}) = \sigma(a) = \frac{1}{1 + \exp(-(\mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}) + w_0))}$$



係数,  $\mathbf{w}$ , 決定のためのコスト関数

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left[ - \sum_{i=1}^N \{y_i \log z(\mathbf{x}) + (1 - y_i) \log(1 - z(\mathbf{x}))\} \right]$$

# これまでの復習

## $K$ クラス識別

入力,  $\mathbf{x}$ , と出力,  $z_k \in (0, 1)$  の関係 ( $k = 1, 2, \dots, K$ ) :

$$z_k(\mathbf{x}) = p(y_k = 1|\mathbf{x}) = \text{SM}_k(a) = \frac{\exp(a_k)}{\sum_{j=1}^K \exp(a_j)},$$

$$a_k = \mathbf{w}_k^\top \boldsymbol{\phi} + w_{k0}$$

学習データ  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i) | i = 1, 2, \dots, N, y_{ik} \in \{0, 1\}, \sum_k y_{ik} = 1\}$   
最尤法による係数  $\mathbf{w}$  の求める:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left[ - \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log \{p(y_{ik} = 1|\mathbf{x})\} \right]$$

$$= \arg \min_{\mathbf{w}} \left[ - \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log \frac{\exp(a_{ik})}{\sum_{j=1}^K \exp(a_{ij})} \right]$$



# ニューラルネットワークによる多クラス識別

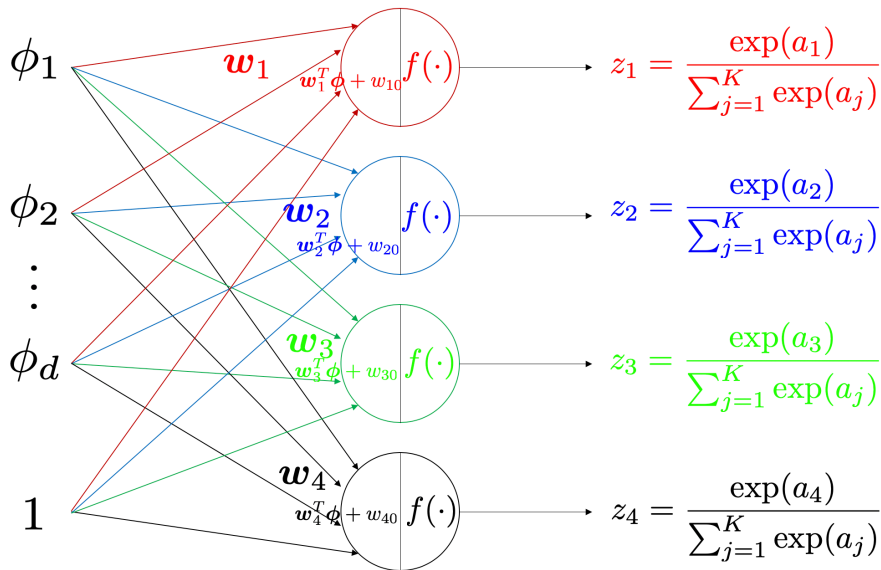
入力,  $\mathbf{x}$ , と出力,  $z_k \in (0, 1)$  の関係 ( $k = 1, 2, \dots, K$ ):

$$z_k(x) = p(y_k = 1 | \mathbf{x}) = \text{SM}_k(a) = \frac{\exp(a_k)}{\sum_{j=1}^K \exp(a_j)}, \quad \text{ただし } a_k = \mathbf{w}_k^\top \boldsymbol{\phi} + w_{k0}$$

最尤法による係数  $\mathbf{w}$  の求める:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left[ - \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log \frac{\exp(a_{ik})}{\sum_{j=1}^K \exp(a_{ij})} \right]$$

