

1. section1_再起型ニューラルネットワークの概念

1. 1. RNN の概念

RNN とは時系列データに対応可能な、ニューラルネットワークである。

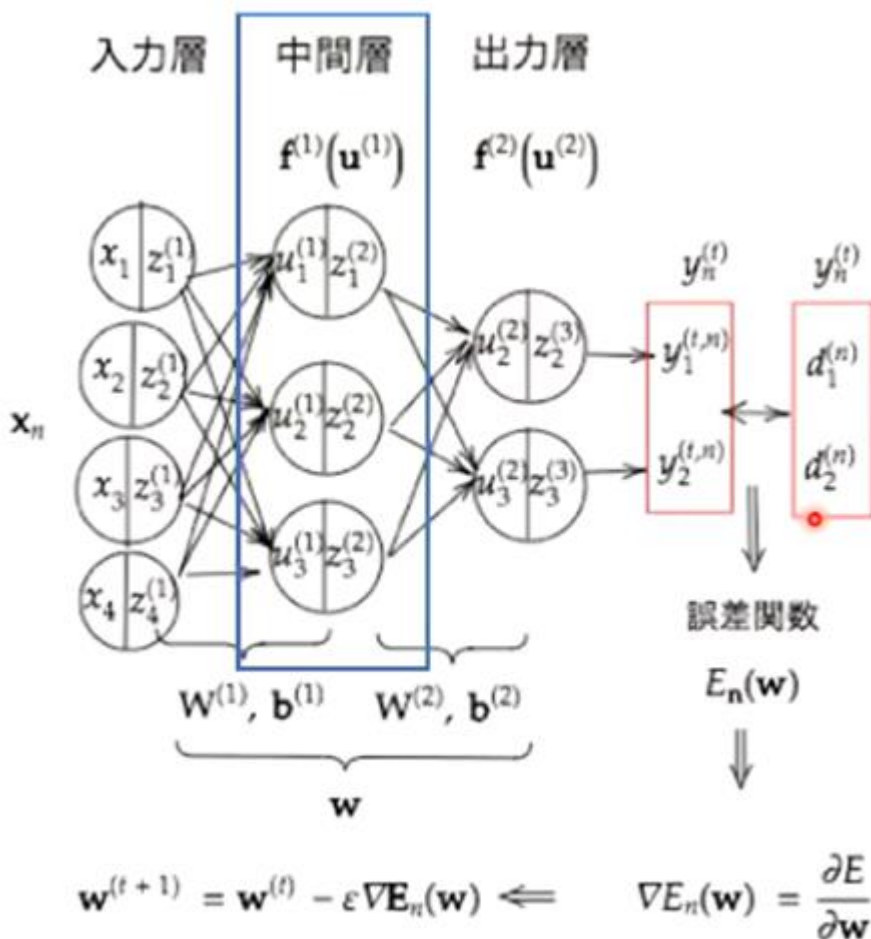
1. 2. 時系列データとは

時間的順序を追って一定間隔毎に観察され、しかも相互に統計的依存関係が認められるようなデータの系列である。

時系列データには

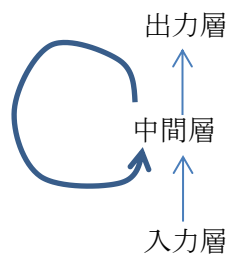
- (1) 音声データ・株価：時系列データ（時間的なつながりがあるデータ）
- (2) テキストデータ：単語 1 つ 1 つを前の単語に繋がって次の単語が出てくる。

1. 3. RNN について



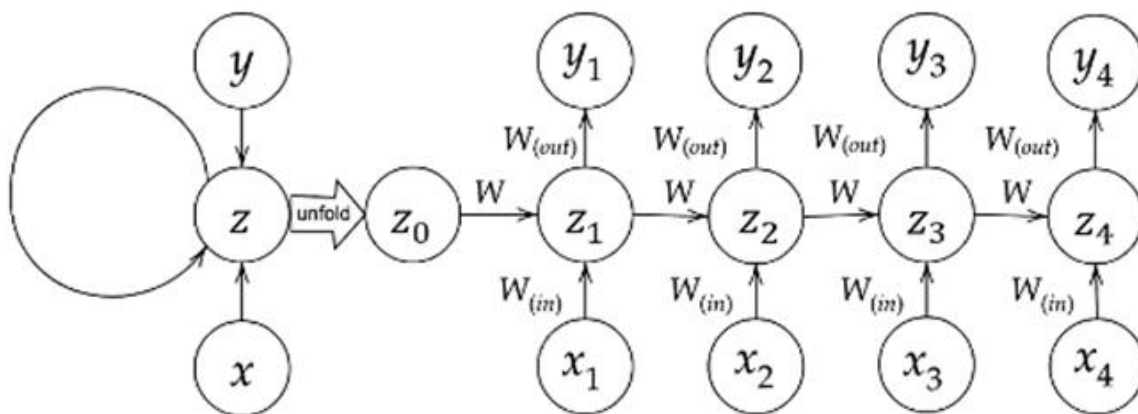
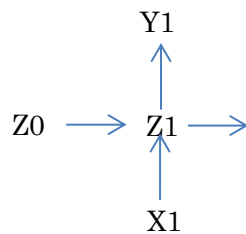
中間層：特徴を取得する。

