

RNN, LSTM

目的

- ・フレームワークを使ってネットワーク(RNN)を構築できる
- ・様々なフレームワークを使える

目次

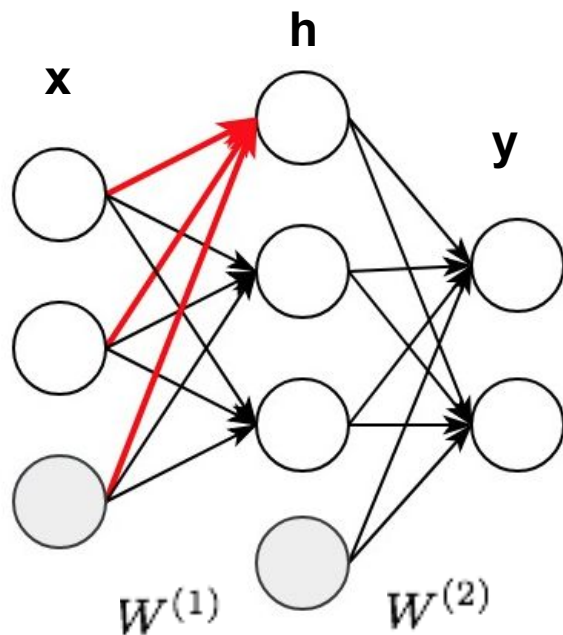
- ・はじめに
 - ・ニューラルネットワーク
 - ・順伝播型ニューラルネットワーク
- ・RNN
 - ・RNNの概要
 - ・RNNの計算
 - ・RNNの問題点
- ・LSTM
 - ・LSTMの概要
 - ・LSTMの計算
- ・まとめ

ニューラルネットワーク

- ・機械学習モデルの一種
 - ・自動特徴量設計
- ・様々なモデル
 - ・CNN(コンピュータビジョン)
 - ・RNN(自然言語処理、時系列予測)
 - ・GAN(生成モデル)

順伝播型ニューラルネットワーク

- 順伝播



$$h^{\ell} = f^{\ell}(W^{\ell}h^{\ell-1} + b^{\ell})$$

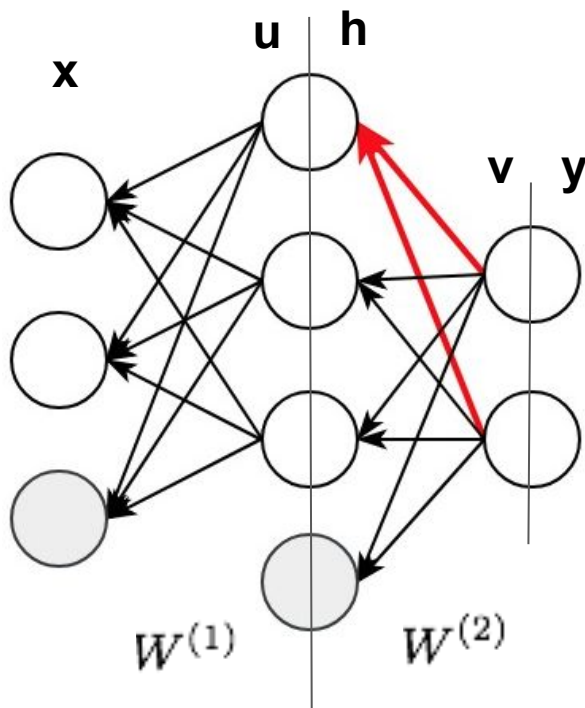


$$h = f(W^{(1)}x + b^{(1)})$$

$$y = g(W^{(2)}h + b^{(2)})$$

順伝播型ニューラルネットワーク

・逆伝播



勾配降下法

$$\theta \leftarrow \theta - \nabla_{\theta} L$$

順伝播の別記法

$$u = W^{(1)}x + b^{(1)}$$

$$h = f(u)$$

$$v = W^{(2)}h + b^{(2)}$$

$$y = g(v)$$

誤差の勾配

$$\delta^{(v)} = \frac{\partial L}{\partial v} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial v} = \frac{\partial L}{\partial y} g'(v)$$

$$\delta^{(u)} = \frac{\partial L}{\partial u} = \frac{\partial L}{\partial v} \frac{\partial v}{\partial h} \frac{\partial h}{\partial u} = (W^{(2)})^T \delta^{(v)} f'(u)$$

順伝播型ニューラルネットワーク

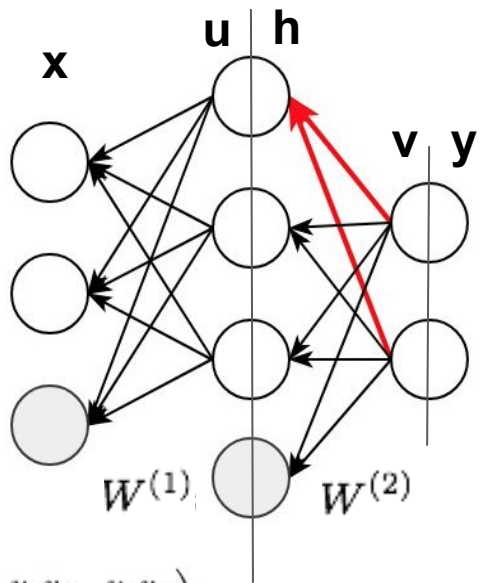
- 各パラメータについての勾配

$$\frac{\partial L}{\partial W^{(2)}} = \frac{\partial L}{\partial v} \frac{\partial v}{\partial W^{(2)}} = \delta^{(v)} \otimes h$$

