

ニコニコAIスクール「脳型人工知能開発者入門コース」#9

再帰的ニューラルネットワーク (1)

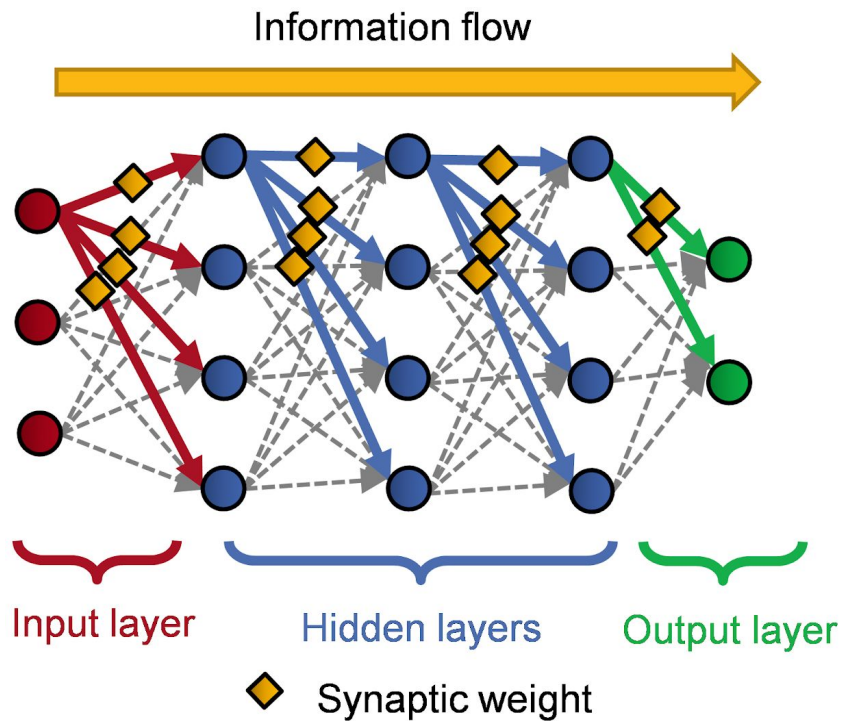
松森 匠哉

2018 / 03 / 11

- BP: Backpropagation
 - 誤差逆伝播法
- FFNN: Feedforward Neural Networks
 - 順方向にしか伝播しないネットワークの総称
- FBNN: Feedback Neural Networks
 - フィードバック構造を持つネットワークの総称
- RNN: Recurrent Neural Networks
 - FBNNの総称
- 線形回帰: Linear Regression
 - データ点に適合する関数を見出すこと

今までのネットワーク

階層型ネットワーク

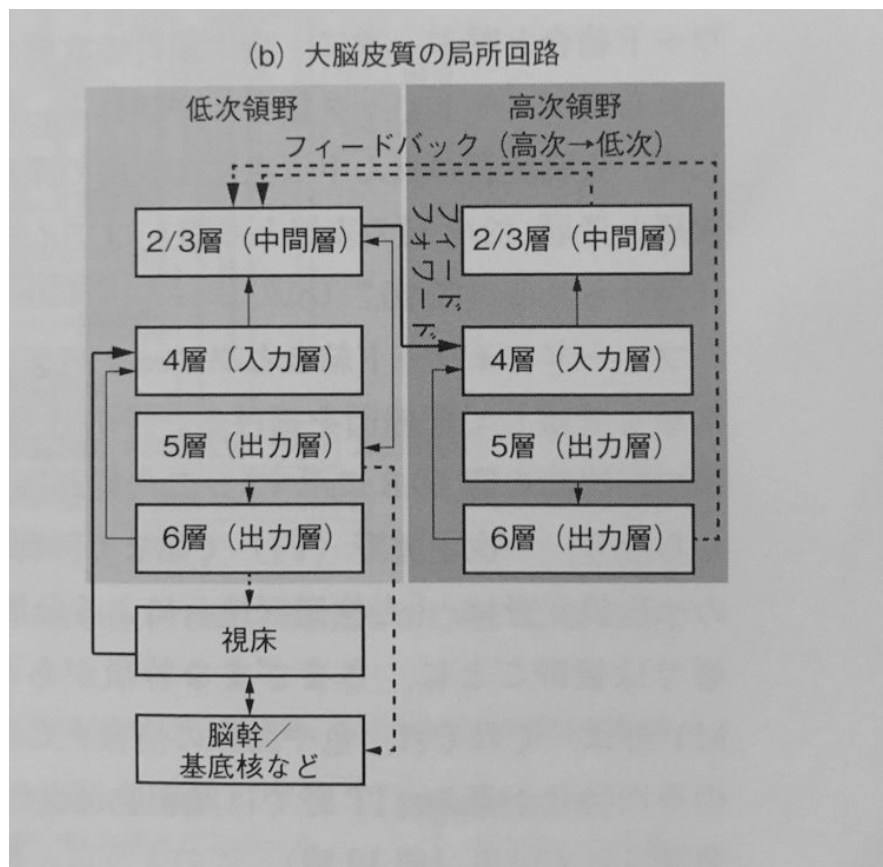


階層型NWの特徴

- 複数のニューロンがグループ化され、それぞれ層を形成している
- 各層は直列に連っており、直前の層のニューロンのみから入力を受ける
- 情報は順方向にしか伝播しない

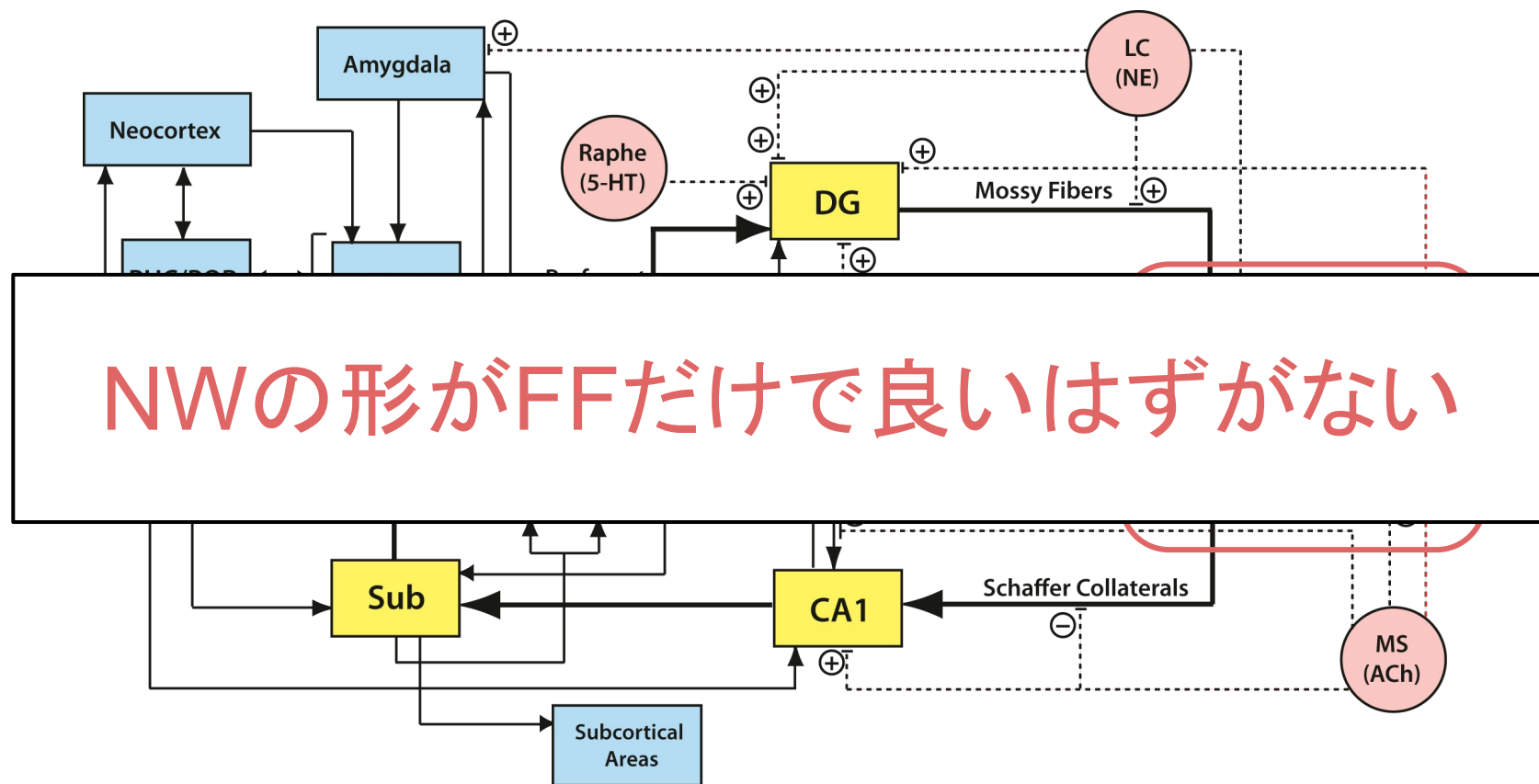
FF vs FB

解剖学的には $FF \ll FB$



FF vs FB

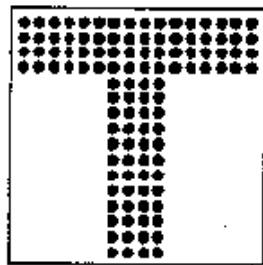
海馬のCA3に存在する再帰的神経投射



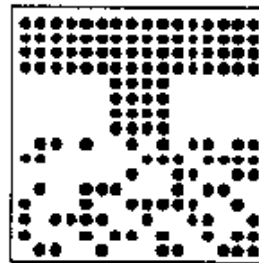
相互結合のモデル化

Hopfield Network

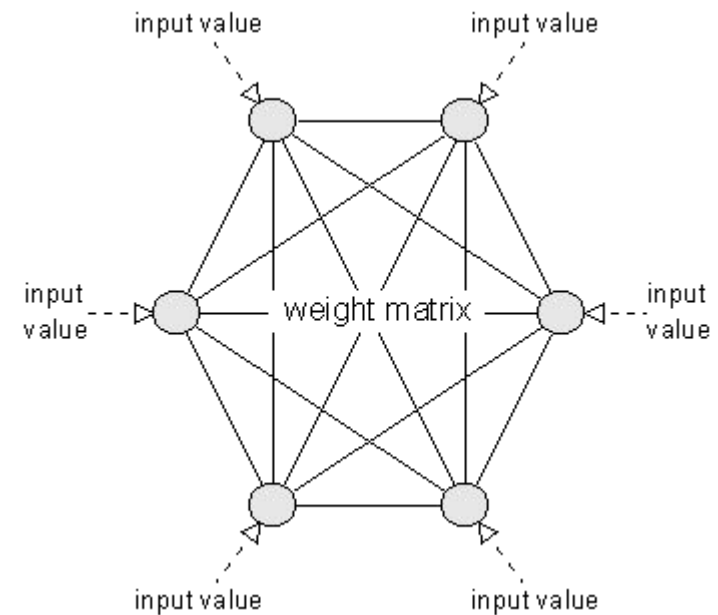
- 物理学者 J.J Hopfield (1982) が提唱
- 部分的情報からネットワークに記憶された情報を再構成することができる
(連想記憶モデル)



Original 'T'



half of image
corrupted by
noise



FF vs FB

様々なネットワークの形

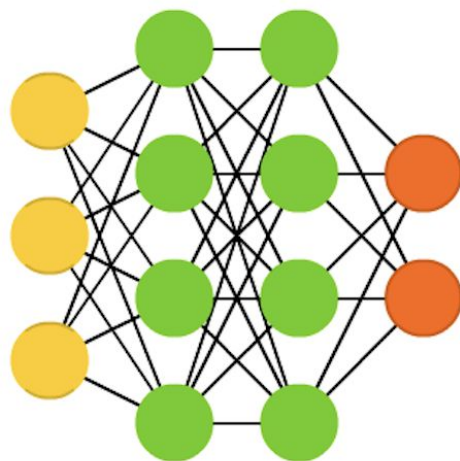
今回はFB結合を含む、様々なネットワークを学んでいきます

● Input Cell

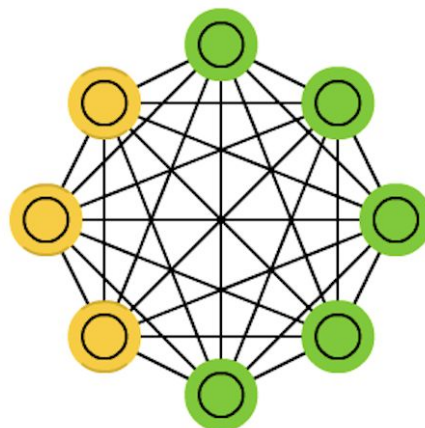
● Hidden Cell

● Output Cell

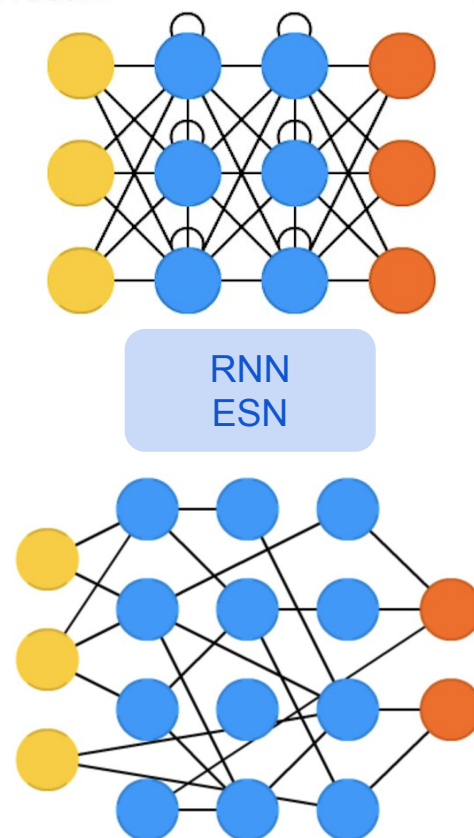
● Recurrent Cell



CNN MLP



Boltzmann Machine



RNN ESN

RNNの位置付け

Neural Network

(§5-7) Feedforward

(§5) 単純パーセプトロン

(§5) 多層パーセプトロン

(§6,7) CNN

(§9-10) Feedback

再帰型

(§9) Simple RNN

(§10) LSTM

(§10) GRU

相互結合型

(§8) RBM

(§9) HFN

(§9) ESN

(§9) LSM

目次

1. Introduction
2. RNNの基礎と学習
 - a. Feedforward Networks v.s. Feedback Networks
 - b. BPTTの解説と問題点
3. 基礎演習: Chainerによるsimple RNNの実装
4. Reservoir Computing
 - a. Echo State Network
 - b. Liquid State Machine
5. 実践演習
 - a. Echo State Networkによる神経活動データの学習