# ニューラルネットワークによる多クラス識別



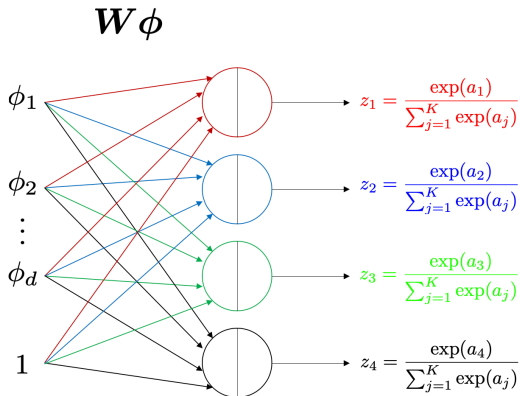
# 誤差の測り方（復習）

基本的に「負の対数尤度」で測る

問題	出力層の活性化関数	典型的な誤差関数
回帰	恒等写像	二乗誤差
2クラス識別	ロジスティック関数	交差エントロピー
多クラス識別	ソフトマックス関数	交差エントロピー

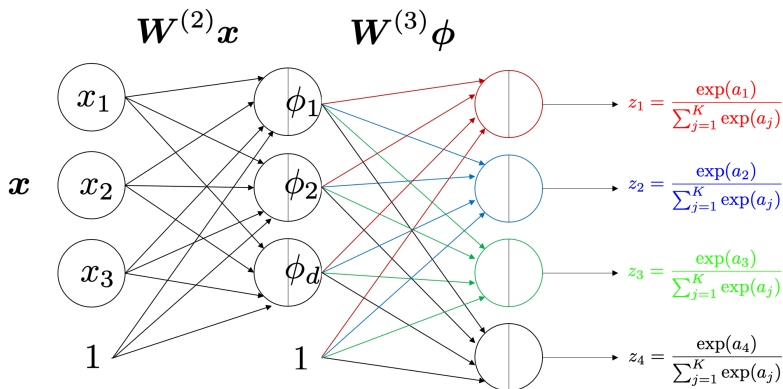
# 特徴 $\phi$ の作り方

入力  $\mathbf{x}$  を変換するために層を増やす



# 特徴 $\phi$ の作り方

入力  $x$  を変換するために層を増やす



# 特徴 $\phi$ の作り方

層を増やした後の関数：

$$\mathbf{z} = \text{SM}(\mathbf{W}^{(3)}\phi(\mathbf{x}))$$

クラス  $k$  に属する事後確率

$$\begin{aligned} z_k = p(y_k = 1|\mathbf{x}) &= \text{SM}_k \left( \mathbf{w}_k^{(3)\top} \phi(\mathbf{x}) \right) = \text{SM}_k \left( \sum_i w_{ki}^{(3)} \phi_i(\mathbf{x}) + w_{k0} \right) \\ &= \text{SM}_k \left( \sum_i w_{ki}^{(3)} \left( \sum_j f(\mathbf{w}_{ij}^{(2)\top} \mathbf{x} + w_{i0}) \right) + w_{k0} \right) \end{aligned}$$

# 特徴 $\phi$ の作り方

## ニューラルネットワークの長所

最尤法による係数  $\mathbf{w}$  の推定により適切な特徴 (写像)  $\phi$  を自動で獲得する

例:  $K$  クラス識別

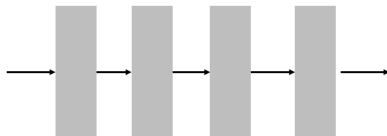
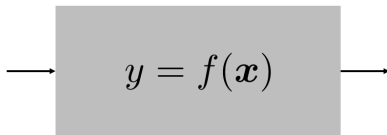
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left[ - \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log \frac{\exp(a_{ik})}{\sum_{j=1}^K \exp(a_{ij})} \right], \quad (a_{ik} = \mathbf{w}_k^\top \phi(\mathbf{x}_i) + w_{k0})$$

特徴  $\phi(\mathbf{x}_i)$  は、前項に沿って書くと、次式のとおり

$$\phi(\mathbf{x}_i) = f(\mathbf{W}^{(2)} \mathbf{x}_i) = \begin{bmatrix} \phi_1(\mathbf{x}_i) \\ \phi_2(\mathbf{x}_i) \\ \phi_3(\mathbf{x}_i) \end{bmatrix} = \begin{bmatrix} f(\mathbf{w}_1^{(2)\top} \mathbf{x}_i + w_{10}) \\ f(\mathbf{w}_2^{(2)\top} \mathbf{x}_i + w_{20}) \\ f(\mathbf{w}_3^{(2)\top} \mathbf{x}_i + w_{30}) \end{bmatrix}$$

