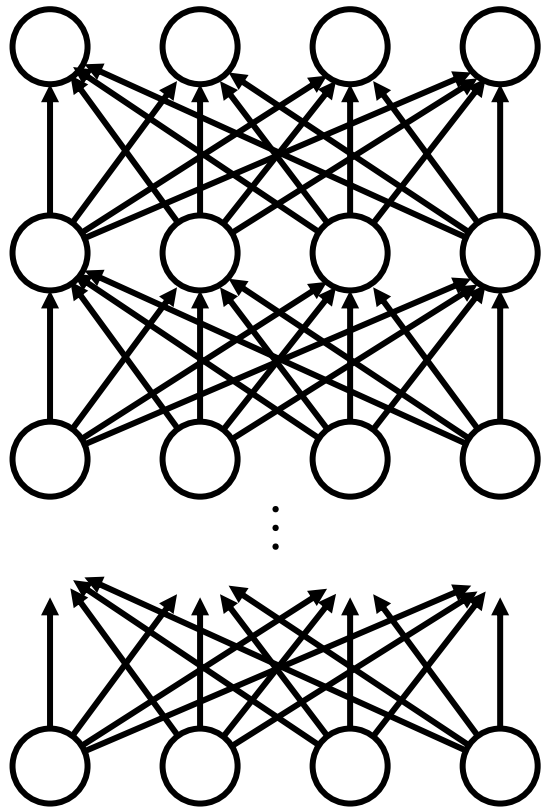


# 誤差逆伝搬法



# MLP のパラメタ推定法: Error back-propagation



$w_{ij}^1$   $o_i^0$  : Output layer

$w_{ij}^2$   $o_i^1$  : Hidden layers

$o_i^2$

$w_{ij}^N$   $o_i^N$  : Input layer

$$o_i^{k-1} = f(x_i^k)$$

$$x_i^k = \sum_j o_j^k \cdot w_{ji}^k = \mathbf{o}^k \cdot \mathbf{w}_i^k$$



$\mathbf{o}^N$  : 入力

$\mathbf{o}^0 = \mathbf{g}(\mathbf{o}^N; \mathbf{w})$  : 出力

$\mathbf{w}$  : モデルパラメタ

$\mathbf{d}$  : 真値

$$E = \frac{1}{2} \sum_i (o_i^0 - d_i)^2 = E(\mathbf{o}^N, \mathbf{d}; \mathbf{w}) \quad : \text{誤差関数}$$

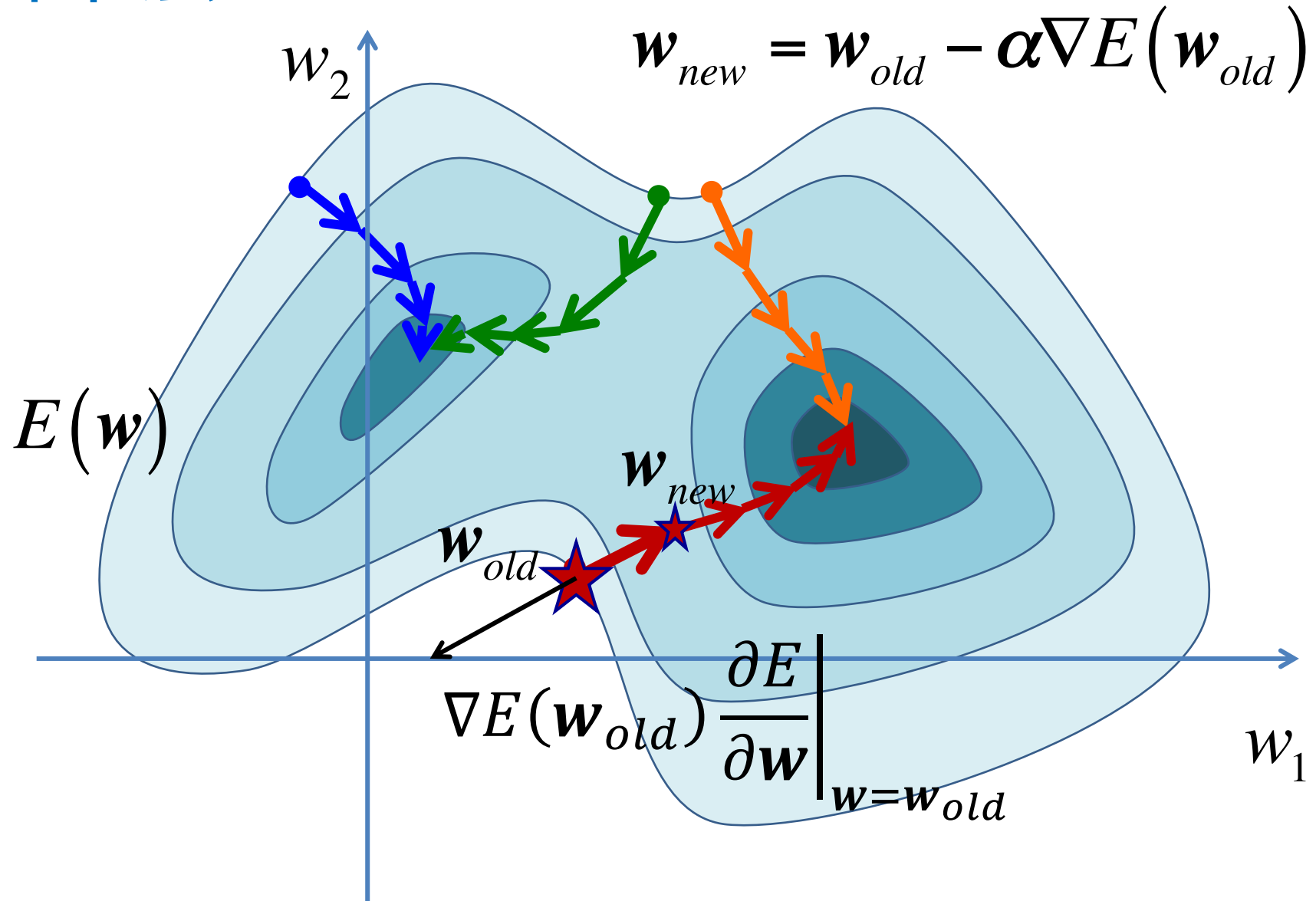
## Error back-propagation

= モデルパラメタの関数であるところの誤差関数を  
モデルパラメタで最小化する

$$w_{ij\_new}^k = w_{ij\_old}^k - \alpha \left. \frac{\partial E}{\partial w_{ij}^k} \right|_{w_{ij}^k = w_{ij\_old}^k}$$



# 最急降下法



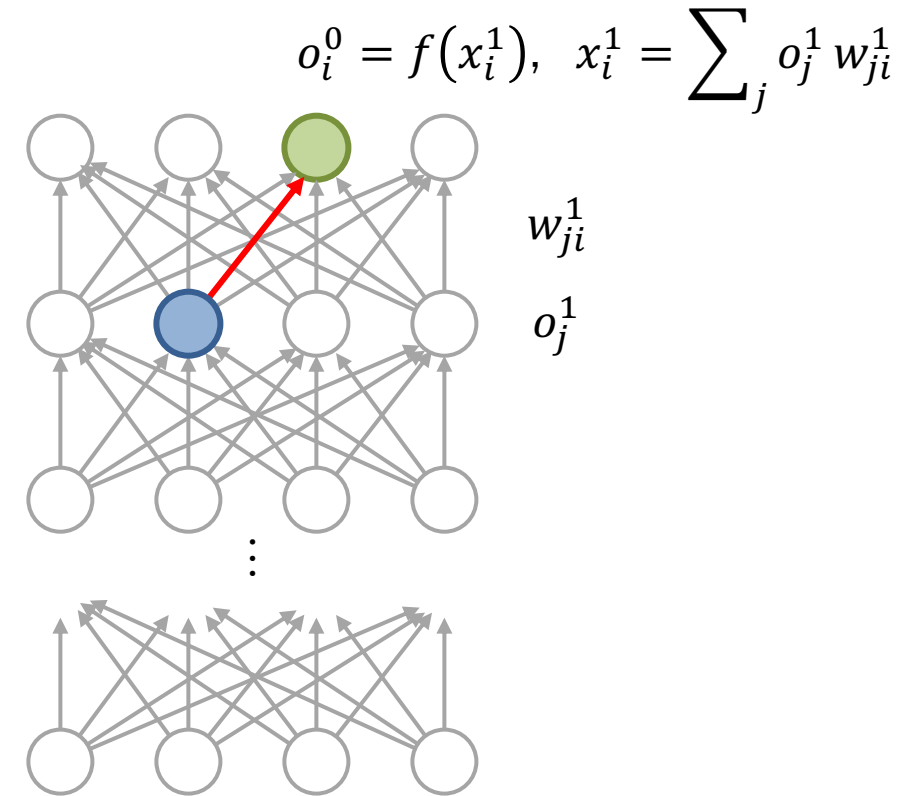
$$E = \frac{1}{2} \sum_i (o_i^0 - d_i)^2 = \frac{1}{2} \sum_i (f(x_i^1) - d_i)^2$$

$$\frac{\partial E}{\partial w_{ji}^1} = \frac{\partial E}{\partial x_i^1} \frac{\partial x_i^1}{\partial w_{ji}^1}$$

$$\frac{\partial E}{\partial x_i^1} = (o_i^0 - d_i) f'(x_i^1) \equiv \delta_i^1$$

$$\frac{\partial x_i^1}{\partial w_{ji}^1} = o_j^1$$

$$\frac{\partial E}{\partial w_{ji}^1} = \delta_i^1 o_j^1, \quad \hat{w}_{ji}^1 = w_{ji}^1 - \varepsilon \delta_i^1 o_j^1$$



$$E = \frac{1}{2} \sum_i (o_i^0 - d_i)^2 = \frac{1}{2} \sum_i (f(x_i^1) - d_i)^2$$

$$\frac{\partial E}{\partial w_{ji}^1} = \frac{\partial E}{\partial x_i^1} \frac{\partial x_i^1}{\partial w_{ji}^1}$$

$$\frac{\partial E}{\partial x_i^1} = (o_i^0 - d_i) f'(x_i^1) \equiv \delta_i^1$$

