

深層学習の基礎

理化学研究所

hataya.tokyo

幡谷龍一郎

目次

- **深層学習以前**
- **深層学習**
- **深層学習の実装**

目次

- **深層学習以前（ニューラルネットワークとは）**
- 深層学習
- 深層学習の実装

今年のノーベル物理学賞

The Nobel Prize in Physics 2024

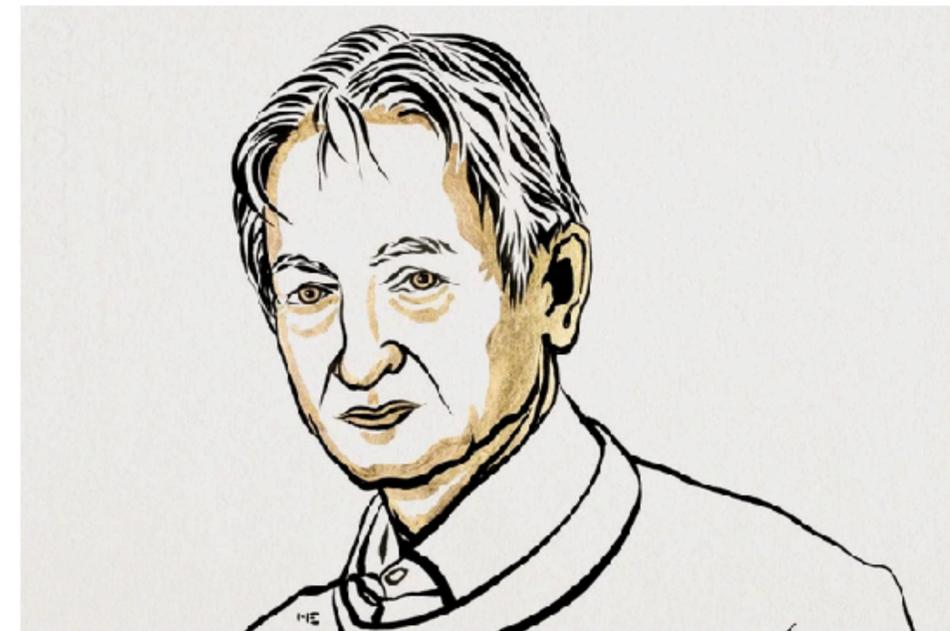
John J. Hopfield

Geoffrey Hinton

*“for foundational discoveries and inventions that enable machine learning with **artificial neural networks**”*



John Hopfield. Ill. Niklas Elmehed © Nobel Prize Outreach

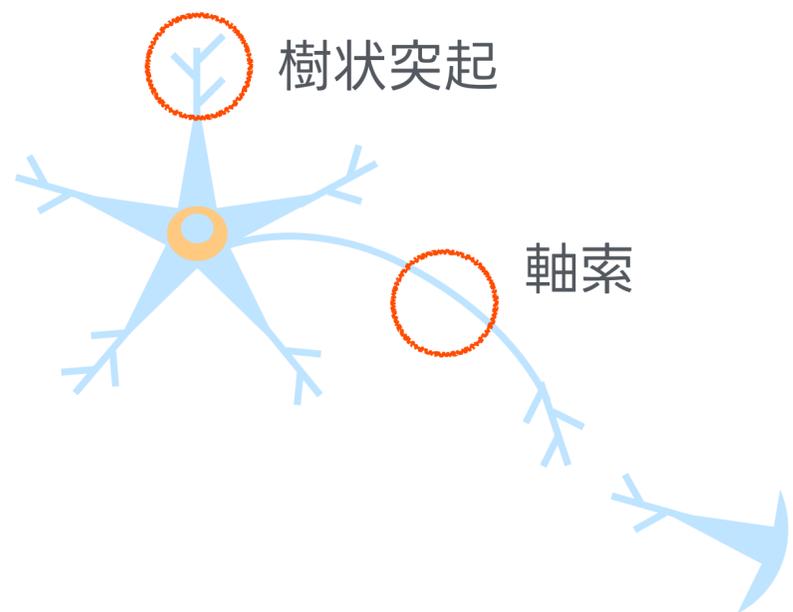


Geoffrey Hinton. Ill. Niklas Elmehed © Nobel Prize Outreach

*ただし今回受賞対象となったタイプのニューラルネットは今回は取り扱わない

人工ニューロン

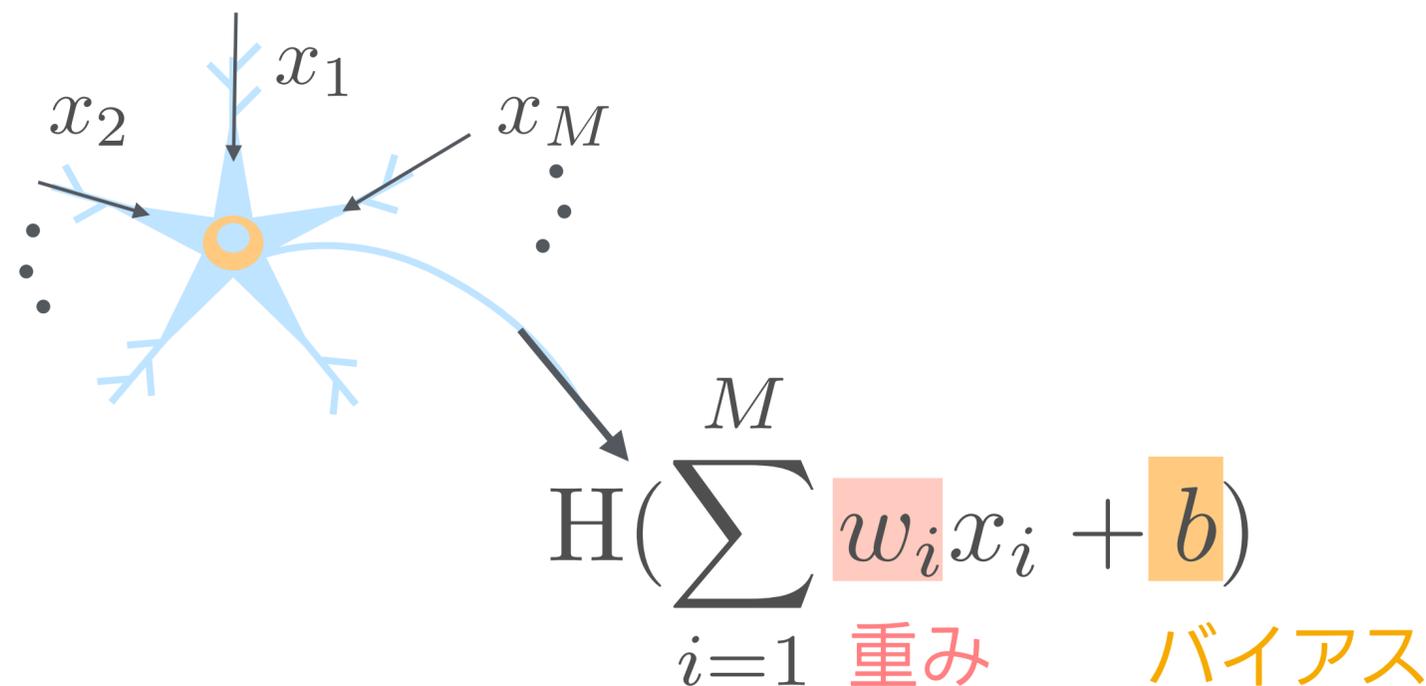
神経細胞



樹状突起への入力が閾値を超えた時だけ
軸索を通じて他の神経細胞に信号が伝達される
(全か無かの法則)

マッカロック・ピッツモデル [McCulloch&Pitts,43]

神経細胞の数理モデル (人工ニューロン)

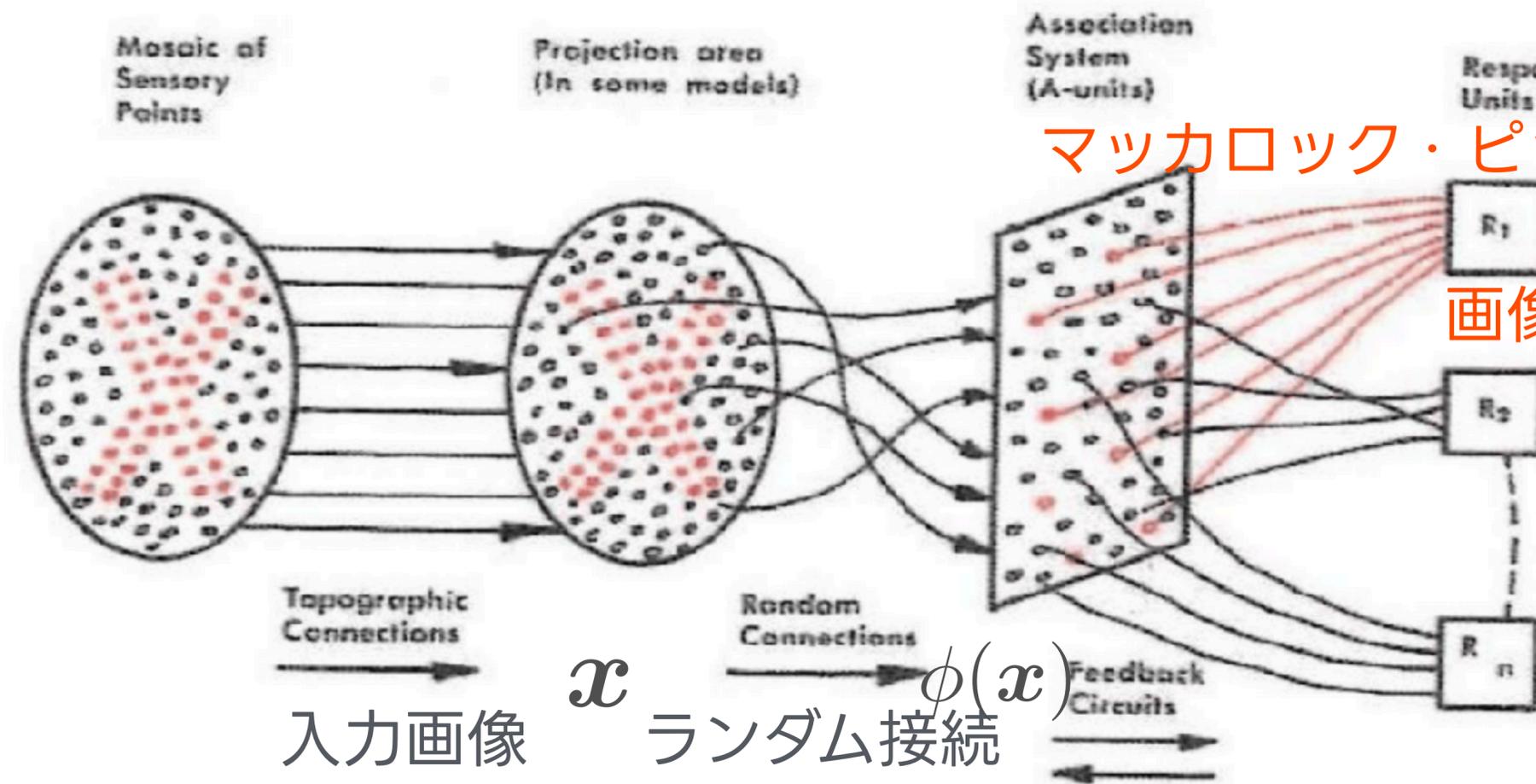


$$H(x) = \begin{cases} 1 & (x \geq 0) \\ 0 & (x < 0) \end{cases}$$

$\sum_{i=1}^M w_i x_i$ が閾値 $-b$ を超えるとき出力は1
そうでなければ0

パーセプトロン [Rosenblatt,57]