# 目的は適切な識別関数の実現

複数の素朴な関数の荷重和と合成で複雑な関数表現



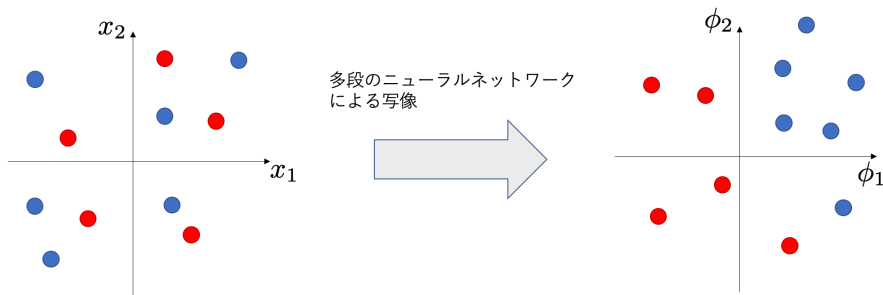
$$y = f(\mathbf{x})$$

$$y = f(\mathbf{x}) = f_L \circ f_{L-1} \circ \dots \circ f_1 \circ \mathbf{x}$$

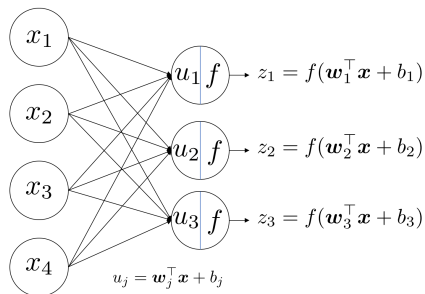


# 複雑な関数による適切な写像

パターン認識であれば線形分離性を高める特徴を得たい



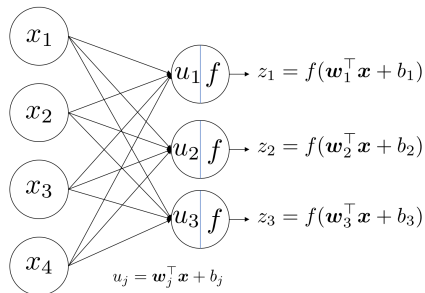




- 入力:  $\mathbf{x} = (x_1, x_2, x_3, x_4)^T$
- 出力:  $\mathbf{z} = (z_1, z_2, z_3)^T$
- $z_j = f(u_j), u_j = \mathbf{w}_j^T \mathbf{x} + b_j$

ただし、

$$\mathbf{w}_j = (w_{j1}, w_{j2}, \dots, w_{jI})^T$$


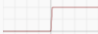









活性化関数も込みで、下記のように表す

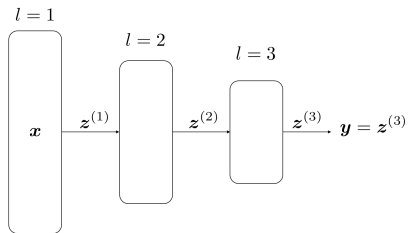
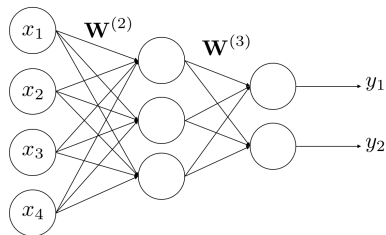
$$\mathbf{z} = \mathbf{f}(\mathbf{u})$$

ただし  $\mathbf{u} = \mathbf{W}\mathbf{x} + \mathbf{b}$

# 活性化関数

Name	Plot	Equation	Derivative
Identity		$f(x) = x$	$f'(x) = 1$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
Tanh		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$
Rectified Linear Unit (ReLU)		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Parametric Rectified Linear Unit (PReLU) [2]		$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Exponential Linear Unit (ELU) [3]		$f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} f(x) + \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
SoftPlus		$f(x) = \log_e(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$