1. 4. RNN の全体像



①入力層→中間層→出力層

②中間層→中間層



重みが以下の3か所に存在します。

- (1) W_{in} : 入力層から中間層への重み
- (2) W_{out} : 中間層から出力層への重み
- (3) W : 前の層からの入力用重み (中間層から中間層)

1. 5. RNN の数学的記述

$$\begin{aligned}u^t &= W_{(in)}x^t + Wz^{t-1} + b \\z^t &= f(W_{(in)}x^t + Wz^{t-1} + b) \\v^t &= W_{(out)}z^t + c \\y^t &= g(W_{(out)}z^t + c)\end{aligned}$$

```
u[:,t+1] = np.dot(X, W_in) + np.dot(z[:,t].reshape(1, -1), W)
z[:,t+1] = functions.sigmoid(u[:,t+1])
```

前の中間層

(1) RNN の順伝播

$f()$ 、 $g()$ が活性化関数

$$z = f(u^t)$$

$$y = g(v^t)$$

```
np.dot(z[:,t+1].reshape(1, -1), W_out)
y[:,t] = functions.sigmoid(np.dot(z[:,t+1].reshape(1, -1), W_out))
```

1. 6. RNN の特徴

時系列モデルを扱うには、初期の状態と過去の時間 ($t-1$) の状態を保持し、そこから次の時間での t を再帰的に求める再起構造が必要になる。

再起構造

Simple RNN

ベクトル加算

右から左に時間的情報がある。

1896

2進数 (A) 1 0 1 1 0 1 1 0
 (B) 0 1 1 1 0 1 0 1
 1 0 1 0 1 1

←

来 (12)

10 → 2
 0 → 0
 1 → 1
 2 → 10
 3 → 11
 4 → 100
 5 → 101

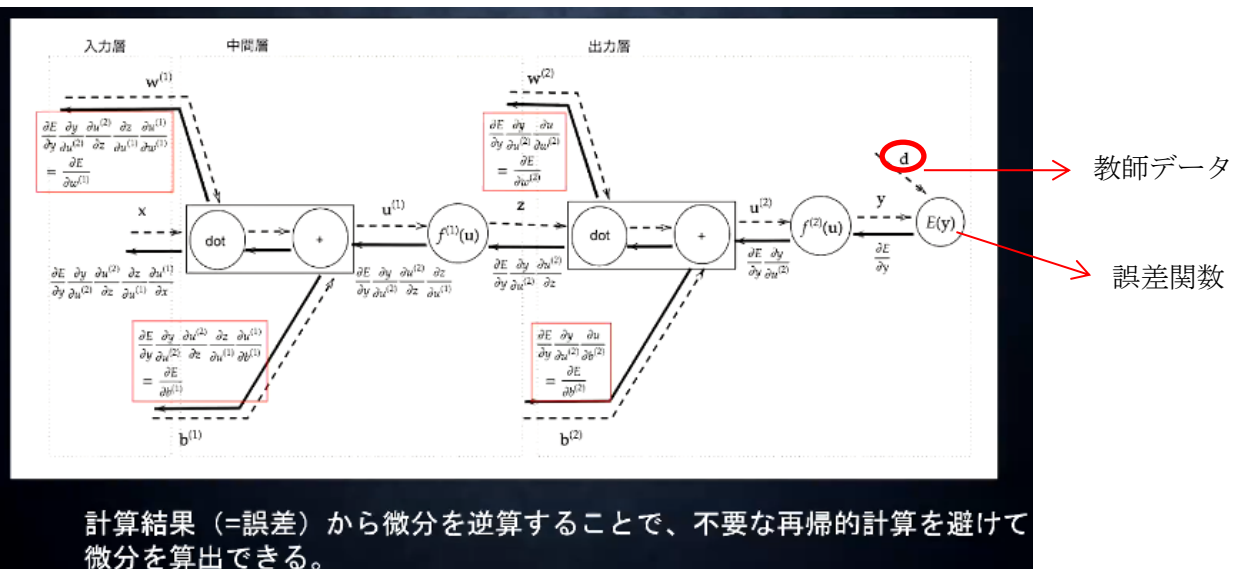
1. 7. BPTT

1. 7. 1 BPTT とは

BackPropagation Through Time

RNN におけるパラメータ調整方法の一種（誤差逆伝播の一種）

1) 誤差逆伝播



微分の連鎖律 → W, B を求める。

$$E(\mathbf{y}) = \frac{1}{2} \sum_{j=1}^J (y_j - d_j)^2 = \frac{1}{2} \|\mathbf{y} - \mathbf{d}\|^2 \quad : \text{誤差関数} = \text{二乗誤差関数}$$

$$\mathbf{y} = \mathbf{u}^{(L)} \quad : \text{出力層の活性化関数} = \text{恒等写像}$$

$$\mathbf{u}^{(l)} = \mathbf{w}^{(l)} \mathbf{z}^{(l-1)} + \mathbf{b}^{(l)} \quad : \text{総入力 of 計算}$$