マッカロック・ピッツモデルによってパターン認識を行う機械学習モデル



マッカロック・ピッツモデル

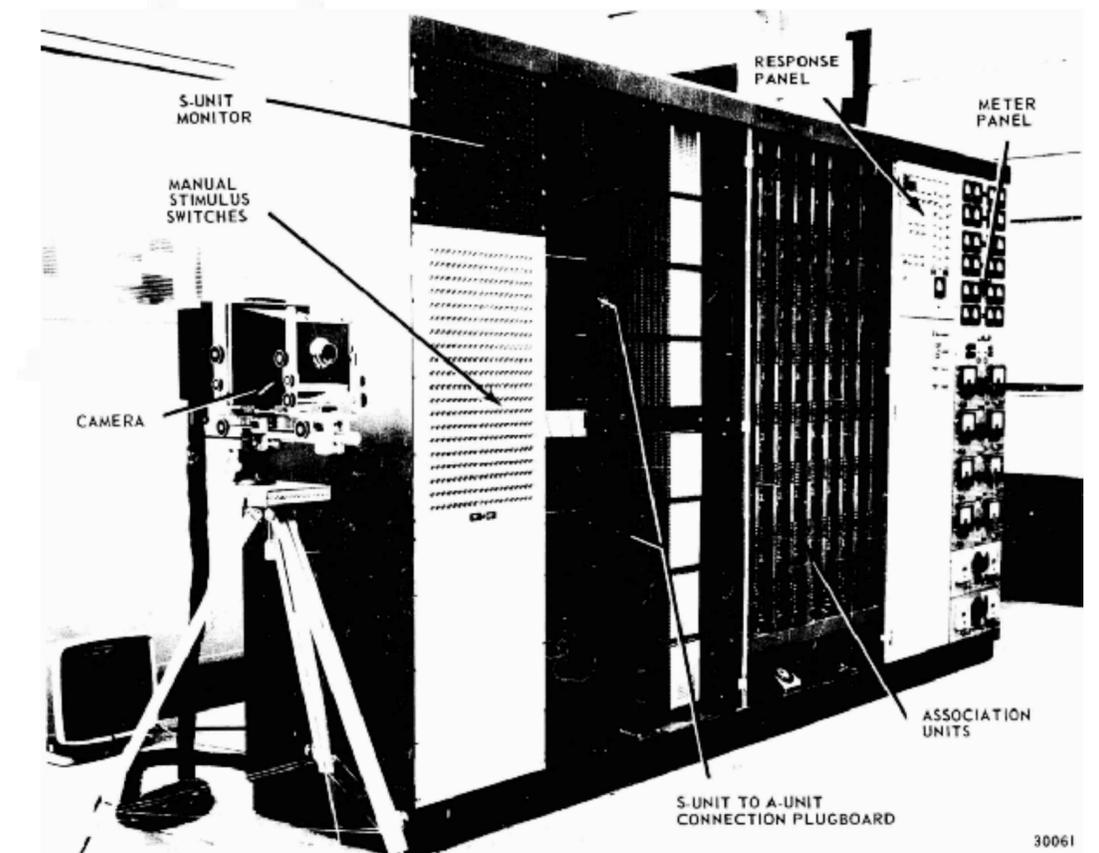
画像はx?

Output Imp. $H(w^T \phi(x) + b)$

w は誤答が少なくなるように更新
ハードウェアとして実装された

人工ニューロンの回路 → 人工ニューラルネットワーク

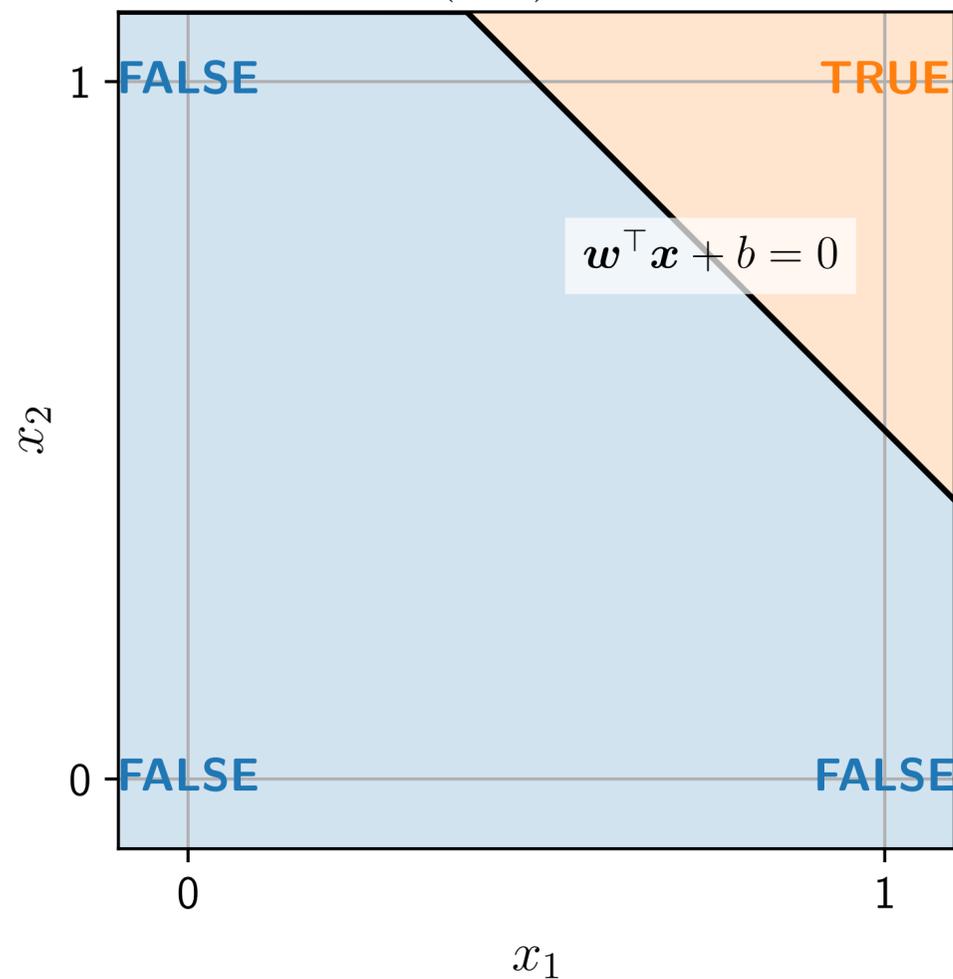
📺 第一次ニューラルネットブームのきっかけ



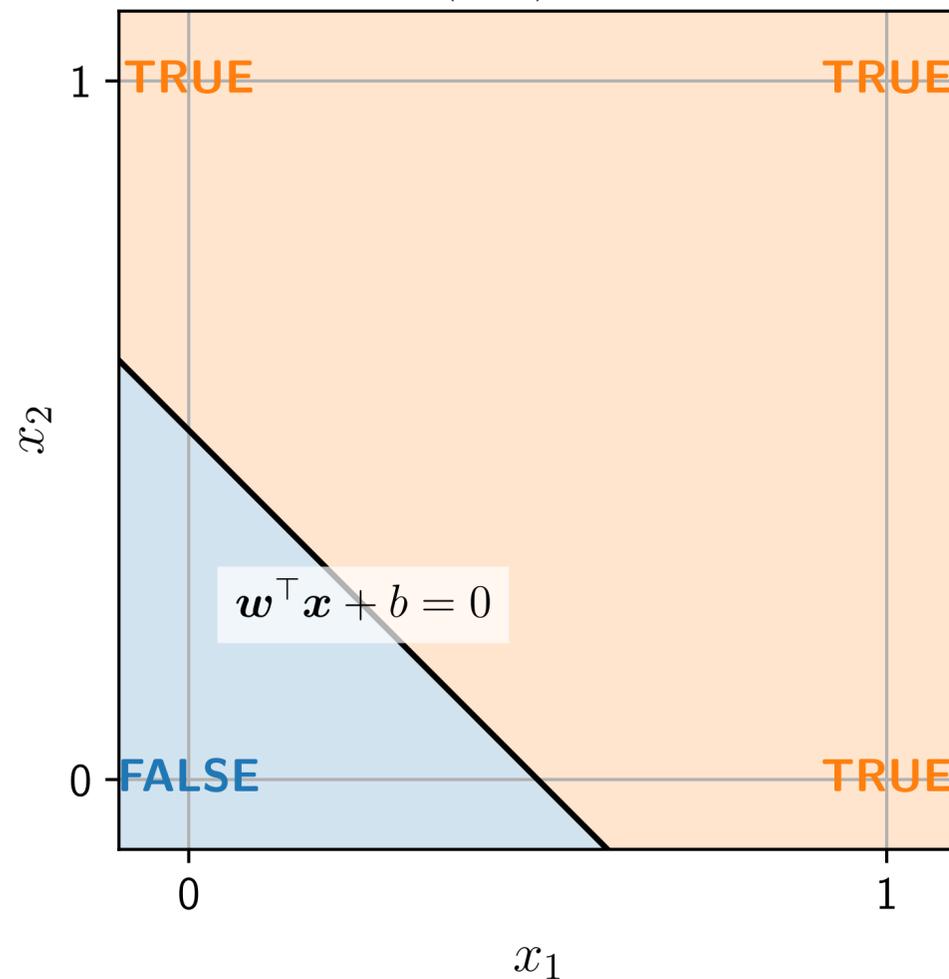
パーセプトロン [Rosenblatt,57]

マッカロック・ピッツモデル（パーセプトロン素子）は論理関数を表現できる

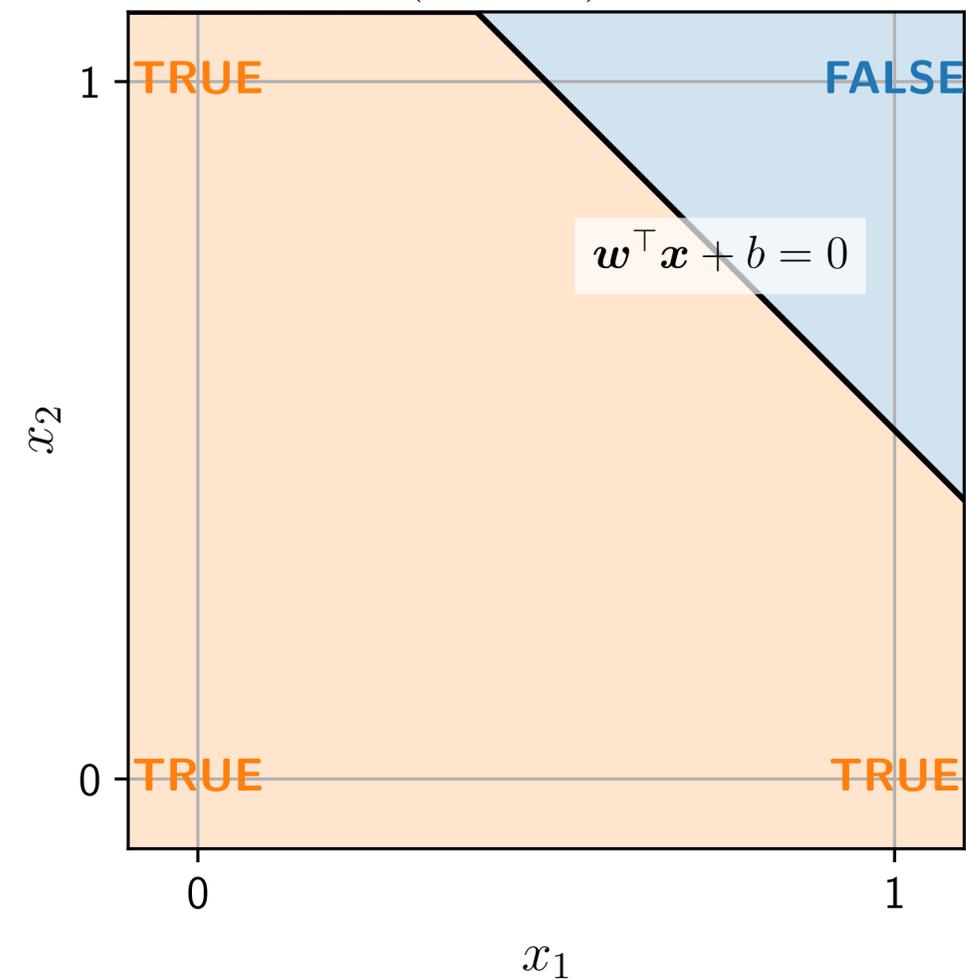
AND $x_1 \wedge x_2$
 $w = (1, 1), b = 1.5$



