

知能プログラミング演習 I

第3回: 深層ニューラルネットワーク

梅津 佑太

2号館 404A: umezu.yuta@nitech.ac.jp

前回作ったディレクトリに移動して今日の課題のダウンロードと解凍

step1: `cd ./DLL`

step2: `wget http://www-als.ics.nitech.ac.jp/~umezu/DLL19/Lec3.zip`

step3: `unzip Lec3.zip`

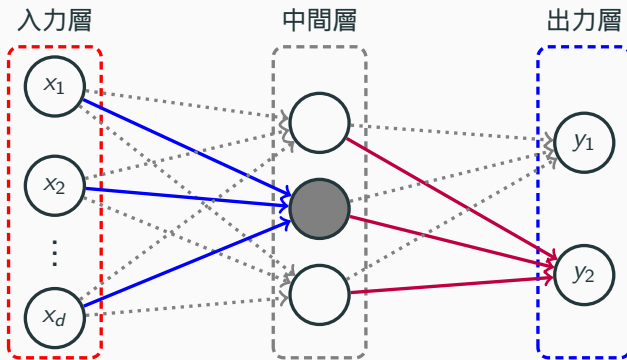
✓ まだ DLL のフォルダを作っていない人は, step1 の前に

`mkdir -p DLL`

でフォルダを作成する

1. 深層ニューラルネットワークの学習

前回の復習: (3層) ニューラルネットワーク



- パーセプトロンの合成によって少し複雑なモデルを記述
 - ✓ 各層のユニットがパーセプトロンの出力を表現している
 - ✓ すべての矢線には, 学習すべき重みがかかっている
 - ✓ 各ユニットを表現するための活性化関数は同じである必要はない

前回の復習: 活性化関数

- 入力層から出力層への活性化関数

- ✓ ReLU¹: $f(x) = \max\{0, x\}$

- ✓ シグモイド関数: $f(x) = 1/(1 + e^{-x})$

- ✓ ハイパボリックタンジェント: $f(x) = \tanh x = (e^x - e^{-x})/(e^x + e^{-x})$

- 中間層から出力層への活性化関数

- ✓ 恒等写像: 回帰問題で用いられる

$$g(x) = x$$

- ✓ ソフトマックス関数: 多クラス分類問題で用いられるもので, 出力されたベクトルは, x の所属確率を表す

$$g(x) = \frac{1}{\sum_{k=1}^m e^{x_k}} (e^{x_1}, \dots, e^{x_m})^T$$

¹Rectified Linear Unit

前回の復習: 誤差関数

教師データ² y_1, \dots, y_m とモデルの出力 $g(\mathbf{v}_1^\top \mathbf{z}), \dots, g(\mathbf{v}_m^\top \mathbf{z})$ に対して³,

- ℓ_2 -誤差: 回帰問題 (g : 恒等写像)

$$E(W, V) = \sum_{j=1}^m (y_j - g(\mathbf{v}_j^\top \mathbf{z}))^2$$

- クロスエントロピー: 分類問題 (g : ソフトマックス関数)

$$E(W, V) = - \sum_{i=1}^m y_i \log g(\mathbf{v}_i^\top \mathbf{z})$$

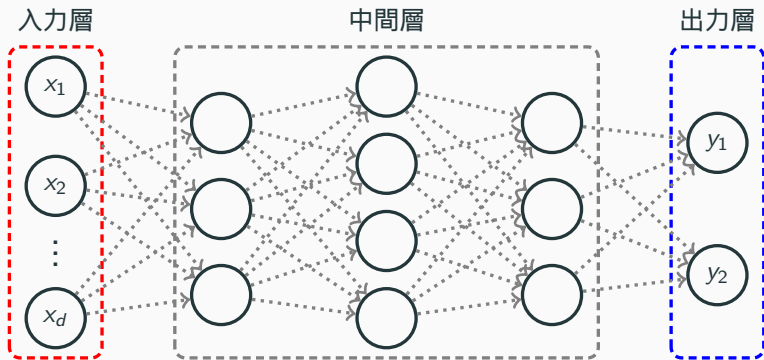
✓ 特に, $m = 2$ で $y \in \{0, 1\}$ の場合は

$$E(\mathbf{w}) = -y \log f(\mathbf{w}^\top \mathbf{x}) - (1 - y) \log(1 - f(\mathbf{w}^\top \mathbf{x}))$$

²正解のラベルや観測した実数値

³各 k に対して, $g(\mathbf{v}_k^\top \mathbf{z}) = g(v_{k0} + v_{k1}f(\mathbf{w}_1^\top \mathbf{x}) + \dots + v_{kq}f(\mathbf{w}_q^\top \mathbf{x}))$

ディープニューラルネットワーク



- 中間層を増やす (深くする) ことで、より複雑なモデルを表現
 - ✓ 深層学習では、入力層と出力層をつなぐ矢線上のすべての重みを学習
 - ✓ 層が深ければ深いほどパラメータ数は増加
 - 上のダイアグラムで $d = 1024$ とした場合、パラメータ数は 3102 個 orz

🕒 2016年10月18日 07時00分 公開

自動運転技術：

ディープニューラルネットワークで自動運転向け画像認識、デンソーと東芝が共同開発

デンソーと東芝は、自動運転や高度運転支援の画像認識システムに向けた人工知能技術を共同開発する。東芝は人工知能技術を搭載した画像認識プロセッサを2018年にサンプル出荷する目標だ。デンソーは2020年以降の実用化を目指す。