微分の連鎖率の式: $\frac{\partial E}{\partial w_{ji}^{(2)}} = \frac{\partial E}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial w_{ji}^{(2)}}$ ここでは出力層の w について微分したい

各関数を微分

$$\frac{\partial \mathbf{y}(\mathbf{u})}{\partial \mathbf{u}} = \frac{\partial \mathbf{u}}{\partial \mathbf{u}} = 1$$

$$\frac{\partial E(\mathbf{y})}{\partial \mathbf{y}} = \frac{\partial}{\partial \mathbf{y}} \frac{1}{2} \|\mathbf{y} - \mathbf{d}\|^2 = \mathbf{y} - \mathbf{d}$$

$$\frac{\partial \mathbf{u}(\mathbf{w})}{\partial w_{ji}^{(2)}} = \frac{\partial}{\partial w_{ji}^{(2)}} (\mathbf{w}^{(l)} \mathbf{z}^{(l-1)} + \mathbf{b}^{(l)}) = \frac{\partial}{\partial w_{ji}^{(2)}} \left(\begin{bmatrix} w_{11}z_1 + \dots + w_{1j}z_j + \dots + w_{1I}z_I \\ \vdots \\ w_{j1}z_1 + \dots + w_{ji}z_j + \dots + w_{jI}z_I \\ \vdots \\ w_{I1}z_1 + \dots + w_{Ij}z_j + \dots + w_{II}z_I \end{bmatrix} + \begin{bmatrix} b_1 \\ \vdots \\ b_j \\ \vdots \\ b_I \end{bmatrix} \right) = \begin{bmatrix} 0 \\ \vdots \\ z_j \\ \vdots \\ 0 \end{bmatrix}$$

微分の連鎖式の3つの式が求められる。

誤差逆伝播法の復習

$$E(\mathbf{y}) = \frac{1}{2} \sum_{j=1}^J (y_j - d_j)^2 = \frac{1}{2} \|\mathbf{y} - \mathbf{d}\|^2 \quad : \text{誤差関数} = \text{二乗誤差関数}$$

$$\mathbf{y} = \mathbf{u}^{(L)} \quad : \text{出力層の活性化関数} = \text{恒等写像}$$

$$\mathbf{u}^{(l)} = \mathbf{w}^{(l)} \mathbf{z}^{(l-1)} + \mathbf{b}^{(l)} \quad : \text{総入力 of 計算}$$

$$\frac{\partial E}{\partial w_{ji}^{(2)}} = \frac{\partial E}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial w_{ji}^{(2)}} = (\mathbf{y} - \mathbf{d}) \cdot \begin{bmatrix} 0 \\ \vdots \\ z_j \\ \vdots \\ 0 \end{bmatrix} = (y_j - d_j) z_j$$

(2) 確認

$$Z=t**2 \quad t=x+y$$

$$Z=(x+y)**2$$

$$Z=x**2+2xy+y**2$$

$$dz/dx=2x+2y$$

1. 7. 2 BPTT の数学的記述

(1) BPTT の数学的記述 1

$$\begin{aligned}\frac{\partial E}{\partial W_{(in)}} &= \frac{\partial E}{\partial u^t} \left[\frac{\partial u^t}{\partial W_{(in)}} \right]^T = \delta^t [x^t]^T \\ \frac{\partial E}{\partial W_{(out)}} &= \frac{\partial E}{\partial v^t} \left[\frac{\partial v^t}{\partial W_{(out)}} \right]^T = \delta^{out,t} [z^t]^T \\ \frac{\partial E}{\partial W} &= \frac{\partial E}{\partial u^t} \left[\frac{\partial u^t}{\partial W} \right]^T = \delta^t [z^{t-1}]^T \\ \frac{\partial E}{\partial b} &= \frac{\partial E}{\partial u^t} \frac{\partial u^t}{\partial b} = \delta^t \\ \frac{\partial E}{\partial c} &= \frac{\partial E}{\partial v^t} \frac{\partial v^t}{\partial c} = \delta^{out,t}\end{aligned}$$