OR $x_1 \vee x_2$
 $w = (1, 1), b = 0.5$



NAND $\neg(x_1 \wedge x_2)$
 $w = (-1, -1), b = 1.5$



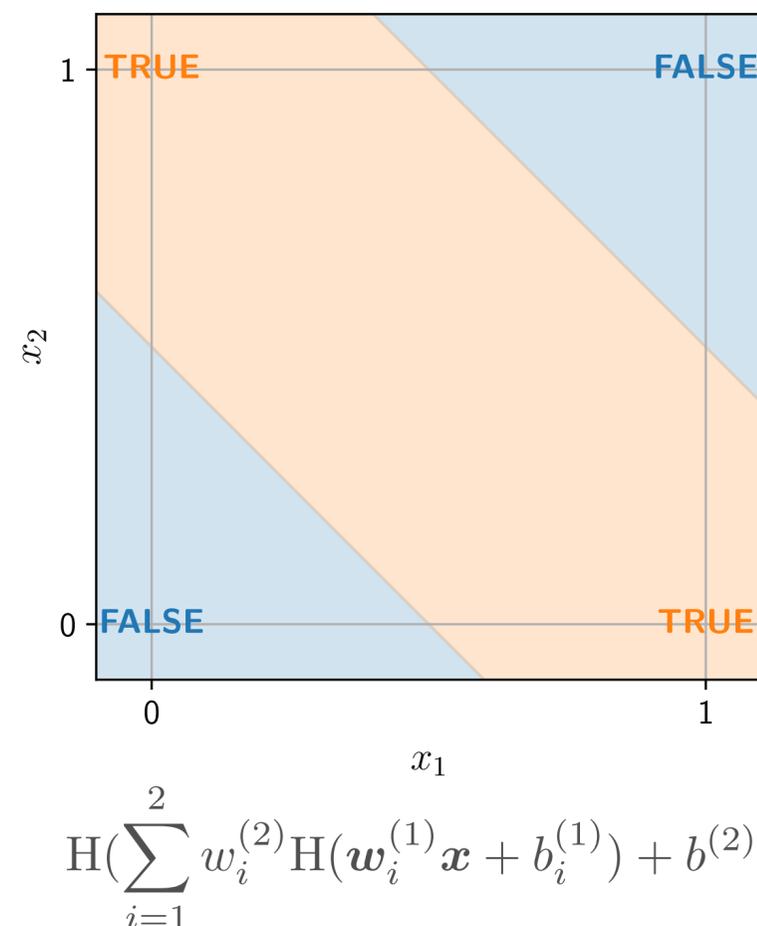
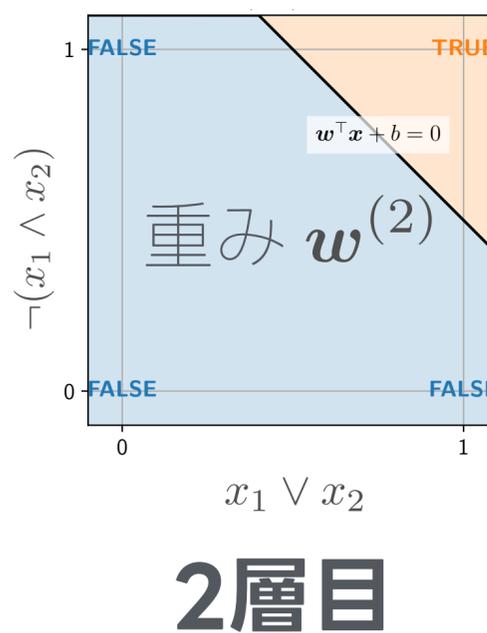
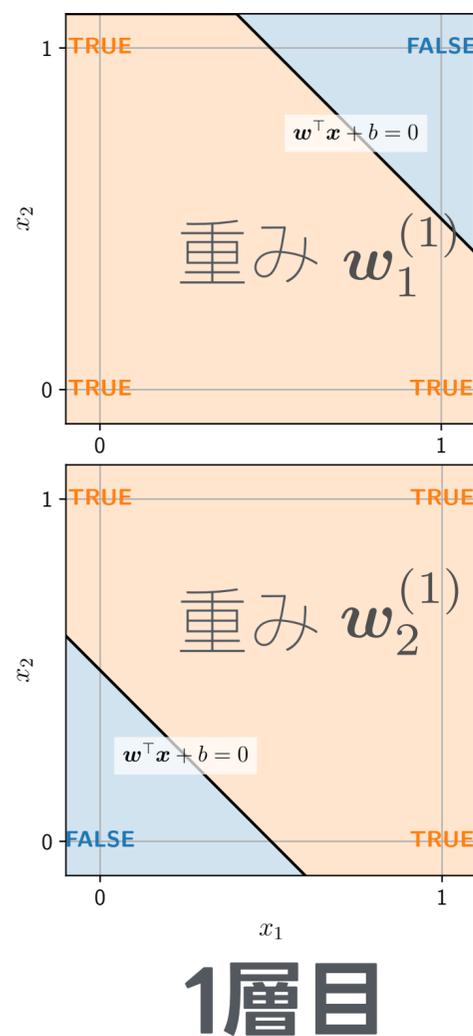
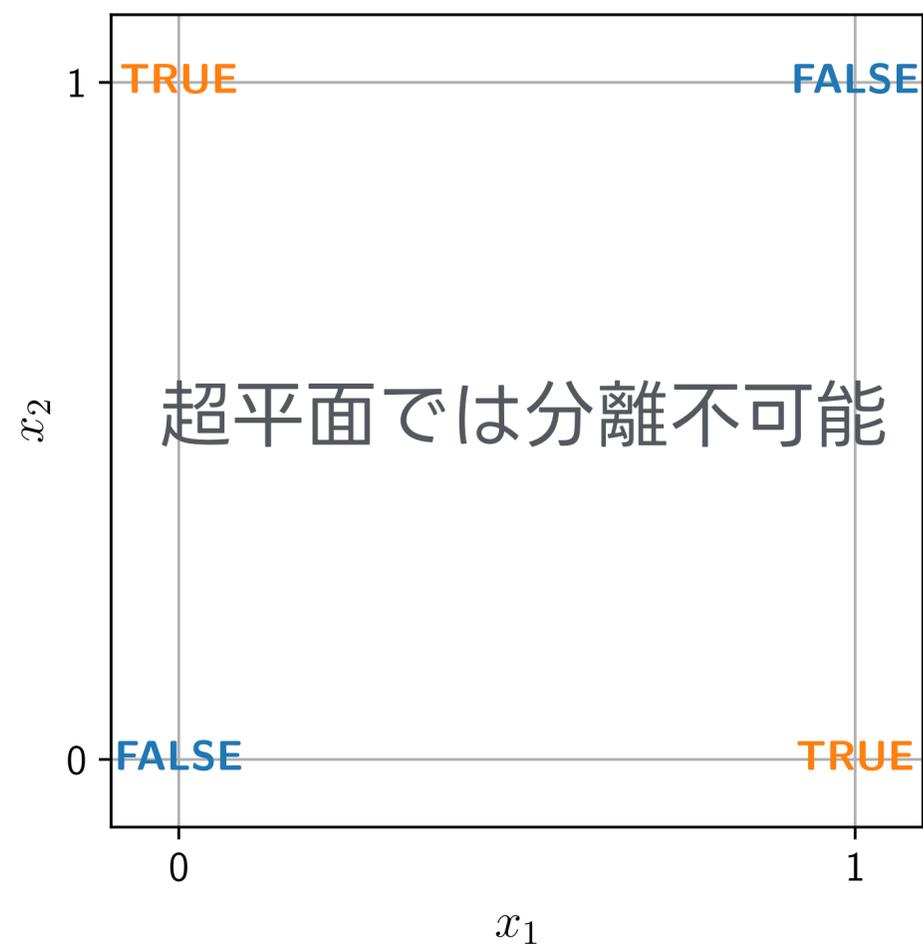
任意の論理関数はNANDの組合せで表現可能
→任意の論理関数は多層化したパーセプトロン素子により表現可能

パーセプトロン [Rosenblatt,57]

マッカロック・ピッツモデル（パーセプトロン素子）は線形分離可能な問題しか扱えない

XOR $x_1 \oplus x_2$

多層にすれば可能 $x_1 \oplus x_2 = (x_1 \vee x_2) \wedge \neg(x_1 \wedge x_2)$



1段だと線形分離不可能なことに注目が集まり第1次ニューラルネットブームは去った（らしい）
 多層の場合，最終層の重み $(w^{(2)}, b^{(2)})$ しか更新できない問題も

誤差逆伝搬法 (backpropagation)

[Amari 67, Rumelhart+86]

Hintonも共著

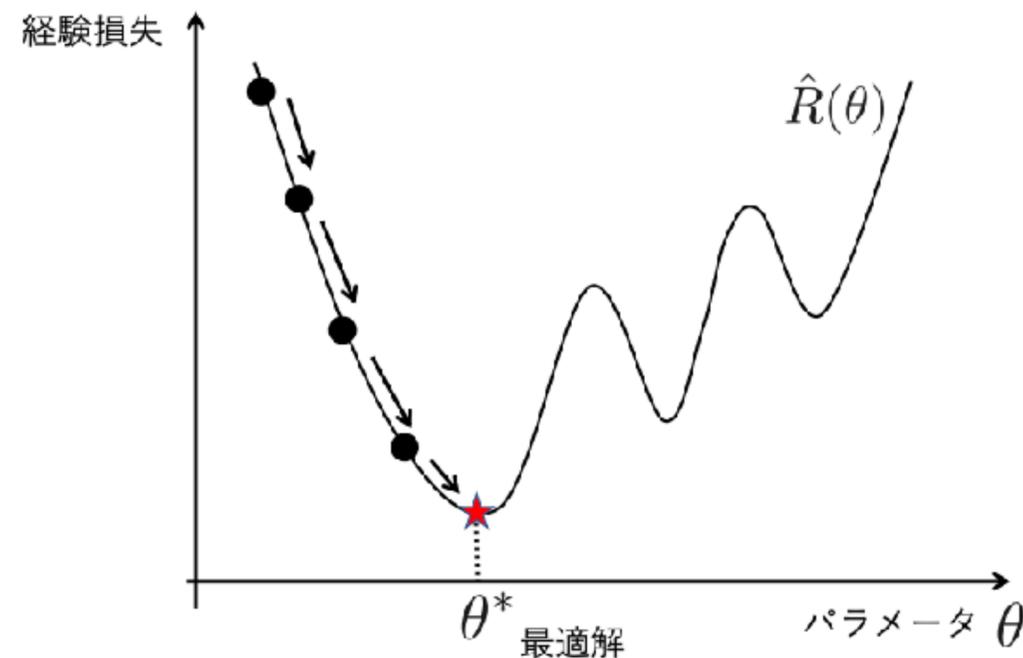


2層のパーセプトロンでは2層目の重みしか更新できない

→ $H\left(\sum_{i=1}^2 w_i^{(2)} H(w_i^{(1)} x + b_i^{(1)}) + b^{(2)}\right)$ の H が微分可能であれば勾配法で更新可能