[齊藤由希, MONOist]


印刷/PDF


通知

見る

 ツイート

52

 シェア

52

 いいね！

0

 B! Bookmark

9

 Pocket



デンソーと東芝は2016年10月17日、自動運転や高度運転支援の画像認識システム向けの人工知能（AI）技術を共同開発すると発表した。両社で車載用プロセッサに実装可能なDNN（Deep Neural Network、ディープニューラルネットワーク）を開発し、GPUを搭載するシステムよりも低消費電力な画像処理を実現する。

<http://monoist.atmarkit.co.jp/mn/articles/1610/18/news021.html>

実社会における深層学習

翻訳

リアルタイム翻訳を無効にする

日本語 英語 韓国語 言語を検出する ▼

↔

英語 日本語 韓国語 ▼

翻訳

深層学習

4/5000

Shinsō gakushū

Depth learning

☆ □ 🔊 <

✎

日本語 英語 韓国語 言語を検出する ▼

↔

英語 日本語 韓国語 ▼

翻訳

機械学習は面白い

8/5000

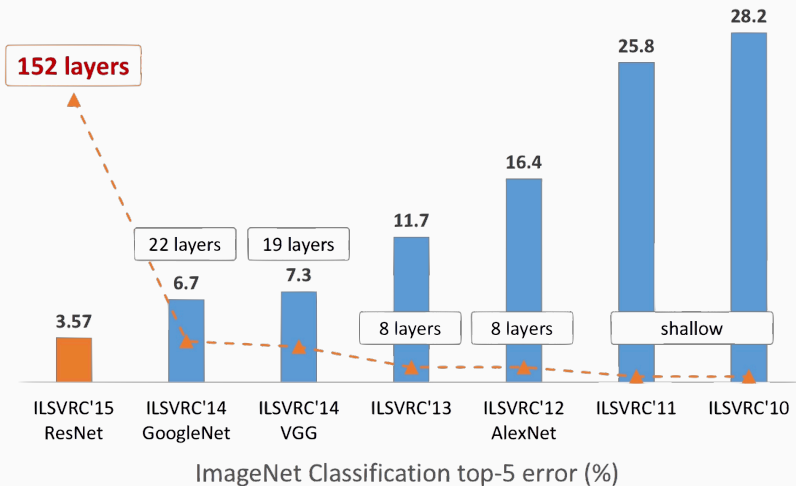
Kikai gakushū wa omoshiroi

Machine learning is interesting

☆ □ 🔊 <

✎

クリックしてその他の翻訳を表示、編集する



http://kaiminghe.com/ilsvrc15/ilsvrc2015_deep_residual_learning_kaiminghe.pdf

なぜ層を深くするのか？

- 3層ニューラルネットワークでも、中間層のユニット数 $q \rightarrow \infty$ で、
任意の関数を任意の精度で近似できる

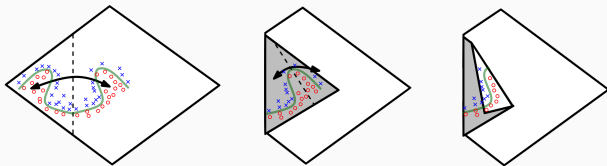
疑問: じゃあ、なんで深い方がいいの？

答え: 深さに対して 表現力が指数的に増大 するから

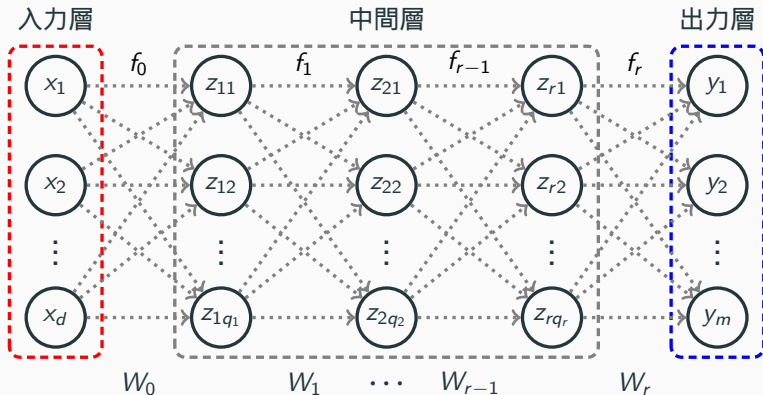
- ニューラルネットの表現力: 領域を何個の多面体に分割できるか
✓ $q (\geq d)$: 中間層の幅 (共通), r : 層の数. ReLU による領域の分割数は

$$\Omega\left(\left(\frac{q}{d}\right)^{r-1} \sum_{j=0}^q \binom{q}{j}\right) \approx \Omega\left(\left(\frac{q}{d}\right)^{(r-1)d} q^d\right)$$