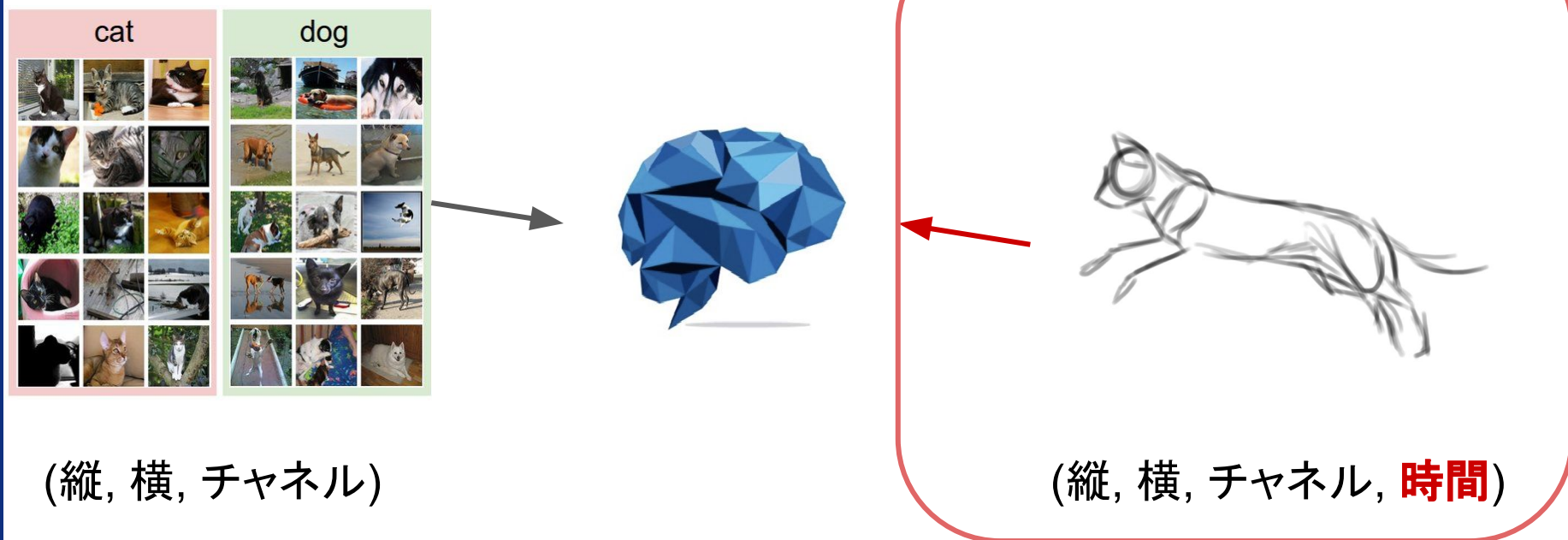
講義 Part1

RNNの基礎と学習



RNN and Back Propagation

人間は静的(static)に外界を捉えていない



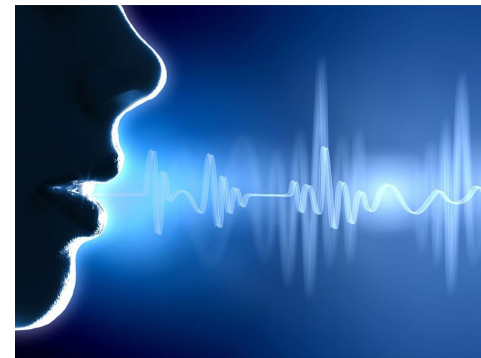
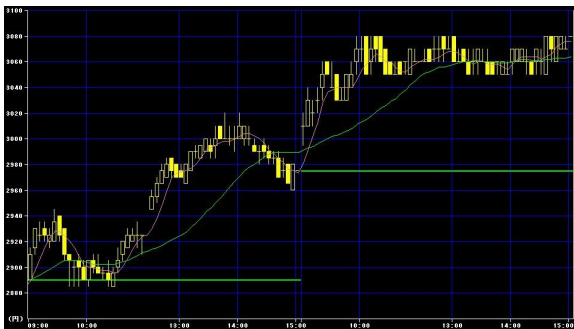
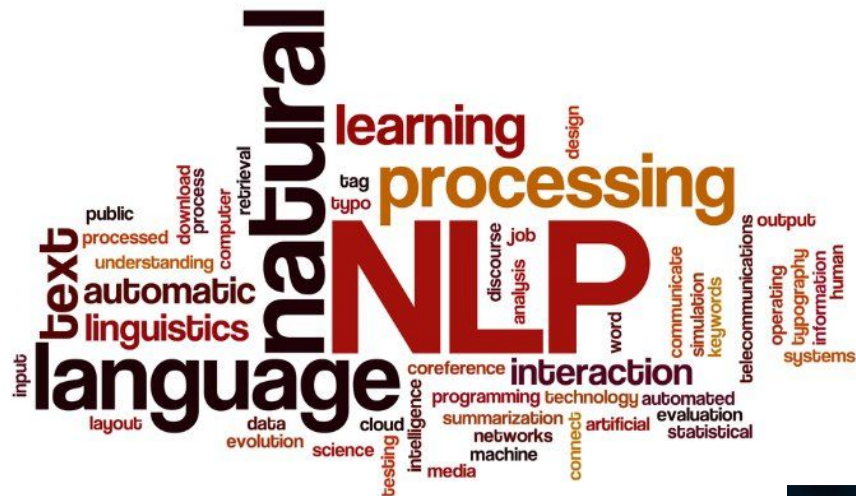
認知のプロセスには時系列が大切