$$\frac{\partial L}{\partial b^{(2)}} = \frac{\partial L}{\partial v} \frac{\partial v}{\partial b^{(2)}} = \delta^{(v)}$$

$$\frac{\partial L}{\partial W^{(1)}} = \frac{\partial L}{\partial u} \frac{\partial u}{\partial W^{(1)}} = \delta^{(u)} \otimes x$$

$$\frac{\partial L}{\partial b^{(1)}} = \frac{\partial L}{\partial u} \frac{\partial u}{\partial b^{(1)}} = \delta^{(u)}$$

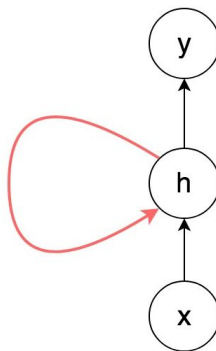


テンソル積(直積)

$$u \otimes v = uv^{\top} = \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{pmatrix} \begin{pmatrix} v_1 & v_2 & v_3 \end{pmatrix} = \begin{pmatrix} u_1 v_1 & u_1 v_2 & u_1 v_3 \\ u_2 v_1 & u_2 v_2 & u_2 v_3 \\ u_3 v_1 & u_3 v_2 & u_3 v_3 \\ u_4 v_1 & u_4 v_2 & u_4 v_3 \end{pmatrix}.$$