- 対称性の高い関数は、多層にすることで得をする



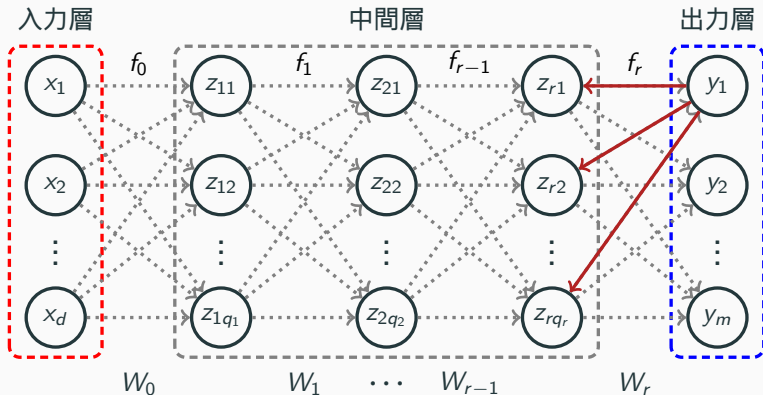
ディープニューラルネットワークにおける誤差逆伝播法



- 最適化すべきパラメータ: $q_0 = d, q_{r+1} = m$ として,

$$W_k = (\mathbf{w}_{k1}, \dots, \mathbf{w}_{kq_{k+1}})^T \in \mathbb{R}^{q_{k+1} \times (q_k + 1)}, \quad k = 0, 1, \dots, r$$

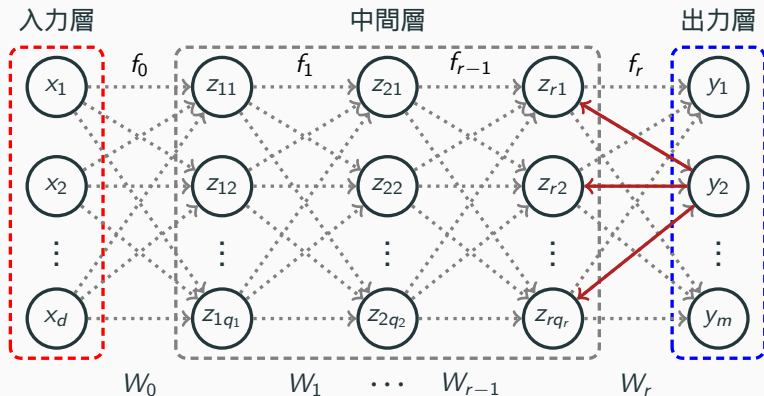
ディープニューラルネットワークにおける誤差逆伝播法



- 最適化すべきパラメータ: $q_0 = d, q_{r+1} = m$ として,

$$W_k = (\mathbf{w}_{k1}, \dots, \mathbf{w}_{kq_{k+1}})^\top \in \mathbb{R}^{q_{k+1} \times (q_k + 1)}, \quad k = 0, 1, \dots, r$$

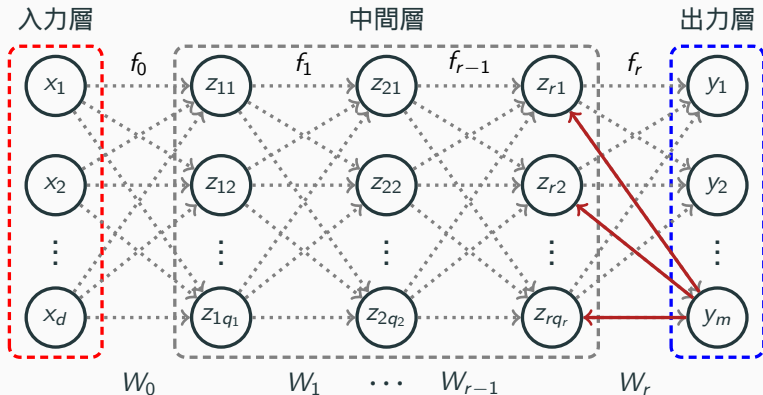
ディープニューラルネットワークにおける誤差逆伝播法



- 最適化すべきパラメータ: $q_0 = d, q_{r+1} = m$ として,

$$W_k = (\mathbf{w}_{k1}, \dots, \mathbf{w}_{kq_{k+1}})^\top \in \mathbb{R}^{q_{k+1} \times (q_k + 1)}, \quad k = 0, 1, \dots, r$$

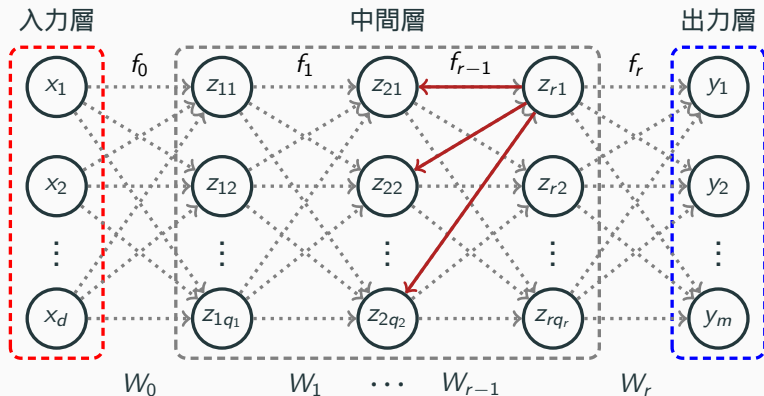
ディープニューラルネットワークにおける誤差逆伝播法



- 最適化すべきパラメータ: $q_0 = d, q_{r+1} = m$ として,

$$W_k = (\mathbf{w}_{k1}, \dots, \mathbf{w}_{kq_{k+1}})^\top \in \mathbb{R}^{q_{k+1} \times (q_k + 1)}, \quad k = 0, 1, \dots, r$$