[動画] 記憶と時系列



逆からいえますか？

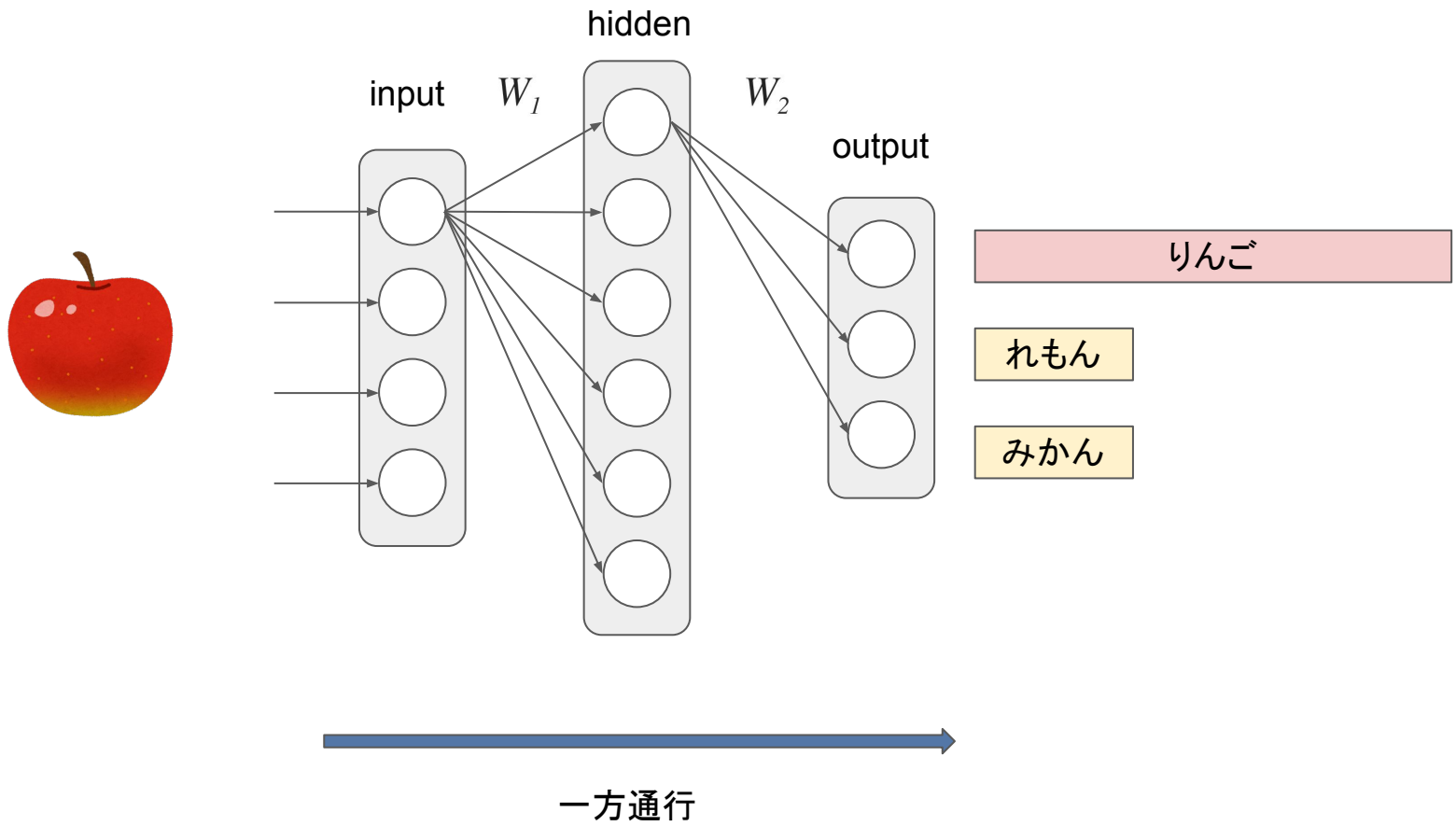
(実は)時系列だらけの世界



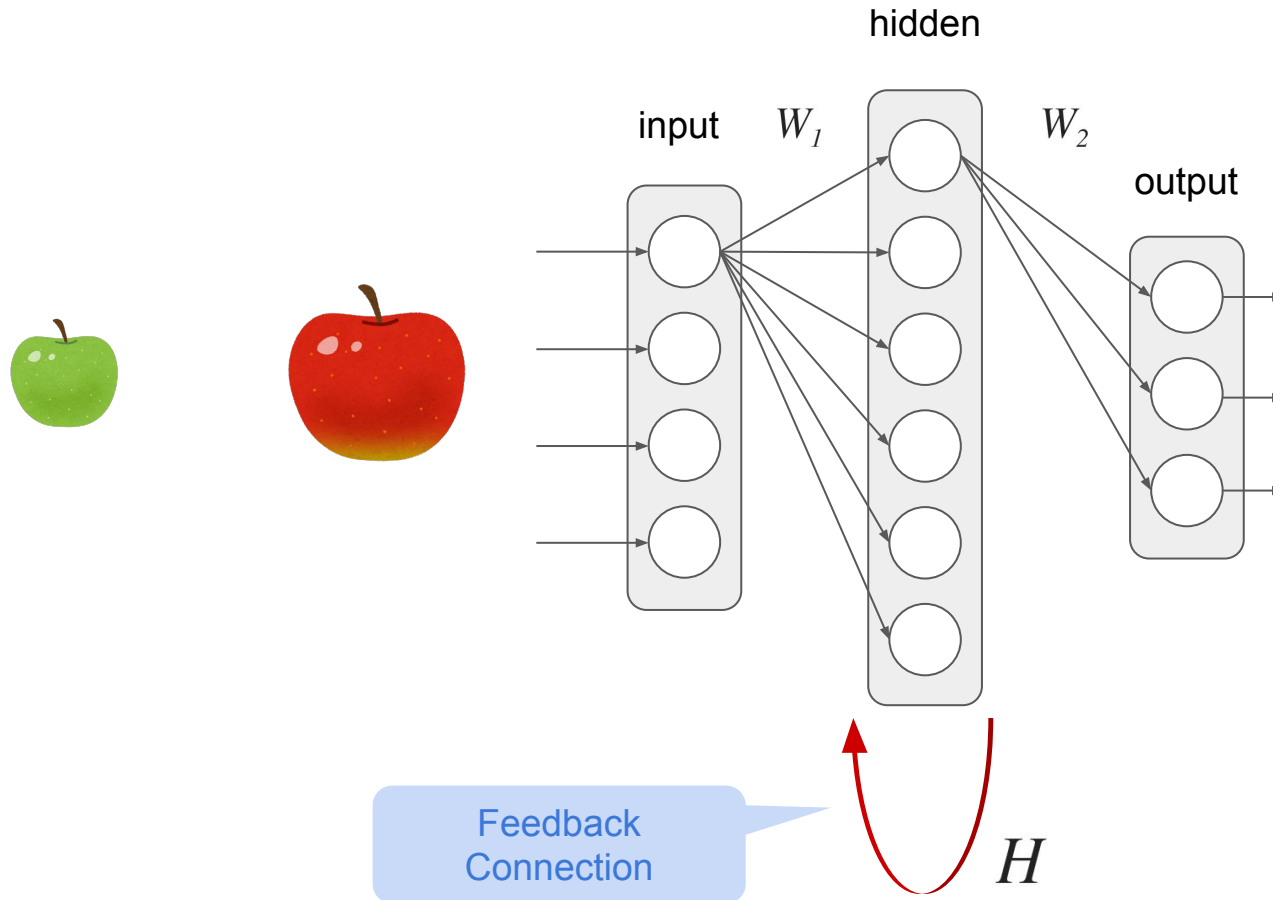
工学的に時系列を扱うことのできるネットワークは非常に重要

Feedforward Network

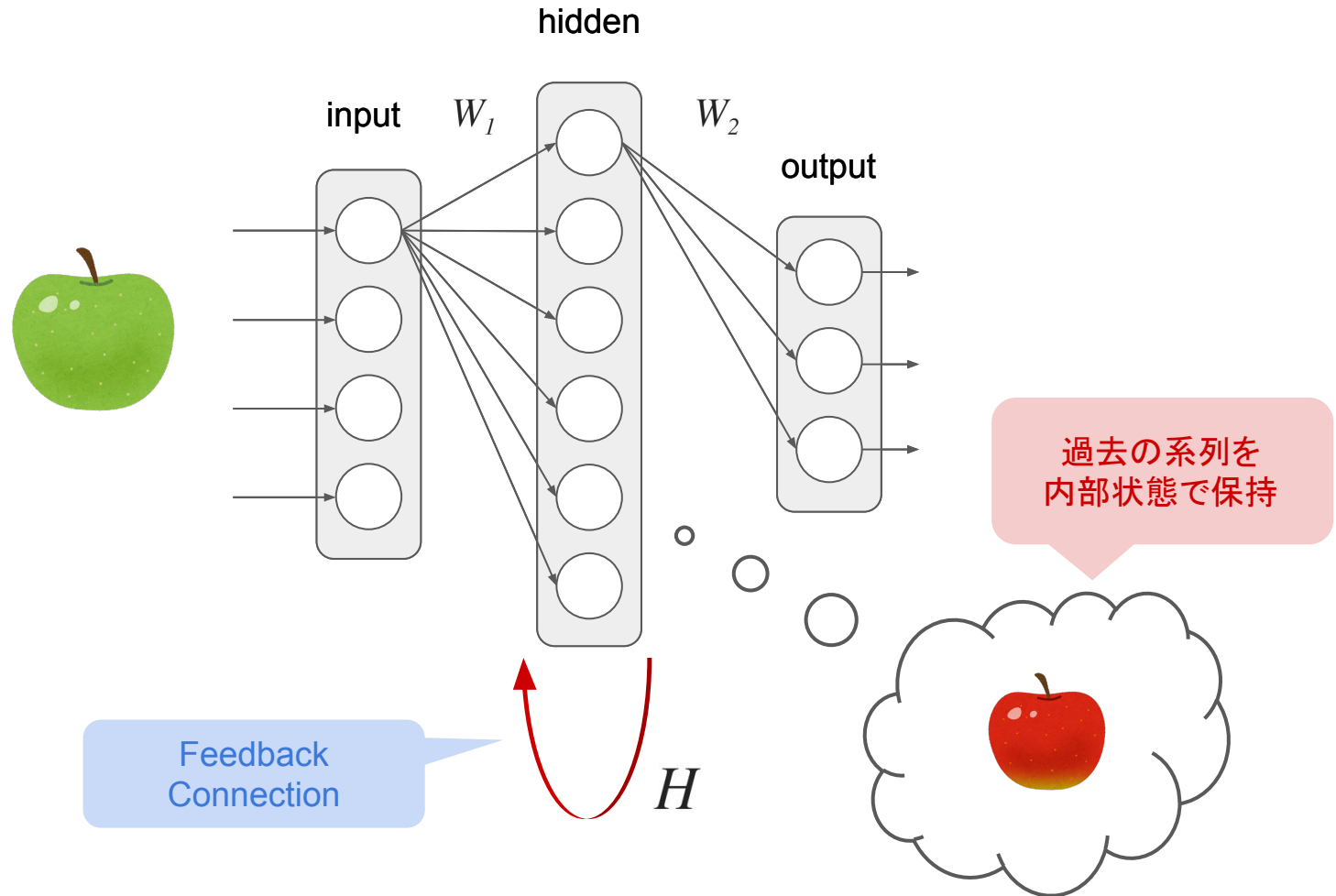
復習



Feedback Network



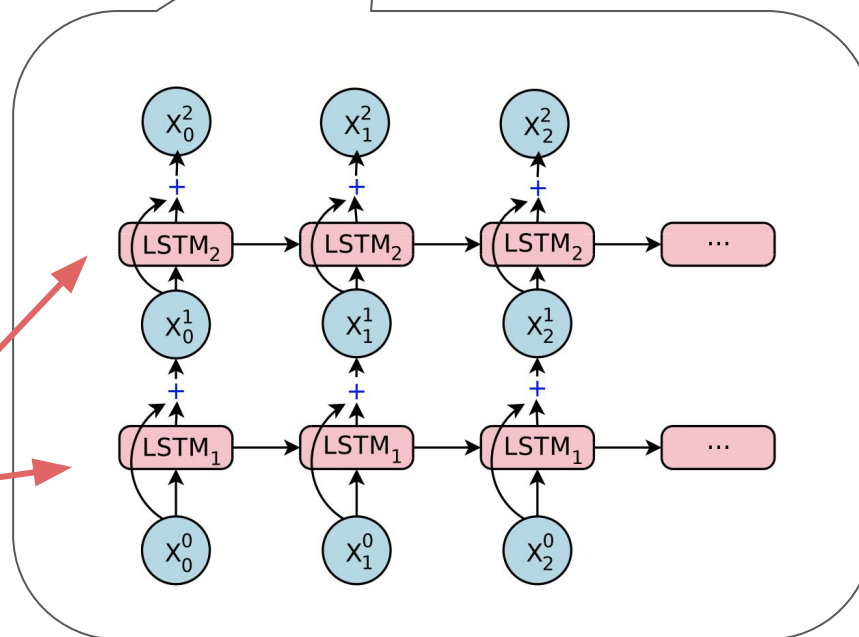
Feedback Network



実際どこに使われているの？



RNN

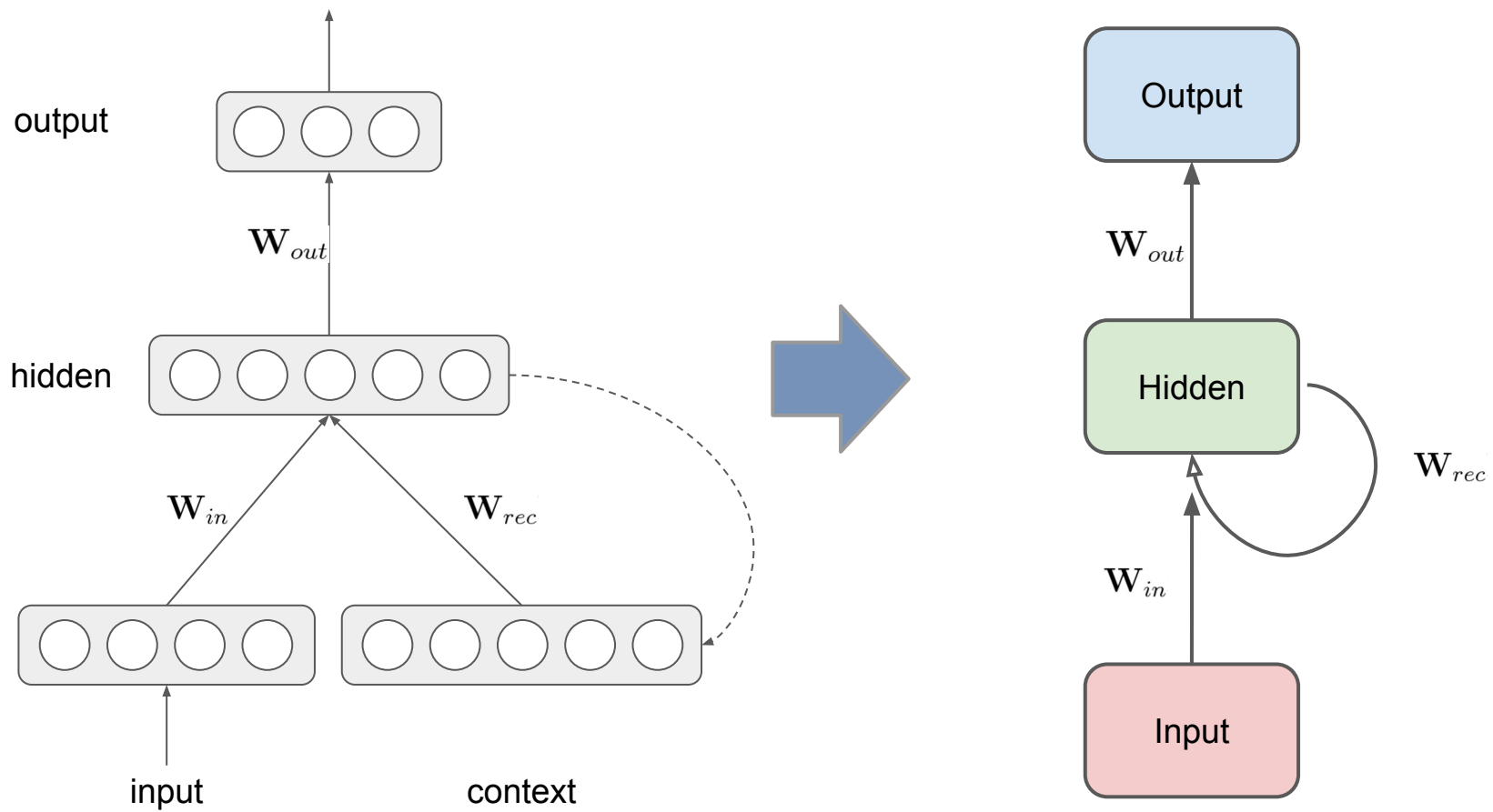


Recurrent Neural Network (RNN)

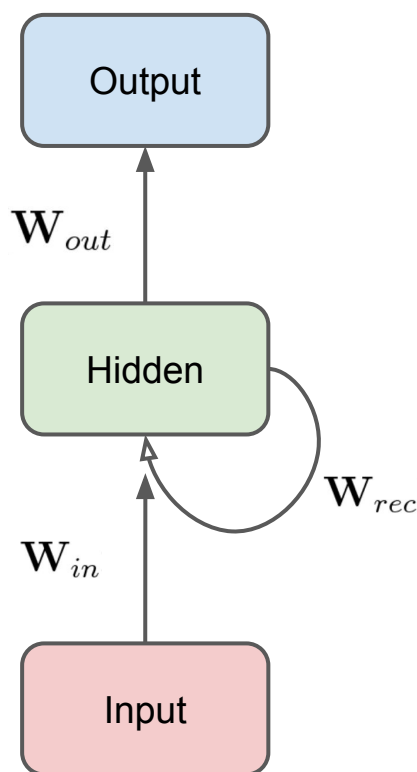
単純再帰型ネットワーク

- 初期モデルはElman(1990)が提唱
 - 文規則学習のシミュレーションに用いられた
- 入力層, 隠れ層, 出力層
直前の隠れ層を次時刻の隠れ層の入力に
 - 時間方向の因果関係を考慮することが可能
- 可変長な入力系列を扱うことが可能
- 誤差逆伝播法 (BP)を用いて学習

RNNの略記法



RNNの定式化



◇ パラメータ集合

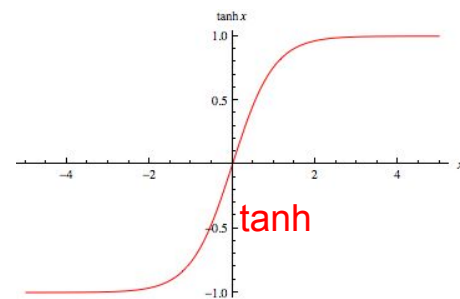
$$\theta = \{\mathbf{W}_{in}, \mathbf{W}_{rec}, \mathbf{W}_{out}, \mathbf{b}, \mathbf{c}\}$$

◇ 隠れ層(Hidden)の状態

$$\mathbf{h}(t) = \Phi_H(\mathbf{W}_{in}\mathbf{x}^{(t)} + \mathbf{W}_{rec}\mathbf{h}^{(t-1)} + \mathbf{b})$$

◇ 出力層(Output)の状態

$$\mathbf{y}(t) = \Phi_O(\mathbf{W}_{out}\mathbf{h}^{(t)} + \mathbf{c})$$

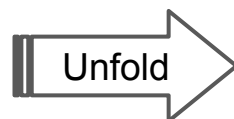


Activation function

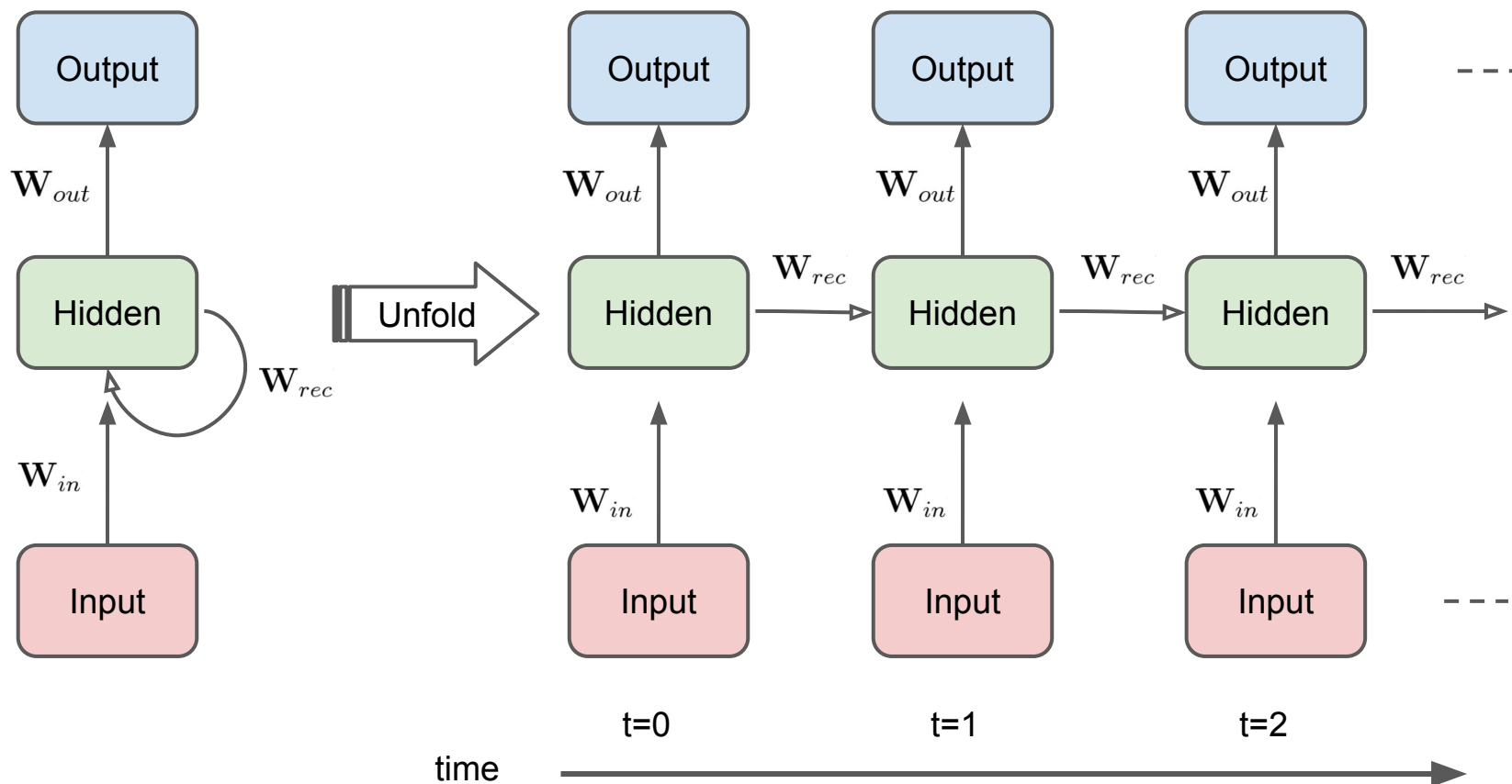
時間軸方向への展開

この変形はよく出てくる

Feedback Network



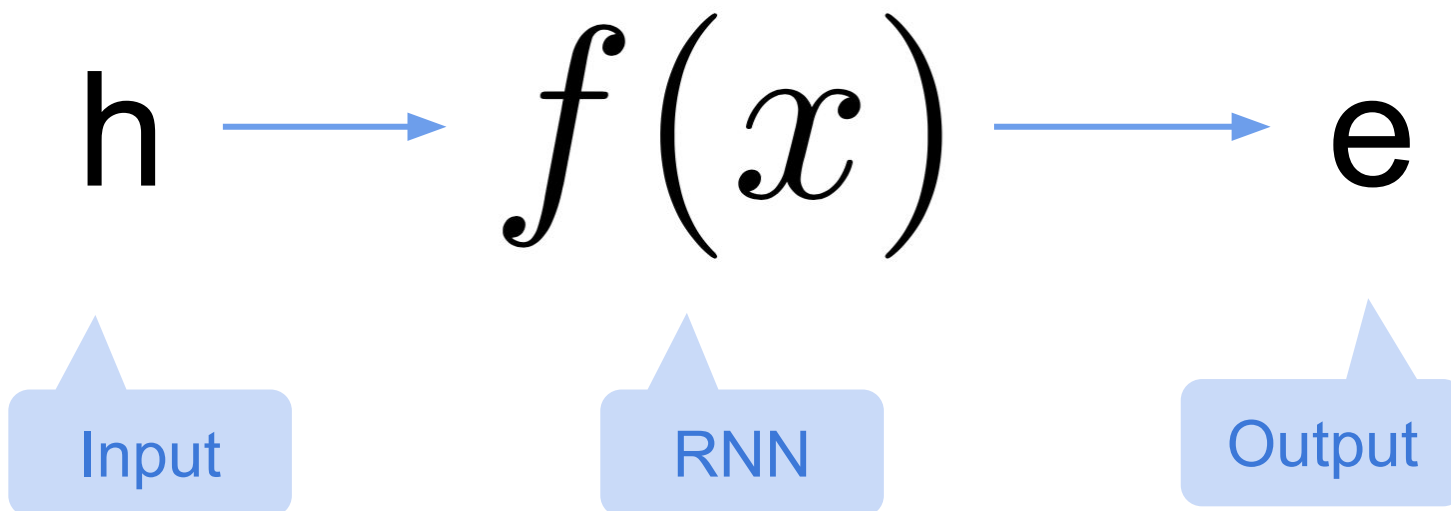
Feedforward Network



(例) 順伝播の計算 (1/3)