参考:

勾配法による経験損失の最小化



最急降下法 (steepest descent) の更新則: $\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \hat{R}(\theta_t)$

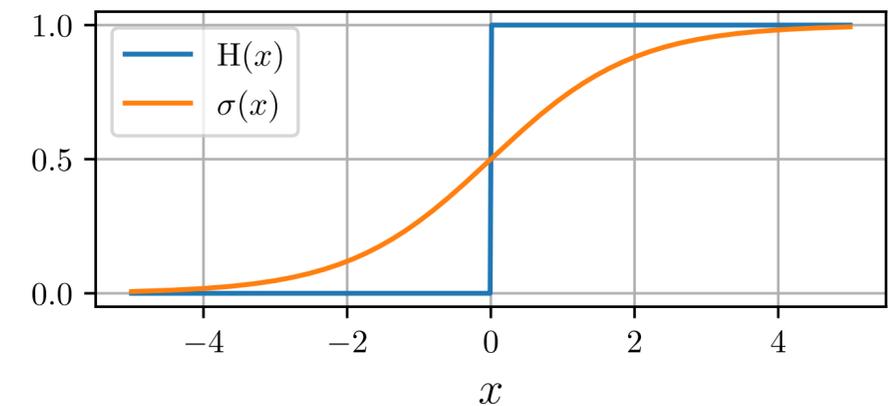
- ▶ 負の勾配 $-\nabla_{\theta} \hat{R}(\theta_t)$ が経験損失の減少方向を表す
- ▶ 学習率 η が更新の大きさを決める

ここでは θ は $w_1^{(1)}, w_2^{(1)}, w^{(2)}$ に対応

誤差逆伝搬法 (backpropagation)

[Amari 67, Rumelhart+86]

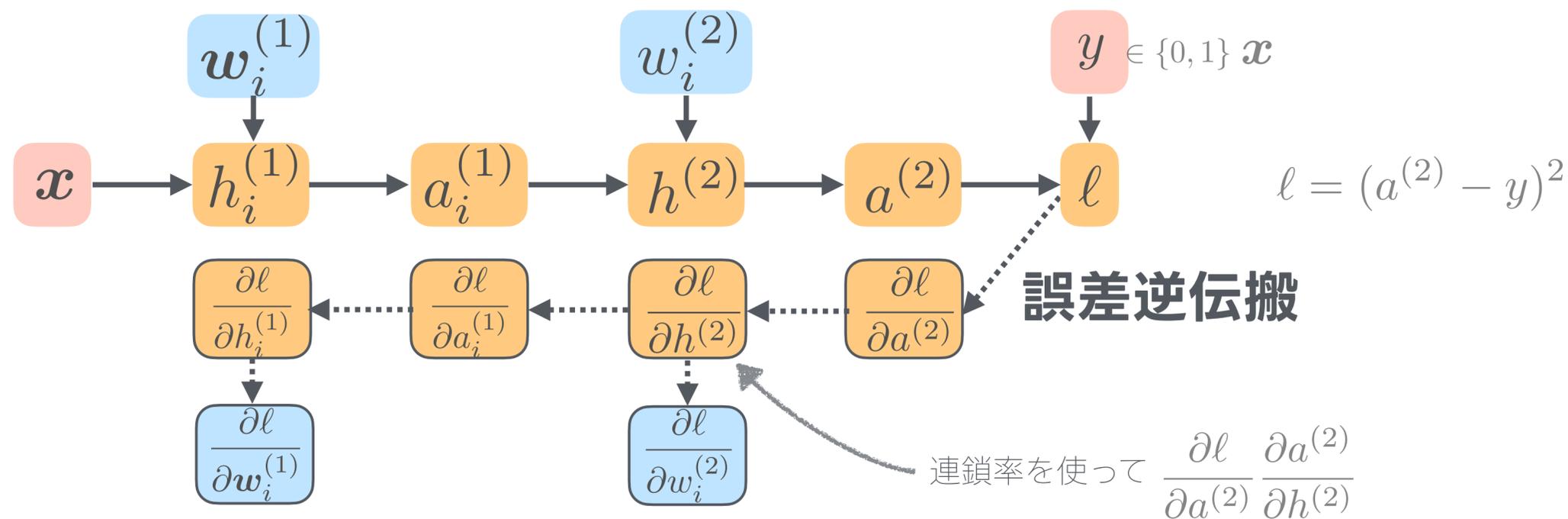
たとえば H をシグモイド関数 $\sigma(x) = \frac{1}{1 + \exp(-x)}$



に置き換え $a_i^{(1)} = \sigma(h_i^{(1)})$, $h_i^{(1)} = \mathbf{w}_i^{(1)\top} \mathbf{x} + b_i^{(1)}$

$$a^{(2)} = \sigma(h^{(2)}), \quad h^{(2)} = \sum_{i=1}^2 w_i^{(2)} a_i^{(1)} + b^{(2)}$$

$$\sigma\left(\sum_{i=1}^2 w_i^{(2)} \sigma(\mathbf{w}_i^{(1)\top} \mathbf{x} + b_i^{(1)}) + b^{(2)}\right)$$



→ 多層のニューラルネットも学習可能に 第二次ニューラルネットブームのきっかけ

ここまでの一般化

1層目 $a_i^{(1)} = \sigma(h_i^{(1)}), \quad h_i^{(1)} = w_i^{(1)\top} \mathbf{x} + b_i^{(1)}$

2層目 $a^{(2)} = \sigma(h^{(2)}), \quad h^{(2)} = \sum_{i=1}^2 w_i^{(2)} a_i^{(1)} + b^{(2)}$

1層目 $a^{(1)} = \eta(h^{(1)}), \quad h^{(1)} = \mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}$

ベクトルの要素ごとに施す
非線形関数 (例: シグモイド関数)
活性化関数 と呼ばれる

$h^{(1)} \in \mathbb{R}^{d_1}, \quad \mathbf{x} \in \mathbb{R}^{d_0}, \quad \mathbf{W}^{(1)} \in \mathbb{R}^{d_1 \times d_0}, \quad \mathbf{b}^{(1)} \in \mathbb{R}^{d_1}$
重み行列 **バイアス**
 まとめて **パラメータ**

2層目 $a^{(2)} = \eta(h^{(2)}), \quad h^{(2)} = \mathbf{W}^{(2)} a^{(1)} + \mathbf{b}^{(2)}$

\vdots
 $\mathbf{W}^{(2)} \in \mathbb{R}^{d_2 \times d_1}, \quad \mathbf{b}^{(2)} \in \mathbb{R}^{d_2}$

l層目 $a^{(l)} = \eta(h^{(l)}), \quad h^{(l)} = \mathbf{W}^{(l)} a^{(l-1)} + \mathbf{b}^{(l)}$

\vdots
 $\mathbf{W}^{(l)} \in \mathbb{R}^{d_l \times d_{l-1}}, \quad \mathbf{b}^{(l)} \in \mathbb{R}^{d_l}$

L層目 $h^{(L)} = \mathbf{W}^{(L)} a^{(L-1)} + \mathbf{b}^{(L)}$

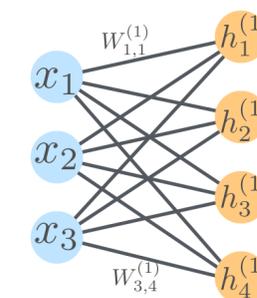
$\mathbf{W}^{(L)} \in \mathbb{R}^{d_L \times d_{L-1}}, \quad \mathbf{b}^{(L)} \in \mathbb{R}^{d_L}$

重み行列が

「普通」の行列のとき、

線形層 (アフィン層)

全結合層 と呼ばれる



線形層+活性化関数を重ねたもの

多層パーセプトロン (MLP)

フィードフォワードネットワーク

と呼ばれる

ここまでの一般化

1層目 $a^{(1)} = \eta(h^{(1)}), \quad h^{(1)} = W^{(1)}x + b^{(1)}$

ベクトルの要素ごとに施す
非線形関数 (例: シグモイド関数)
活性化関数と呼ばれる

$h^{(1)} \in \mathbb{R}^{d_1}$ $x \in \mathbb{R}^{d_0}$ $W^{(1)} \in \mathbb{R}^{d_1 \times d_0}, b^{(1)} \in \mathbb{R}^{d_1}$
重み行列 バイアス
まとめてパラメータ

2層目 $a^{(2)} = \eta(h^{(2)}), \quad h^{(2)} = W^{(2)}a^{(1)} + b^{(2)}$

$W^{(2)} \in \mathbb{R}^{d_2 \times d_1}, b^{(2)} \in \mathbb{R}^{d_2}$

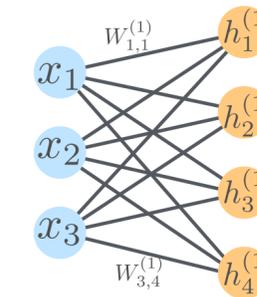
l 層目 $a^{(l)} = \eta(h^{(l)}), \quad h^{(l)} = W^{(l)}a^{(l-1)} + b^{(l)}$

$W^{(l)} \in \mathbb{R}^{d_l \times d_{l-1}}, b^{(l)} \in \mathbb{R}^{d_l}$

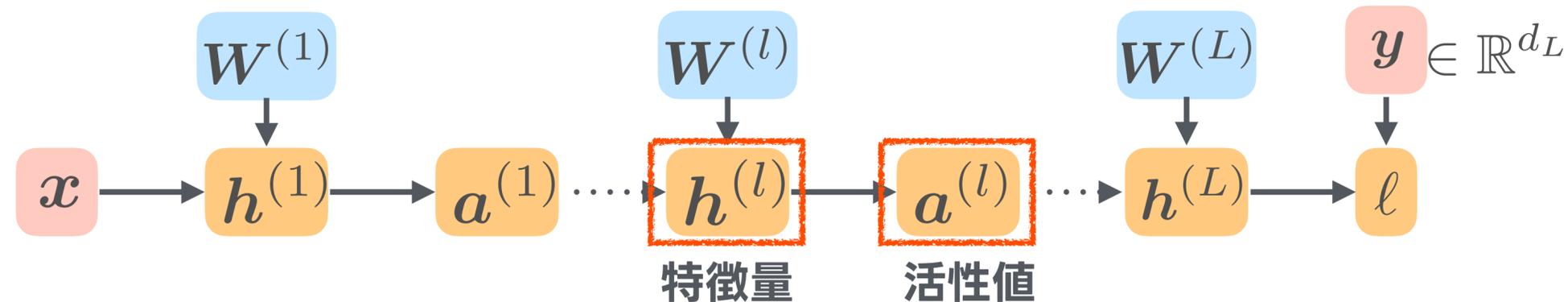
L 層目 $h^{(L)} = W^{(L)}a^{(L-1)} + b^{(L)}$