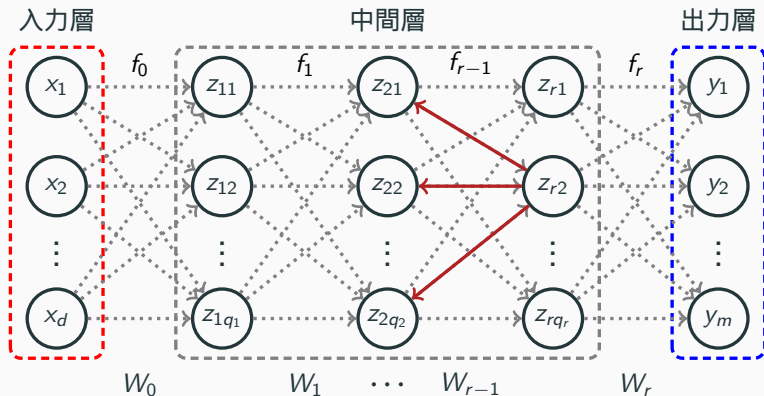
ディープニューラルネットワークにおける誤差逆伝播法



- 最適化すべきパラメータ: $q_0 = d, q_{r+1} = m$ として,

$$W_k = (\mathbf{w}_{k1}, \dots, \mathbf{w}_{kq_{k+1}})^\top \in \mathbb{R}^{q_{k+1} \times (q_k + 1)}, \quad k = 0, 1, \dots, r$$

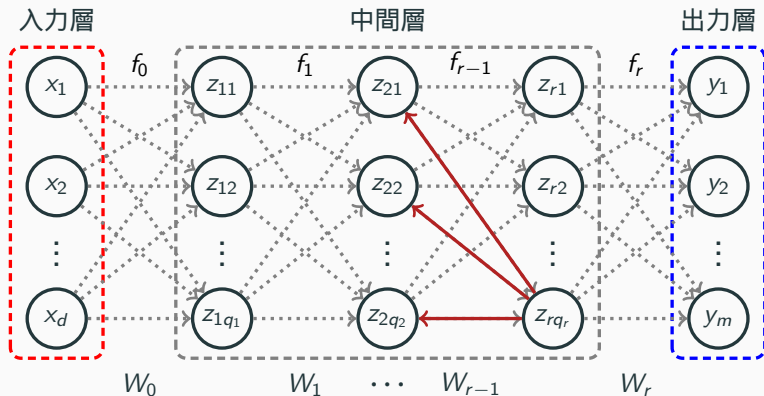
ディープニューラルネットワークにおける誤差逆伝播法



- 最適化すべきパラメータ: $q_0 = d, q_{r+1} = m$ として,

$$W_k = (\mathbf{w}_{k1}, \dots, \mathbf{w}_{kq_{k+1}})^\top \in \mathbb{R}^{q_{k+1} \times (q_k + 1)}, \quad k = 0, 1, \dots, r$$

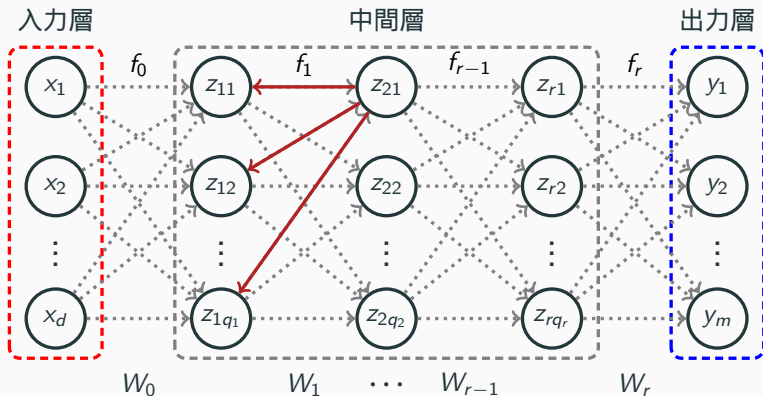
ディープニューラルネットワークにおける誤差逆伝播法



- 最適化すべきパラメータ: $q_0 = d, q_{r+1} = m$ として,

$$W_k = (\mathbf{w}_{k1}, \dots, \mathbf{w}_{kq_{k+1}})^T \in \mathbb{R}^{q_{k+1} \times (q_k + 1)}, \quad k = 0, 1, \dots, r$$

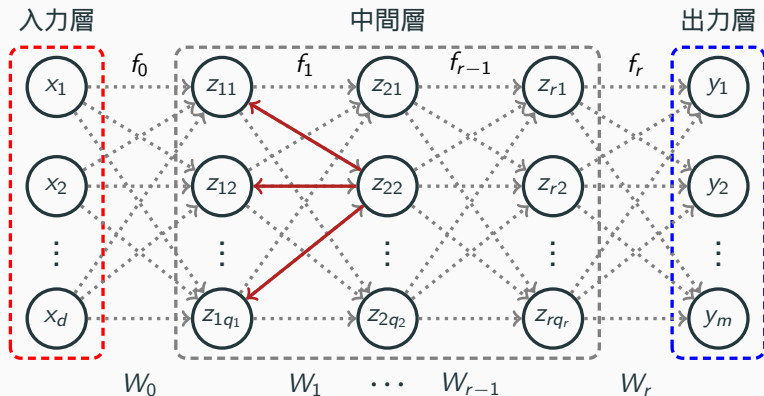
ディープニューラルネットワークにおける誤差逆伝播法



- 最適化すべきパラメータ: $q_0 = d, q_{r+1} = m$ として,

$$W_k = (\mathbf{w}_{k1}, \dots, \mathbf{w}_{kq_{k+1}})^\top \in \mathbb{R}^{q_{k+1} \times (q_k + 1)}, \quad k = 0, 1, \dots, r$$