“hello”を学習した
Character Level RNN

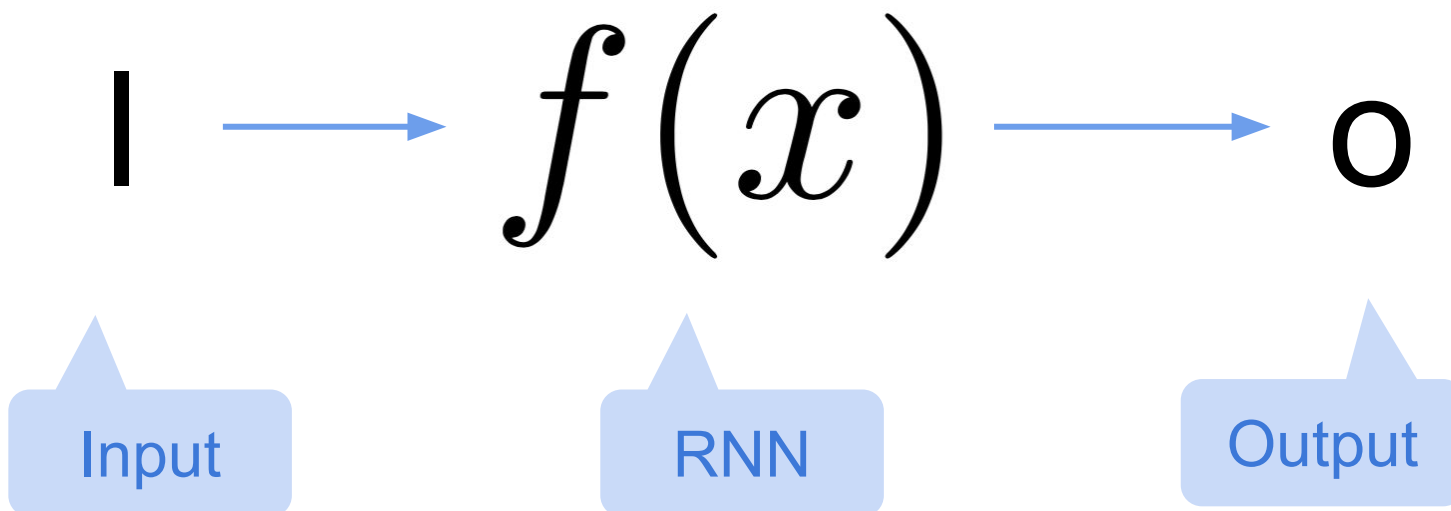
“hello”



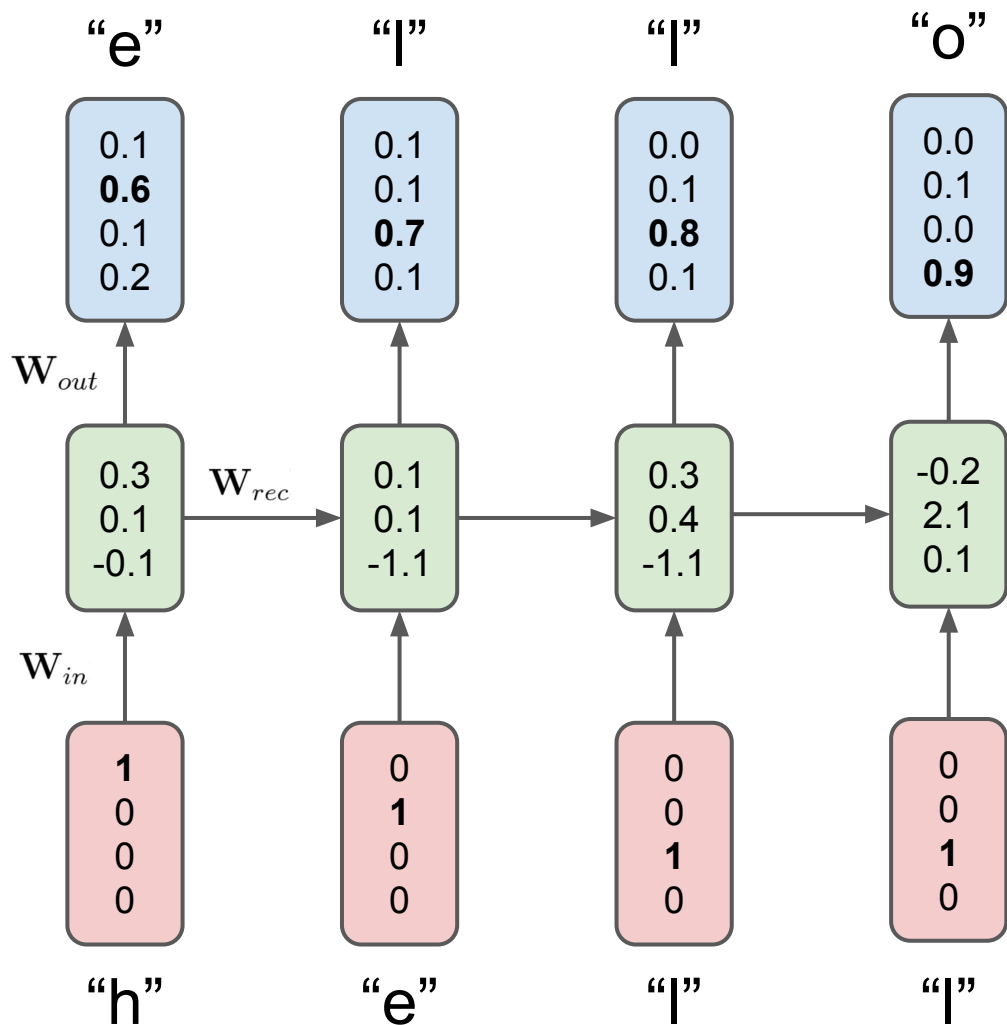
(例) 順伝播の計算 (2/3)

“hello”を学習した
Character Level RNN

“hello”

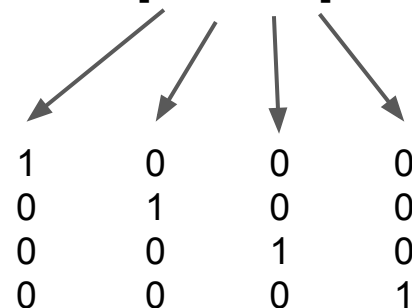


(例) 順伝播の計算 (3/3)



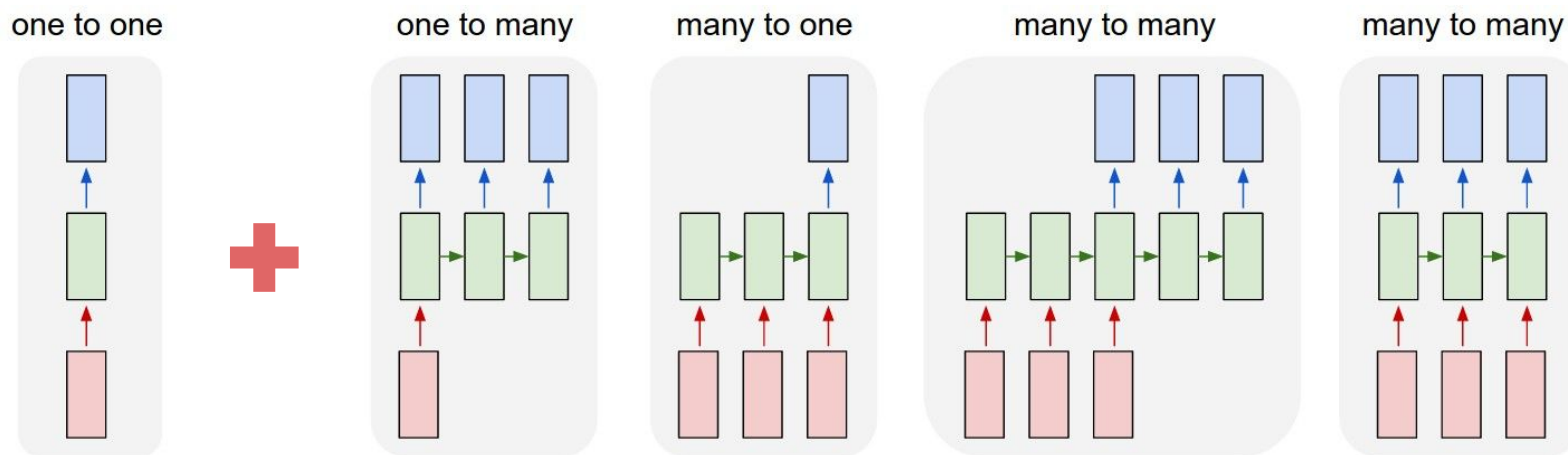
“hello”を学習した
Character Level RNN

文字: [h, e, l, o]



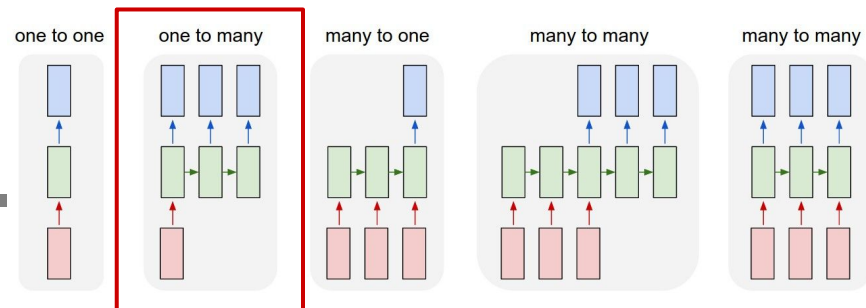
One-hot vector表現
分散表現の一種

RNNの柔軟性 (1/5)

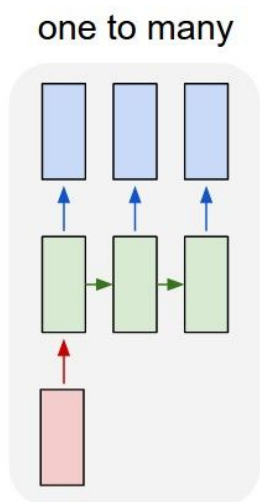


入出力の対応関係を変えることで様々な用途に対応することが可能

RNNの柔軟性 (2/5)

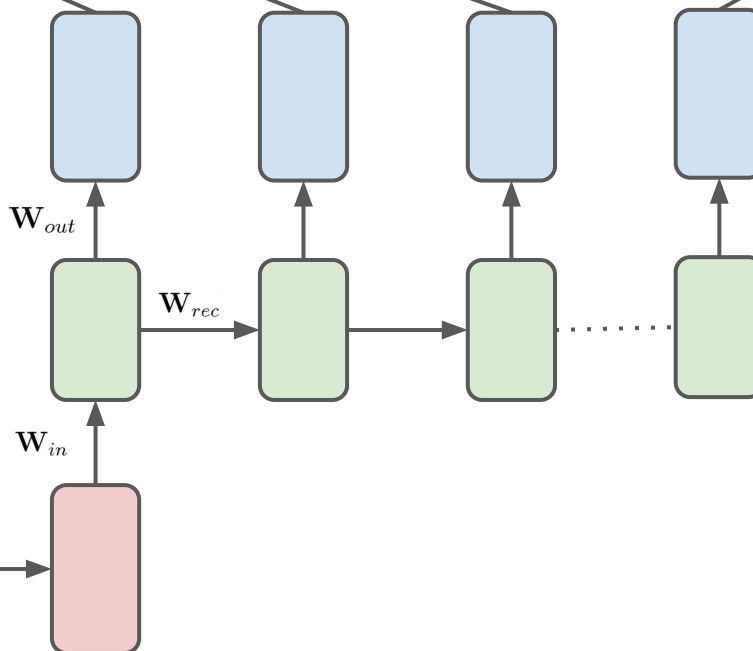
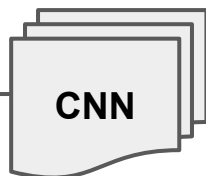


画像キャプションニング

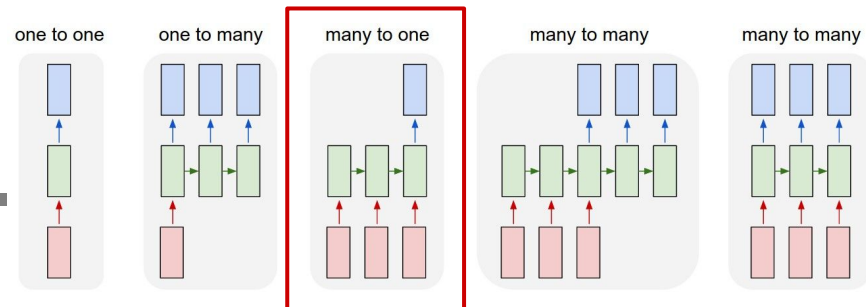


画像

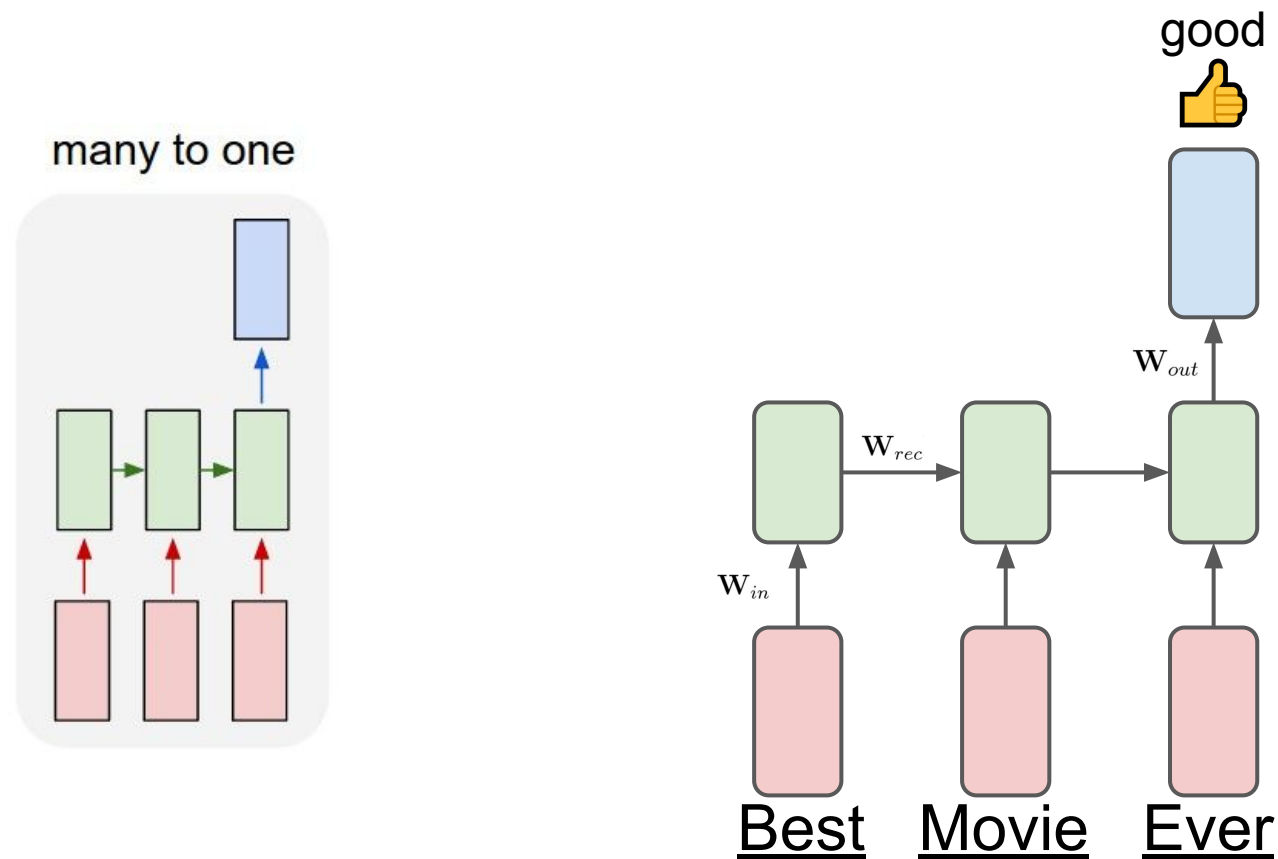
A person riding a motorcycle on a dirt road.