3つに対して微分する

$$\frac{\partial E}{\partial W_{(in)}} = \frac{\partial E}{\partial u^t} \left[\frac{\partial u^t}{\partial W_{(in)}} \right]^T = \delta^t [x^t]^T$$

```
np.dot(X.T, delta[:,t].reshape(1,-1))
```

$$\frac{\partial E}{\partial W_{(out)}} = \frac{\partial E}{\partial v^t} \left[\frac{\partial v^t}{\partial W_{(out)}} \right]^T = \delta^{out,t} [z^t]^T$$

```
np.dot(z[:,t+1].reshape(-1,1), delta_out[:,t].reshape(-1,1))
```

$$\frac{\partial E}{\partial W} = \frac{\partial E}{\partial u^t} \left[\frac{\partial u^t}{\partial W} \right]^T = \delta^t [z^{t-1}]^T$$

```
np.dot(z[:,t].reshape(-1,1), delta[:,t].reshape(1,-1))
```

$$\delta = \nabla E / \nabla u \ast \ast t$$

X.T は時間の意味（転置とは違う）

2) BPTT の数学的記述 2

$$u^t = W_{(in)}x^t + W z^{t-1} + b$$
$$z^t = f(W_{(in)}x^t + W z^{t-1} + b)$$
$$v^t = W_{(out)}z^t + c$$
$$y^t = g(W_{(out)}z^t + c)$$

順伝播の式

$$u^t = W_{(in)}x^t + W z^{t-1} + b$$

```
u[:,t+1] = np.dot(X, W_in) + np.dot(z[:,t].reshape(1, -1), W)
```

$$z^t = f(W_{(in)}x^t + W z^{t-1} + b)$$

```
z[:,t+1] = functions.sigmoid(u[:,t+1])
```

$$v^t = W_{(out)}z^t + c$$

```
np.dot(z[:,t+1].reshape(1, -1), W_out)
```

$$y^t = g(W_{(out)}z^t + c)$$

```
y[:,t] = functions.sigmoid(np.dot(z[:,t+1].reshape(1, -1), W_out))
```

$$\frac{\partial E}{\partial b} = \frac{\partial E}{\partial u^t} \frac{\partial u^t}{\partial b} = \delta^t \quad \frac{\partial u^t}{\partial b} = 1$$

コード記載なし（本講座のsimple RNNコードでは簡略化のため省略）

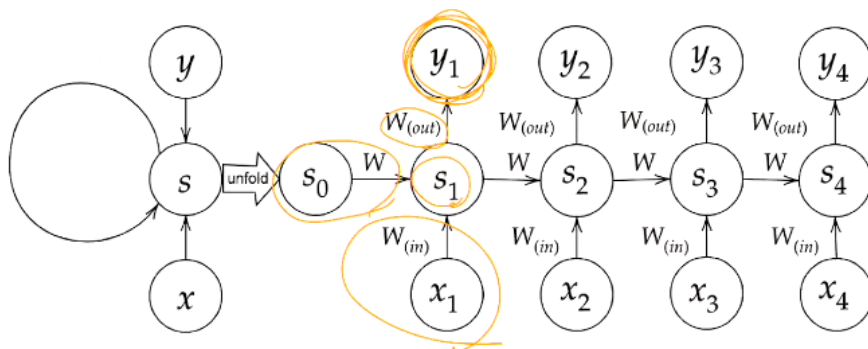
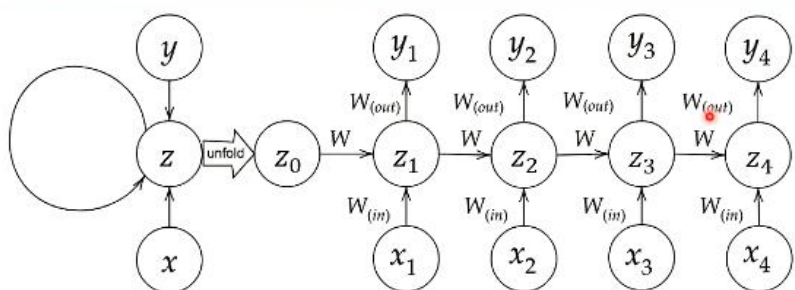
$$\frac{\partial E}{\partial c} = \frac{\partial E}{\partial v^t} \frac{\partial v^t}{\partial c} = \delta^{out,t} \quad \frac{\partial v^t}{\partial c} = 1$$

コード記載なし（本講座のsimple RNNコードでは簡略化のため省略）

バイアスの微分式

●確認

下図の y_1 を $x \cdot z_0 \cdot z_1 \cdot w_{in} \cdot w \cdot w_{out}$ を用いて数式で表せ。
 ※バイアスは任意の文字で定義せよ。
 ※また中間層の出力にシグモイド関数 $g(x)$ を作用させよ。
 (7分)



$$\begin{cases} s_1 = x_1 \cdot w_{in} + s_0 \cdot w + b \\ y_1 = g(w_{out} \cdot s_1) \end{cases}$$

$$y_1 = g(w_{out} \cdot s_1 + c)$$

$$s_1 = f(w_{in} \cdot x_1 + w \cdot s_0 + b)$$

(3) BPTT の数学的記述 3

$$\frac{\partial E}{\partial u^t} = \frac{\partial E}{\partial v^t} \frac{\partial v^t}{\partial u^t} = \frac{\partial E}{\partial v^t} \frac{\partial \{W_{(out)} f(u^t) + c\}}{\partial u^t} = f'(u^t) W_{(out)}^T \delta^{out,t} = \delta^t$$