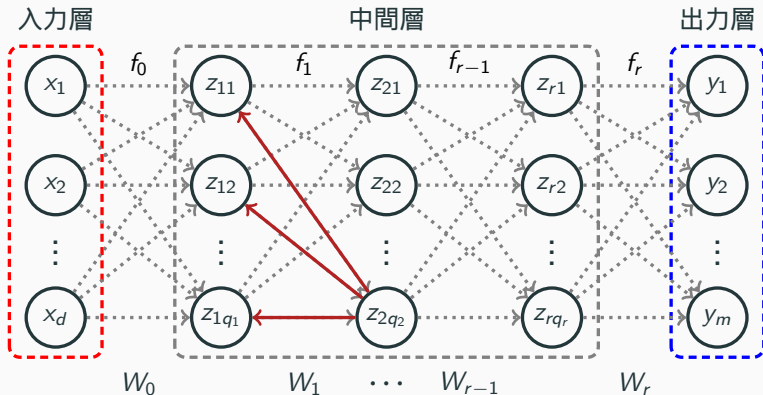
ディープニューラルネットワークにおける誤差逆伝播法



- 最適化すべきパラメータ: $q_0 = d, q_{r+1} = m$ として,

$$W_k = (\mathbf{w}_{k1}, \dots, \mathbf{w}_{kq_{k+1}})^\top \in \mathbb{R}^{q_{k+1} \times (q_k + 1)}, \quad k = 0, 1, \dots, r$$

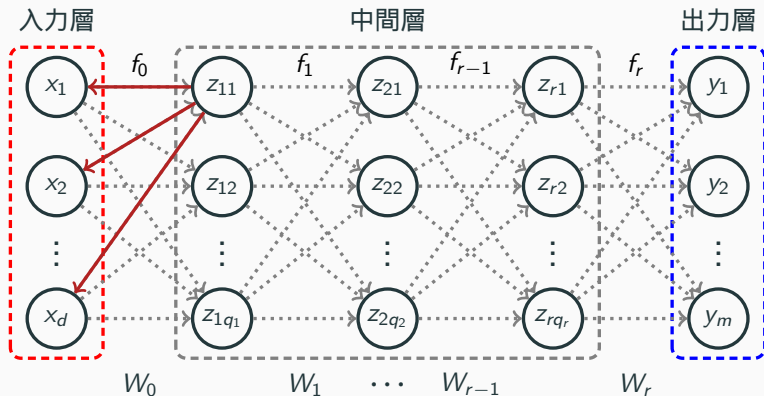
ディープニューラルネットワークにおける誤差逆伝播法



- 最適化すべきパラメータ: $q_0 = d, q_{r+1} = m$ として,

$$W_k = (\mathbf{w}_{k1}, \dots, \mathbf{w}_{kq_{k+1}})^\top \in \mathbb{R}^{q_{k+1} \times (q_k + 1)}, \quad k = 0, 1, \dots, r$$

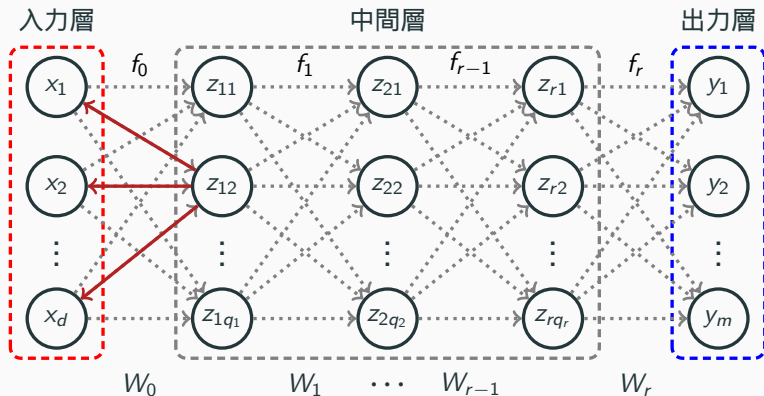
ディープニューラルネットワークにおける誤差逆伝播法



- 最適化すべきパラメータ: $q_0 = d, q_{r+1} = m$ として,

$$W_k = (\mathbf{w}_{k1}, \dots, \mathbf{w}_{kq_{k+1}})^\top \in \mathbb{R}^{q_{k+1} \times (q_k + 1)}, \quad k = 0, 1, \dots, r$$