$W^{(L)} \in \mathbb{R}^{d_L \times d_{L-1}}, b^{(L)} \in \mathbb{R}^{d_L}$

重み行列・バイアスが
普通の行列のとき,
線形層 (アフィン層)
全結合層と呼ばれる



線形層+活性化関数を重ねたもの
多層パーセプトロン (MLP)
フィードフォワードネットワーク
と呼ばれる

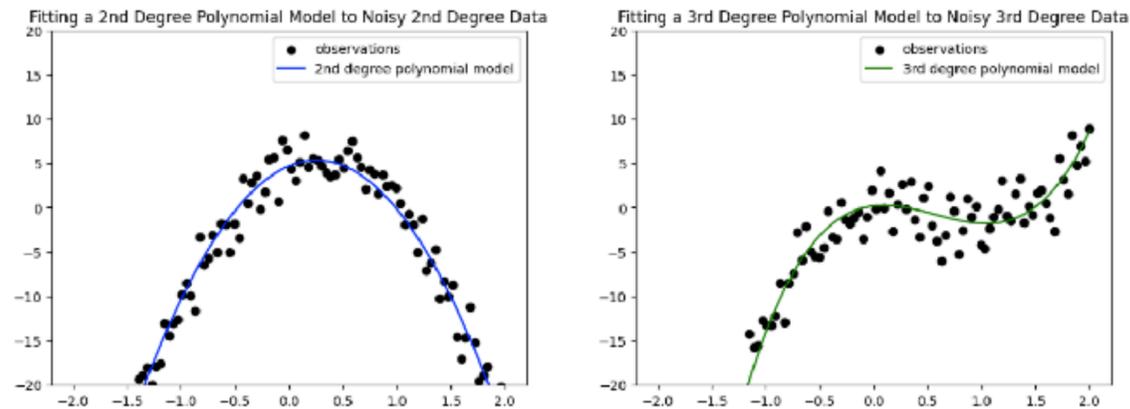


誤差逆伝播が可能

基底関数モデルとの関係

参考：

線形基底関数モデルの効果



$$\phi(x) = (1, x, x^2)$$

を用いた線形モデル

▶ 入力に関して非線形なモデルになるので、複雑なデータも説明

$$\phi(x) = (1, x, x^2, x^3)$$

を用いた線形モデル

松井先生の前回の講義

基底関数モデル：

$$f(x) = \sum_{i=1}^M w_i \phi_i(x)$$

基底関数は**固定**

2層ニューラルネット：

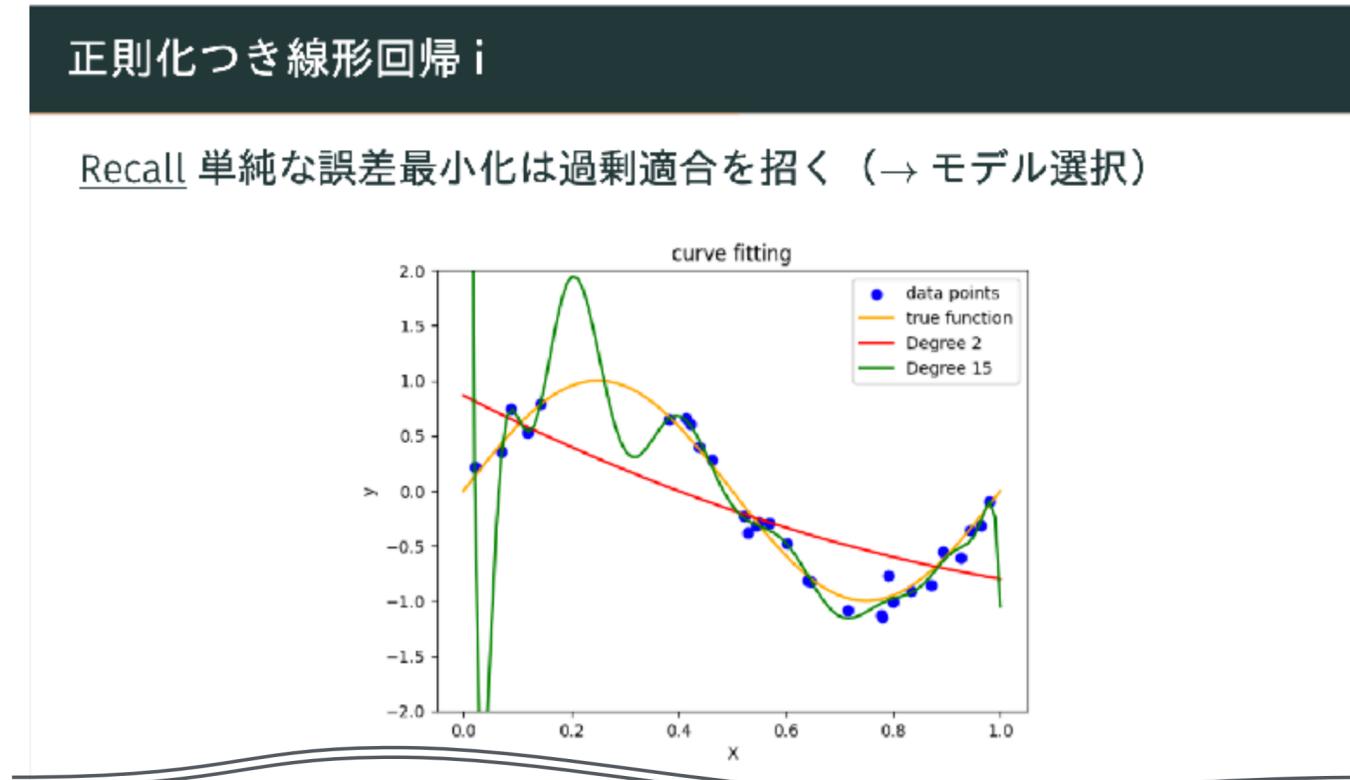
$$f(x) = \sum_{i=1}^M w_i^{(2)} \eta(w_i^{(1)\top} x + b_i^{(1)})$$

基底関数がデータに**適応可能**

重み共有によるパラメータ削減

パラメータが多いと過剰適合が生じる可能性がある

参考：



松井先生の前回の講義

ニューラルネットはパラメータが多いため過剰適合しやすい

容易にデータセットを「丸暗記」してしまう

→ データの構造を利用し重みを共有することでパラメータ数を削減

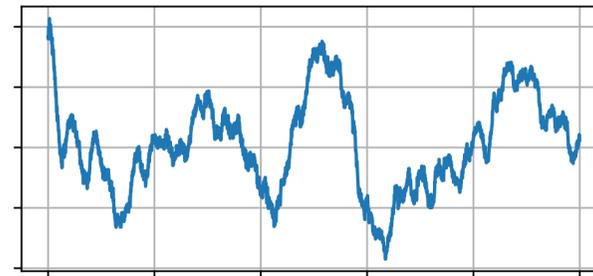
$$f(\mathbf{x}) = \sum_{i=1}^M w_i \phi_i(\mathbf{x})$$
$$f(\mathbf{x}) = \sum_{i=1}^M w_i^{(2)} \eta(\mathbf{w}_i^{(1)\top} \mathbf{x} + b_i^{(1)})$$

畳み込みニューラルネット [LeCun+89]

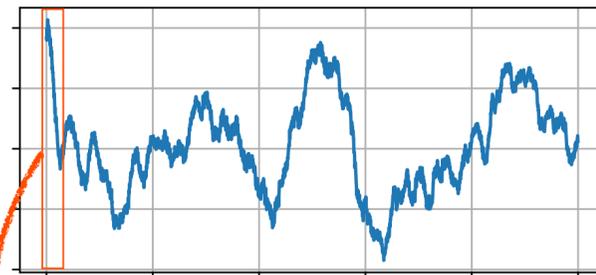
画像：



音声：

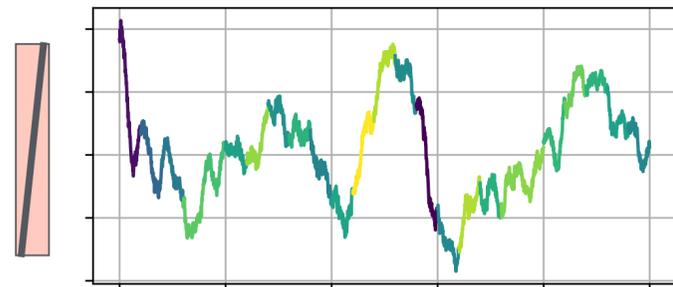
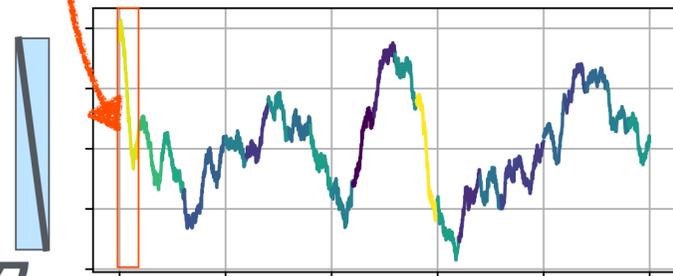


畳み込み

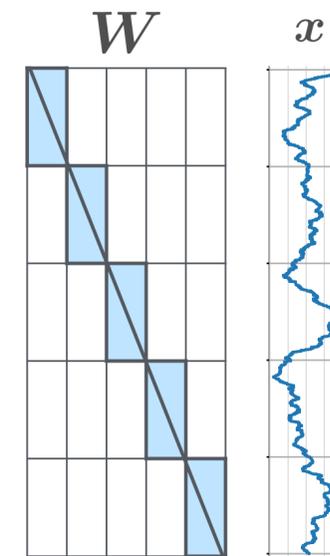


内積

フィルター



行列積で表現可能



共有された \downarrow のみ学習
(ほかはゼロ)

→ パラメータ削減

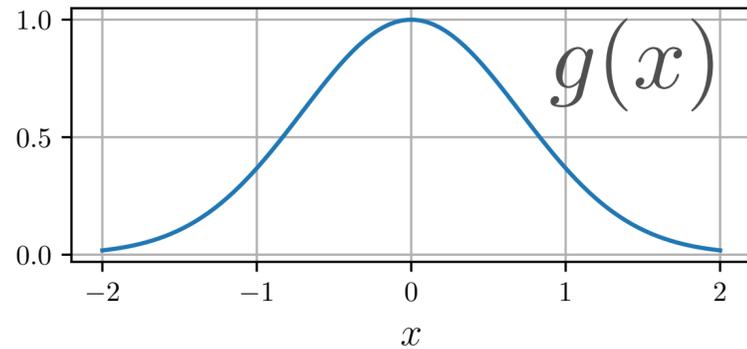
(データの構造を捉えているので
学習もしやすい)

(詳細は次回以降)

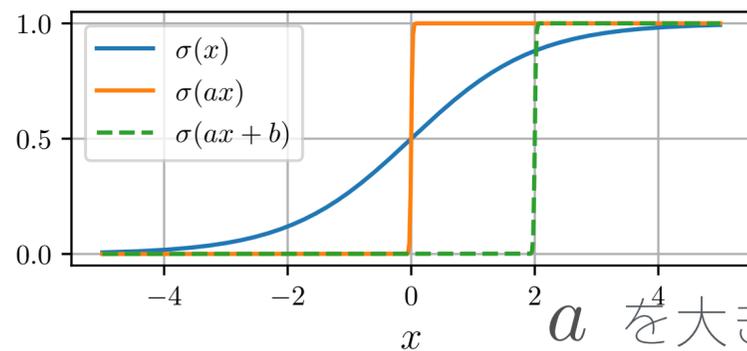
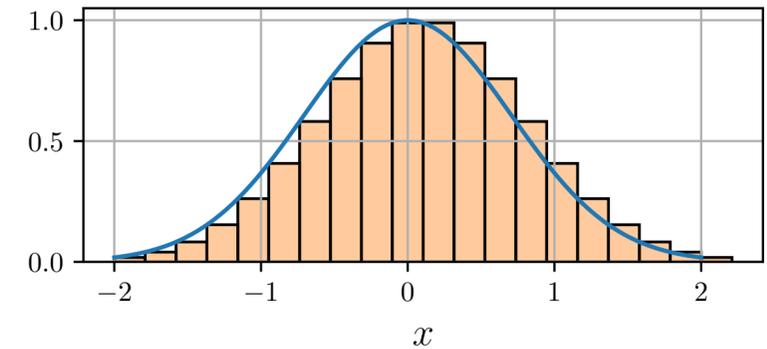
ニューラルネットの万能近似性 [Cybenko89] etc.,