Eを減じるために  
どれだけ  $w_{ji}^1$  を変えるべきか

は

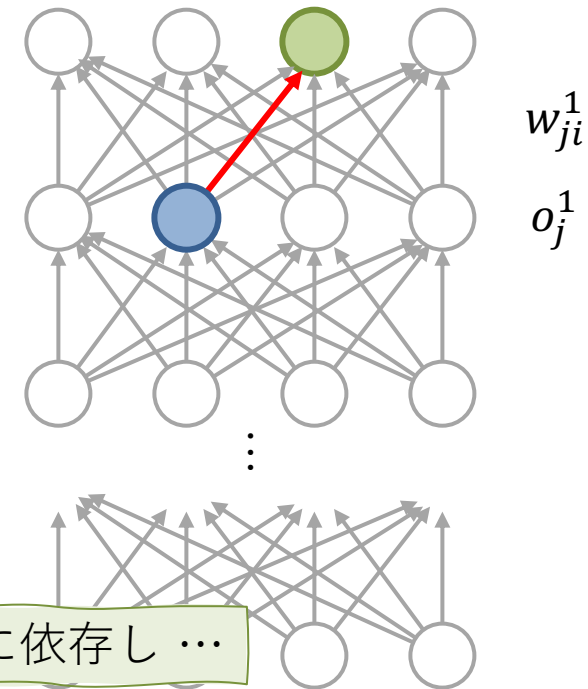
Eを減じるために  
どれだけ  $x_i^1$  を変えるべきか

に依存し …

$$\frac{\partial E}{\partial w_{ji}^1} = \delta_i^1 o_j^1,$$

$$\hat{w}_{ji}^1 = w_{ji}^1 - \varepsilon \delta_i^1 o_j^1$$

$$o_i^0 = f(x_i^1), \quad x_i^1 = \sum_j o_j^1 w_{ji}^1$$



$$E = \frac{1}{2} \sum_i (o_i^0 - d_i)^2 = \frac{1}{2} \sum_i (f(x_i^1) - d_i)^2$$

$$\frac{\partial E}{\partial w_{ji}^1} = \frac{\partial E}{\partial x_i^1} \frac{\partial x_i^1}{\partial w_{ji}^1}$$

$$\frac{\partial E}{\partial x_i^1} = (o_i^0 - d_i) f'(x_i^1) \equiv \delta_i^1$$

Eを減じるために  
どれだけ  $w_{ji}^1$  を変えるべきか

Eを減じるために  
どれだけ  $x_i^1$  を変えるべきか

Eを減じるために  
どれだけ  $x_i^1$  を変えるべきか

に依存して決まり...

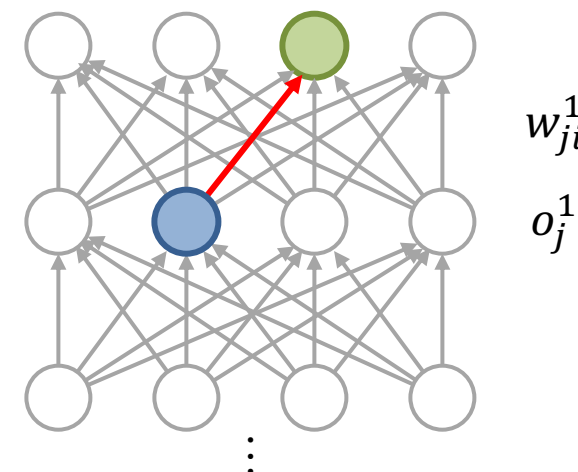
$$\frac{\partial E}{\partial w_{ji}^1} = \delta_i^1 o_j^1,$$

は

Error

に依存して決まる

$$o_i^0 = f(x_i^1), \quad x_i^1 = \sum_j o_j^1 w_{ji}^1$$



$$E = \frac{1}{2} \sum_i (o_i^0 - d_i)^2 = \frac{1}{2} \sum_i (f(x_i^1) - d_i)^2$$

$$\frac{\partial E}{\partial w_{ji}^1} = \frac{\partial E}{\partial x_i^1} \frac{\partial x_i^1}{\partial w_{ji}^1}$$

$$\frac{\partial E}{\partial x_i^1} = (o_i^0 - d_i) f'(x_i^1) \equiv \delta_i^1$$

$$\frac{\partial x_i^1}{\partial w_{ji}^1} = o_j^1$$

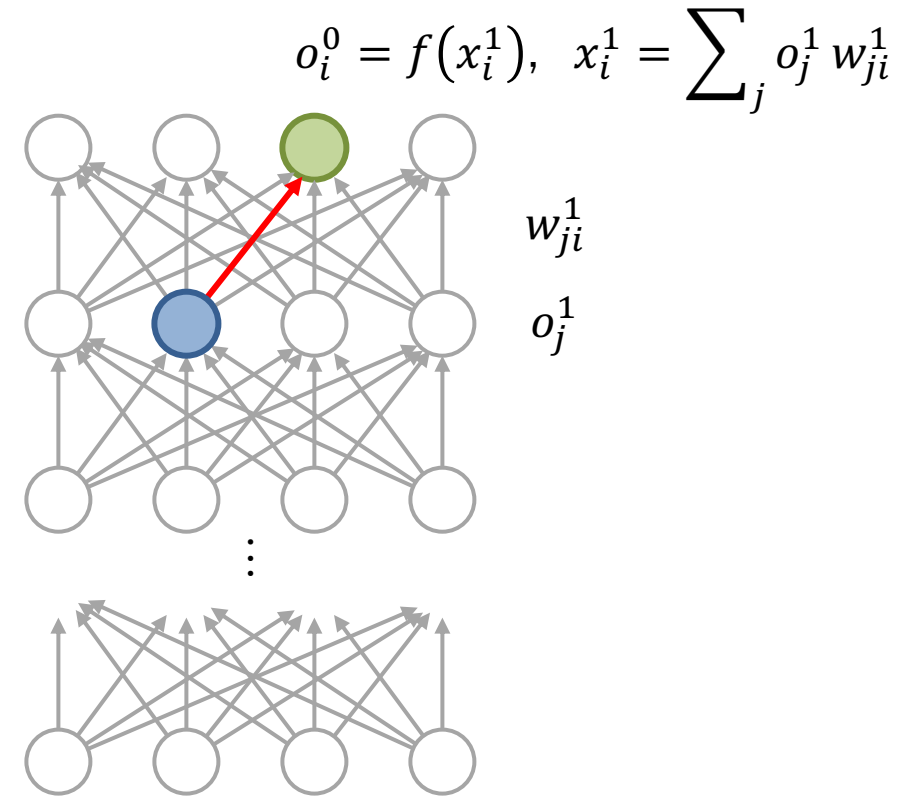
$$\frac{\partial E}{\partial w_{ji}^1} = \delta_i^1 o_j^1,$$