RNNの柔軟性 (3/5)

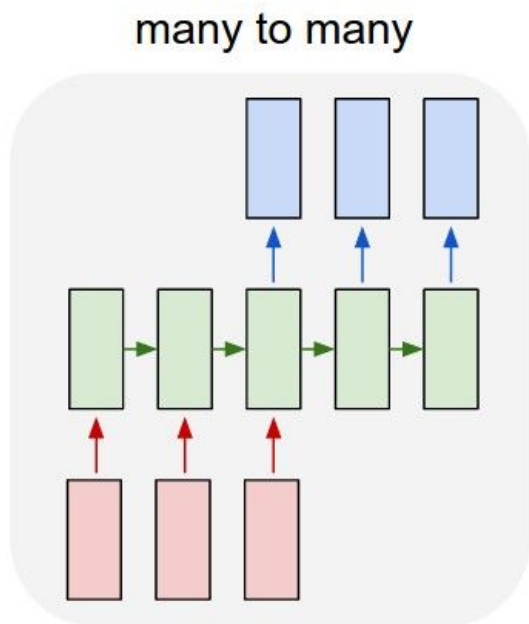
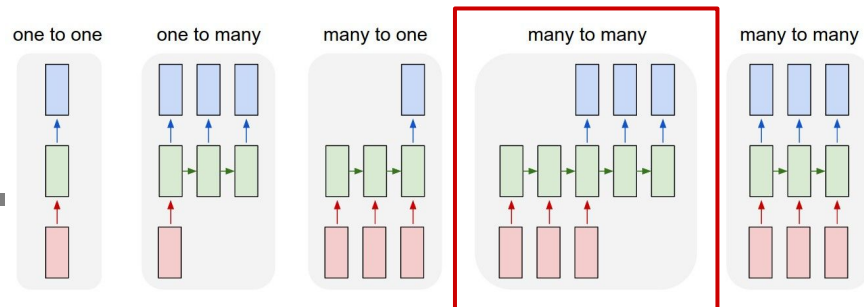


要約



RNNの柔軟性 (4/5)

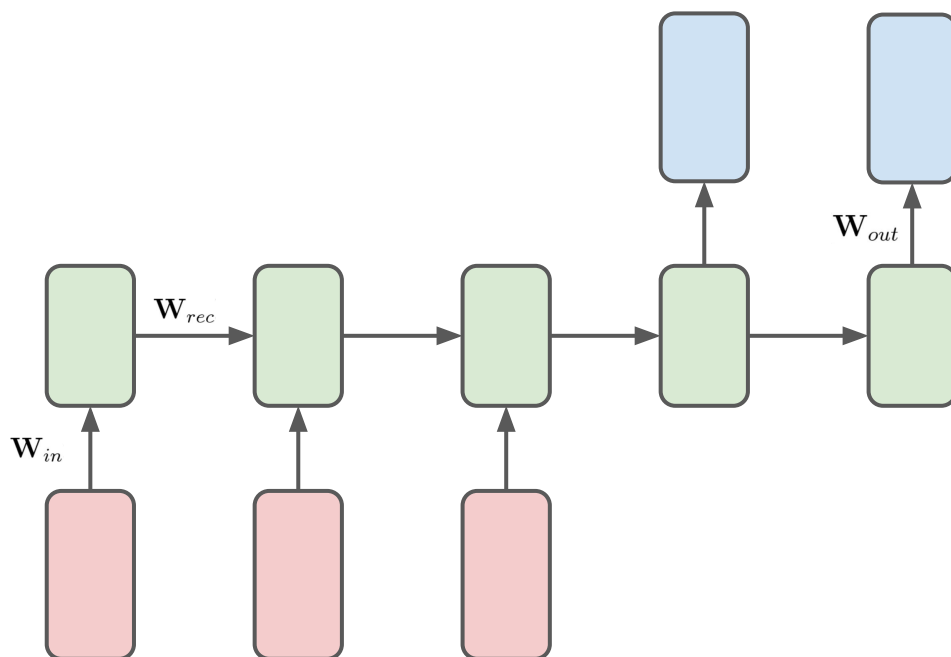
機械翻訳



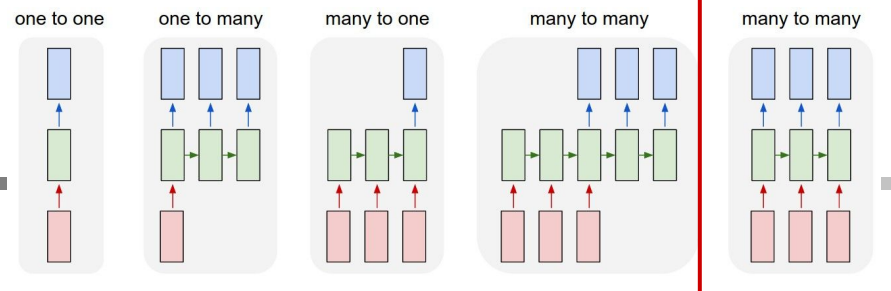
Thank you <EOS>



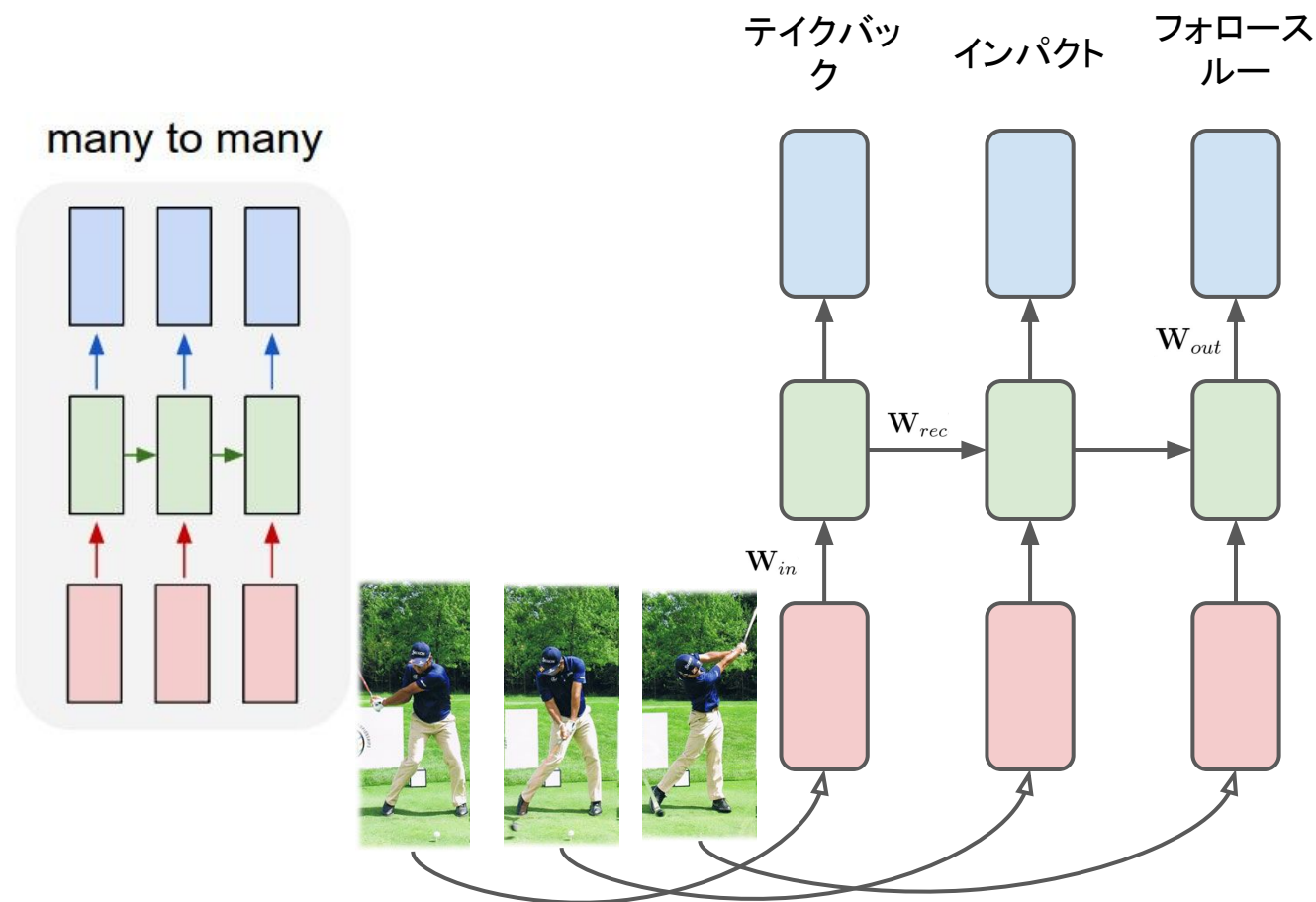
Danke <EOS>



RNNの柔軟性 (5/5)



動画キャプションング



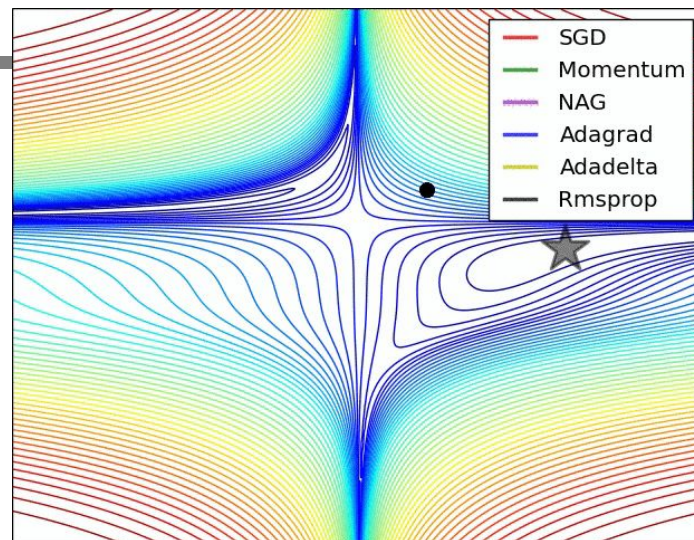
RNNの学習法

Gradient Descent (GD)

- 最急降下法

$$\theta^{k+1} = \theta^k - \alpha \frac{\partial J}{\partial \theta}$$

勾配を求めたい

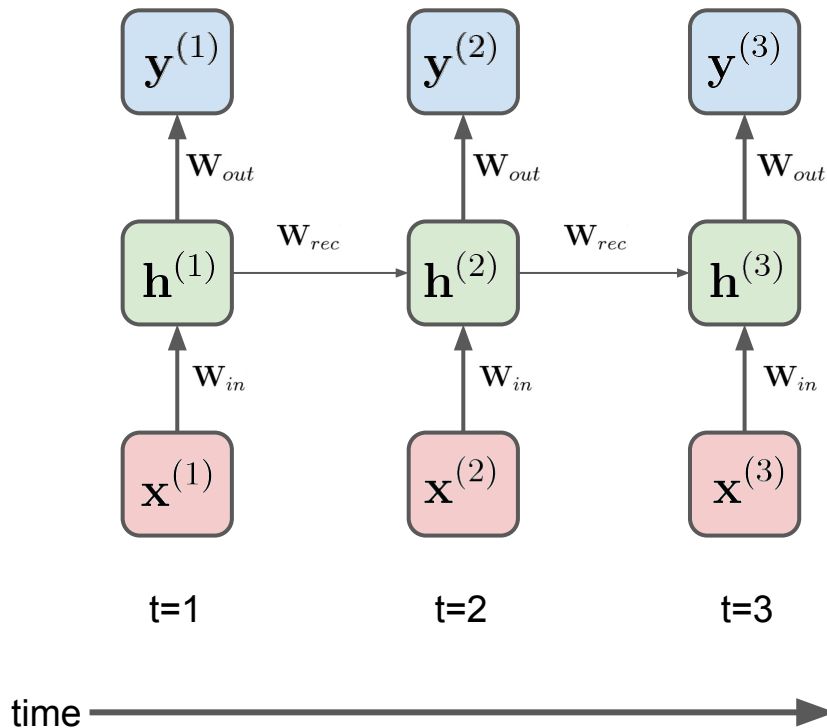


Backpropagation Through Time (BPTT)

- 勾配を計算
- RNNを展開して、FFNNとしてBPを行う
- 時間軸方向(Through Time)のBackpropagation

BPTTによる学習: 具体例(1/6)

順伝搬計算



◇ 時刻 $t = 1$

$$h^{(1)} = \Phi_H(W_{in}x^{(1)} + W_{rec}h^{(0)} + b)$$

$$y^{(1)} = \Phi_O(W_{out}h^{(1)} + c)$$

定数

◇ 時刻 $t = 2$

$$h^{(2)} = \Phi_H(W_{in}x^{(2)} + W_{rec}h^{(1)} + b)$$

$$y^{(2)} = \Phi_O(W_{out}h^{(2)} + c)$$

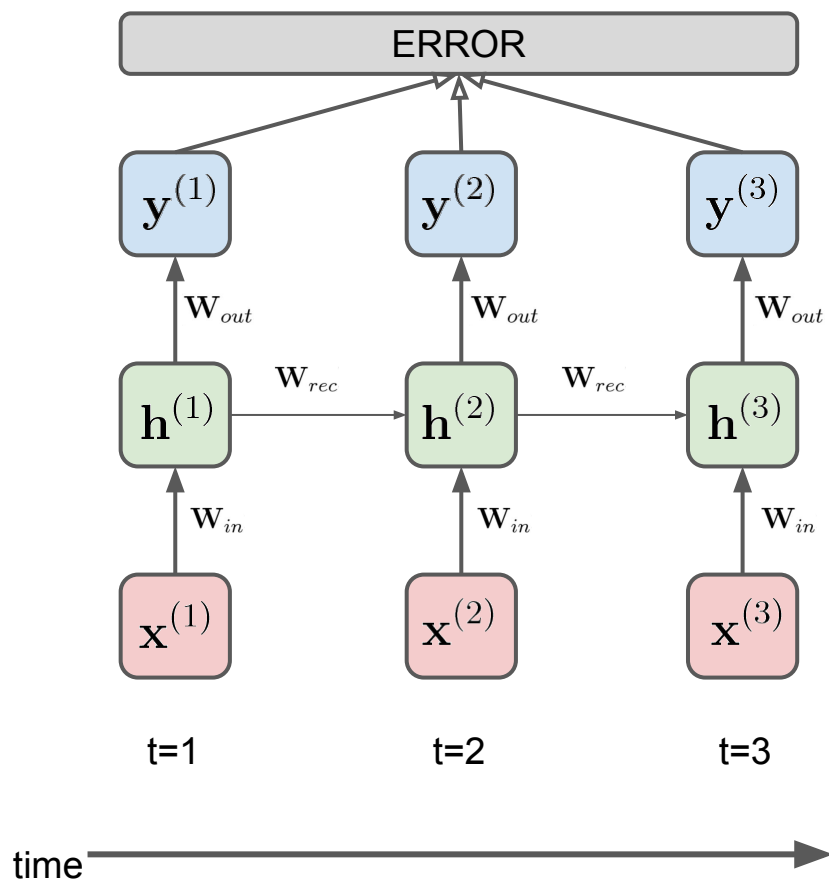
◇ 時刻 $t = 3$

$$h^{(3)} = \Phi_H(W_{in}x^{(3)} + W_{rec}h^{(2)} + b)$$

$$y^{(3)} = \Phi_O(W_{out}h^{(3)} + c)$$

BPTTによる学習: 具体例 (2/6)

最急降下法による最小値問題



◇ 誤差の計算 (L2ノルム)