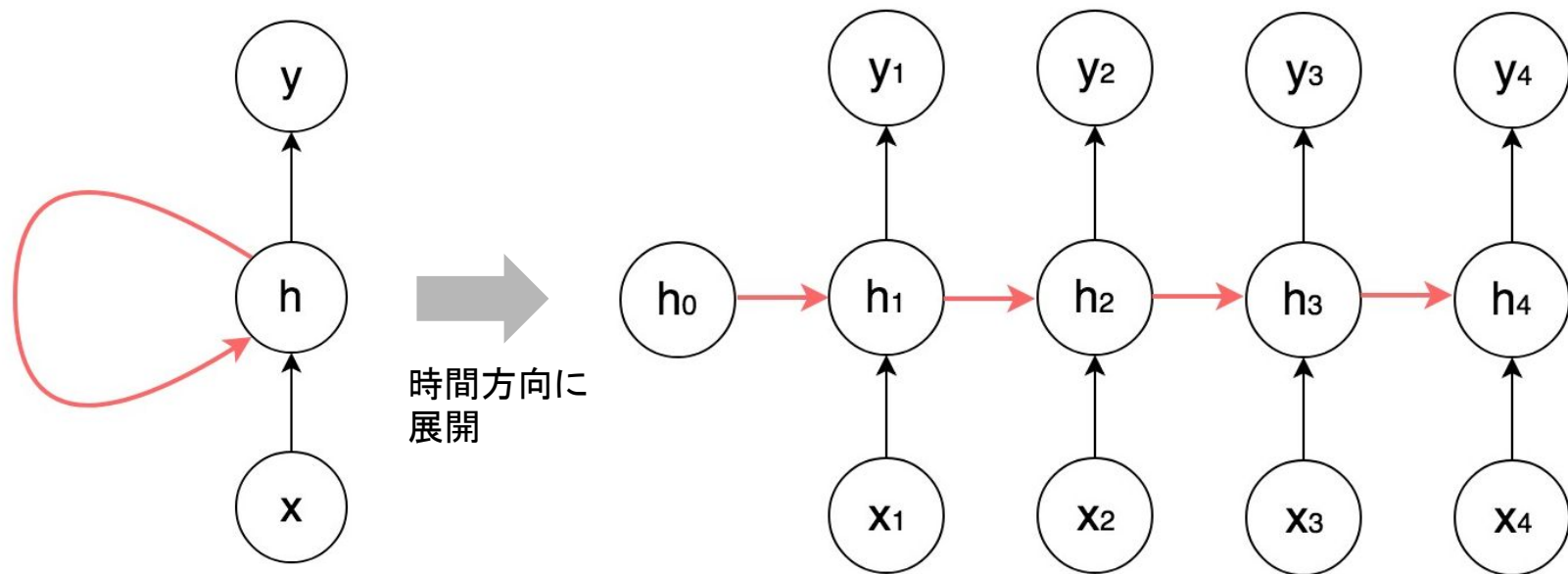
RNN(Recurrent Neural Network)とは

- ・過去の隠れ層の状態も入力に

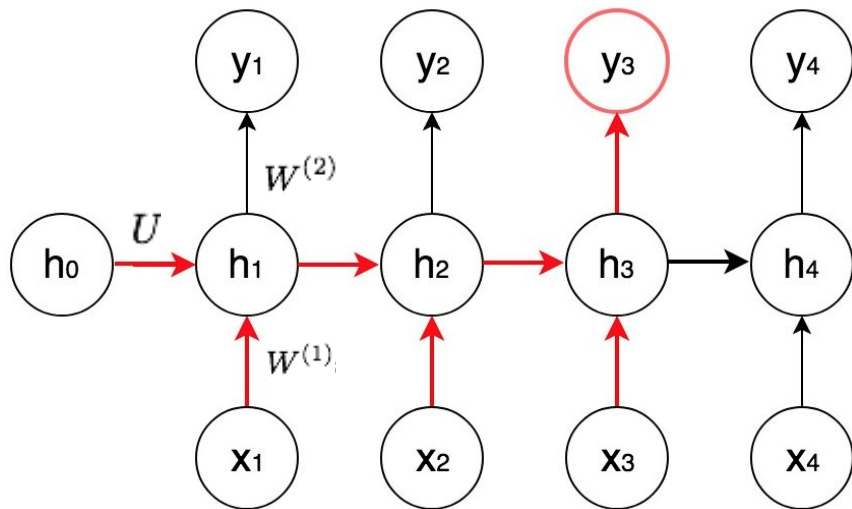


- ・系列データを扱うのに適したニューラルネットワーク
 - ・過去の情報を保持
 - ・可変長入力
 - ・例)「I like sports.」, 「I got up early this morning.」

RNNユニット



RNNの順伝播



$$h_t = f(W^{(1)}x_t + \underline{Uh_{t-1}} + b^{(1)})$$

$$y_t = g(W^{(2)}h_t + b^{(2)})$$

RNNの誤差逆伝播

- ・BPTT(BackPropagation Through Time)

- ・時間展開したうえでの誤差逆伝播

- ・Truncated BPTT

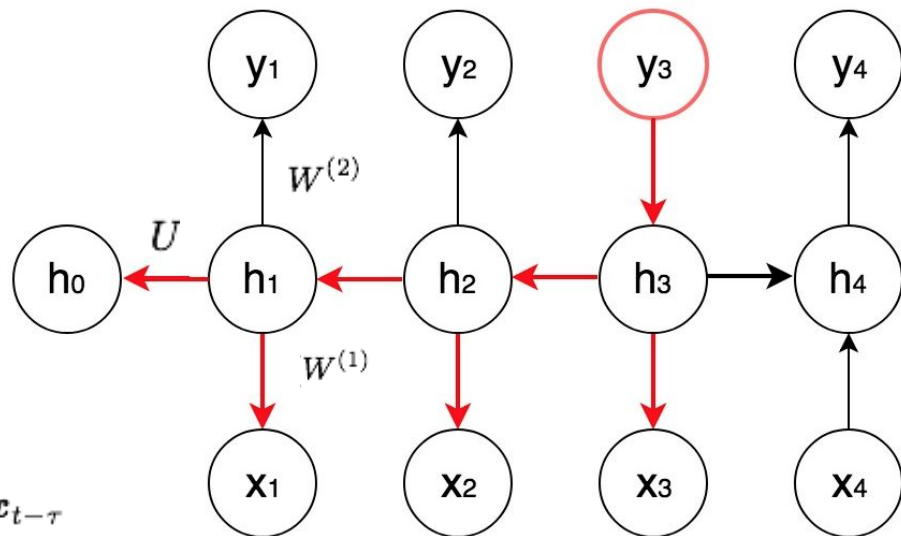
- ・遡るステップ数を限定したBPTT

例)

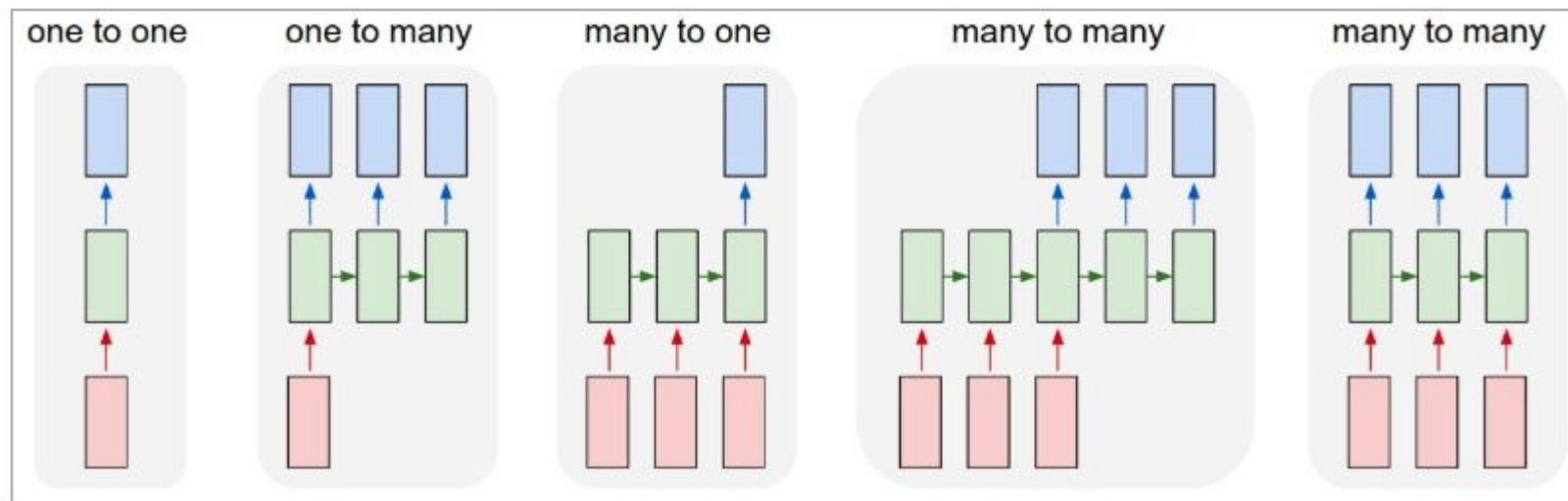
$$\frac{\partial L_t}{\partial W^{(2)}} = \frac{\partial L_t}{\partial \mathbf{v}_t} \frac{\partial \mathbf{v}_t}{\partial W^{(2)}} = \boldsymbol{\delta}_t^{(v)} \otimes \mathbf{h}_t$$