1変数入力1変数出力の2層ニューラルネット $f(x) = \sum_{i=1}^M w_i^{(2)} \sigma(w_i^{(1)} x + b_i^{(1)})$ は
 任意の1変数関数 $g(x)$ を近似可能

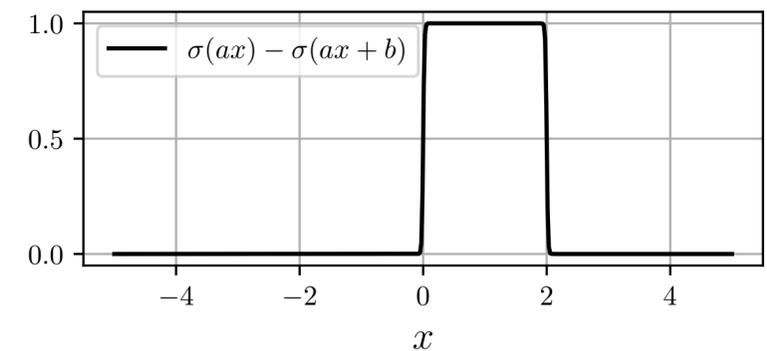
証明の方針



$g(x)$ は区分定数関数
 によって近似可能



区分定数関数はシグモイド関数
 2つで近似可能



a を大きく取る

区分定数関数の高さを決める

$$f(x) = \sum_{i=1}^M w_i^{(2)} \sigma(w_i^{(1)} x + b_i^{(1)})$$

区分定数関数を作る

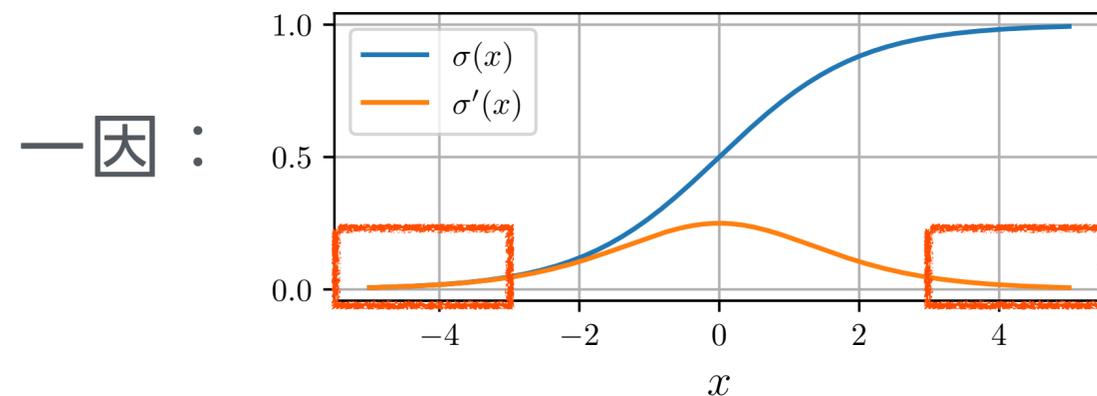
任意の1変数関数
 $g(x)$ を近似可能

ニューラルネットの層を増やす

ニューラルネットの層を増やす（深くする）と
データを階層的に処理できるようになる

より複雑なタスクを行うためには深くする必要

深くすると誤差逆伝搬した勾配が消失・発散して
学習が失敗してしまう

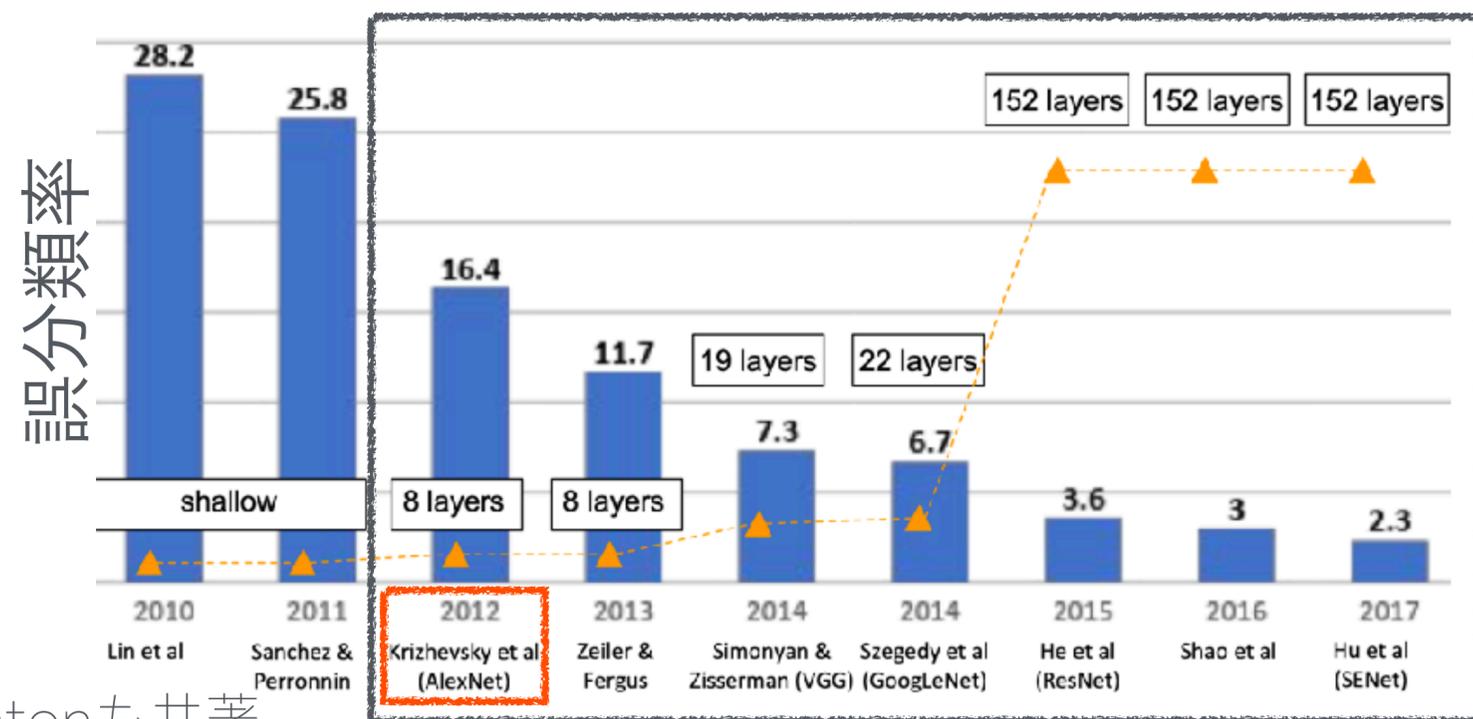


目次

- 深層学習以前
- **深層学習（深層学習を可能にする技術）**
- 深層学習の実装

深層学習の勃興

ImageNetデータセットの
1000カテゴリ 120万枚の画像を分類するコンペ



Hintonも共著



深層畳み込みニューラルネットワーク

AlexNetの成功理由

活性化関数をReLUへの変更

過剰適合を抑制する「ドロップアウト」の導入

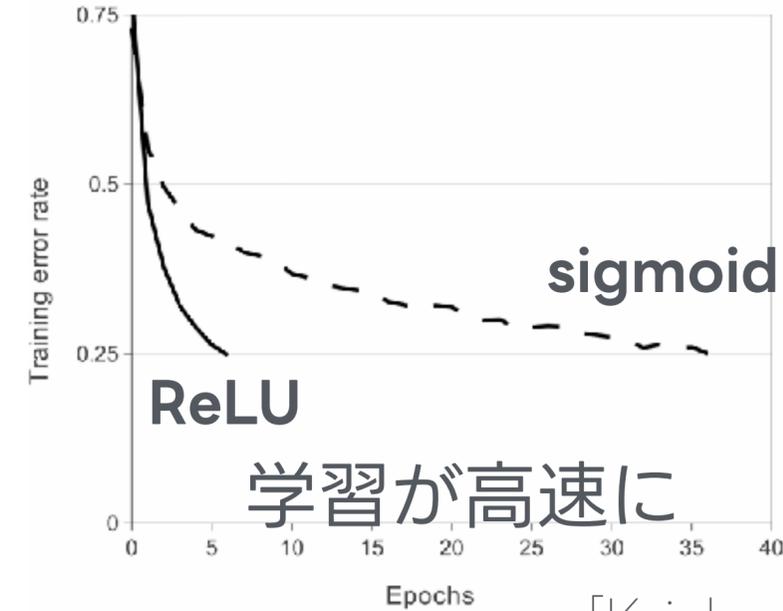
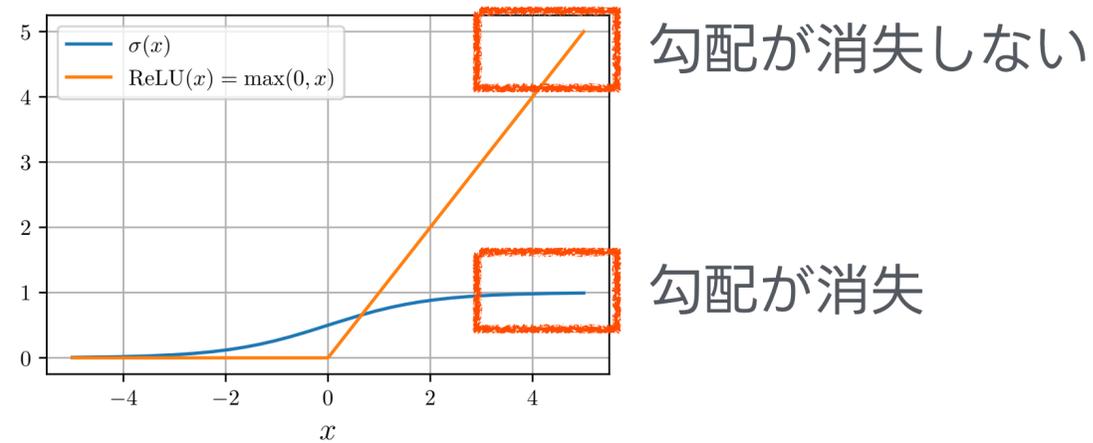
データ量の増加+計算機の進化 (GPU)

ImageNet中の画像



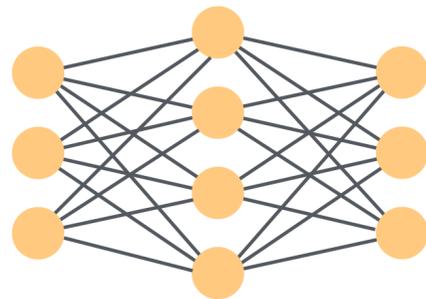
AlexNetの成功理由

ReLU活性化関数 $\text{ReLU}(x) = \max(0, x)$



[Krizhevsky+12]