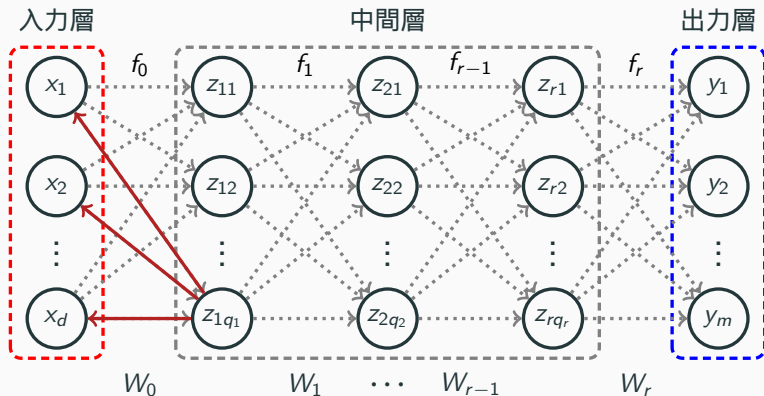
ディープニューラルネットワークにおける誤差逆伝播法



- 最適化すべきパラメータ: $q_0 = d, q_{r+1} = m$ として,

$$W_k = (\mathbf{w}_{k1}, \dots, \mathbf{w}_{kq_{k+1}})^\top \in \mathbb{R}^{q_{k+1} \times (q_k + 1)}, \quad k = 0, 1, \dots, r$$

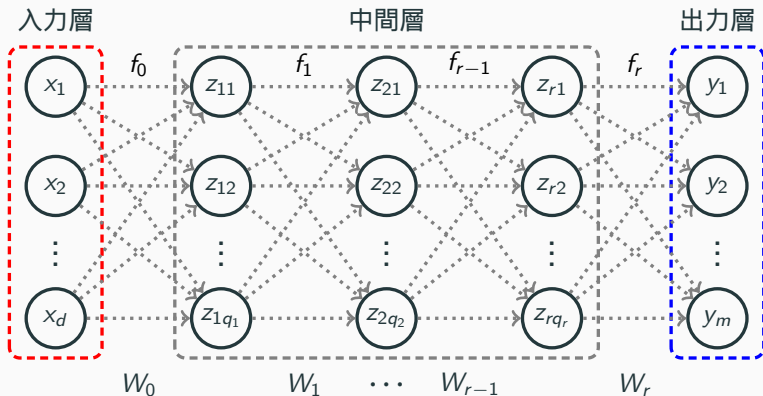
ディープニューラルネットワークにおける誤差逆伝播法



- 最適化すべきパラメータ: $q_0 = d, q_{r+1} = m$ として,

$$W_k = (\mathbf{w}_{k1}, \dots, \mathbf{w}_{kq_{k+1}})^{\top} \in \mathbb{R}^{q_{k+1} \times (q_k + 1)}, \quad k = 0, 1, \dots, r$$

誤差逆伝播法



- モデルの関数表現: $k-1$ 層から k 層への活性化関数を f_k として,

$$\mathbf{z}_1 = f_0(W_0 \mathbf{x}) \left(= (f_0(\mathbf{w}_{01}^\top \mathbf{x}), \dots, f_0(\mathbf{w}_{0q_1}^\top \mathbf{x}))^\top \right),$$

$$\mathbf{z}_k = f_{k-1}(W_{k-1} \mathbf{z}_{k-1}), \quad k = 2, \dots, r,$$

$$\mathbf{y} = f_r(W_r \mathbf{z}_r)$$

誤差関数とその微分

切片項も考慮して, $(1, \mathbf{z}_j^\top)^\top$ を改めて \mathbf{z}_j とする.

- 回帰問題: 最後の間層から出力層への活性化関数を恒等写像 $f_r(W_r \mathbf{z}_r) = W_r \mathbf{z}_r$ とし, 誤差関数は最小二乗誤差を用いる:

$$E(W_0, \dots, W_r) = \frac{1}{2} \|\mathbf{y} - W_r \mathbf{z}_r\|^2 = \frac{1}{2} \sum_{j=1}^m (y_j - \mathbf{w}_{rj}^\top \mathbf{z}_r)^2$$

- 分類問題: 最後の間層から出力層への活性化関数をソフトマックス関数 $f_r(W_r \mathbf{z}_r) = e^{W_r \mathbf{z}_r} / (\sum_{i=1}^m e^{\mathbf{w}_{ri}^\top \mathbf{z}_r})$ とし, 誤差関数はクロスエントロピーを用いる:

$$E(W_0, \dots, W_r) = - \sum_{j=1}^m y_j \log f_r(\mathbf{w}_{rj}^\top \mathbf{z}_r)$$

誤差関数とその微分

切片項も考慮して, $(1, \mathbf{z}_j^\top)^\top$ を改めて \mathbf{z}_j とする.

- 回帰問題: 最後の間層から出力層への活性化関数を恒等写像 $f_r(W_r \mathbf{z}_r) = W_r \mathbf{z}_r$ とし, 誤差関数は最小二乗誤差を用いる:

$$E(W_0, \dots, W_r) = \frac{1}{2} \|\mathbf{y} - W_r \mathbf{z}_r\|^2 = \frac{1}{2} \sum_{j=1}^m (y_j - \mathbf{w}_{rj}^\top \mathbf{z}_r)^2$$

- 分類問題: 最後の間層から出力層への活性化関数をソフトマックス関数 $f_r(W_r \mathbf{z}_r) = e^{W_r \mathbf{z}_r} / (\sum_{i=1}^m e^{\mathbf{w}_{ri}^\top \mathbf{z}_r})$ とし, 誤差関数はクロスエントロピーを用いる:

$$E(W_0, \dots, W_r) = - \sum_{j=1}^m y_j \log f_r(\mathbf{w}_{rj}^\top \mathbf{z}_r)$$

$$\Rightarrow \frac{\partial E(W_0, \dots, W_r)}{\partial W_r} = (f_r(W_r \mathbf{z}_r) - \mathbf{y}) \mathbf{z}_r^\top$$

W_r に関する誤差関数の微分の導出 I

回帰問題の場合は簡単なので、分類問題だと思って計算する.