$$\delta^{t-1} = \frac{\partial E}{\partial u^{t-1}} = \frac{\partial E}{\partial u^t} \frac{\partial u^t}{\partial u^{t-1}} = \delta^t \left\{ \frac{\partial u^t}{\partial z^{t-1}} \frac{\partial z^{t-1}}{\partial u^{t-1}} \right\} = \delta^t \{ W f'(u^{t-1}) \}$$

$$\delta^{t-z-1} = \delta^{t-z} \{ W f'(u^{t-z-1}) \}$$

```
delta[:,t] = (np.dot(delta[:,t+1].T, W.T) + np.dot(delta_out[:,t].T, W_out.T)) *
functions.d_sigmoid(u[:,t+1])
```

時間的つながりがある

t 時間目と t - 1 時間目とは関係がある。

(4) BPTT の数学的記述 4

パラメータの更新式

$$W_{(in)}^{t+1} = W_{(in)}^t - \epsilon \frac{\partial E}{\partial W_{(in)}} = W_{(in)}^t - \epsilon \sum_{z=0}^{T_t} \delta^{t-z} [x^{t-z}]^T$$

$$W_{(out)}^{t+1} = W_{(out)}^t - \epsilon \frac{\partial E}{\partial W_{(out)}} = W_{(out)}^t - \epsilon \delta^{out,t} [z^t]^T$$

$$W^{t+1} = W^t - \epsilon \frac{\partial E}{\partial W} = W_{(in)}^t - \epsilon \sum_{z=0}^{T_t} \delta^{t-z} [z^{t-z-1}]^T$$

$$b^{t+1} = b^t - \epsilon \frac{\partial E}{\partial b} = b^t - \epsilon \sum_{z=0}^{T_t} \delta^{t-z}$$

$$c^{t+1} = c^t - \epsilon \frac{\partial E}{\partial c} = c^t - \epsilon \delta^{out,t}$$

重み

バイアス

W o u t の更新式には時間的に遡る処理はない

W i n には時間的遡りがあるので Σ がある。

ϵ : 学習率

$$W_{(in)}^{t+1} = W_{(in)}^t - \epsilon \frac{\partial E}{\partial W_{(in)}} = W_{(in)}^t - \epsilon \sum_{z=0}^{T_t} \delta^{t-z} [x^{t-z}]^T$$

W_in -= learning_rate * W_in_grad

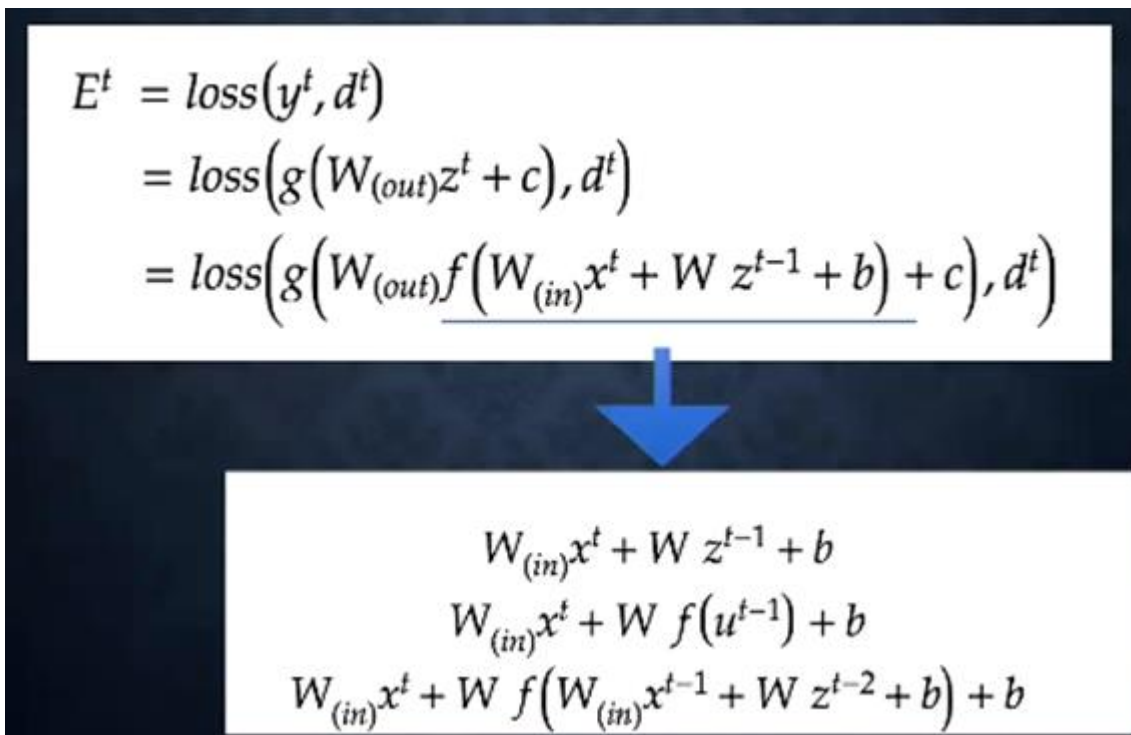
$$W_{(out)}^{t+1} = W_{(out)}^t - \epsilon \frac{\partial E}{\partial W_{(out)}} = W_{(out)}^t - \epsilon \delta^{out,t} [z^t]^T$$