E を減じるために  
どれだけ  $x_i^1$  を変えるべきか

を

$\delta_i^1$

とおく



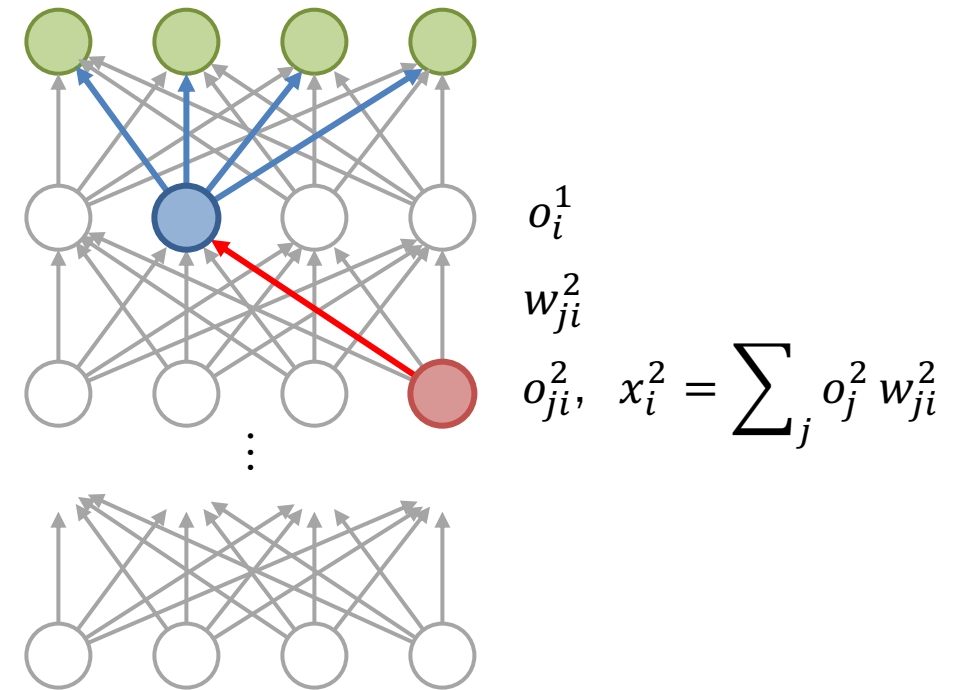


$$\frac{\partial E}{\partial w_{ji}^2} = \frac{\partial E}{\partial x_i^2} \frac{\partial x_i^2}{\partial w_{ji}^2}$$

$$\begin{aligned} \frac{\partial E}{\partial x_i^2} &= \sum_k \frac{\partial E}{\partial x_k^1} \frac{\partial x_k^1}{\partial o_i^1} \frac{\partial o_i^1}{\partial x_i^2} \\ &= \sum_k \delta_k^1 w_{ik}^1 f'(x_i^2) \equiv \delta_i^2 \end{aligned}$$

$$\frac{\partial E}{\partial w_{ji}^2} = \delta_i^2 o_j^2$$

$$\hat{w}_{ji}^2 = w_{ji}^2 - \varepsilon \delta_i^2 o_j^2$$



$$\frac{\partial E}{\partial w_{ji}^2} = \frac{\partial E}{\partial x_i^2} \frac{\partial x_i^2}{\partial w_{ji}^2}$$

$$\frac{\partial E}{\partial x_i^2} = \sum_k \frac{\partial E}{\partial x_k^1} \frac{\partial x_k^1}{\partial o_i^1} \frac{\partial o_i^1}{\partial x_i^2}$$

$$= \sum_k \delta_k^1 w_{ik}^1 f'(x_i^2) \equiv \delta_i^2$$

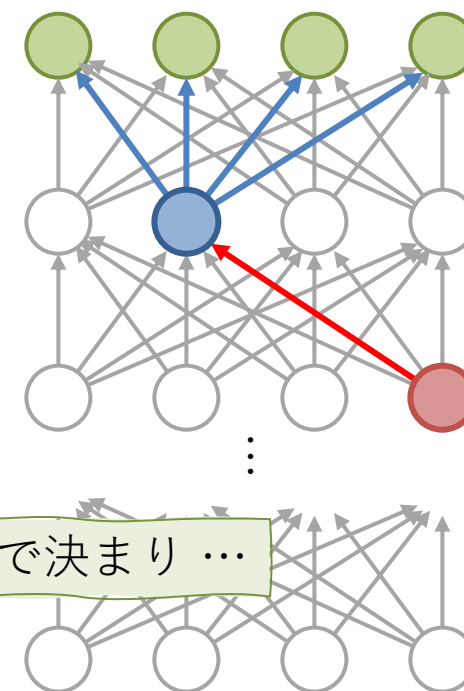
Eを減じるために  
どれだけ  $w_{ji}^2$  を変えるべきか

は

Eを減じるために  
どれだけ  $x_i^2$  を変えるべきか

$$\frac{\partial E}{\partial w_{ji}^2} = \delta_i^2 o_j^2$$

$$\hat{w}_{ji}^2 = w_{ji}^2 - \varepsilon \delta_i^2 o_j^2$$



$$o_i^1$$

$$w_{ji}^2$$

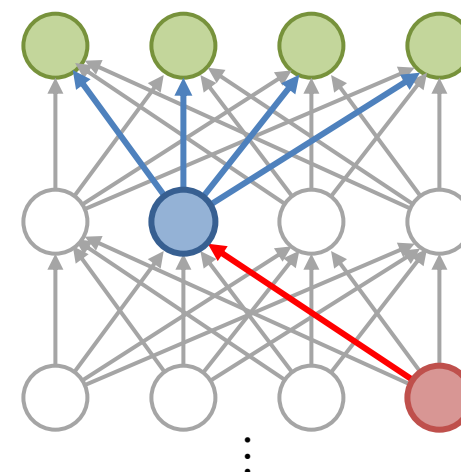
$$o_{ji}^2, \quad x_i^2 = \sum_j o_j^2 w_{ji}^2$$



$$\frac{\partial E}{\partial w_{ji}^2} = \frac{\partial E}{\partial x_i^2} \frac{\partial x_i^2}{\partial w_{ji}^2}$$

$$\frac{\partial E}{\partial x_i^2} = \sum_k \frac{\partial E}{\partial x_k^1} \frac{\partial x_k^1}{\partial o_i^1} \frac{\partial o_i^1}{\partial x_i^2}$$

$$= \sum_k \delta_k^1 w_{ik}^1 f'(x_i^2) \equiv \delta_i^2$$



$$o_i^1, w_{ji}^2, x_i^2 = \sum_j o_j^2 w_{ji}^2$$

Eを減じるために  
どれだけ  $w_{ji}^2$  を変えるべきか

は

Eを減じるために  
どれだけ  $x_i^2$  を変えるべきか

で決まり ...

$$\frac{\partial E}{\partial w_{ji}^2} = \delta_i^2 o_j^2$$

Eを減じるために  
どれだけ  $x_i^2$  を変えるべきか

は

Eを減じるために  
どれだけ  $x_i^1$  を変えるべきか

で決まる

$$\hat{w}_{ji}^2 = w_{ji}^2 - \varepsilon \delta_i^2 o_j^2$$

Eを減じるために  
どれだけ  $x_i^1$  を変えるべきか

は

Error

で決まる



$$\frac{\partial E}{\partial w_{ji}^2} = \frac{\partial E}{\partial x_i^2} \frac{\partial x_i^2}{\partial w_{ji}^2}$$

$$\frac{\partial E}{\partial x_i^2} = \sum_k \frac{\partial E}{\partial x_k^1} \frac{\partial x_k^1}{\partial o_i^1} \frac{\partial o_i^1}{\partial x_i^2}$$

$$= \sum_k \delta_k^1 w_{ik}^1 f'(x_i^2) \equiv \delta_i^2$$

$$\frac{\partial E}{\partial w_{ji}^2} = \delta_i^2 o_j^2$$

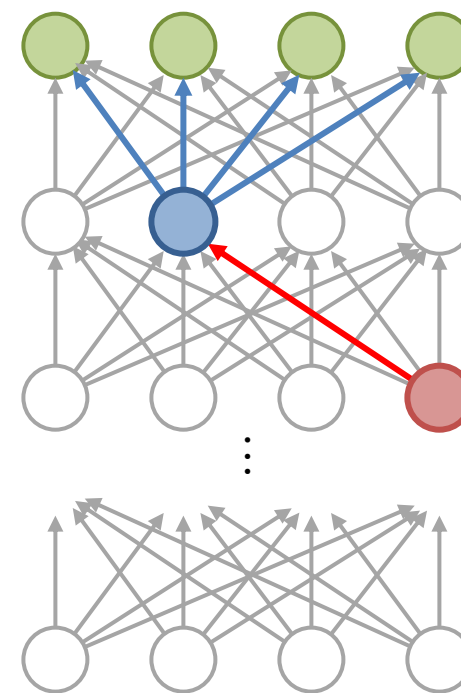
Eを減じるために  
どれだけ  $x_i^2$  を変えるべきか

は

$\delta_i^1$

で決まる

$$\hat{w}_{ji}^2 = w_{ji}^2 - \varepsilon \delta_i^2 o_j^2$$



$o_i^1$

$w_{ji}^2$

$$o_{ji}^2, x_i^2 = \sum_j o_j^2 w_{ji}^2$$

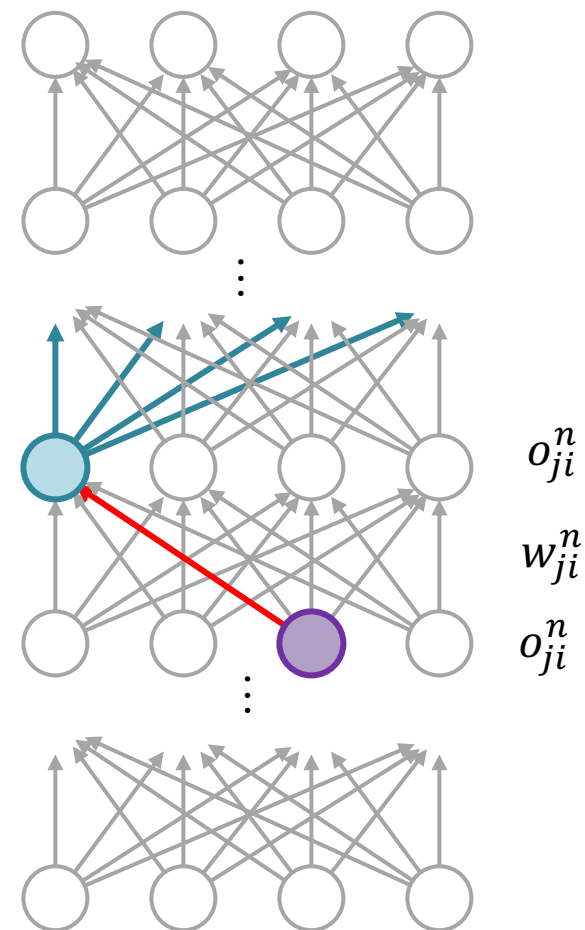


$$\frac{\partial E}{\partial w_{ji}^n} = \frac{\partial E}{\partial x_i^n} \frac{\partial x_i^n}{\partial w_{ji}^n}$$