- 準備: ソフトマックス関数の第 j 成分を \mathbf{w}_{rk} で微分する.

$$f_r(\mathbf{w}_{rj}^\top \mathbf{z}_r) = \frac{1}{\sum_{i=1}^m e^{\mathbf{w}_{ri}^\top \mathbf{z}_r}} e^{\mathbf{w}_{rj}^\top \mathbf{z}_r}$$

より、分数関数の微分と、 $\frac{\partial e^{\mathbf{w}_{rj}^\top \mathbf{z}_r}}{\partial \mathbf{w}_{rk}} = e^{\mathbf{w}_{rk}^\top \mathbf{z}_r} \mathbf{z}_r$ ($j = k$) & 0 ($j \neq k$) などを用いて、

$$j = k \Rightarrow \frac{\partial f_r(\mathbf{w}_{rk}^\top \mathbf{z}_r)}{\partial \mathbf{w}_{rk}} = (1 - f_r(\mathbf{w}_{rk}^\top \mathbf{z}_r)) f_r(\mathbf{w}_{rk}^\top \mathbf{z}_r) \mathbf{z}_r$$

$$j \neq k \Rightarrow \frac{\partial f_r(\mathbf{w}_{rj}^\top \mathbf{z}_r)}{\partial \mathbf{w}_{rk}} = -f_r(\mathbf{w}_{rj}^\top \mathbf{z}_r) f_r(\mathbf{w}_{rk}^\top \mathbf{z}_r) \mathbf{z}_r$$

W_r に関する誤差関数の微分の導出 II

ソフトマックス関数の微分を用いて, 誤差関数 (クロスエントロピー) を \mathbf{w}_{rk} で微分すると,

$$\begin{aligned}\frac{\partial E(W_0, \dots, W_r)}{\partial \mathbf{w}_{rk}} &= - \sum_{j=1}^m y_j \frac{\partial \log f_r(\mathbf{w}_{rj}^\top \mathbf{z}_r)}{\partial \mathbf{w}_{rk}} \\ &= - \sum_{j=1}^m y_j \frac{1}{f_r(\mathbf{w}_{rj}^\top \mathbf{z}_r)} \frac{\partial f_r(\mathbf{w}_{rj}^\top \mathbf{z}_r)}{\partial \mathbf{w}_{rk}} = (f_r(\mathbf{w}_{rk}^\top \mathbf{z}_r) - y_k) \mathbf{z}_r\end{aligned}$$

となる⁴. $W_r = (\mathbf{w}_{r1}, \dots, \mathbf{w}_{rm})^\top$ であったから,

$$\begin{aligned}\frac{\partial E(W_0, \dots, W_r)}{\partial W_r} &= \left(\frac{\partial E(W_0, \dots, W_r)}{\partial \mathbf{w}_{r1}}, \dots, \frac{\partial E(W_0, \dots, W_r)}{\partial \mathbf{w}_{rm}} \right)^\top \\ &= (f_r(W_r \mathbf{z}_r) - \mathbf{y}) \mathbf{z}^\top\end{aligned}$$

⁴ \mathbf{y} は正解ラベルで 1, それ以外は 0 なので, \mathbf{y} の要素をすべて足すと 1 になることに注意!

W_{r-1} に関する誤差関数の微分 I

同様に, 誤差関数を $\mathbf{w}_{r-1,k}$ ($k = 1, \dots, q_r$) で微分すると,

$$\begin{aligned}\frac{\partial E(W_0, \dots, W_r)}{\partial \mathbf{w}_{r-1,k}} &= \sum_{j=1}^m \frac{\partial E(W_0, \dots, W_r)}{\partial \mathbf{w}_{rj}^\top \mathbf{z}_r} \frac{\partial \mathbf{w}_{rj}^\top \mathbf{z}_r}{\partial \mathbf{w}_{r-1,k}} \\ &= \sum_{j=1}^m (f_r(\mathbf{w}_{rj}^\top \mathbf{z}_r) - y_j) \frac{\partial \mathbf{w}_{rj}^\top (1, f_{r-1}(W_{r-1} \mathbf{z}_{r-1}))^\top}{\partial \mathbf{w}_{r-1,k}} \\ &= \sum_{j=1}^m (f_r(\mathbf{w}_{rj}^\top \mathbf{z}_r) - y_j) w_{rjk} \nabla f_{r-1}(\mathbf{w}_{r-1,k}^\top \mathbf{z}_{r-1}) \mathbf{z}_{r-1}\end{aligned}$$

ここで, $\tilde{\mathbf{w}}_{rk} = (w_{r1k}, \dots, w_{rmk})^\top$ は W_r の k 列目なので,

$$\frac{\partial E(W_0, \dots, W_r)}{\partial \mathbf{w}_{r-1,k}} = f_r(W_r \mathbf{z}_r - \mathbf{y})^\top \tilde{\mathbf{w}}_{rk} \nabla f_{r-1}(\mathbf{w}_{r-1,k}^\top \mathbf{z}_{r-1}) \mathbf{z}_{r-1}$$