ドロップアウト



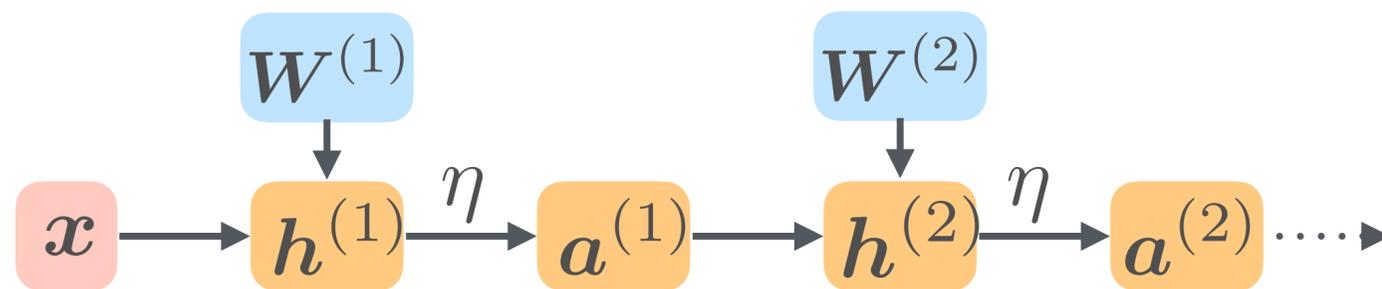
学習時にランダムに各層の出力ベクトルの要素を0に
→ 些細な特徴を学習しにくくなる (過剰適合の抑制)

データ量の増加+計算機の進化 (GPU)

インターネットの普及によりデータの収集が容易に
深層学習に現れる行列積の並列計算がGPUにより高速に

パラメータの初期化

パラメータはランダムに初期化



$h^{(1)}, h^{(2)}, \dots$, の分散が大きく変化しないように初期化したい

Kaiming法 (He法とも) 活性化関数がReLUのとき用いる

$$W_{ij}^{(l)} \sim \mathcal{N}\left(0, \frac{2}{d_{l-1}}\right) \text{ あるいは } \mathcal{U}\left(-\sqrt{\frac{6}{d_{l-1}}}, \sqrt{\frac{6}{d_{l-1}}}\right)$$

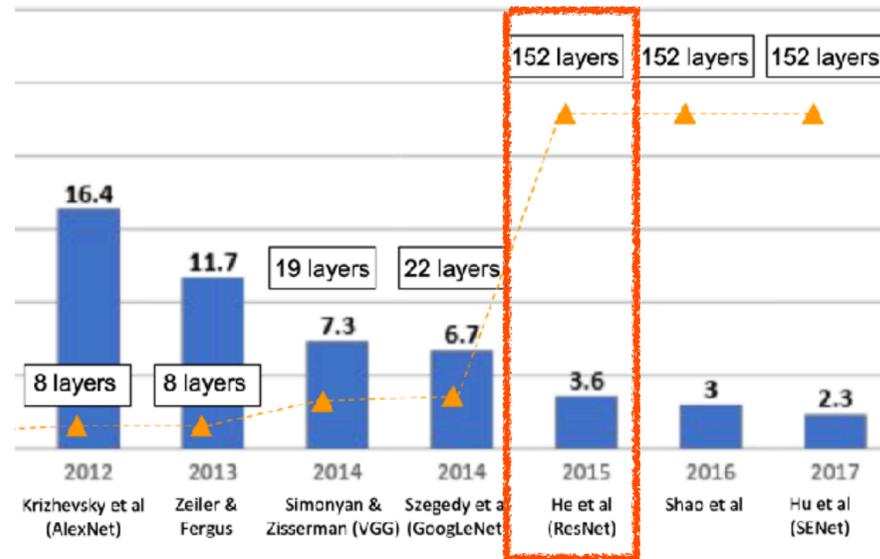
d_{l-1} : 層への入力ベクトルの次元

実際
$$h_j^{(l)} = \sum_{i=1}^{d_{l-1}} W_{ji}^{(l)} a_i^{(l-1)}$$

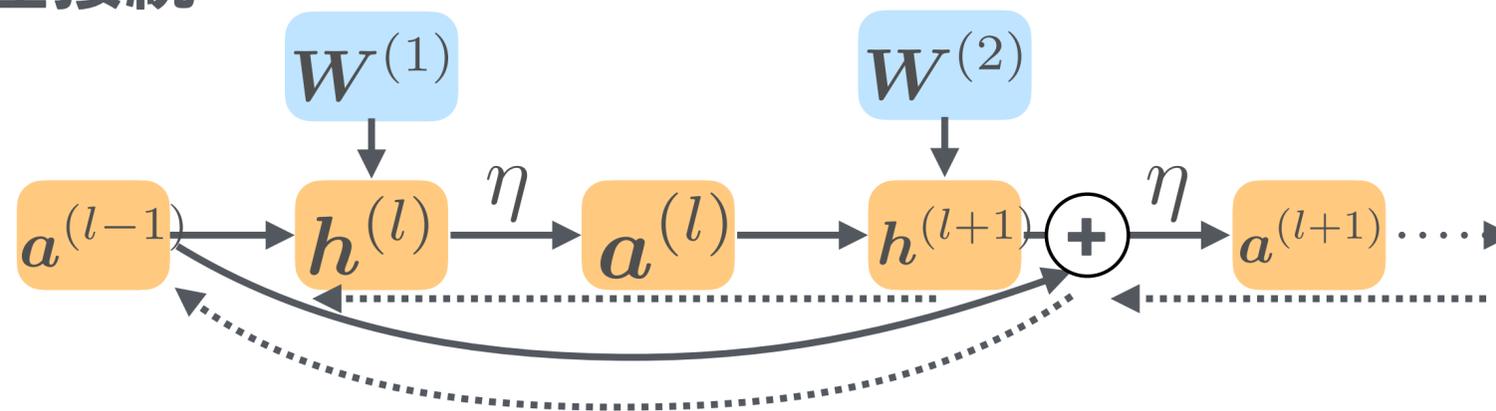
$$\begin{aligned} \mathbb{V}(h_j^{(l)}) &= d_{l-1} \mathbb{V}(W_{ji}^{(l)} a_i^{(l-1)}) \\ &= d_{l-1} \mathbb{E}(W_{ji}^{(l)2}) \mathbb{E}(a_i^{(l-1)2}) \\ &= d_{l-1} \mathbb{V}(W_{ji}^{(l)}) \frac{\mathbb{V}(h_i^{(l-1)})}{2} \end{aligned} \left. \begin{array}{l} \right\} W_{ji}^{(l)}, a_i^{(l)} \text{ は独立, } \mathbb{V}(x) = \mathbb{E}(x^2) - \mathbb{E}(x)^2 \\ \left. \right\} \text{ReLUなので右側だけ}$$

ResNet

ResNetは100層以上の深さを実現

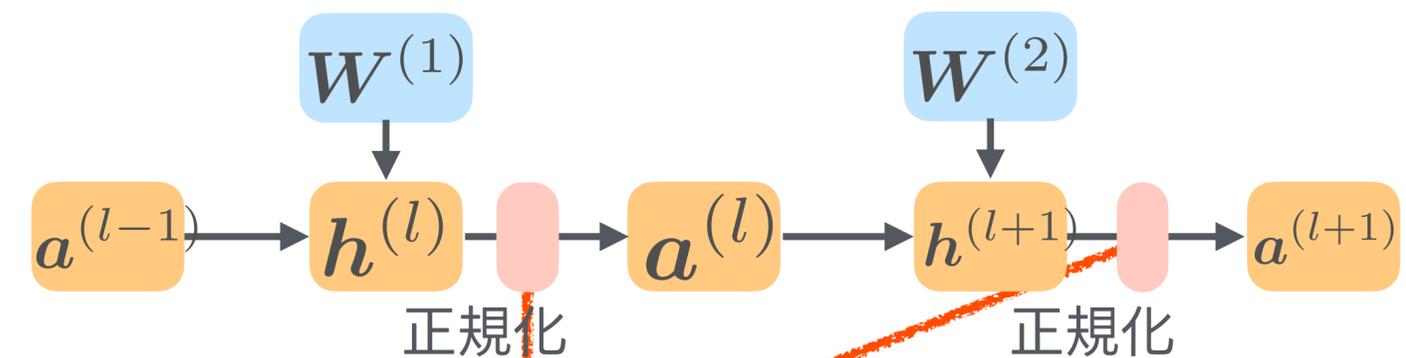


残差接続



勾配の消失を防ぐ
→ 100~1000層も可能に

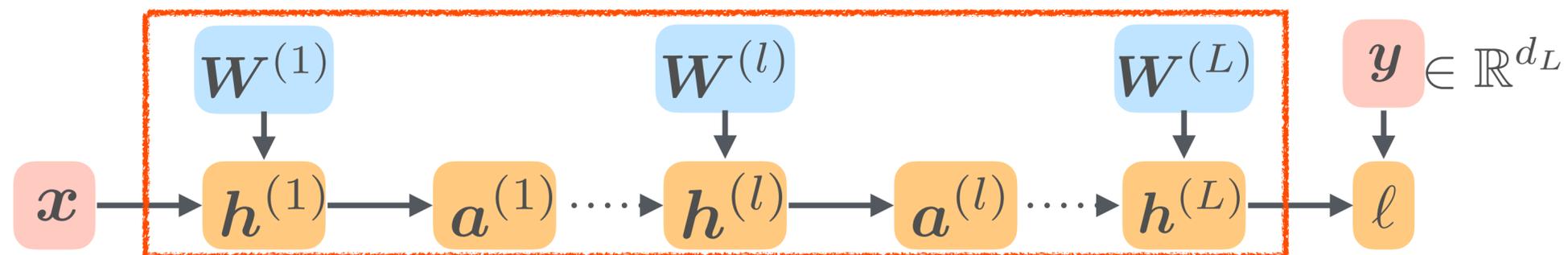
特徴量の正規化



同じ平均・分散
→ 学習の安定

→ 残差接続・特徴量の正規化は最新のニューラルネットでも利用されている技術

確率的勾配降下法 (SGD)



$\theta = (W^{(1)}, b^{(1)}, \dots, W^{(L)}, b^{(L)})$ として $f_\theta : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{d_L}$ と表記

深層学習では一般的

バッチ学習

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \frac{1}{N} \sum_{\substack{(\mathbf{x}_n, \mathbf{y}_n) \\ \in \mathcal{D}}} L(f_{\theta}(\mathbf{x}_n), \mathbf{y}_n)$$

(ミニバッチを用いた)
確率的勾配降下法

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \frac{1}{M} \sum_{\substack{(\mathbf{x}_n, \mathbf{y}_n) \\ \in \mathcal{B}}} L(f_{\theta}(\mathbf{x}_n), \mathbf{y}_n)$$

$\mathcal{B} \subset \mathcal{D}, |\mathcal{B}| = M$ は大きさ M のミニバッチ
でランダムに選択

確率的勾配降下法

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} L(f_{\theta}(\mathbf{x}_n), \mathbf{y}_n)$$

n は毎回ランダムに選択

全データにアクセスせずに済む
並列化しやすい
ランダム性が正則化に

α : 学習率

$\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$: データセット

確率的勾配降下法 (SGD)