$$\begin{aligned} \frac{\partial E}{\partial x_i^n} &= \sum_k \frac{\partial E}{\partial x_k^{n-1}} \frac{\partial x_k^{n-1}}{\partial o_i^{n-1}} \frac{\partial o_i^{n-1}}{\partial x_i^n} \\ &= \sum_k \delta_k^{n-1} w_{ik}^{n-1} f'(x_i^n) \equiv \delta_i^n \end{aligned}$$

$$\frac{\partial E}{\partial w_{ji}^n} = \delta_i^n o_j^n$$

$$\hat{w}_{ji}^n = w_{ji}^n - \varepsilon \delta_i^n o_j^n$$



$$\frac{\partial E}{\partial w_{ji}^n} = \frac{\partial E}{\partial x_i^n} \frac{\partial x_i^n}{\partial w_{ji}^n}$$

$$\frac{\partial E}{\partial x_i^n} = \sum_k \frac{\partial E}{\partial x_k^{n-1}} \frac{\partial x_k^{n-1}}{\partial o_i^{n-1}} \frac{\partial o_i^{n-1}}{\partial x_i^n}$$

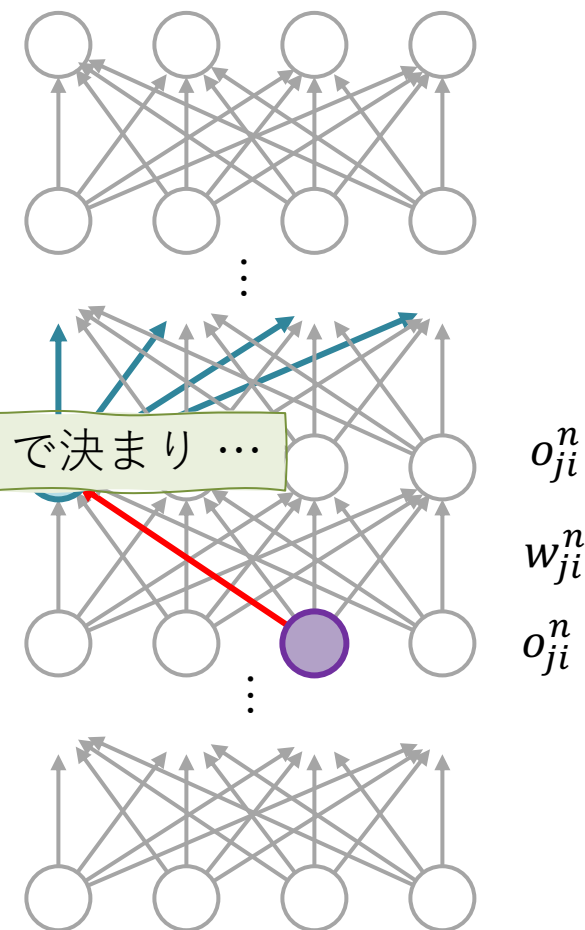
Eを減じるために  
どれだけ  $w_{ji}^n$  を変えるべきか

は

Eを減じるために  
どれだけ  $x_i^n$  を変えるべきか

$$\frac{\partial E}{\partial w_{ji}^n} = \delta_i^n o_j^n$$

$$\hat{w}_{ji}^n = w_{ji}^n - \varepsilon \delta_i^n o_j^n$$



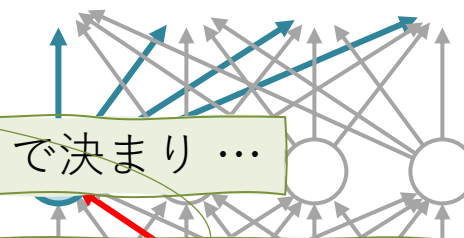
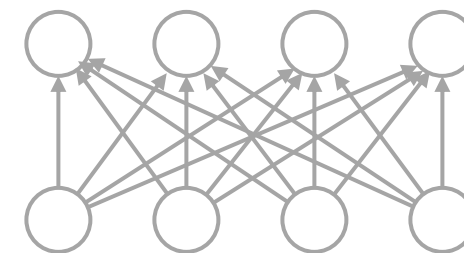
$$\frac{\partial E}{\partial w_{ji}^n} = \frac{\partial E}{\partial x_i^n} \frac{\partial x_i^n}{\partial w_{ji}^n}$$

$$\frac{\partial E}{\partial x_i^n} = \sum_k \frac{\partial E}{\partial x_k^{n-1}} \frac{\partial x_k^{n-1}}{\partial o_i^{n-1}} \frac{\partial o_i^{n-1}}{\partial x_i^n}$$

Eを減じるために  
どれだけ  $w_{ji}^n$  を変えるべきか

は

Eを減じるために  
どれだけ  $x_i^n$  を変えるべきか



で決まり ...

$o_{ji}^n$

$$\frac{\partial E}{\partial w_{ji}^n} = \delta_i^n o_j^n$$

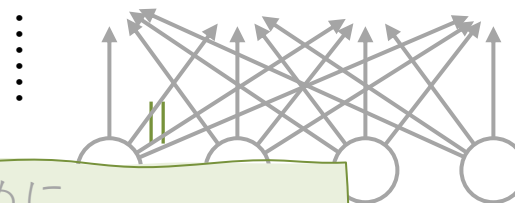
Eを減じるために  
どれだけ  $x_i^n$  を変えるべきか

は

Eを減じるために  
どれだけ  $x_i^{n-1}$  を変えるべきか

で決まり ...

$$\hat{w}_{ji}^n = w_{ji}^n - \varepsilon \delta_i^n o_j^n$$



Eを減じるために  
どれだけ  $x_i^1$  を変えるべきか

は

Error

で決まる

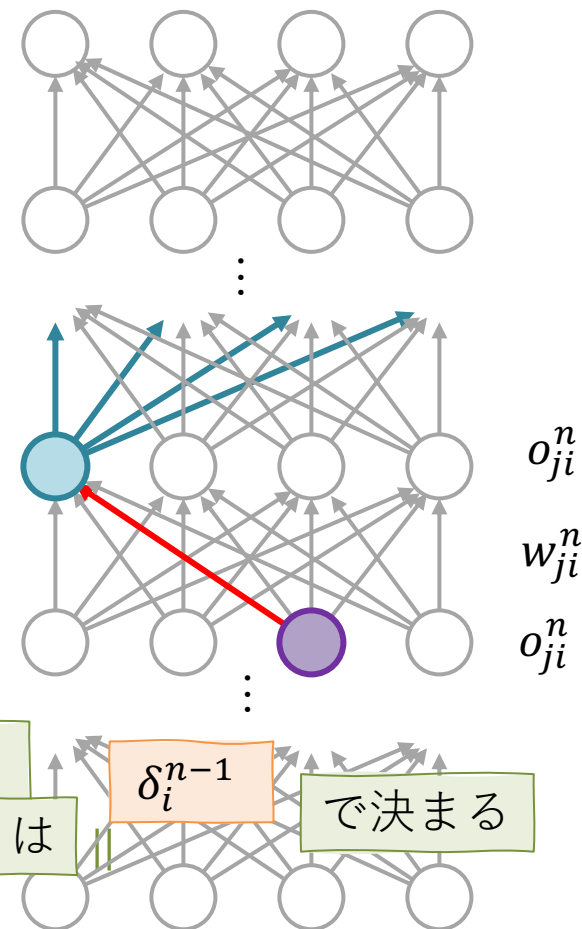


$$\frac{\partial E}{\partial w_{ji}^n} = \frac{\partial E}{\partial x_i^n} \frac{\partial x_i^n}{\partial w_{ji}^n}$$

$$\begin{aligned} \frac{\partial E}{\partial x_i^n} &= \sum_k \frac{\partial E}{\partial x_k^{n-1}} \frac{\partial x_k^{n-1}}{\partial o_i^{n-1}} \frac{\partial o_i^{n-1}}{\partial x_i^n} \\ &= \sum_k \delta_k^{n-1} w_{ik}^{n-1} f'(x_i^n) \equiv \delta_i^n \end{aligned}$$