W_{r-1} に関する誤差関数の微分 II

したがって,

$$\begin{aligned}\frac{\partial E(W_0, \dots, W_r)}{\partial W_{r-1}} &= \left(\frac{\partial E(W_0, \dots, W_r)}{\partial \mathbf{w}_{r-1,1}}, \dots, \frac{\partial E(W_0, \dots, W_r)}{\partial \mathbf{w}_{r-1,q_r}} \right)^\top \\ &= \begin{pmatrix} \tilde{\mathbf{w}}_{r1}^\top f_r(W_r \mathbf{z}_r - \mathbf{y}) \nabla f_{r-1}(\mathbf{w}_{r-1,1}^\top \mathbf{z}_{r-1}) \mathbf{z}_{r-1}^\top \\ \vdots \\ \tilde{\mathbf{w}}_{rq_r}^\top f_r(W_r \mathbf{z}_r - \mathbf{y}) \nabla f_{r-1}(\mathbf{w}_{r-1,q_r}^\top \mathbf{z}_{r-1}) \mathbf{z}_{r-1}^\top \end{pmatrix} \\ &= \left[\tilde{W}_r^\top (f_r(W_r \mathbf{z}_r) - \mathbf{y}) \odot \nabla f_{r-1}(W_{r-1} \mathbf{z}_{r-1}) \right] \mathbf{z}_{r-1}^\top\end{aligned}$$

ただし, \tilde{W}_r は W_r の第 1 列 (切片項に対応する部分) を除いた $q_{r+1} \times q_r$ 次元の行列.

より一般に, $u_{kl} = \mathbf{w}_{kl}^\top \mathbf{z}_k$ ($l = 1, \dots, q_{k+1}$) とすれば,

$$\begin{aligned} \frac{\partial E(W_0, \dots, W_r)}{\partial u_{kl}} &= \sum_{j=1}^{q_{k+2}} \frac{\partial E(W_0, \dots, W_r)}{\partial u_{k+1,j}} \frac{\partial u_{k+1,j}}{\partial u_{kl}} \\ &= \sum_{j=1}^{q_{k+2}} \frac{\partial E(W_0, \dots, W_r)}{\partial u_{k+1,j}} w_{k+1,jl} \nabla f_k(u_{kl}) \end{aligned}$$

なので, $\delta_{kl} = \partial E(W_0, \dots, W_r) / \partial u_{kl}$ として, 漸化式

$$\delta_{kl} = \sum_{j=1}^{q_{k+2}} \delta_{k+1,j} w_{k+1,jl} \nabla f_k(u_{kl})$$

が得られる. このとき, $\delta_k = (\delta_{k1}, \dots, \delta_{kq_{k+1}})^\top$ は, W_{k+1} の第 1 列を除いた $q_{k+2} \times q_{k+1}$ 次元の行列を用いて, 以下のように書ける:

$$\delta_k = \tilde{W}_{k+1}^\top \delta_{k+1} \odot \nabla f_k(W_k \mathbf{z}_k), \quad (k = 0, 1, \dots, r-1)$$

したがって、第 k 層において、パラメータに関する微分は、

$$\frac{\partial E(W_0, \dots, W_r)}{\partial \mathbf{w}_{kl}} = \frac{\partial E(W_0, \dots, W_r)}{\partial u_{kl}} \frac{\partial u_{kl}}{\partial \mathbf{w}_{kl}} = \delta_{kl} \mathbf{z}_k$$

より、

$$\frac{\partial E(W_0, \dots, W_r)}{\partial W_k} = \begin{pmatrix} \delta_{k1} \mathbf{z}_k^\top \\ \vdots \\ \delta_{kq_{k+1}} \mathbf{z}_k^\top \end{pmatrix} = \left\{ \tilde{W}_{k+1}^\top \boldsymbol{\delta}_{k+1} \odot \nabla f_k(W_k \mathbf{z}_k) \right\} \mathbf{z}_k^\top$$

したがって、パラメータの更新式は

$$W_k \leftarrow W_k - \eta \left\{ \tilde{W}_{k+1}^\top \boldsymbol{\delta}_{k+1} \odot \nabla f_k(W_k \mathbf{z}_k) \right\} \mathbf{z}_k^\top$$

エポック内で誤差逆伝播法のアルゴリズムは以下の通り. ただし, η は適当な学習率.

1. ランダムにデータを読み込む: $(\mathbf{y}, \mathbf{x}) \in \{(\mathbf{y}_i, \mathbf{x}_i) \mid i = 1, \dots, n\}$
2. 順伝播:
 - 2.1 $\mathbf{z}_1 \leftarrow f_0(W_0 \mathbf{x})$
 - 2.2 $\mathbf{z}_{k+1} \leftarrow f_k(W_k \mathbf{z}_k) \ (k = 1, \dots, r-1)$
3. 逆伝播:
 - 3.1 $\delta_r \leftarrow f_r(W_r \mathbf{z}_r) - \mathbf{y}$
 - 3.2 $\delta_k \leftarrow \tilde{W}_{k+1}^\top \delta_{k+1} \odot \nabla f_k(W_k \mathbf{z}_k) \ (k = 0, 1, \dots, r-1; \mathbf{z}_0 = \mathbf{x})$
4. 重みの更新: $W_k \leftarrow W_k - \eta \delta_k \mathbf{z}_k^\top$