教師信号

$$J = \frac{1}{2} \sum_{1 \leq t \leq 3} (y^{(t)} - d^{(t)})^2$$

最小化したい

FFNETと同じように最急降下法で
最小値問題を解く

◇ 最急降下法

$$\theta^{k+1} = \theta^k - \alpha \frac{\partial J}{\partial \theta}$$

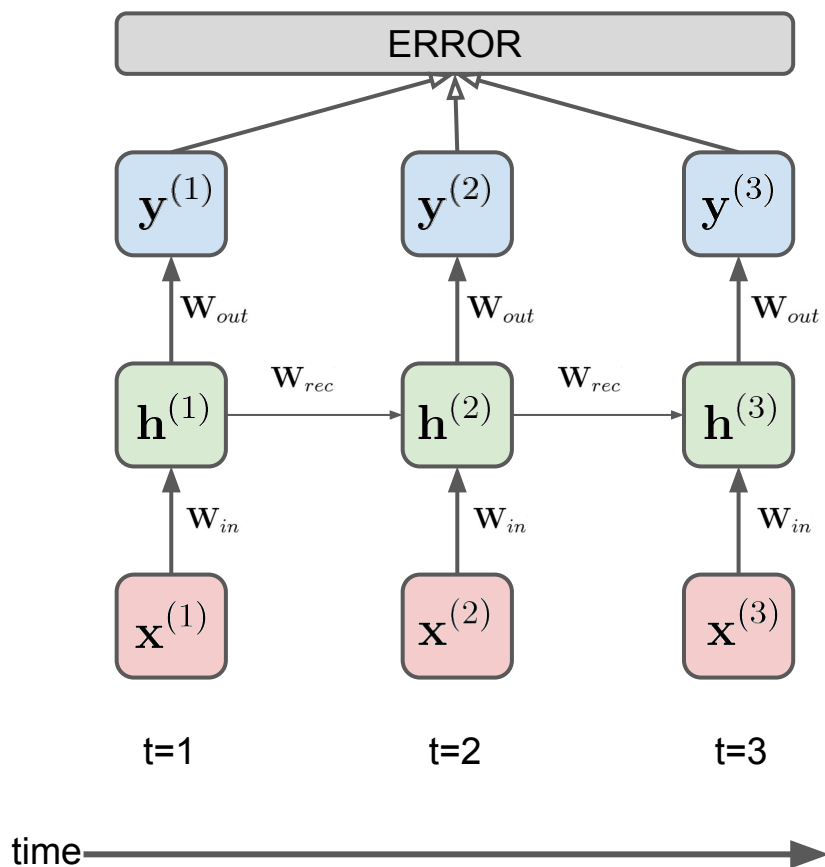
BPTTで
求める

□ パラメータ

$$\theta = \{W_{in}, W_{rec}, W_{out}, b, c\}$$

BPTTによる学習: 具体例 (3/6)

誤差逆伝播法



◇ 誤差逆伝播

□ 置き換え

$$J = \sum_{1 \leq t \leq 3} C^{(t)}$$

時刻tの出力誤差

□ 偏微分計算

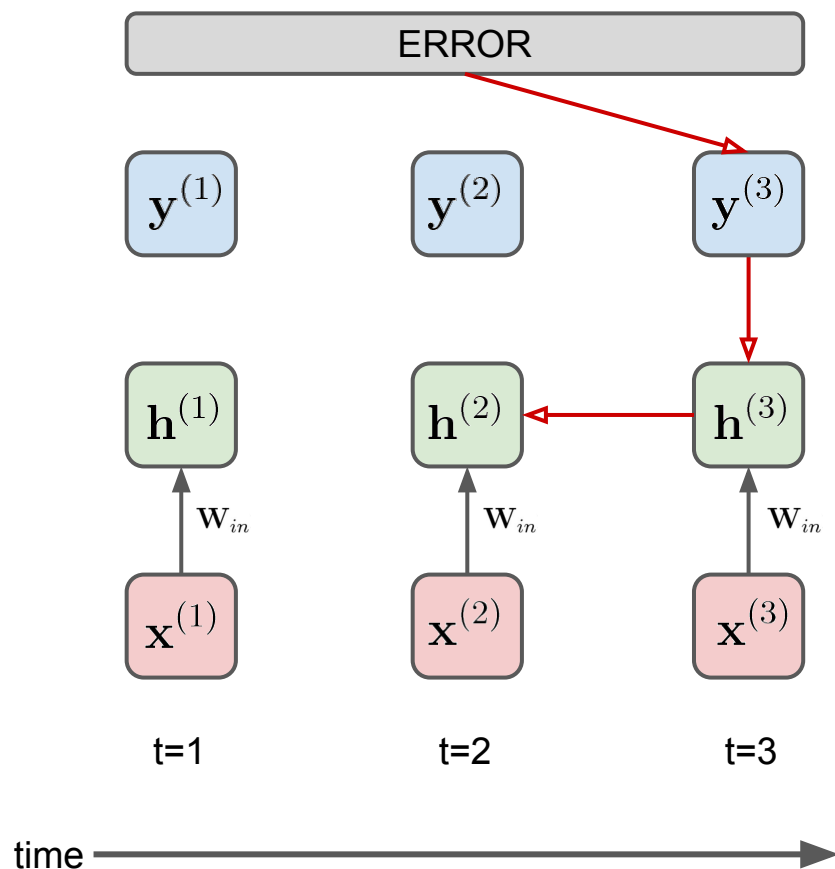
$$\begin{aligned} \frac{\partial J}{\partial \theta} &= \frac{\partial}{\partial \theta} \sum_{1 \leq t \leq 3} C^{(t)} \\ &= \sum_{1 \leq t \leq 3} \boxed{\frac{\partial C^{(t)}}{\partial \theta}} \end{aligned}$$

t=3の場合を求めてみる

$$C^{(k)} \triangleq \frac{1}{2} (y^{(k)} - d^{(k)})$$

BPTTによる学習: 具体例 (4/6)

誤差逆伝播法



◇ 誤差逆伝播

t=3の場合を求めている

□ C3の偏微分計算

$$\frac{\partial C^{(3)}}{\partial \theta} = \frac{\partial C^{(3)}}{\partial \mathbf{y}^{(3)}} \frac{\partial \mathbf{y}^{(3)}}{\partial \mathbf{h}^{(3)}} \quad \text{h3の関数}$$

y3の関数

$$\frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}} \quad \text{h2の関数}$$

◇ 時刻 t = 3

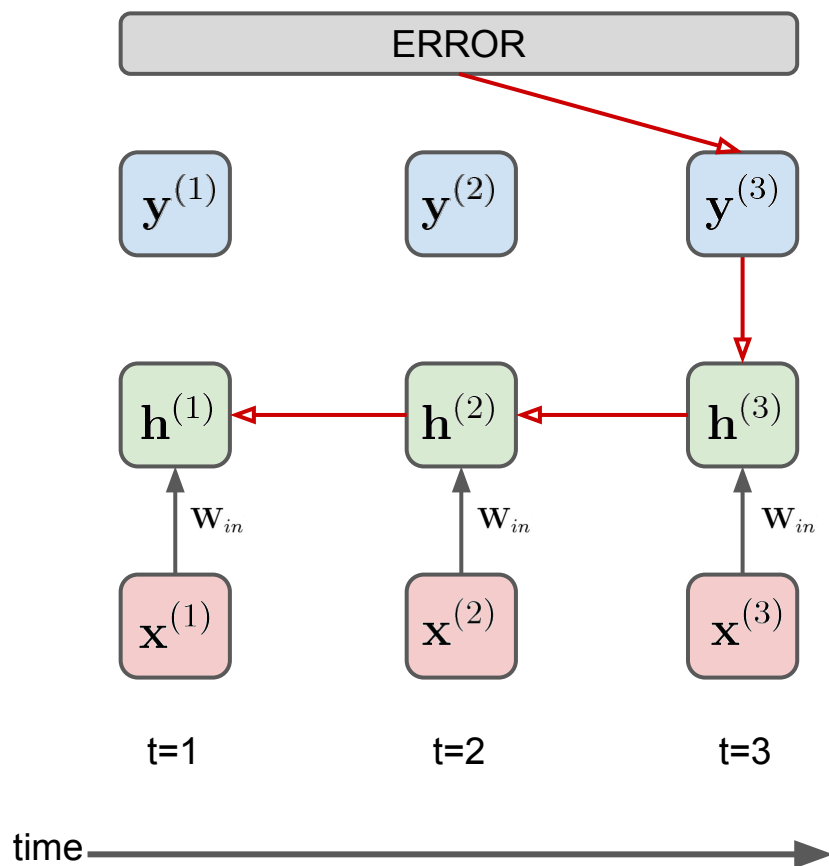
$$J = \frac{1}{2} \sum_{1 \leq t \leq 3} (\mathbf{y}^{(t)} - \mathbf{d}^{(t)})^2$$

$$\mathbf{y}^{(3)} = \Phi_O(\mathbf{W}_{\text{out}} \mathbf{h}^{(3)} + \mathbf{c})$$

$$\mathbf{h}^{(3)} = \Phi_H(\mathbf{W}_{\text{in}} \mathbf{x}^{(3)} + \mathbf{W}_{\text{rec}} \mathbf{h}^{(2)} + \mathbf{b})$$

BPTTによる学習: 具体例 (5/6)

誤差逆伝播法



◇ 誤差逆伝播

t=3の場合を求めている

□ C3の偏微分計算

$$\frac{\partial C^{(3)}}{\partial \theta} = \frac{\partial C^{(3)}}{\partial y^{(3)}} \frac{\partial y^{(3)}}{\partial h^{(3)}} \quad \text{h3の関数}$$

y3の関数

$$\frac{\partial h^{(3)}}{\partial h^{(2)}} \quad \text{h2の関数}$$

$$\frac{\partial h^{(2)}}{\partial h^{(1)}} \frac{\partial h^{(1)}}{\partial \theta} \quad \text{h1の関数}$$

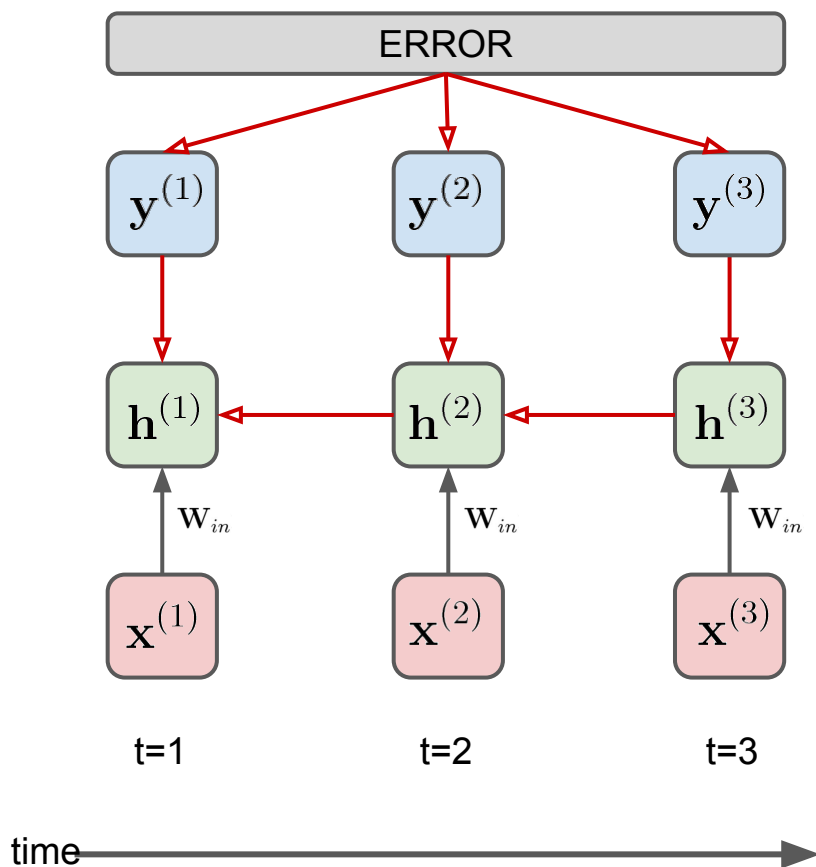
◇ 時刻 t = 2, 1

$$h^{(2)} = \Phi_H(W_{in}x^{(2)} + W_{rec}h^{(1)} + b)$$

$$h^{(1)} = \Phi_H(W_{in}x^{(1)} + W_{rec}h^{(0)} + b) \quad \text{定数}$$

BPTTによる学習: 具体例 (6/6)

誤差逆伝播法



◇ 誤差逆伝播

$t=1, 2, 3$ の場合

□ C3の偏微分計算

$$\frac{\partial C^{(3)}}{\partial \theta} = \frac{\partial C^{(3)}}{\partial \mathbf{y}^{(3)}} \frac{\partial \mathbf{y}^{(3)}}{\partial \mathbf{h}^{(3)}} \frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}} \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}} \frac{\partial \mathbf{h}^{(1)}}{\partial \theta}$$

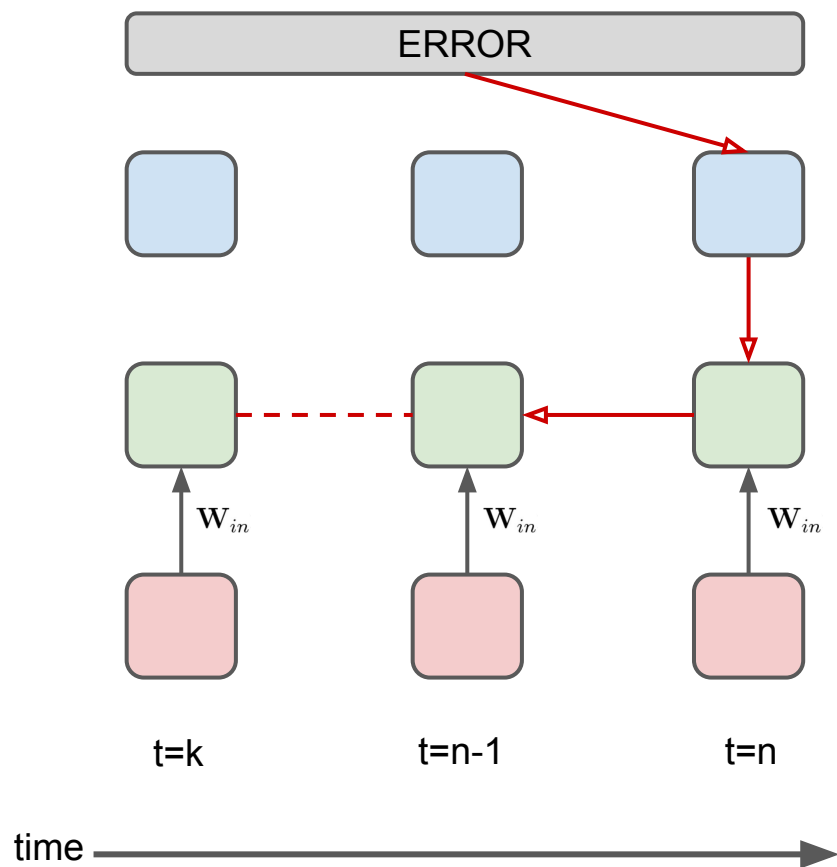
□ C2の偏微分計算

$$\frac{\partial C^{(2)}}{\partial \theta} = \frac{\partial C^{(2)}}{\partial \mathbf{y}^{(2)}} \frac{\partial \mathbf{y}^{(2)}}{\partial \mathbf{h}^{(2)}} \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}} \frac{\partial \mathbf{h}^{(1)}}{\partial \theta}$$

□ C1の偏微分計算

$$\frac{\partial C^{(1)}}{\partial \theta} = \frac{\partial C^{(1)}}{\partial \mathbf{y}^{(1)}} \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}} \frac{\partial \mathbf{h}^{(1)}}{\partial \theta}$$