# 層を増やす



# 層を増やす

$$\mathbf{z}^{(1)} = \mathbf{x}$$

$$\mathbf{u}^{(2)} = \mathbf{W}^{(2)} \mathbf{z}^{(1)} + \mathbf{b}^{(2)}$$

$$\mathbf{z}^{(2)} = \mathbf{f}(\mathbf{u}^{(2)})$$

前ページの図では、 $\mathbf{W}^{(2)}$  は  $3 \times 4$  の行列。

$$\mathbf{u}^{(3)} = \mathbf{W}^{(3)} \mathbf{z}^{(2)} + \mathbf{b}^{(3)}$$

$$\mathbf{z}^{(3)} = \mathbf{f}(\mathbf{u}^{(3)})$$

前ページの図では、 $\mathbf{W}^{(3)}$  は  $2 \times 3$  の行列。

多層ニューラルネットワークの  $(l+1)$  番目の層は、直前の  $l$  層の出力  $\mathbf{z}^{(l)}$  を受け取り次式を計算する

$$\mathbf{u}^{(l+1)} = \mathbf{W}^{(l+1)} \mathbf{z}^{(l)} + \mathbf{b}^{(l+1)}$$

$$\mathbf{z}^{(l+1)} = \mathbf{f}(\mathbf{u}^{(l+1)})$$

以下、バイアス項  $b$  も重みに組み込んで記述を見やすくする

- $\boldsymbol{x} \leftarrow (\boldsymbol{x}^\top, 1)$
- $\boldsymbol{w} \leftarrow (\boldsymbol{w}^\top, b)$
- $\boldsymbol{w}^\top \boldsymbol{x} \leftarrow \boldsymbol{w}^\top \boldsymbol{x} + b$



# ニューラルネットワークの学習

## 目的

学習データの集合,  $\mathcal{X} = \{(\mathbf{x}_1, \mathbf{d}_1), (\mathbf{x}_2, \mathbf{d}_2), \dots, (\mathbf{x}_N, \mathbf{d}_N)\}$  を用いて、ニューラルネットワークが望ましい関数を表すように、各層の重み係数  $\mathbf{W}^{(2)}, \mathbf{W}^{(3)}, \dots$  を求めること

最後の層  $L$  からの出力  $\mathbf{y} = \mathbf{z}^{(L)}$  は入力  $\mathbf{x}$  と重み係数  $\mathbf{W}^{(2)}, \mathbf{W}^{(3)}, \dots, \mathbf{W}^{(L)}$  の関数：

$$\mathbf{y} = \mathbf{y}(\mathbf{x}; \mathbf{w}) = \mathbf{y}(\mathbf{x}; \mathbf{W}^{(2)}, \mathbf{W}^{(3)}, \dots, \mathbf{W}^{(L)})$$

(すべての重み係数を  $\mathbf{w}$  で略記:  $\mathbf{w} = \{\mathbf{W}^{(2)}, \mathbf{W}^{(3)}, \dots, \mathbf{W}^{(L)}\}$ )

## 方針

データ  $\mathbf{x}_n$  を入力したときの出力,  $\mathbf{y}(\mathbf{x}_n; \mathbf{w})$ , と  $\mathbf{d}_n$  の誤差を  $L_n(\mathbf{w})$  で表す。全データに対する誤差の総和を小さくするように重み係数  $\mathbf{w}$  を求める

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_n L_n(\mathbf{w})$$

誤差の測り方は問題に応じて変える

# 係数の求め方

コスト関数の勾配

$$\nabla L(\mathbf{w}) = \frac{\partial L}{\partial \mathbf{w}} = \left[ \frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_M} \right]^\top$$

勾配降下法によるコストの更新

- ①  $\mathbf{w} = \mathbf{w}^0$  で初期化する
- ②  $\mathbf{w}^{t+1} = \mathbf{w}^t - \epsilon \nabla L$  による更新を繰り返す。