W_out -= learning_rate * W_out_grad

$$W^{t+1} = W^t - \epsilon \frac{\partial E}{\partial W} = W_{(in)}^t - \epsilon \sum_{z=0}^{T_t} \delta^{t-z} [z^{t-z-1}]^T$$

W -= learning_rate * W_grad

1. 7. 3 BPTT の全体像



誤差関数は y と d から求まる。

$Z \ t \rightarrow X \ t-1 \rightarrow Z \ t-2 \dots$ 過去の繋がりがある

(1) プログラム 3_1_simple_RNN_after

(2) コードの演習問題

コード演習問題

```
def bptt(xs, ys, W, U, V):
    """
    ys: labels, (batch_size, n_seq, output_size)
    """
    # 前向き
    hiddens, outputs = rnn_net(xs, W, U, V)
    # 損失関数のW, U, Vに関する偏微分
    dW = np.zeros_like(W) # dL/dW
    dU = np.zeros_like(U) # dL/dU
    dV = np.zeros_like(V) # dL/dV
    # 損失関数の出力値に関する偏微分
    do = _calculate_do(outputs, ys) # dL/do, (batch_size, n_seq, output_size)

    batch_size, n_seq = ys.shape[:2]
    # 時間を逆方向にたどり、パラメータの偏微分を計算 (バックプロパゲーション)
    # dL/dV = do/dV * dL/do = h * dL/do
    # dL/dU = do/dU * dL/do
    # dL/dW = do/dW * dL/do
    # do/dW = (dh_t/dW * d/dh_t + dh_{t-1}/dW * d/dh_{t-1} + ...) o_t(t)
    # = (x_t * x_{t-1} * U * x_{t-2} * U^2 + ...) * V
    # do/dU = (dh_t/dU * d/dh_t + dh_{t-1}/dU * d/dh_{t-1} + ...) o_t(t)
    # = (h_{t-1} + h_{t-2} * U + h_{t-3} * U^2 + ...) * V
    for t in reversed(range(n_seq)):
        dv += np.dot(do[:, t], hiddens[:, t]) / batch_size
        delta_t = do[:, t].dot(V)
        # 時刻tの出力は時刻t以前の中間層すべてに依存するため
        # W, Uはさらに遡って計算
        for bptt_step in reversed(range(t+1)):
            dw += np.dot(delta_t.T, xs[:, bptt_step]) / batch_size
            du += np.dot(delta_t.T, hiddens[:, bptt_step]) / batch_size
            delta_t = (お)
    return dW, dU, dV
```

中間層の過力の出力

(回答: (2) $\text{delta_t} \cdot \text{dot}(U)$)

【解説】

RNNでは中間層出力 $h_{\{t\}}$ が過去の中間層出力 $h_{\{t-1\}}, \dots, h_{\{1\}}$ に依存する。

RNNにおいて損失関数を重み W や U に関して偏微分するときは、それを考慮する必要があり、 $dh_{\{t\}}/dh_{\{t-1\}} = U$ 。

であることに注意すると、過去に遡るたびに U が掛けられる。つまり、 $\text{delta_t} = \text{delta_t} \cdot \text{dot}(U)$ となる。

- (1) $\text{delta_t} \cdot \text{dot}(W)$
- (2) $\text{delta_t} \cdot \text{dot}(U)$
- (3) $\text{delta_t} \cdot \text{dot}(V)$
- (4) $\text{delta_t} * V$

2. section2_LSTM

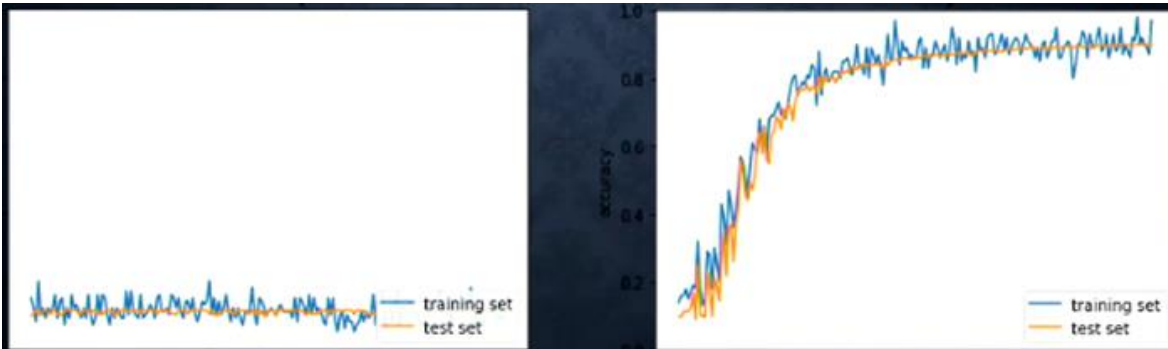
2. 1 RNN の課題

時系列を遡れば遡るほど、勾配が消失していくーー>長い時系列の学習が困難
(解決策)