$$\frac{\partial L_t}{\partial W^{(1)}} = \sum_{\tau=0}^{t-1} \frac{\partial L_t}{\partial \mathbf{u}_{t-\tau}} \frac{\partial \mathbf{u}_{t-\tau}}{\partial W^{(1)}_{(t-\tau)}} = \sum_{\tau=0}^{t-1} \boldsymbol{\delta}_{t-\tau}^{(u)} \otimes \mathbf{x}_{t-\tau}$$

誤差 $L = \sum_{t=1}^T L_t(\mathbf{y}_t, \mathbf{t}_t)$



RNNモデルの応用例



画像認識

画像キャプション
生成

感情分析

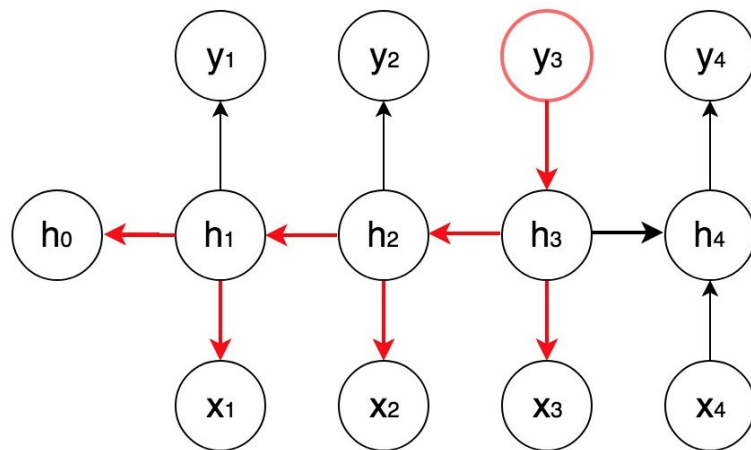
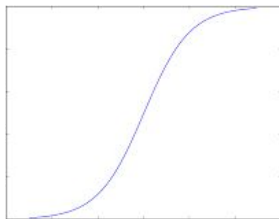
機械翻訳

動画分類

RNNの問題

- ・長期的な依存関係を学習できない。
- ・勾配消失(層を遡るごとに指数関数的に勾配が小さくなる)

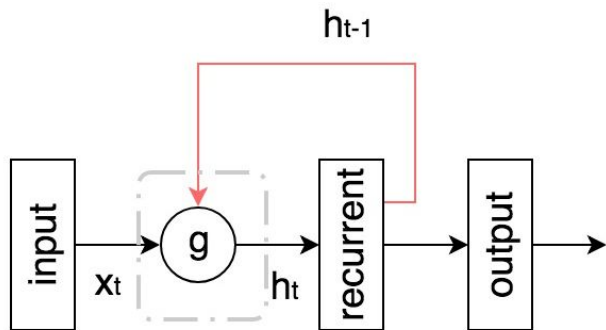
$$\begin{aligned}\delta_{t-1}^{(u)} &= \frac{\partial L_t}{\partial \mathbf{u}_{t-1}} = \frac{\partial L_t}{\partial \mathbf{u}_t} \frac{\partial \mathbf{u}_t}{\partial \mathbf{h}_{t-1}} \frac{\partial \mathbf{h}_{t-1}}{\partial \mathbf{u}_{t-1}} \\ &= U^T \delta_t^{(u)} \underline{f'(\mathbf{h}_{t-1})}\end{aligned}$$