鞍点でよく停まる

# コスト関数の鞍点

学習データの一部のみで勾配を求めて係数  $w$  を更新する

## Stochastic Gradient Descent: SGD

- ①  $N$  個の学習データの内の 1 つ,  $x_{n^*}$ , をランダムに選択
- ②  $w^{t+1} = w^t - \epsilon \nabla L_{n^*}$  により係数を更新。

鞍点を回避しやすい

## 注：学習率

$$\boldsymbol{w}^{t+1} = \boldsymbol{w}^t - \epsilon \nabla L_{n^*}$$

係数  $\epsilon$  を学習率と呼ぶ。

- 小さすぎると収束に時間がかかる
- 大きすぎると局所解を通り過ぎる

実験により適切な値を定める必要がある

## 最適化法について

Neural Network の学習で得られる解は（ほぼ全て）局所解

- 局所解：近傍でコスト最小
  - 複数の局所解が存在
  - 初期値に依存して得られる解は変化
- 大域解：定義域全体の中でコスト最小



# 勾配降下法

- 確率的勾配降下法 (SGD): 訓練データをランダムに一つずつ利用
- ミニバッチ勾配降下法: 訓練データの一部を複数個利用
- 勾配降下法: 訓練データの全てを毎回利用

# ミニバッチとエポック

- ミニバッチ：訓練データを複数の「ミニバッチ」に分割
  - 訓練データをランダムに並べ替えて  $M$  分割 ( $B_1, B_2, \dots, B_M$ )
- エポック：全てのデータを1回ずつ更新に使うこと
  - 全データを  $M$  個のバッチに分割して、全てのバッチ  $B_1, B_2, \dots, B_M$  を1回ずつ更新に使うと1エポック

各種勾配降下法：更新時のパラメータの値における勾配のみ利用

- 収束に時間がかかることが多い

モメンタムを利用する更新法：過去の履歴を参照

- momentum = 「慣性」
  - 過去の更新傾向を今回の更新に反映
  - 振動成分を無視して「傾向」を求める

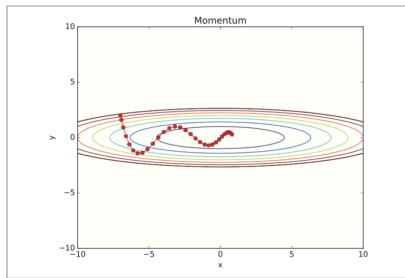
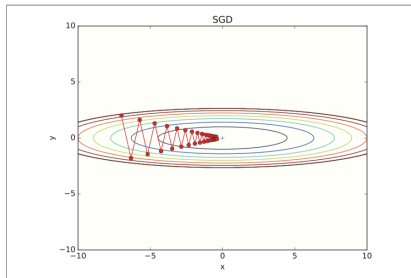
モメンタム項を加えた SGD の更新式

$$\boldsymbol{w}^{t+1} = \boldsymbol{w}^t - \epsilon \nabla L_{n^*} + \mu \boldsymbol{v}^t$$