$$\frac{\partial E}{\partial w_{ji}^n} = \delta_i^n o_j^n$$

$$\hat{w}_{ji}^n = w_{ji}^n - \varepsilon \delta_i^n o_j^n$$





# MLPのパラメタ推定が 抱える問題



# パラメタ推定における問題

## □ 勾配法が抱える一般的問題

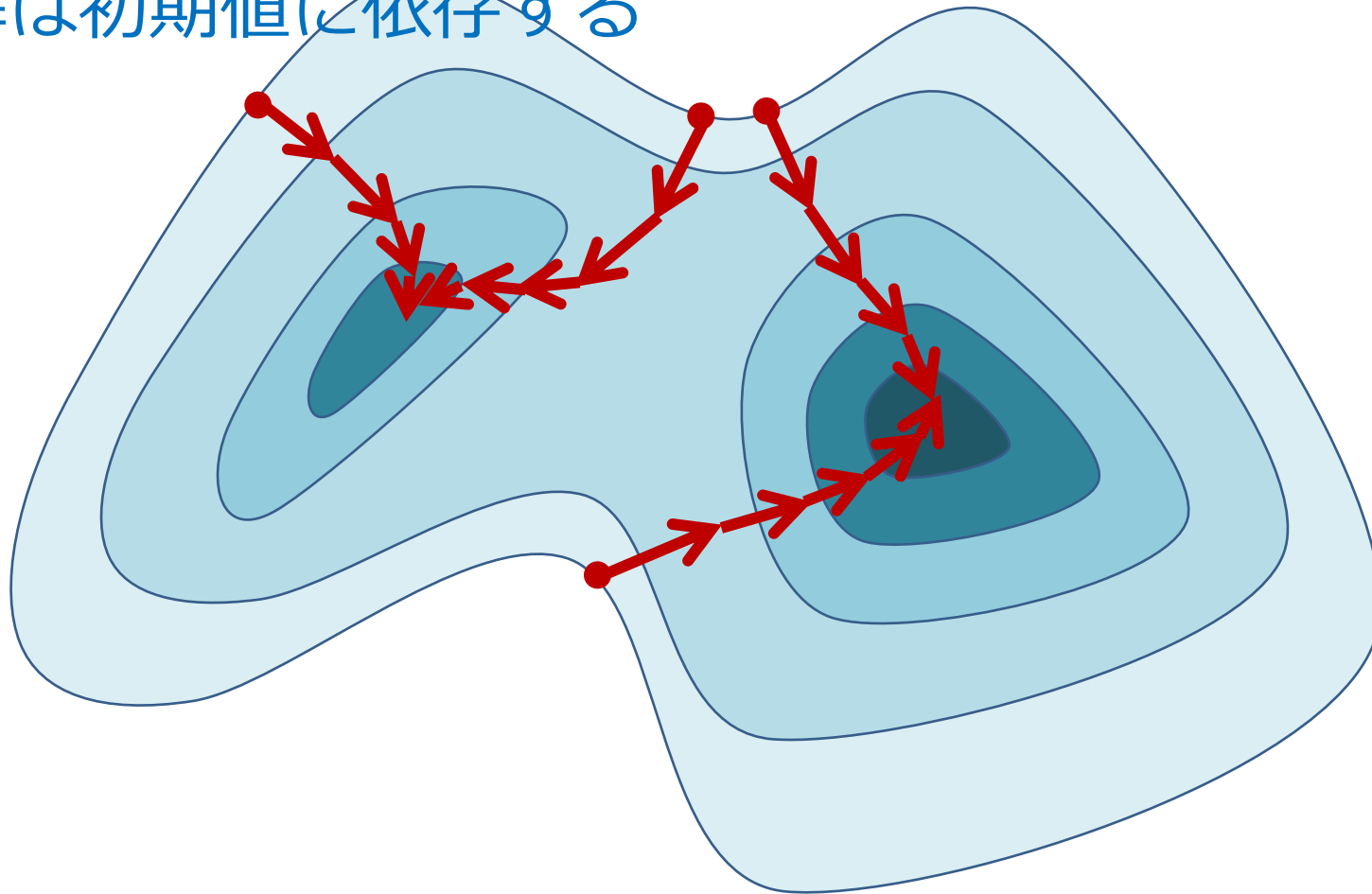
- 解の初期値依存性
  - 局所最適解にトラップされる
  - 勾配消失問題
- 不定解を生じやすい