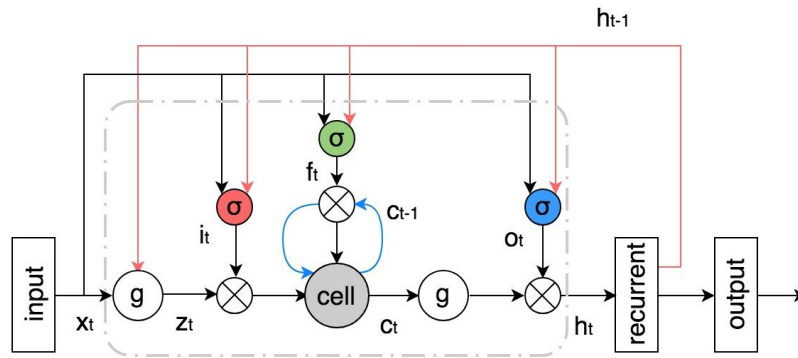
LSTM(Long Short Term Memory)とは

- ・RNNの派生モデル

- ・長期的な依存関係を学習できる(勾配消失を避ける)
- ・セルとゲートを導入



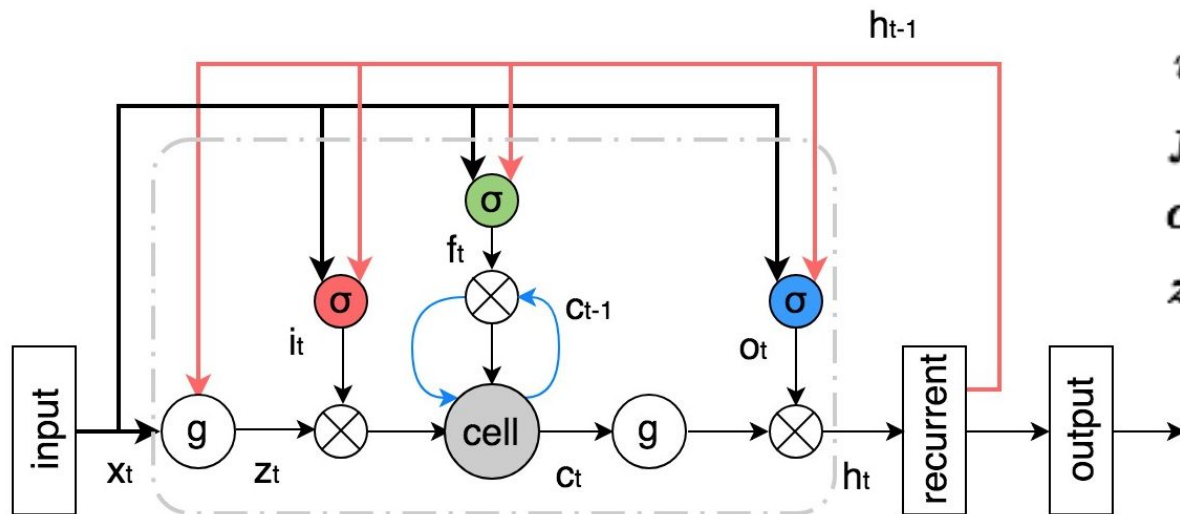
RNN



LSTM

LSTMの順伝播①

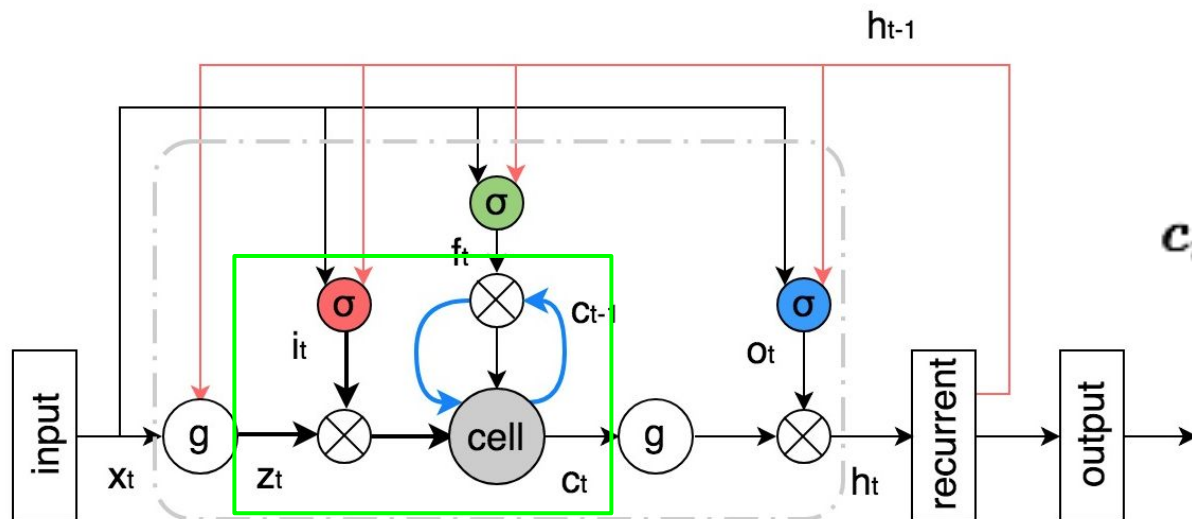
- ・ゲート(入力ゲート・忘却ゲート・出力ゲート)とセルへの入力計算



$$\begin{aligned}i_t &= \sigma(W_i x_t + U_i h_{t-1} + b_i) \\f_t &= \sigma(W_f x_t + U_f h_{t-1} + b_f) \\o_t &= \sigma(W_o x_t + U_o h_{t-1} + b_o) \\z_t &= g(W_z x_t + U_z h_{t-1} + b_z)\end{aligned}$$