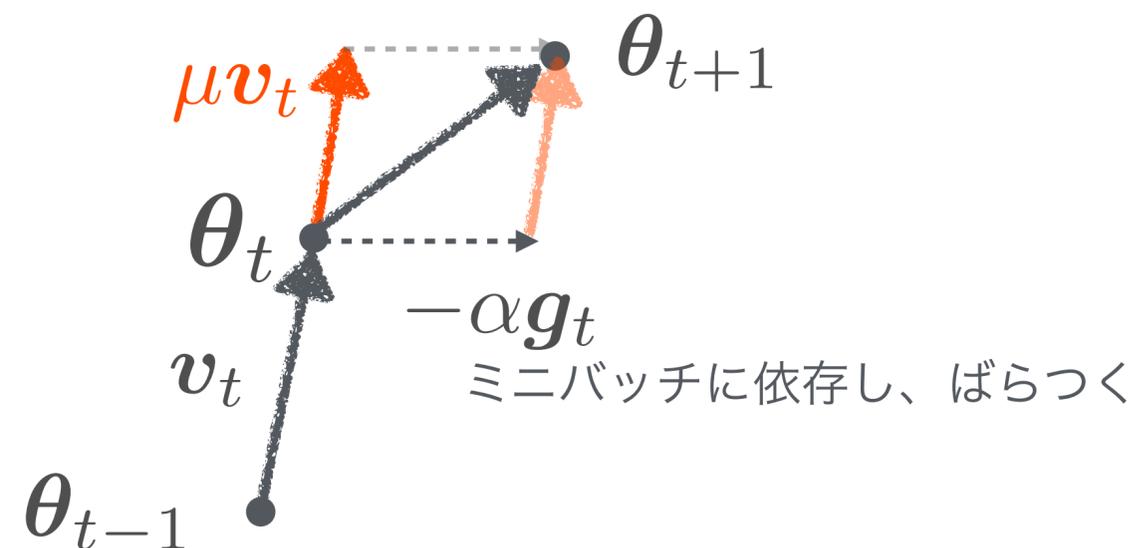
通常SGDは以下とともに用いる $g_t = \nabla_{\theta} \frac{1}{M} \sum_{\substack{(\mathbf{x}_n, \mathbf{y}_n) \\ \in \mathcal{B}}} L(f_{\theta_t}(\mathbf{x}_n), \mathbf{y}_n)$ とする

モメンタム

$$\theta_{t+1} = \theta_t - \alpha g_t + \mu v_t$$
$$v_t = \theta_t - \theta_{t-1}$$

SGDを安定化させ収束性を改善する補正
実質的な学習率を大きくする効果も

$\mu = 0.9$ 程度に設定される



重み減衰

$$\theta_{t+1} = \theta_t - \alpha(g_t + \lambda \theta_t)$$

以下の正則化付きの損失を用いていることと等価

$$L'(f_{\theta}(\mathbf{x}_n), \mathbf{y}_n) = L(f_{\theta}(\mathbf{x}_n), \mathbf{y}_n) + \frac{\lambda}{2} \|\theta\|^2$$

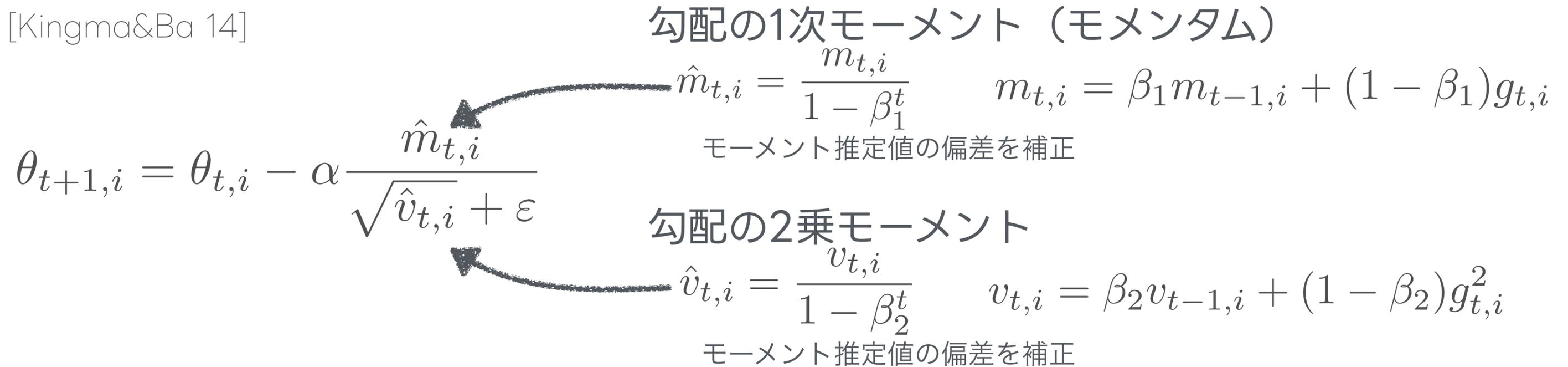
より発展的な確率的勾配降下法

Adam [Kingma&Ba 14]

$$\theta_{t+1,i} = \theta_{t,i} - \alpha \frac{\hat{m}_{t,i}}{\sqrt{\hat{v}_{t,i} + \epsilon}}$$

勾配の1次モーメント (モメンタム)
 $\hat{m}_{t,i} = \frac{m_{t,i}}{1 - \beta_1^t}$ $m_{t,i} = \beta_1 m_{t-1,i} + (1 - \beta_1) g_{t,i}$
モーメント推定値の偏差を補正

勾配の2乗モーメント
 $\hat{v}_{t,i} = \frac{v_{t,i}}{1 - \beta_2^t}$ $v_{t,i} = \beta_2 v_{t-1,i} + (1 - \beta_2) g_{t,i}^2$
モーメント推定値の偏差を補正



これまで勾配が大きかった成分 i は控えめに更新
小さかった成分 i は積極的に更新

→ 学習率をパラメータごとに適応的に調整している

AdamはSGDよりも安定した学習が可能で、広く用いられている

現代の深層学習はAdamに過剰適合しているとも…

なぜ深層学習を使うのか

(データが十分にあれば) 様々な関数を表現できる

言語モデル

ChatGPT

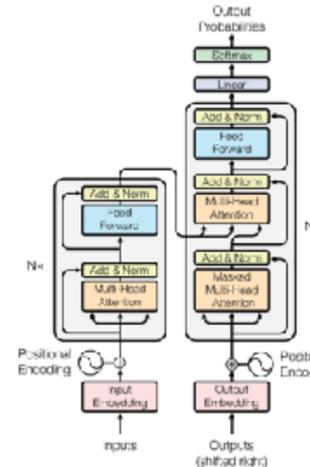
共有する

「深層学習の基礎」という講義をすることになりました。講義のシラバスを考えてください。

以下は、「深層学習の基礎」という講義に適したシラバス案です。学生が深層学習の概念を理解し、実際に簡単なモデルを構築・活用できるレベルを目指す構成になっています。

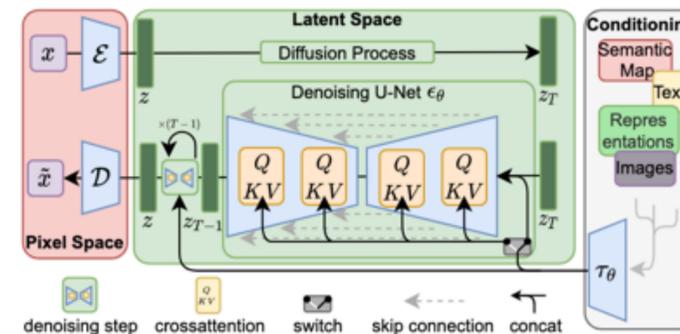
講義名

深層学習の基礎



これまでの単語列を入力、次の単語を出力とする関数

画像生成モデル



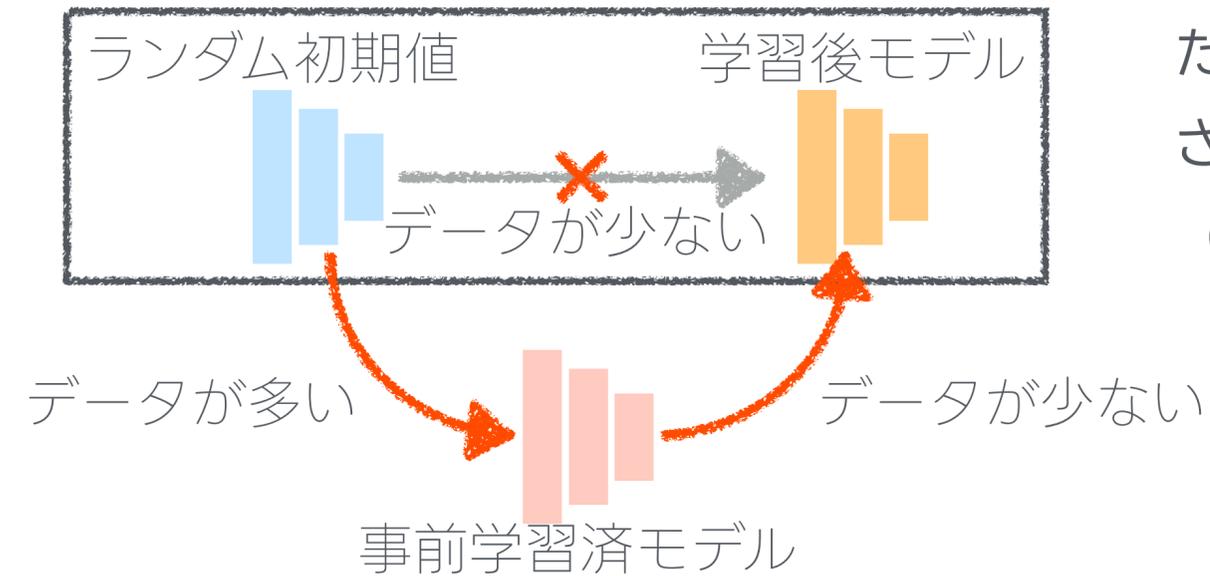
画像のノイズを少し除去する関数
(何度も繰り返してノイズから画像をつくる)

ニューラルネットは各層は簡単な計算なので大規模化しやすい
層をブロックのように組み合わせることで様々な関数をモデル化できる

なぜ深層学習を使うのか

転移学習

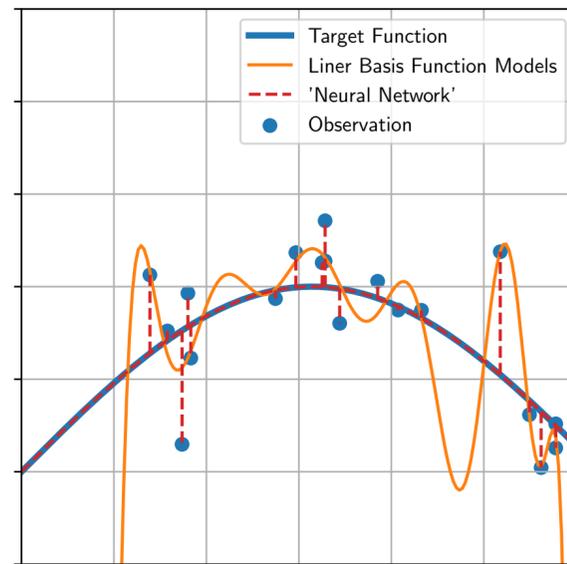
事前学習済みのパラメータを初期値とすることで、少ないデータでも学習可能



たとえばImageNetで事前学習した画像認識モデルは
さまざまな画像認識タスクに転移されている

(ただし転移学習がネガティブに働く場合もあるので注意)

良性過剰適合



ニューラルネットは学習データを「丸暗記」してもよく汎化する
→ 十分にパラメータがあるので、目標の関数をよく近似しつつ
観測ノイズはスパイクで表現している

図はイメージです

目次

- 深層学習以前
- 深層学習
- **深層学習の実装**

実装

- ウェブ上で実行できるGoogle colabで共有しました
- <https://bit.ly/3DhALFH>

文献紹介

- 岡谷貴之「深層学習」講談社, 2022.
 - 深層学習を概観する和書
- 斎藤康毅「ゼロから作るDeep Learning」2016.
 - Pythonを使って文字通りゼロから深層学習の理解を目指すシリーズ
- Zhang, Lipton, Li, and Smola, “Dive into Deep Learning,” 2023.
 - オンライン版にはPyTorchなどのコードも付属 <https://d2l.ai>
- PyTorch tutorial <https://pytorch.org/tutorials/index.html>
 - 日本語翻訳（更新は数年前に止まっているが基本的には問題ないはず） https://yutarogawa.github.io/pytorch_tutorials_jp/