## □ 勾配法が一般的に抱える問題は、DNNの構造が複雑なだけに顕在化しやすい。



# 勾配法と初期値問題

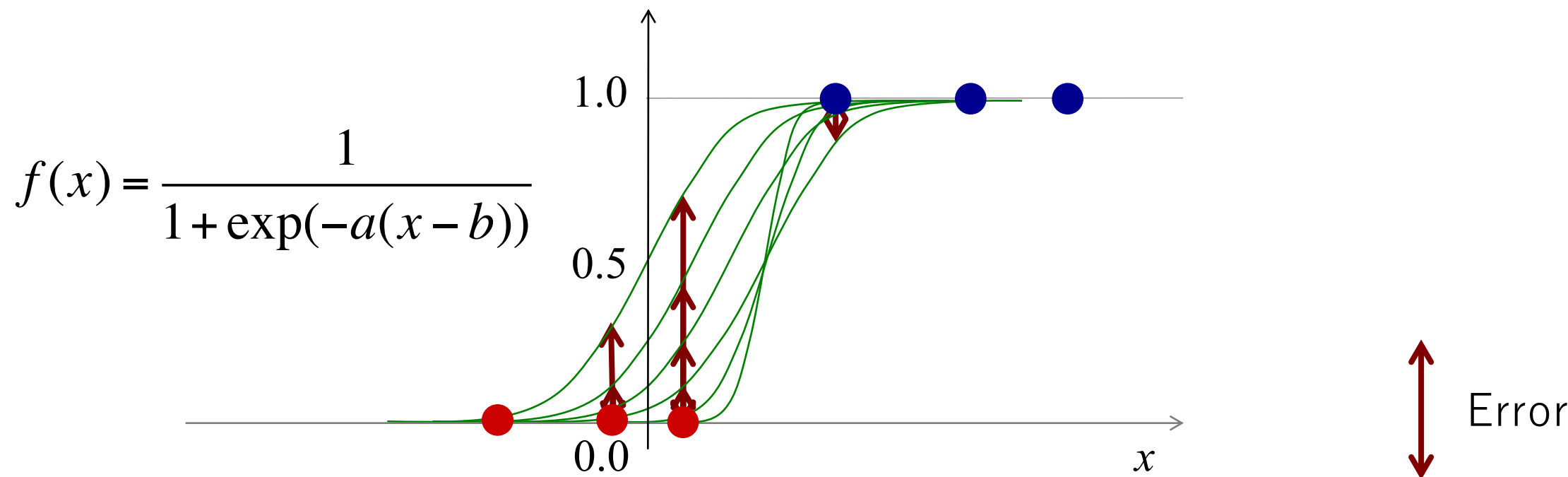
勾配法においては、  
得られる解は初期値に依存する



# 勾配消失問題

## 1-dimensional view

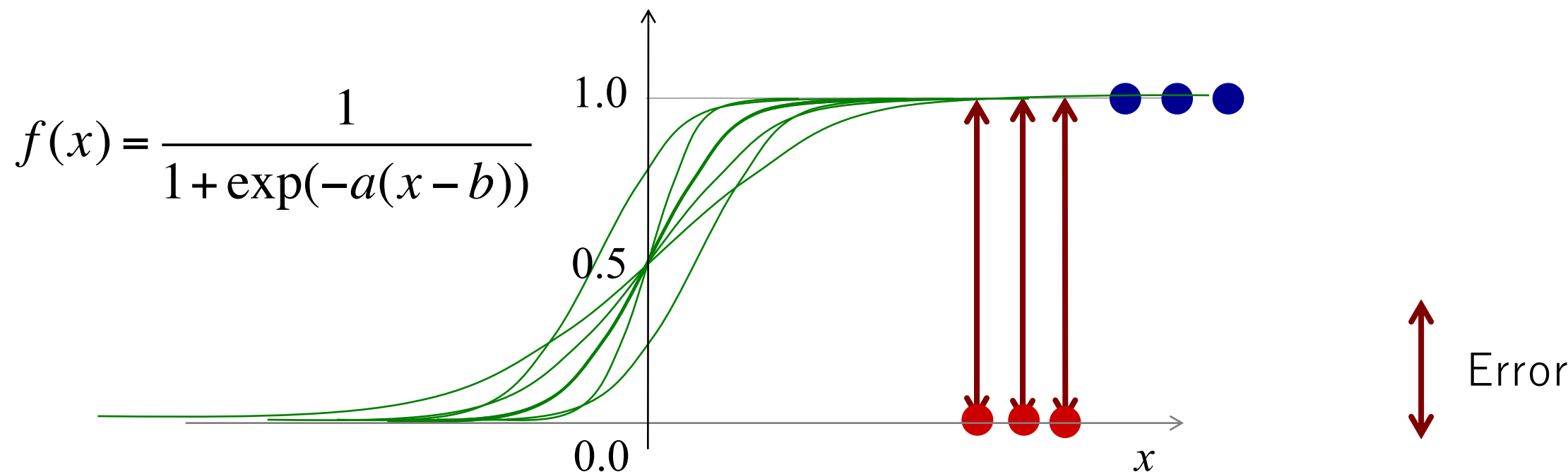
初期解が最適解に近ければ，最適解をみつけて収束する。



# 勾配消失問題

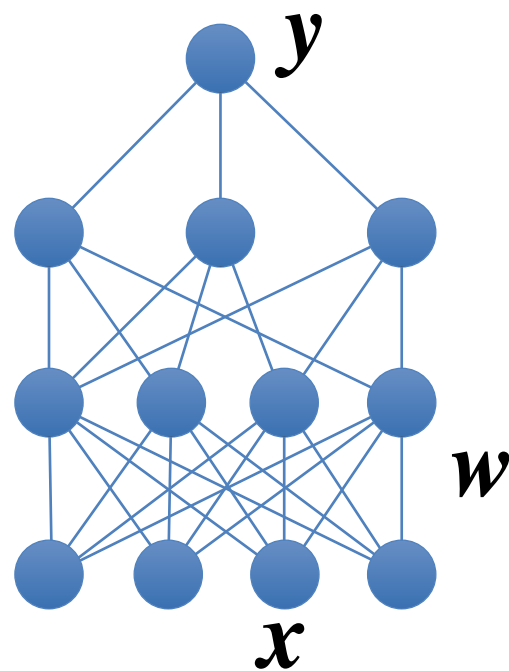
## 1-dimensional view

初期解が最適解から離れているとき、  
解を見つけることなくパラメタ更新が止まる。

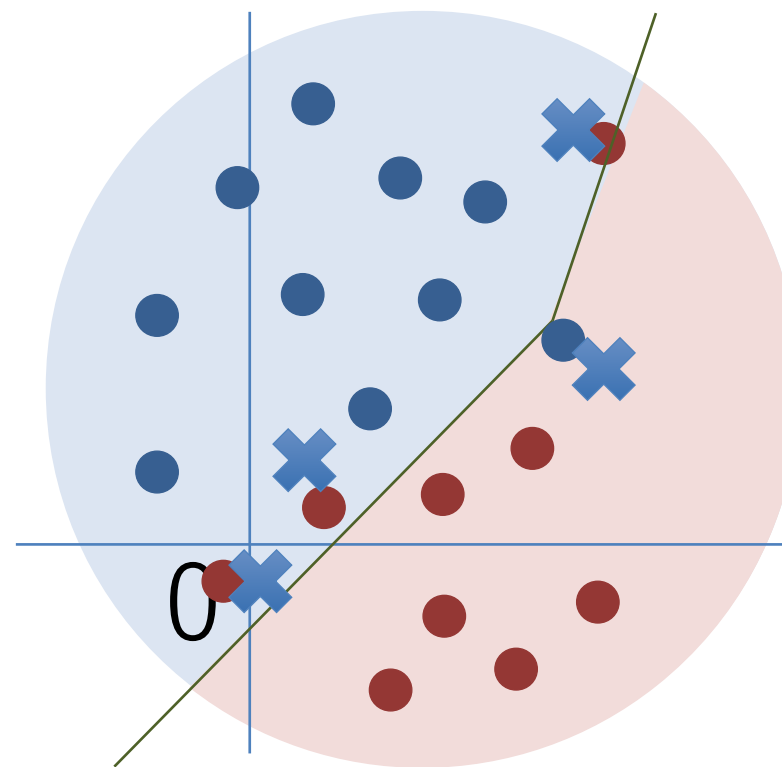


# 勾配消失問題

初期解が最適解に近ければ,



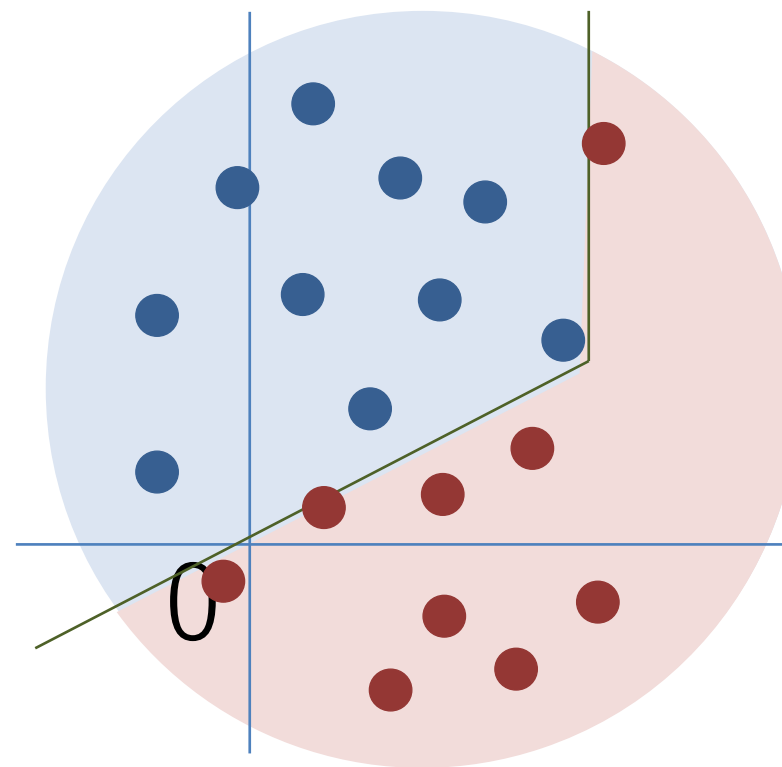
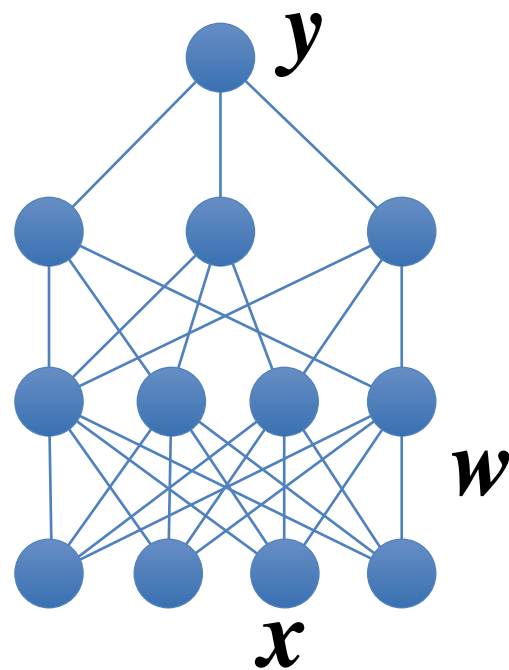
2-dimensional view