LSTMの順伝播②

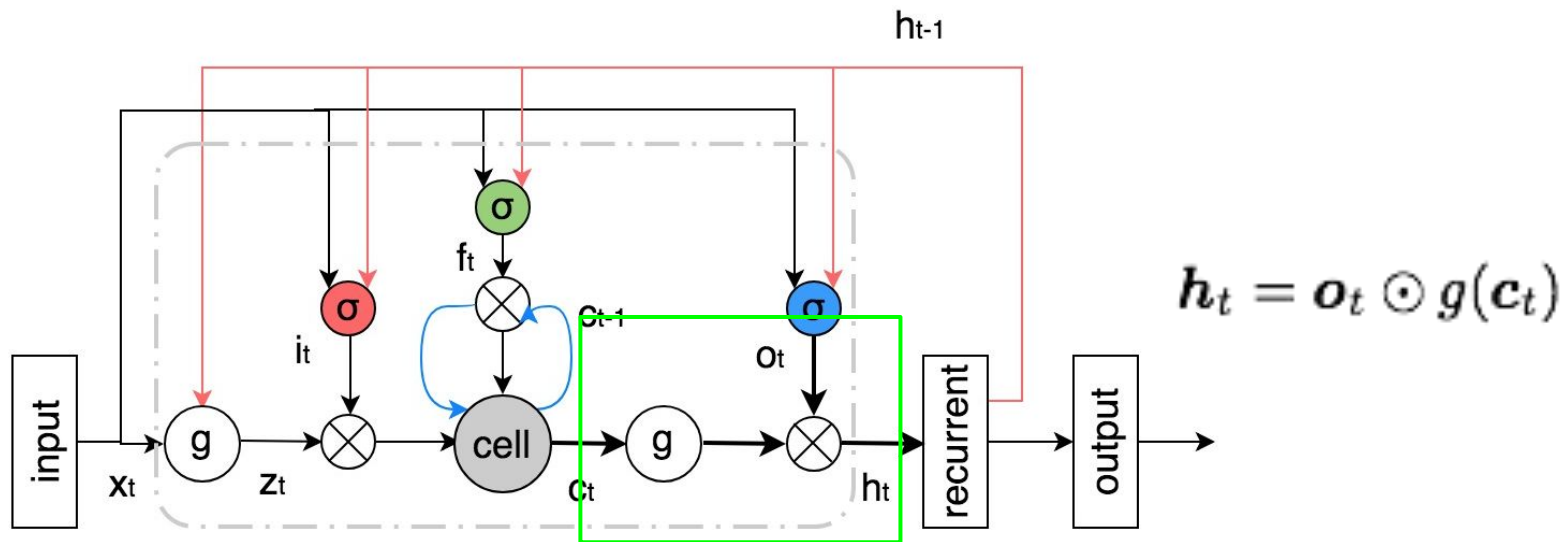
- ・セルの値を計算



$$c_t = f_t \odot c_{t-1} + i_t \odot z_t$$

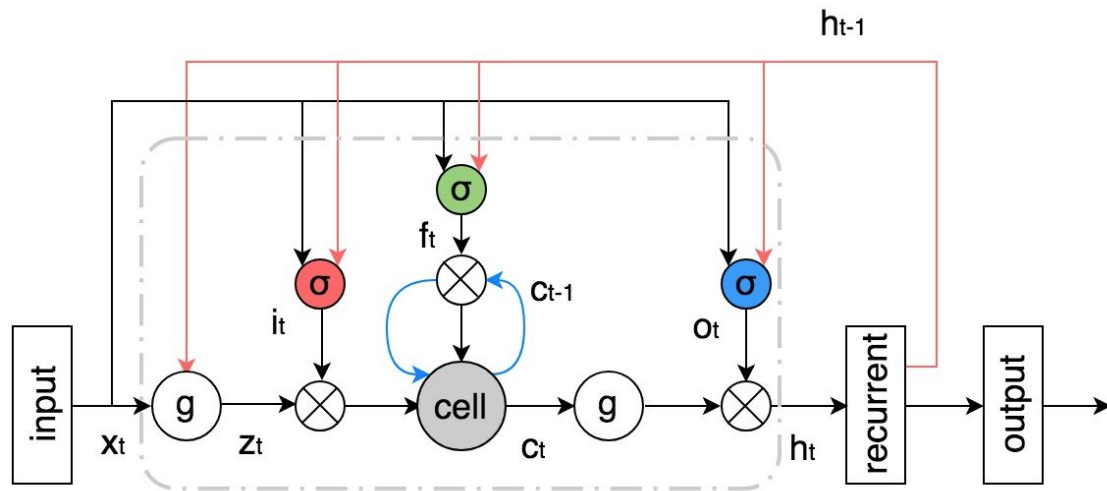
LSTMの順伝播③

- 出力値を計算



LSTM例

・空欄補充



「映画おもしろかったね。ところで、とてもお腹が空いたから、何か_____。」

食べよう
見に行こう
話そう

LSTMの誤差逆伝播

- ・RNNと同じくBPTT
- ・RNNに比べてパラメータ数が多い(4倍)
- ・勾配爆発を避けるために勾配クリッピングを使うことがある
 - ・大きい勾配に対してノルムをしきい値以下に正規化