右辺第三項がモメンタム項（慣性項）

$$\boldsymbol{v}^t = \boldsymbol{w}^t - \boldsymbol{w}^{t-1}$$

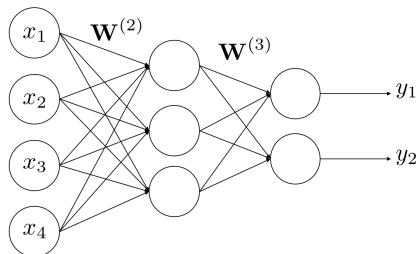
# SGD vs Momentum SGD



<https://www.dragonarrow.work/articles/195>

## 誤差逆伝搬

# 誤差逆伝播法（準備）



2層ネットワークを題材に

入力は  $\mathbf{x} = (x_1, x_2, x_3, x_4)^\top$ 。  $z_j^{(1)} = x_j$  とする。

次の層の出力は次のとおり

$$z_j^{(2)} = f(u_j^{(2)}) = f\left(\sum_i w_{ji}^{(2)} z_i^{(1)}\right)$$

# 誤差逆伝播法（準備）

出力層は次のとおり。回帰を想定して活性化関数は恒等写像

$$y_j(\mathbf{x}) = z_j^{(3)} = u_j^{(3)} = \sum_i w_{ji}^{(3)} z_i^{(2)}$$

誤差関数は二乗誤差とする

$$L_n = \frac{1}{2} \|\mathbf{y}(\mathbf{x}_n) - \mathbf{d}_n\|^2 = \frac{1}{2} \sum_j (y_j(\mathbf{x}_n) - d_{nj})^2$$



まず出力層の重みで微分

$$\frac{\partial L_n}{\partial w_{ji}^{(3)}} = (\mathbf{y}(\mathbf{x}_n) - \mathbf{d}_n)^\top \frac{\partial \mathbf{y}}{\partial w_{ji}^{(3)}}$$

$\partial \mathbf{y} / \partial w_{ji}^{(3)}$  は第  $j$  成分のみ  $z_i^{(2)}$ 。

$$\frac{\partial \mathbf{y}}{\partial w_{ji}^{(3)}} = [0, \dots, 0, z_i^{(2)}, 0, \dots, 0]$$

よって

$$\frac{\partial L_n}{\partial w_{ji}^{(3)}} = (y_j(\mathbf{x}_n) - d_j) z_i^{(2)}$$

# 誤差逆伝播法

中間層の重みによる微分  $\partial L_n / \partial w_{ji}^{(2)}$  の計算。重み  $w_{ji}^{(2)}$  は  $u_j^{(2)} = \sum_i w_{ji}^{(2)} z_i^{(1)}$  の中でのみ現れる。

$$\frac{\partial L_n}{\partial w_{ji}^{(2)}} = \frac{\partial L_n}{\partial u_j^{(2)}} \frac{\partial u_j^{(2)}}{\partial w_{ji}^{(2)}}$$

右辺第二項は  $u_j^{(2)} = \sum_i w_{ji}^{(2)} z_i^{(1)}$  より次の通り

$$\frac{\partial u_j^{(2)}}{\partial w_{ji}^{(2)}} = z_i^{(1)}$$

右辺第一項  $\partial L_n / \partial u_j^{(2)}$  は、 $u_j^{(2)}$  の変化による  $L_n$  の変化を表現。 $u_j^{(2)}$  の変化は活性化関数を介して出力層の  $u_k^{(3)}$  を変化させる。 $L_n$  は  $u_k^{(3)}$  の変化に伴い変化する。

$$\frac{\partial L_n}{\partial u_j^{(2)}} = \sum_k \frac{\partial L_n}{\partial u_k^{(3)}} \frac{\partial u_k^{(3)}}{\partial u_j^{(2)}}$$

# 誤差逆伝播法

$$L_n = 1/2 \sum_k (y_k(\mathbf{x}_n) - d_k)^2 = 1/2 \sum_k (u_k^{(3)} - d_k)^2$$

$$\frac{\partial L_n}{\partial u_k^{(3)}} = u_k^{(3)} - d_k$$

$$u_k^{(3)} = \sum_i w_{ki}^{(3)} f(u_i^{(2)}) \text{ より}$$

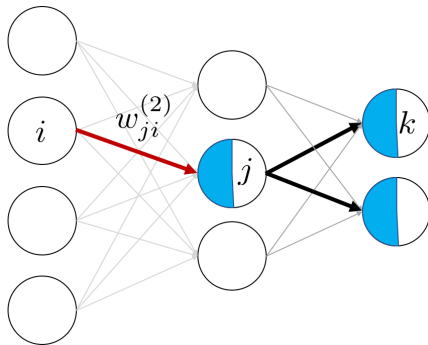
$$\frac{\partial u_k^{(3)}}{\partial u_j^{(2)}} = w_{kj}^{(3)} f'(u_j^{(2)})$$