# 勾配消失問題

## 2-dimensional view

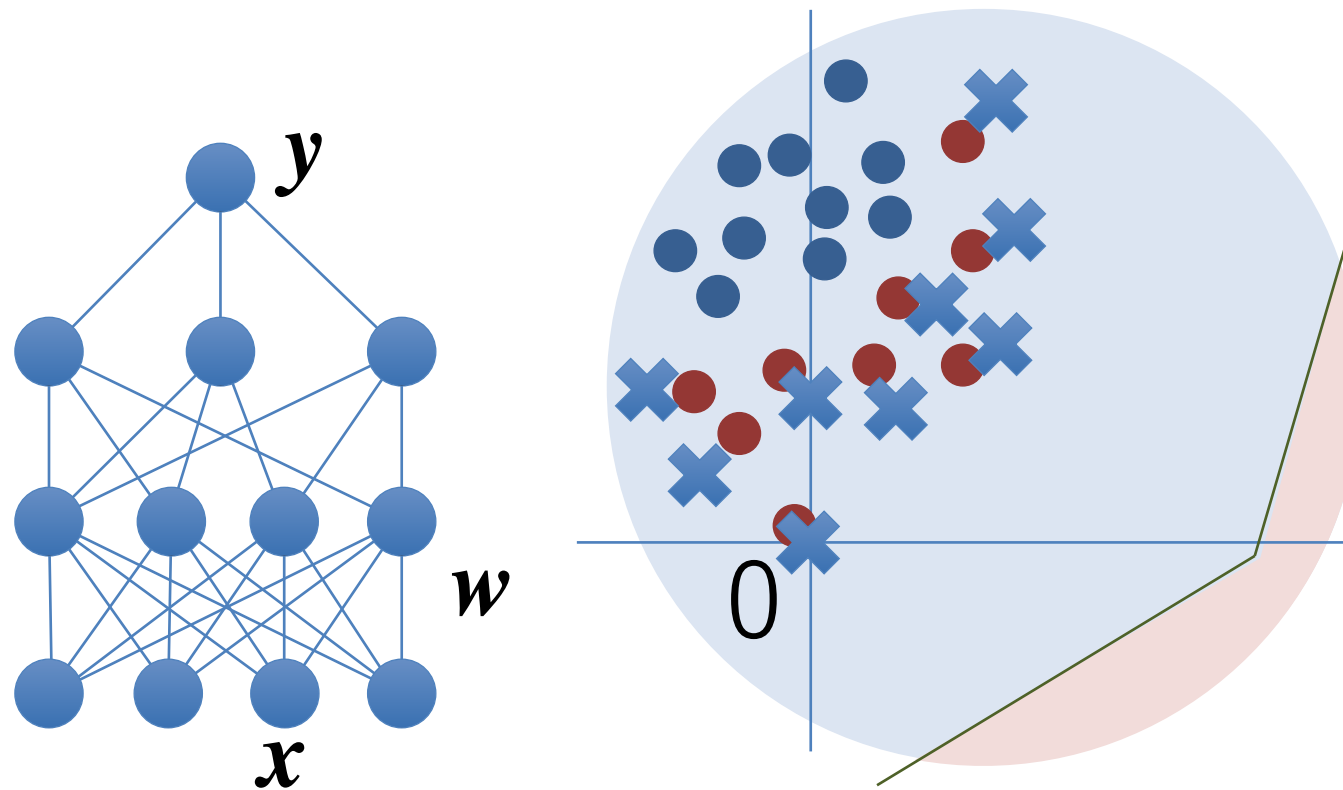
初期解が最適解に近ければ，最適解をみつけて収束する。



# 勾配消失問題

2-dimensional view

初期解が最適解から離れているとき,

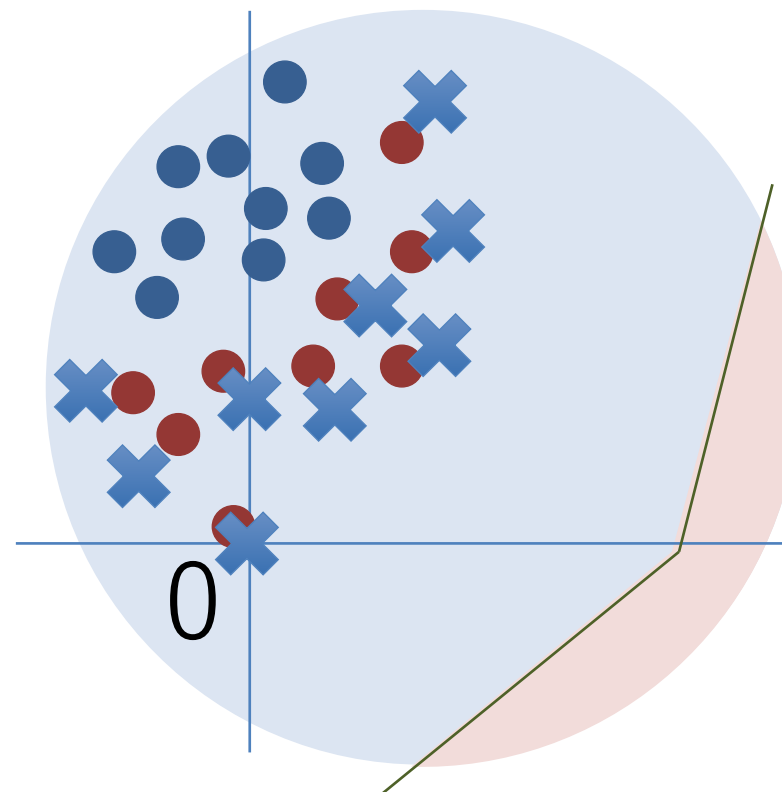
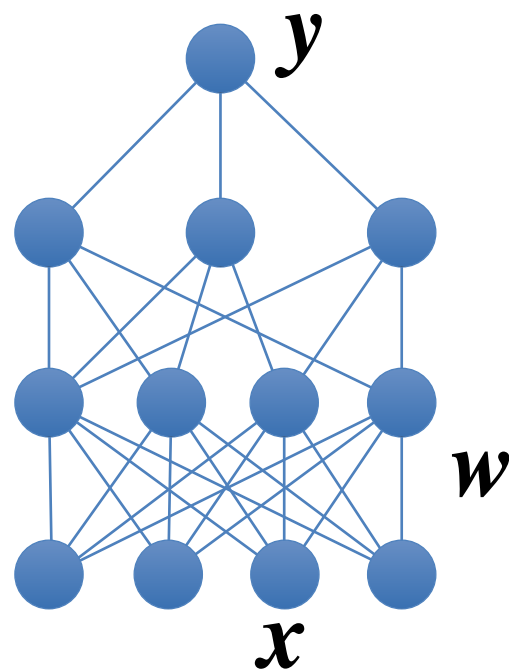




# 勾配消失問題

## 2-dimensional view

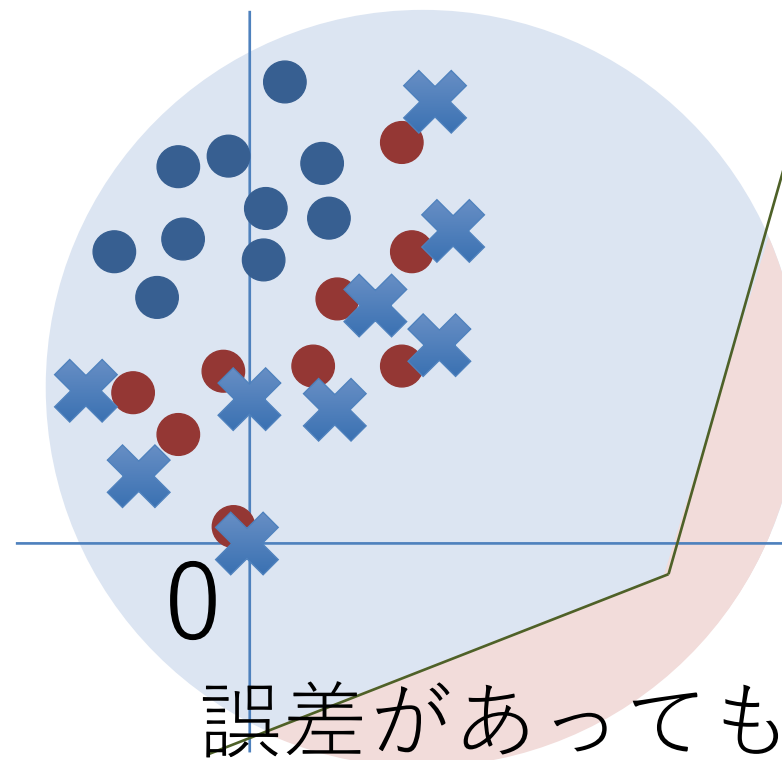
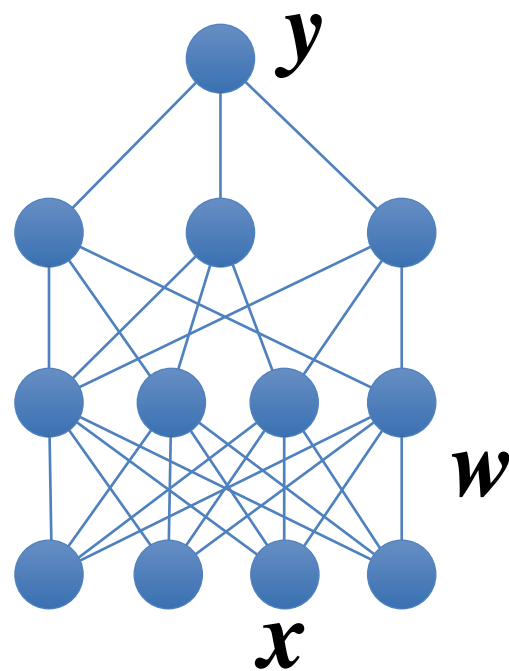
初期解が最適解から離れているとき、  
解を見つけることなくパラメタ更新が止まる。



# 勾配消失問題

2-dimensional view

初期解が最適解から離れているとき、  
解を見つけることなくパラメタ更新が止まる。



$$\left. \frac{\partial E}{\partial \mathbf{w}} \right|_{old} = 0$$

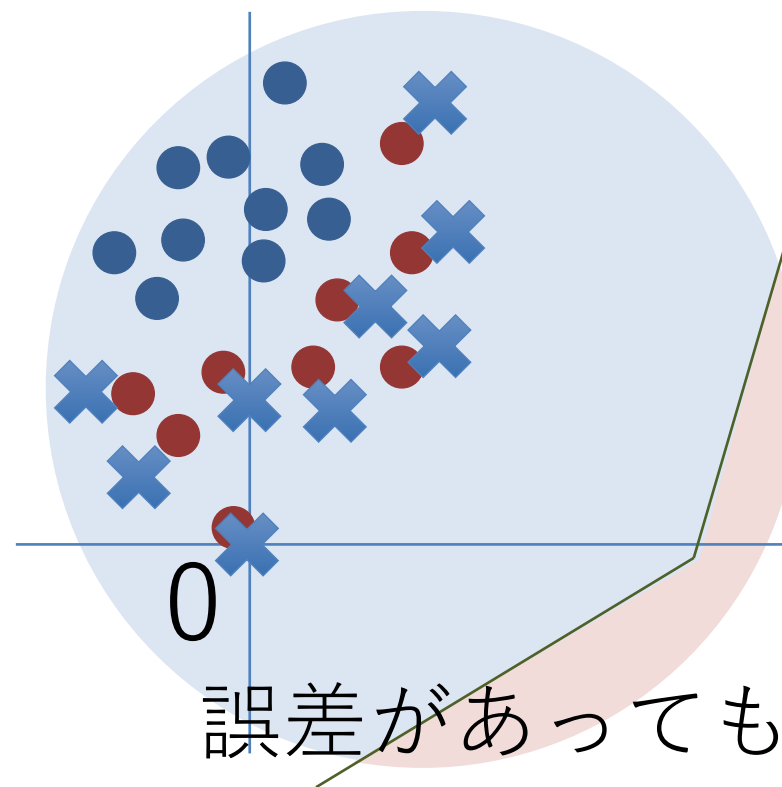
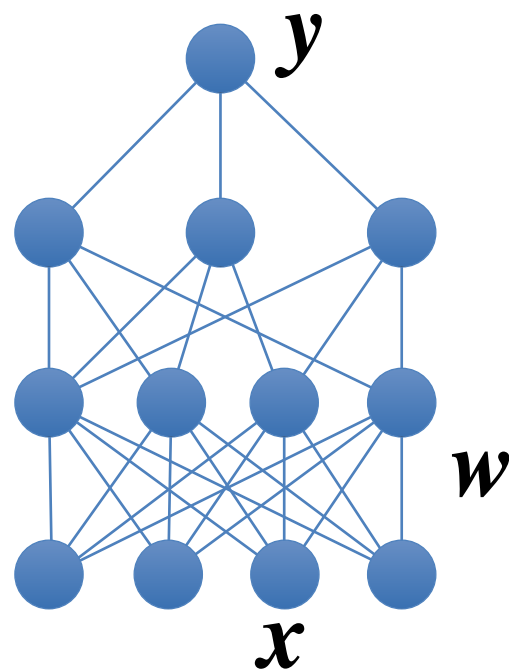
$$\mathbf{w}_{new} = \mathbf{w}_{old}$$