以前の勾配消失とは別で、構造自体を変えて解決したものが **LSTM** である。

2. 1. 1 勾配消失問題の復習

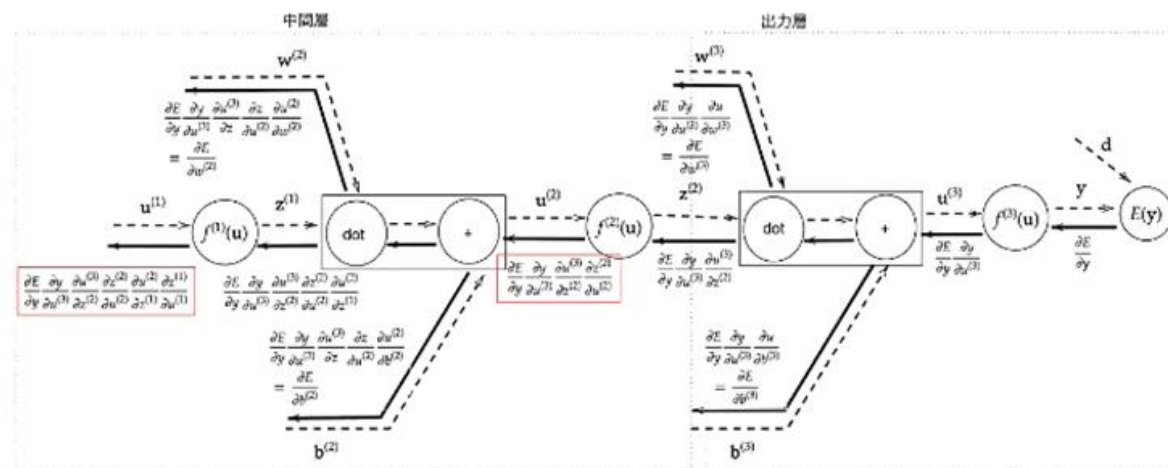
誤差逆伝播法が階層に進んでいくに連れて、勾配がどんどん緩やかになっていく。そのため、勾配降下法による更新では、階層のパラメータはほとんど変わらず、訓練は最適値に収束しなくなる。



うまく進まない

うまく進む

2. 1. 2 勾配消失問題のビジョン



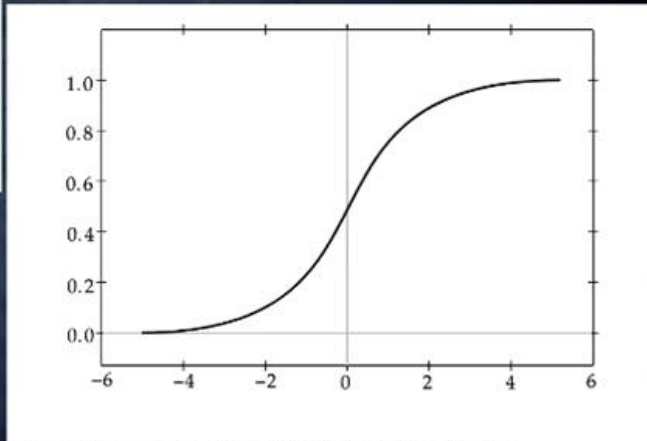
(1) 活性化関数：シグモイド関数

サンプルコード

```
def sigmoid(x):  
    return 1/(1 + np.exp(-x))
```

数式

$$f(u) = \frac{1}{1 + e^{-u}}$$



0 ~ 1の間を緩やかに変化する関数で、ステップ関数ではON/OFFしかない状態に対し、信号の強弱を伝えられるようになり、予想ニューラルネットワーク普及のきっかけとなった。

課題

大きな値では出力の変化が微小なため、勾配消失問題を引き起こす事があった。

$$\frac{\partial E}{\partial y} \frac{\partial y}{\partial u^{(3)}} \frac{\partial u^{(3)}}{\partial z^{(2)}} \boxed{\frac{\partial z^{(2)}}{\partial u^{(2)}}}$$