以上より

$$\frac{\partial L_n}{\partial w_{ji}^{(2)}} = \left( f'(u_j^{(2)}) \sum_k w_{kj}^{(3)} (u_k^{(3)} - d_k) \right) z_i^{(1)}$$

# 誤差逆伝播法

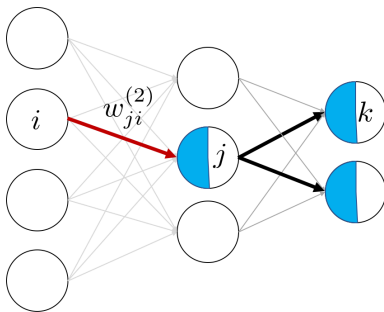
$$\frac{\partial L_n}{\partial w_{ji}^{(2)}} = \left( f'(u_j^{(2)}) \sum_k w_{kj}^{(3)} (u_k^{(3)} - d_k) \right) z_i^{(1)}$$



# 誤差逆伝播法

より一般的に

$$\frac{\partial L_n}{\partial w_{ji}^{(l)}} = \frac{\partial L_n}{\partial u_j^{(l)}} \frac{\partial u_j^{(l)}}{\partial w_{ji}^{(l)}}$$



右辺第一項は  $u_j^{(l)}$  の変化により  $L_n$  に生じる変化。 $u_j^{(l)}$  の変化は、次の  $(l+1)$  層の  $u_k^{(l+1)}$  を介して出力を変化させて、コスト関数に影響を与える。

$$\frac{\partial L_n}{\partial u_j^{(l)}} = \sum_k \frac{\partial L_n}{\partial u_k^{(l+1)}} \frac{\partial u_k^{(j+1)}}{\partial u_j^{(l)}}$$

ここで、第  $l$  層の第  $j$  番目のユニットについて、次式を定義する

$$\delta_j^{(l)} = \frac{\partial L_n}{\partial u_j^{(l)}}$$

$$u_k^{(l+1)} = \sum_i w_{ki}^{(l+1)} f(u_i^{(l)}) \text{ より}$$

$$\delta_j^{(l)} = \sum_k \delta_k^{(l+1)} \left( w_{kj}^{(l+1)} f'(u_j^{(l)}) \right)$$

$\delta_j^{(l)}$  は、ひとつ上の層の  $\delta_k^{(l+1)}$  から計算できる。

出力層から入力層に向かって  $\delta_j^{(l)}$  の計算を「伝播」させる。

$\partial u_j^{(l)} / \partial w_{ji}^{(l)} = z_i^{(l-1)}$  より結局次式を得る

$$\frac{\partial L_n}{\partial w_{ji}^{(l)}} = \delta_j^{(l)} z_i^{(l-1)}$$

各計算にネットワーク全体を見る必要はない。局所的に計算可能。出力層での微分は採用したコスト関数に依存して変わる。

$$\delta_j^{(L)} = \frac{\partial L_n}{\partial u_j^{(L)}}$$

## 誤差逆伝播による勾配の計算方法

- 入力：学習データ  $(\mathbf{x}_n, \mathbf{d}_n)$
- 出力：コスト関数  $L_n(\mathbf{w})$  の、各重み係数による微分,  $\partial L_n / \partial w_{ji}^{(l)}$
- ①  $\mathbf{z}^{(1)} = \mathbf{x}_n$ 。各層への入力  $\mathbf{u}^{(l)}$  と出力  $\mathbf{z}^{(l)}$  を順に計算する
- ② 出力層での  $\delta_j^{(L)}$  を計算する。普通は  $\delta_j^{(L)} = z_j - d_j$
- ③ 中間層での  $\delta_j^{(l)}$  を次式に従って逆向きに計算する ( $l = L, L-1, \dots$ )

$$\delta_j^{(l)} = \sum_k \delta^{(l+1)} \left( w_{kj}^{(l+1)} f'(u_j^{(l)}) \right)$$

- ④ 各層  $(l)$  のパラメータによる微分を次式で計算する

$$\frac{\partial L_n}{\partial w_{ji}^{(l)}} = \delta_j^{(l)} z_i^{(l-1)}$$