まとめ

- ・RNN

- ・系列データを扱うのに適したニューラルネットワーク

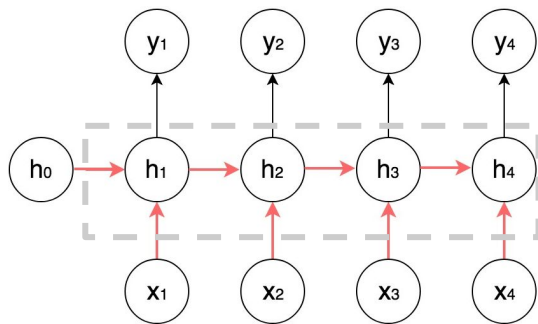
- ・LSTM

- ・RNNの派生
 - ・長期的な依存関係を学習できる
 - ・ほとんどのRNNモデルはLSTMを使用

Appendix

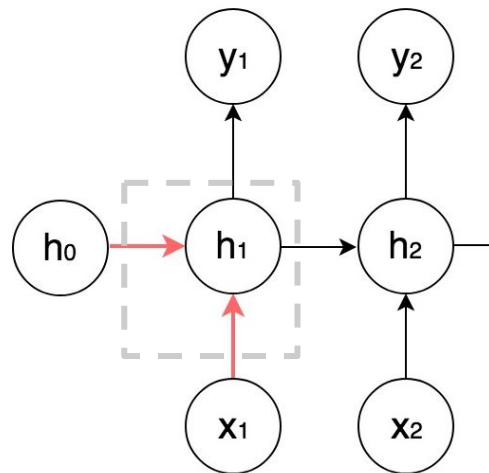
RNNのAPI(TensorFlow)

- `tf.contrib.rnn.BasicRNNCell`
(`num_units`, `activation=tanh`)
 - RNNセルを生成
- `tf.contrib.rnn.BasicLSTMCell`
(`num_units`, `forget_bias=1.0`, `activation=tanh`)
 - LSTMセルを生成
- `tf.contrib.rnn.static_rnn`
(`cell`, `inputs`, `initial_state=None`)
 - 時間方向に展開
 - 出力と最後の内部状態を返す



LSTMのAPI(Chainer)

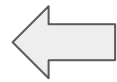
- `chainer.links.LSTM`
(`in_size`, `out_size`, `forget_bias_init=1.0`)
- LSTM層を生成
- `__call__(x)`
 - 内部状態を更新
 - LSTMの出力を返す



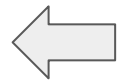
順伝播型ニューラルネットワークの誤差勾配

・テンソル積・転置となる過程

$$\frac{\partial L}{\partial W^{(2)}} = \frac{\partial L}{\partial \mathbf{v}} \frac{\partial \mathbf{v}}{\partial W^{(2)}} = \boldsymbol{\delta}^{(v)} \otimes \mathbf{h}$$


$$\frac{\partial L}{\partial W_{ij}^{(2)}} = \frac{\partial L}{\partial v_i} \frac{\partial v_i}{\partial W_{ij}^{(2)}} = \delta_i^{(v)} h_j$$

$$\boldsymbol{\delta}^{(u)} = \frac{\partial L}{\partial \mathbf{u}} = \frac{\partial L}{\partial \mathbf{v}} \frac{\partial \mathbf{v}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{u}} = (W^{(2)})^T \boldsymbol{\delta}^{(v)} f'(\mathbf{u})$$


$$\begin{aligned} \delta_i^{(u)} &= \frac{\partial L}{\partial u_i} = \sum_j \left(\frac{\partial L}{\partial v_j} \frac{\partial v_j}{\partial h_i} \right) \frac{\partial h_i}{\partial u_i} \\ &= \sum_j (\delta_j^{(v)} W_{ji}^{(2)}) f'(u_i) \end{aligned}$$

RNNのパラメータについての誤差勾配

順伝播

$$\mathbf{u}_t = W^{(1)} \mathbf{x}_t + U \mathbf{h}_{t-1} + \mathbf{b}^{(1)}$$

$$\mathbf{h}_t = f(\mathbf{u}_t)$$

$$\mathbf{v}_t = W^{(2)} \mathbf{h}_t + \mathbf{b}^{(2)}$$

$$\mathbf{y}_t = g(\mathbf{v}_t)$$

全体誤差は各時刻での誤差の和

$$L = \sum_{t=1}^T L_t(\mathbf{y}_t, \mathbf{t}_t)$$

u, vについての
誤差勾配

$$\delta_t^{(v)} = \frac{\partial L_t}{\partial \mathbf{v}_t} = \frac{\partial L_t}{\partial \mathbf{y}_t} \frac{\partial \mathbf{y}_t}{\partial \mathbf{v}_t} = \frac{\partial L_t}{\partial \mathbf{y}_t} g'(\mathbf{v}_t)$$

$$\delta_t^{(u)} = \frac{\partial L_t}{\partial \mathbf{u}_t} = \frac{\partial L_t}{\partial \mathbf{v}_t} \frac{\partial \mathbf{v}_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{u}_t} = (W^{(2)})^T \delta_t^{(v)} f'(\mathbf{u}_t)$$