→パラメタは更新されない

# 勾配消失問題

## 2-dimensional view

初期解が最適解から離れているとき、  
解を見つけることなくパラメタ更新が止まる。



$$\left. \frac{\partial E}{\partial \mathbf{w}} \right|_{old} = 0$$

$$\mathbf{w}_{new} = \mathbf{w}_{old}$$

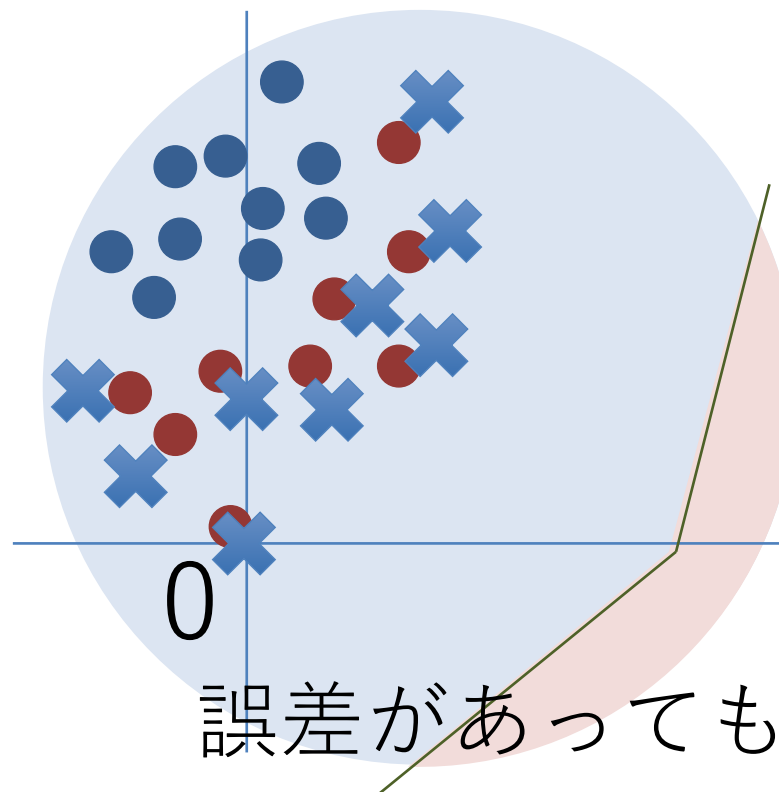
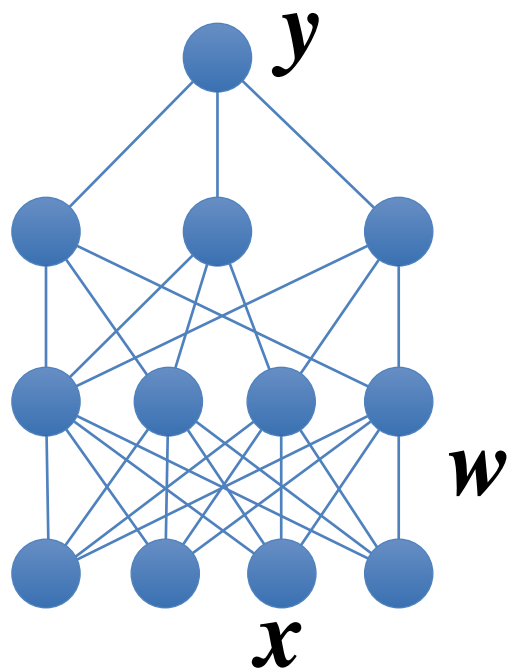
誤差があっても、傾きがゼロ

→パラメタは更新されない

# 勾配消失問題

## 2-dimensional view

初期解が最適解から離れているとき、  
解を見つけることなくパラメタ更新が止まる。



$$\left. \frac{\partial E}{\partial \mathbf{w}} \right|_{old} = 0$$

$$\mathbf{w}_{new} = \mathbf{w}_{old}$$

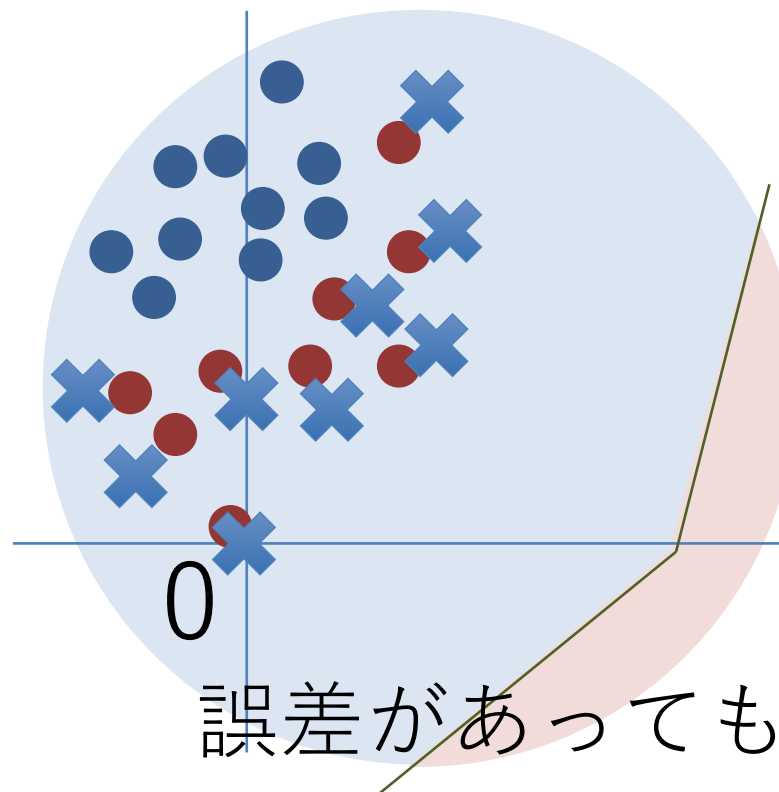
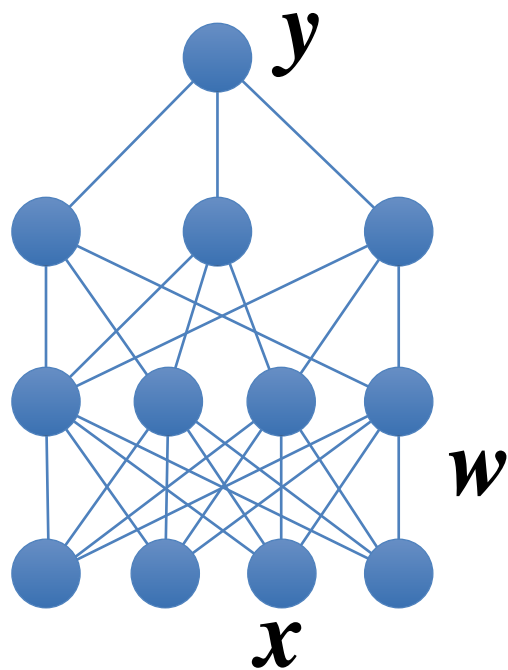
誤差があっても、傾きがゼロ

→パラメタは更新されない

# 勾配消失問題

2-dimensional view

初期解が最適解から離れているとき、  
解を見つけることなくパラメタ更新が止まる。



$$\left. \frac{\partial E}{\partial \mathbf{w}} \right|_{old} = 0$$

$$\mathbf{w}_{new} = \mathbf{w}_{old}$$

誤差があっても、傾きがゼロ

→パラメタは更新されない

Vanishing gradient problem

# まとめ

- MLPのパラメタ推定には，誤差逆伝播法が用いられる。
- 誤差逆伝播法は，勾配法を採用する。
- 誤差逆伝播法は，初期値問題，勾配消失問題などの問題を持つ。