BPTTによる学習: 一般化



◇ 誤差逆伝播

t=n の場合

□ Cnの偏微分計算

k=1

$$\frac{\partial C^{(n)}}{\partial \theta} = \frac{\partial C^{(n)}}{\partial \mathbf{y}^{(n)}} \frac{\partial \mathbf{y}^{(n)}}{\partial \mathbf{h}^{(n)}} \frac{\partial \mathbf{h}^{(n)}}{\partial \mathbf{h}^{(k)}} \frac{\partial \mathbf{h}^{(k)}}{\partial \theta}$$

$$\frac{\partial \mathbf{h}^{(n)}}{\partial \mathbf{h}^{(k)}} = \prod_{k < i \leq n} \frac{\partial \mathbf{h}^{(i)}}{\partial \mathbf{h}^{(i-1)}}$$

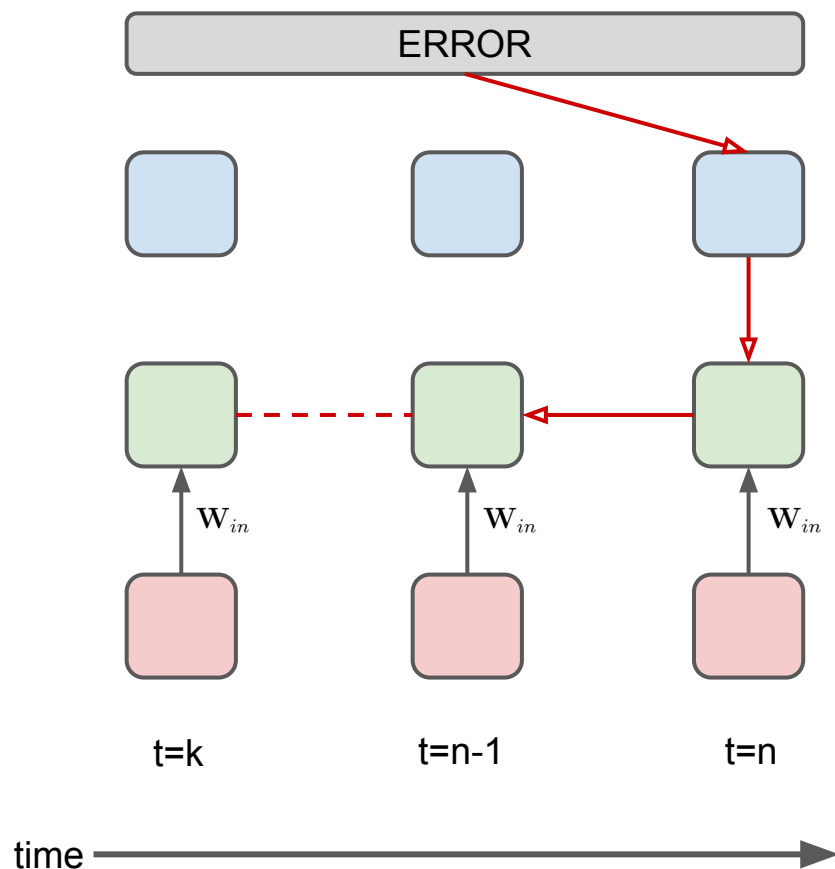
h(z)とする

$$= \prod_{k < i \leq n} \frac{\partial \mathbf{h}^{(i)}}{\partial \mathbf{z}^{(i)}} \frac{\partial \mathbf{z}^{(i)}}{\partial \mathbf{h}^{(i-1)}}$$

$$= \prod_{k < i \leq n} \Phi'_H(\mathbf{z}^{(i)}) W_{\text{rec}}^T$$

$$\triangleq \prod_{k < i \leq n} \Phi'_H(\mathbf{z}^{(i)}) W_{\text{rec}}^T$$

BPTTによる学習: 勾配爆発・消失 (1/2)



◇ 誤差逆伝播

t=n の場合

□ Cnの偏微分計算

$$\frac{\partial C^{(n)}}{\partial \theta} = \frac{\partial C^{(n)}}{\partial \mathbf{y}^{(n)}} \frac{\partial \mathbf{y}^{(n)}}{\partial \mathbf{h}^{(n)}} \left(\prod_{k < i \leq n} \Phi_H'^{(i)} W_{\text{rec}}^T \right) \frac{\partial \mathbf{h}^{(k)}}{\partial \theta}$$

時系列の長い
データ

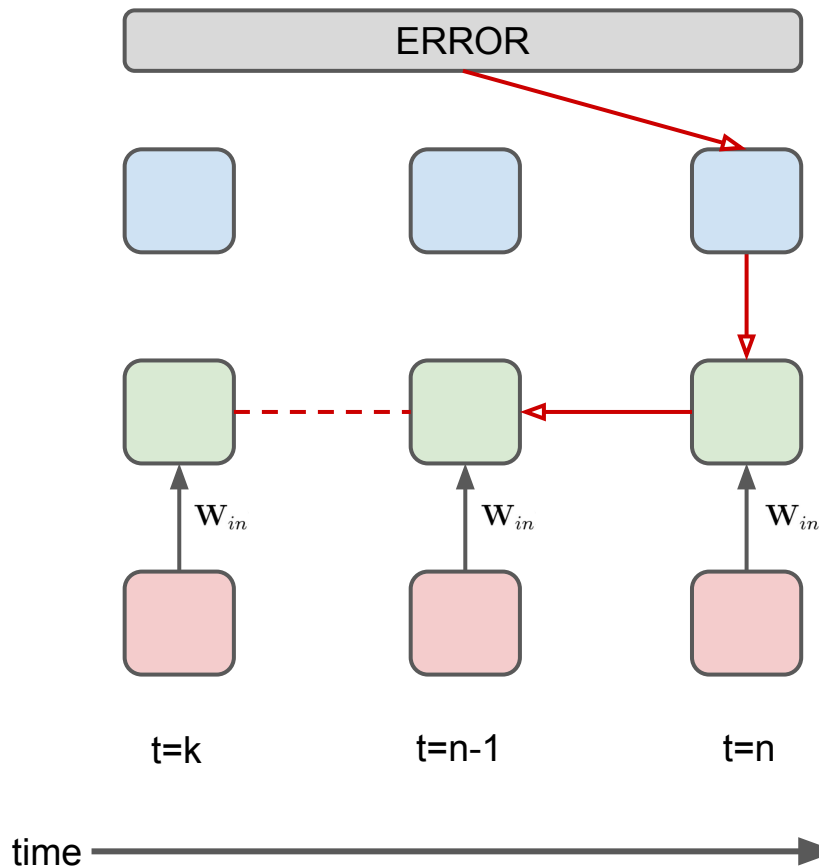
n大, k小 → chainが長くなる

□ n = 100, k = 1の場合

$$\begin{aligned} \frac{\partial C^{(100)}}{\partial \theta} &= \frac{\partial C^{(100)}}{\partial \mathbf{y}^{(100)}} \frac{\partial \mathbf{y}^{(100)}}{\partial \mathbf{h}^{(100)}} \|\Phi_H'\|^{99} \|W_{\text{rec}}^T\|^{99} \frac{\partial \mathbf{h}^{(1)}}{\partial \theta} \\ &\approx \|\Phi_H'\|^{99} \|W_{\text{rec}}^T\|^{99} \end{aligned}$$

$\|x\| < 1$: で勾配が消失 ($\rightarrow 0$)
 $\|x\| > 1$: で勾配が爆発 ($\rightarrow \infty$)

BPTTによる学習: 勾配爆発・消失 (2/2)



◇ 誤差逆伝播

t=n の場合

□ Cnの偏微分計算 (まとめ)

$$\frac{\partial C^{(n)}}{\partial \theta} = \frac{\partial C^{(n)}}{\partial \mathbf{y}^{(n)}} \frac{\partial \mathbf{y}^{(n)}}{\partial \mathbf{h}^{(n)}} \left(\prod_{k < i \leq n} \Phi_H'^{(i)} W_{\text{rec}}^T \right) \frac{\partial \mathbf{h}^{(k)}}{\partial \theta}$$

$$\approx ||\Phi_H'||^{n-k} ||W_{\text{rec}}^T||^{n-k}$$

系列の長いデータを扱う時(n大の時)は
時間の展開を小さくする(kを大きくする)
ことで**勾配爆発・消失**を回避する

※ 勾配爆発・消失に対する詳しいアプローチは第10回の講義で解説

RNNまとめ

- RNNでは過去の隠れ層の状態を入力に加えることで、時系列を保持している
- ネットワークを時間軸方向に展開して、BPを行う(BPTT)
- BPTTでは勾配爆発・消失問題がある
- 長期時系列データを扱うことは難しい
 - ⇒ LSTM, GRU, Gradient Clipping (第10章)



Work #9-1

人工データを用いたRNNの学習

講義 Part 2

BPを使わないRNN



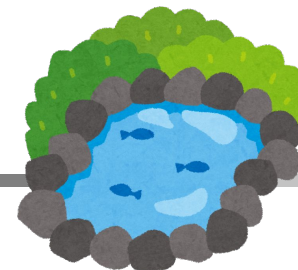
Echo State Network

一般的なRNNの問題点

- RNNは計算コストが高すぎる
 - 勾配消失・爆発問題
 - LSTM, GRU, Gradient Clipping などの発明
 - それでも、学習が難しいのが現状
- そもそも、脳はBPTTなんてしているのか
 - 過去の入出力や結合荷重を保存して
時間を遡って誤差を計算する...アヤシイ

Reservoir Computing

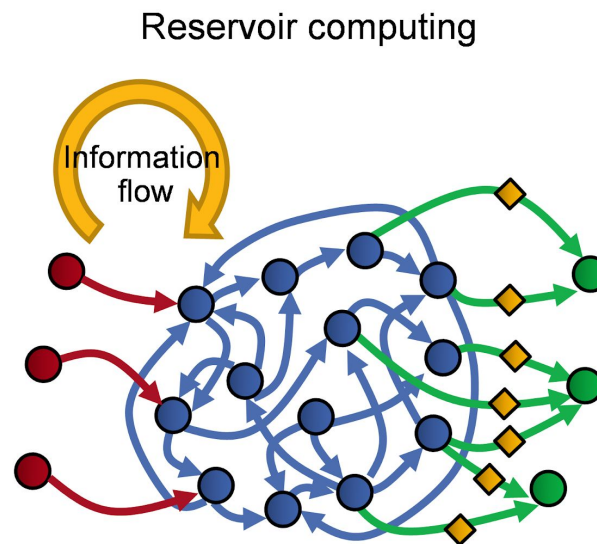
Reservoir: 貯水池 →



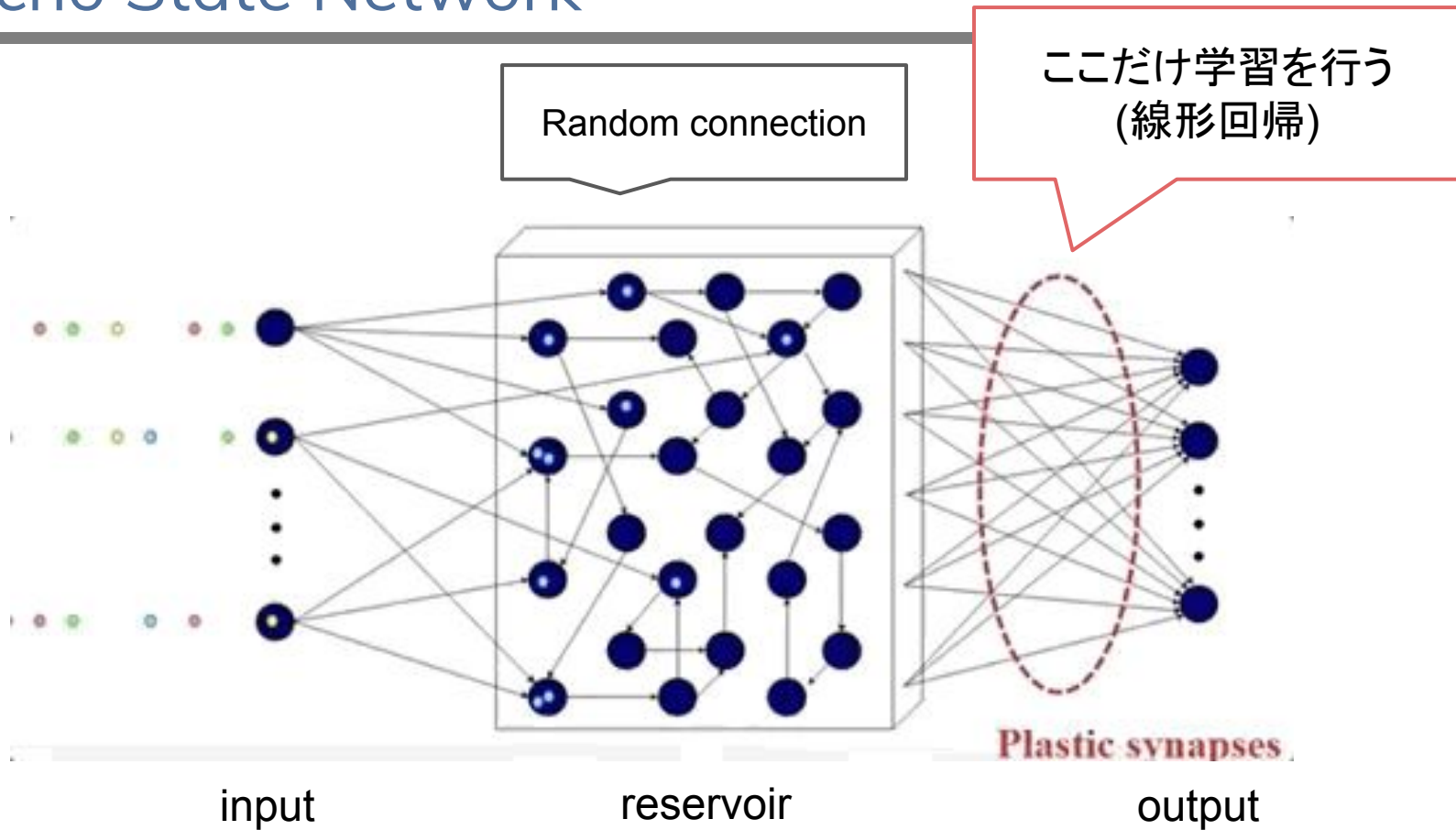
- 入力層, Reservoir層, 出力層からなる
- ランダムに結合したReservoir層が時系列を保持
- 学習は出力層の線形回帰だけ
→ 計算量が少ないため爆速
- 脳の神経モデル

代表的なネットワーク:

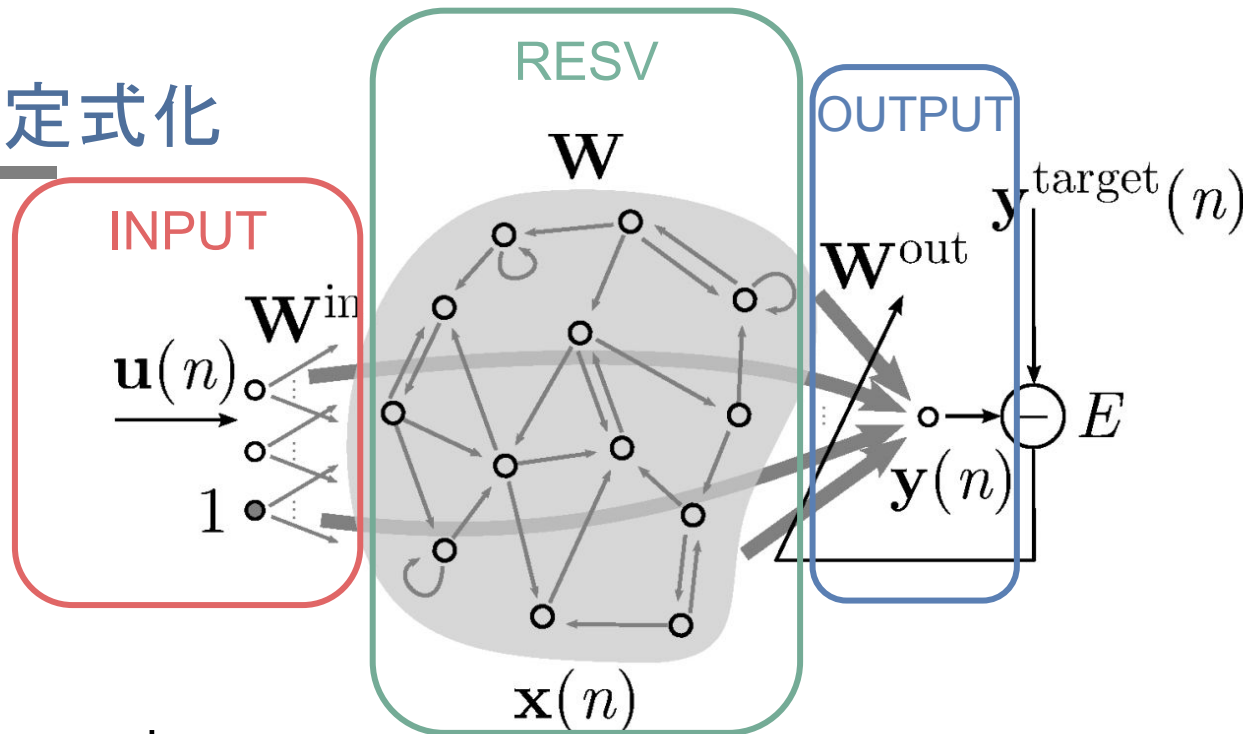
- Echo State Networks (ESN)
- Liquid State Machine (LSM)



Echo State Network



ESNの定式化



Update: reservoir

$$\tilde{\mathbf{x}}(n) = \tanh(\mathbf{W}^{\text{in}}[1; \mathbf{u}(n)] + \mathbf{W}\mathbf{x}(n-1)),$$

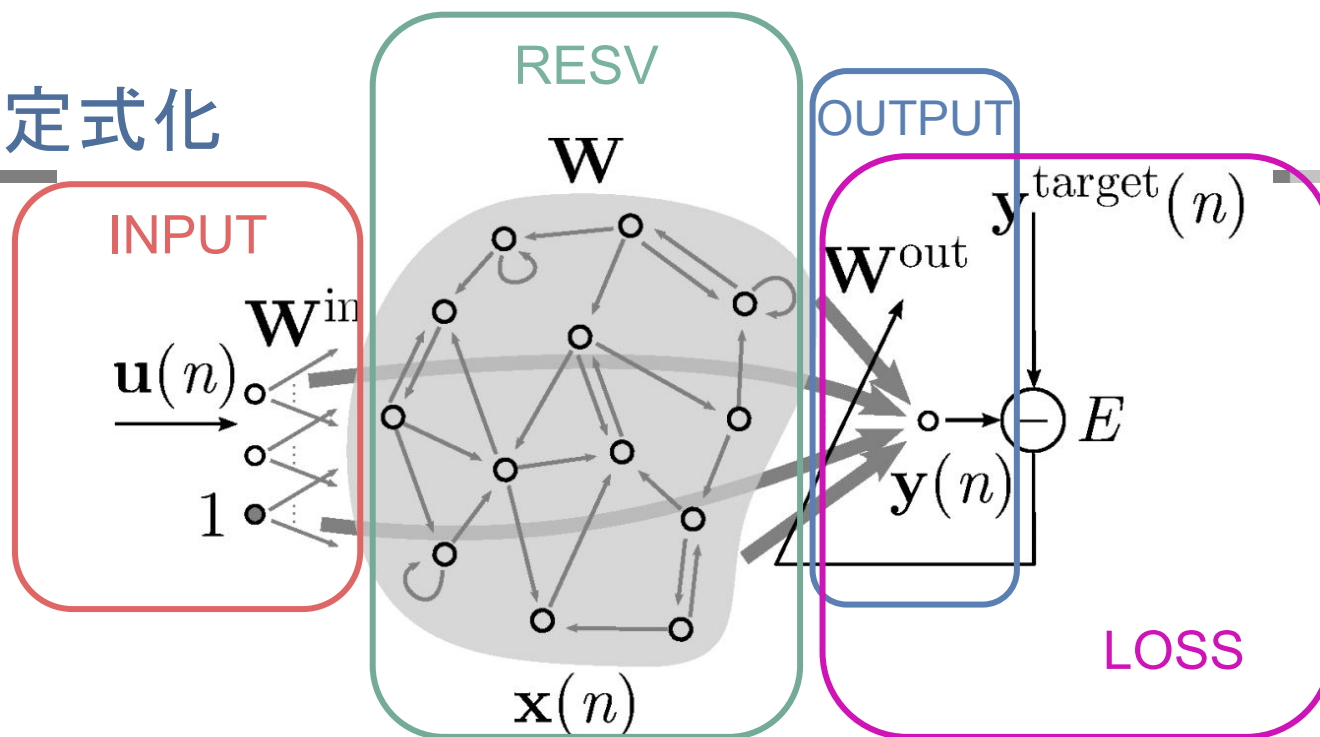
$$\mathbf{x}(n) = (1 - \alpha)\mathbf{x}(n-1) + \alpha\tilde{\mathbf{x}}(n),$$

α : leaking rate
記憶の強さを調整

Compute: output

$$\mathbf{y}(n) = \mathbf{W}^{\text{out}}[1; \mathbf{u}(n); \mathbf{x}(n)],$$

ESNの定式化



Compute: output

$$\mathbf{y}(n) = \mathbf{W}^{\text{out}}[1; \mathbf{u}(n); \mathbf{x}(n)],$$

Compute: loss

$$E(\mathbf{y}, \mathbf{y}^{\text{target}}) = \frac{1}{N_y} \sum_{i=1}^{N_y} \sqrt{\frac{1}{T} \sum_{n=1}^T (y_i(n) - y_i^{\text{target}}(n))^2}$$

RMSEを最小化
線形回帰で学習

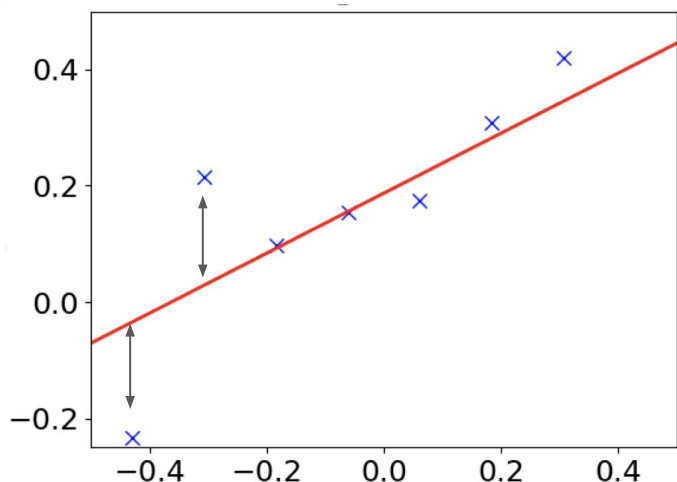
ESNの学習: 線形回帰

L2正則化最小二乗回帰

最小二乗法の復習

- ESNの学習は通常リッジ回帰を使って行われる
 - 最も単純な方法では、すべての出力を用意して回帰計算を行う
 - リッジ回帰については『第三回: 線形回帰』で解説済み

今回は簡単のため、単純線形回帰で求めている



モデルの予測と正解の二乗和誤差が最小になるようなパラメータ θ を探す

$$L_{LS}(\mathbf{W}) = \frac{1}{2} \sum_{i=1}^T (f(y_i; \mathbf{W}) - y_i^{tgt})^2$$

$$\mathbf{W}_{LS} = \operatorname{argmin}_{\mathbf{W}} L_{LS}(\mathbf{W})$$

ESNの学習: 単純線形回帰

求めるWは自乗誤差Lを最小化するW $L_{LS}(W)$

$$W_{LS}^{out} = \arg \min_{W^{out}} \overbrace{\frac{1}{2} ||W^{out} X - Y^{tgt}||^2}$$

Wに関して2次関数の形 (凸関数) をしているので、 $L_{LS}(W)$
の偏微分=0を与えるWが求値

$$W_{LS}^{out} = Y^{tgt} X^T (X X^T)^{-1}$$

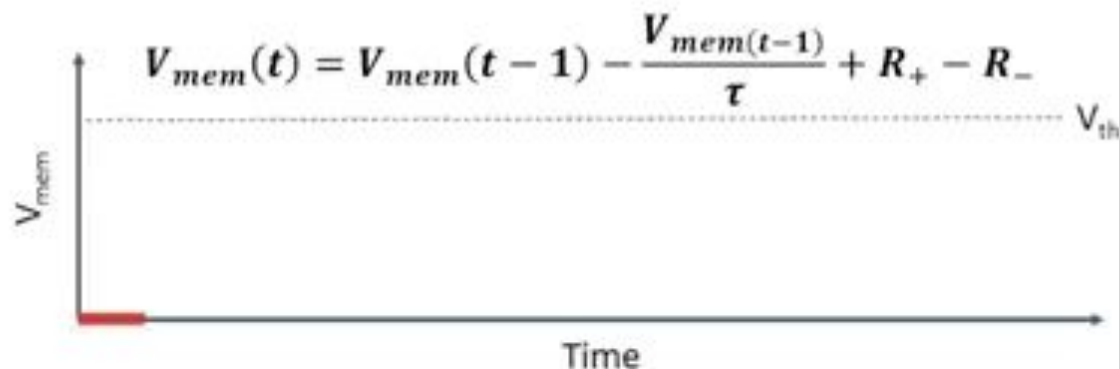
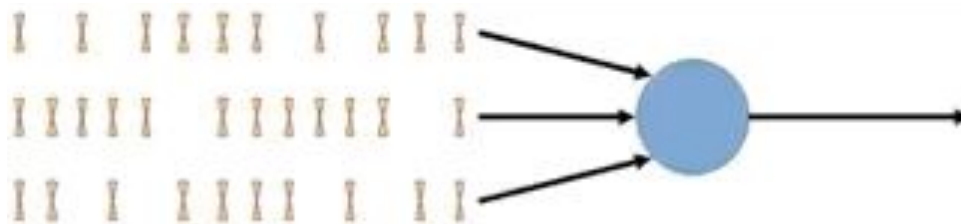
Liquid State Machine

特徴

- ESNに類似したRCモデルの一つ
- 脳っぽさを謳っている
- ニューロン発火モデルとして、
IAF model(積分発火モデル)を使用 (c.f. tanh for ESN)
- 連続的に変化するデータを扱うことができる
(\Leftrightarrow ESN: 離散的)
- ハードウェアと親和性が高く、VLSIなどでしばし実装される

LSM: IAFモデル

Integrate and Fire



時間を通してニューロンに電位が溜まっていき、
しきい値を超えた際に発火する (⇒積分発火モデル)
(c.f. ESNではtanhでactivationしていた)

なぜRCが上手くいくか

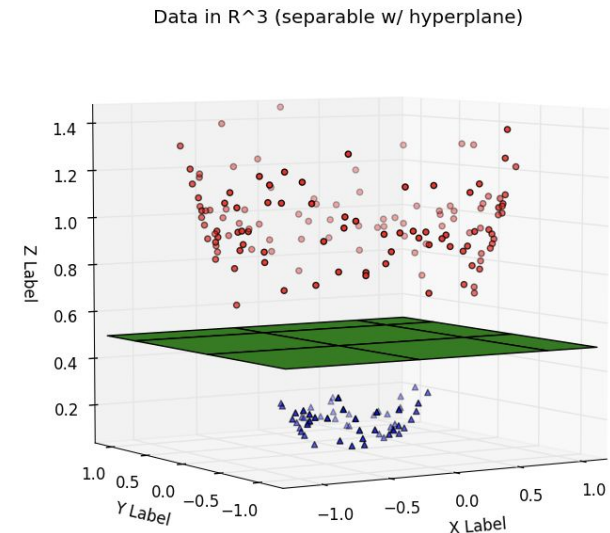
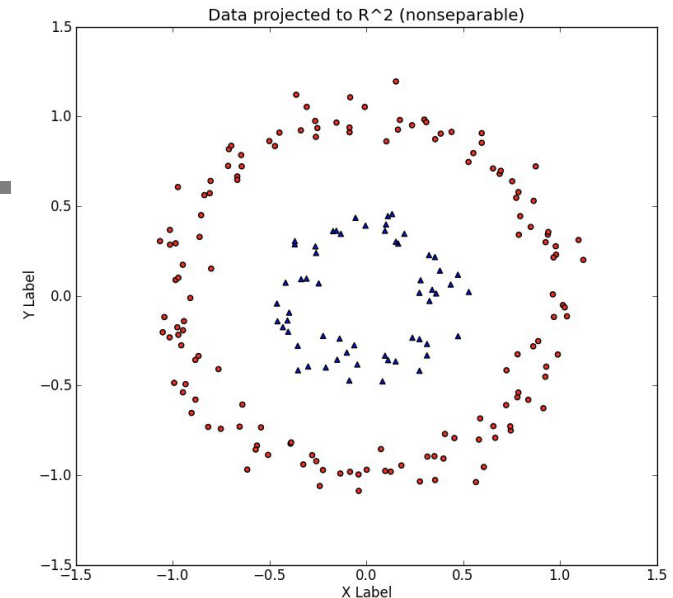
Kernel Trickとの類似性

Kernel Trick:

サポートベクターマシン(SVM)で使用
線形分離不能なデータの次元を上げて
頑張って線形分離する手法

Reservoir Computing:

- reservoir層で次元を拡張
→ Kernel Trickに相当



RCまとめ

- RCではランダムに結合したreservoir層を使用
- Reader(output)では単純な線形回帰を使用
 - ⇒ 計算コスト低 (c.f. RNN)
- 2つのモデルを紹介
 - ESN: tanhを使用
 - LSM: スパイクモデルを使用
- 計算がうまくいくのはカーネルトリックと似た原理
- Reservoirの重みは基本的に学習しない
 - 学習するモデルも存在するが、あまりうまく行っていない模様

Appendix

参考文献の表記

参考文献は以下のように表記されています

- 論文の場合: [Hinton+ 12]
 - APPENDIXのReferencesはAPA形式
- 本の場合: [Asakawa: 16]
 - APPENDIXのReferencesはAPA形式
- WEBページの場合: [Karpathy:: 17]
 - APPENDIXのReferencesはAPA形式

References

- [Reno911::18] https://www.youtube.com/watch?v=Ukgii7Yd_cU&index=4&list=PLdtFmJC5lu-AbZvwnTX7pQtsb6TpJqbiV&t=336s
- [Karpathy::18] (n.d.). Retrieved January 28, 2018, from <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
- [Jaeger 07] Jaeger, H. (2007). Echo state network. *Scholarpedia*, 2(9), 2330.
- [Lukoševičius 12] Lukoševičius, M. (2012). A practical guide to applying echo state networks. In *Neural networks: tricks of the trade* (pp. 659-686). Springer Berlin Heidelberg.
- [Maass+ 10] Maass, W. (2010). Liquid state machines: motivation, theory, and applications. *Computability in context: computation and logic in the real world*, 275-296.
- [Lukoševičius & Jaeger 09] Lukoševičius, M., & Jaeger, H. (2009). Reservoir computing approaches to recurrent neural network training. *Computer Science Review*, 3(3), 127-149.

References

[Lecture] Evolution: from vanilla RNN to GRU & LSTMs

<https://towardsdatascience.com/lecture-evolution-from-vanilla-rnn-to-gru-lstms-58688f1da83a>

[Wu+ 16] Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., ... & Klingner, J. (2016). Google's neural machine translation system: Bridging the gap between human and machine translation. arXiv preprint arXiv:1609.08144.