1ステップ前のu
の誤差勾配

$$\begin{aligned} \delta_{t-1}^{(u)} &= \frac{\partial L_t}{\partial \mathbf{u}_{t-1}} = \frac{\partial L_t}{\partial \mathbf{u}_t} \frac{\partial \mathbf{u}_t}{\partial \mathbf{h}_{t-1}} \frac{\partial \mathbf{h}_{t-1}}{\partial \mathbf{u}_{t-1}} \\ &= U^T \delta_t^{(u)} f'(\mathbf{h}_{t-1}) \end{aligned}$$

RNNの各パラメータについての誤差勾配

$$\begin{aligned}\frac{\partial L_t}{\partial W^{(2)}} &= \frac{\partial L_t}{\partial \mathbf{v}_t} \frac{\partial \mathbf{v}_t}{\partial W^{(2)}} = \boldsymbol{\delta}_t^{(v)} \otimes \mathbf{h}_t \\ \frac{\partial L_t}{\partial \mathbf{b}^{(2)}} &= \frac{\partial L_t}{\partial \mathbf{v}_t} \frac{\partial \mathbf{v}_t}{\partial \mathbf{b}^{(2)}} = \boldsymbol{\delta}_t^{(v)}\end{aligned}$$

$$\frac{\partial L_t}{\partial W^{(1)}} = \sum_{\tau=0}^{t-1} \frac{\partial L_t}{\partial \mathbf{u}_{t-\tau}} \frac{\partial \mathbf{u}_{t-\tau}}{\partial W_{(t-\tau)}^{(1)}} = \sum_{\tau=0}^{t-1} \boldsymbol{\delta}_{t-\tau}^{(u)} \otimes \mathbf{x}_{t-\tau}$$

$$\frac{\partial L_t}{\partial U^{(1)}} = \sum_{\tau=0}^{t-1} \frac{\partial L_t}{\partial \mathbf{u}_{t-\tau}} \frac{\partial \mathbf{u}_{t-\tau}}{\partial U_{(t-\tau)}^{(1)}} = \sum_{\tau=0}^{t-1} \boldsymbol{\delta}_{t-\tau}^{(u)} \otimes \mathbf{h}_{t-\tau-1}$$

$$\frac{\partial L_t}{\partial \mathbf{b}^{(1)}} = \sum_{\tau=0}^{t-1} \frac{\partial L_t}{\partial \mathbf{u}_{t-\tau}} \frac{\partial \mathbf{u}_{t-\tau}}{\partial \mathbf{b}_{(t-\tau)}^{(1)}} = \sum_{\tau=0}^{t-1} \boldsymbol{\delta}_{t-\tau}^{(u)}$$

LSTMが勾配消失を避ける理由

cについての勾配は、1ステップ前においても、活性化関数の微分を通さず、忘却ゲートの値倍となる(忘却ゲートの値が1に近ければ誤差は消失しない)

$$\delta_t^{(c)} = \frac{\partial L_t}{\partial \mathbf{c}_t}$$
$$\delta_{t-1}^{(c)} = \frac{\partial L_t}{\partial \mathbf{c}_{t-1}} = \frac{\partial L_t}{\partial \mathbf{c}_t} \frac{\partial \mathbf{c}_t}{\partial \mathbf{c}_{t-1}} = \delta_t^{(c)} \mathbf{f}_t$$

各パラメータについての勾配は、cについての勾配を元に計算できるため、cについての勾配が消失しないなら、より以前の情報も考慮できる。(以下は例として W_i と W_o について具体的に計算)

$$\frac{\partial L_t}{\partial W_i} = \sum_{\tau=0}^{t-1} \frac{\partial L_t}{\partial \mathbf{c}_{t-\tau}} \frac{\partial \mathbf{c}_{t-\tau}}{\partial \mathbf{i}_{t-\tau}} \frac{\partial \mathbf{i}_{t-\tau}}{\partial W_{i,(t-\tau)}}$$
$$= \sum_{\tau=0}^{t-1} (\delta_{t-\tau}^{(c)} \mathbf{z}_{t-\tau}) \otimes \mathbf{x}_{t-\tau}$$

$$\delta_t^{(h)} = \frac{\partial L_t}{\partial \mathbf{h}_t} = \frac{\partial L_t}{\partial \mathbf{c}_t} \frac{\partial \mathbf{c}_t}{\partial \mathbf{h}_t} = \delta_t^{(c)} \mathbf{o}_t \frac{1}{\cosh^2(\mathbf{c}_t)}$$

$$\frac{\partial L_t}{\partial W_o} = \sum_{\tau=0}^{t-1} \frac{\partial L_t}{\partial \mathbf{h}_{t-\tau}} \frac{\partial \mathbf{h}_{t-\tau}}{\partial \mathbf{o}_{t-\tau}} \frac{\partial \mathbf{o}_{t-\tau}}{\partial W_{o,(t-\tau)}}$$
$$= \sum_{\tau=0}^{t-1} (\delta_{t-\tau}^{(h)} \tanh(\mathbf{c}_{t-\tau})) \otimes \mathbf{x}_{t-\tau}$$