$$(1 - \text{sigmoid}(x)) \cdot \text{sigmoid}(x)$$

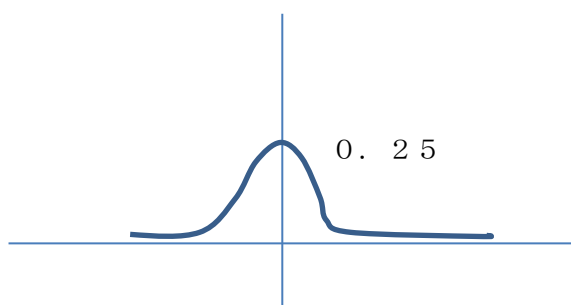
$$\frac{\partial E}{\partial y} \frac{\partial y}{\partial u^{(3)}} \frac{\partial u^{(3)}}{\partial z^{(2)}} \boxed{\frac{\partial z^{(2)}}{\partial u^{(2)}}} \frac{\partial u^{(2)}}{\partial z^{(1)}} \boxed{\frac{\partial z^{(1)}}{\partial u^{(1)}}}$$

$$(1 - \text{sigmoid}(x)) \cdot \text{sigmoid}(x)$$

$$(1 - \text{sigmoid}(x)) \cdot \text{sigmoid}(x)$$

(2) 確認

シグモイド関数を微分した時最大値が0.25となる。

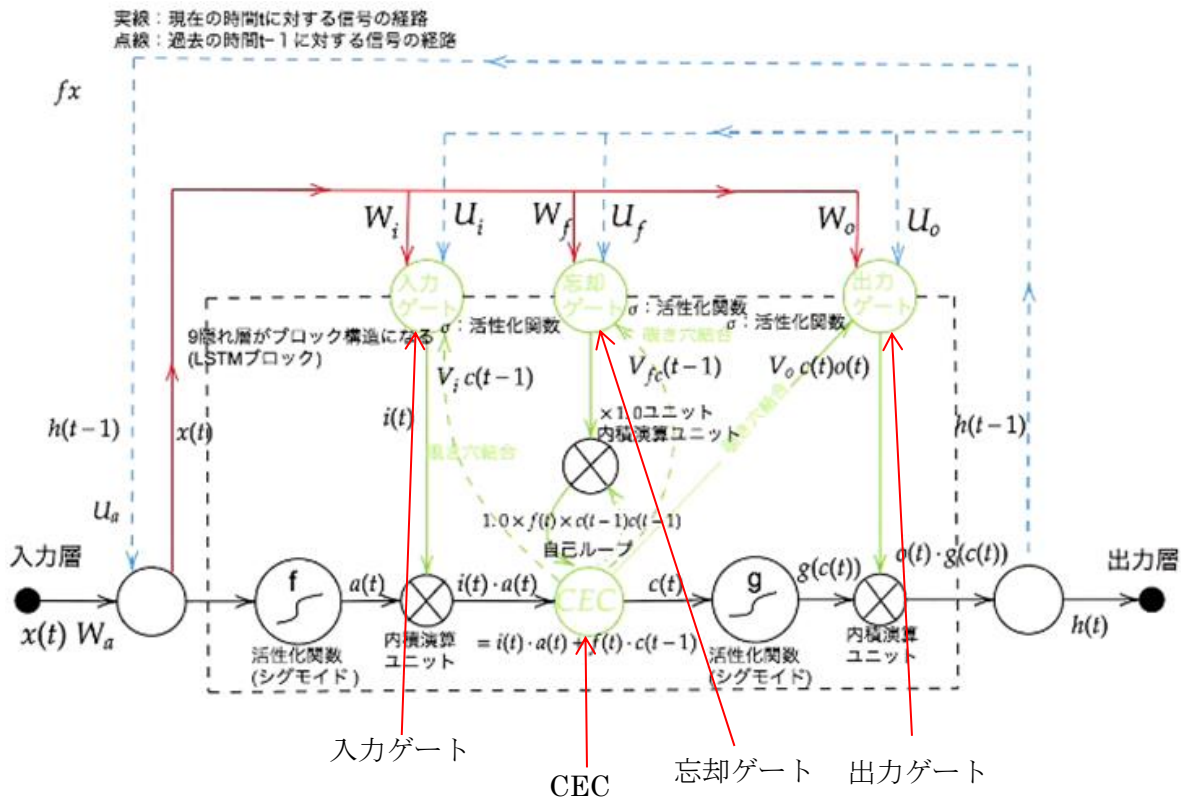


2. 1. 3 勾配爆発

勾配が、層を逆伝播するごとに指数関数的に大きくなっていくこと。
逆伝播で勾配が大きくなる場合、恒等関数等を使用する。

2. 1. 4 LSTM の全体図

Long Short-Time Memory



RNN の一種

青い点線：時間的なループ

黒い点線：中間層

赤い矢印：今回の入力値 (W)

$a(t) \cdot i(t) = c(t)$ ：今回の状態

2. 2 CEC

入力値や中間層の記憶を指せる。

勾配消失および勾配爆発の解決方法として、勾配が、1であれば解決できる。

$$\delta^{t-z-1} = \delta^{t-z} \{ W f'(u^{t-z-1}) \} = 1$$

$$\boxed{\frac{\partial E}{\partial c^{t-1}}} = \frac{\partial E}{\partial c^t} \frac{\partial c^t}{\partial c^{t-1}} = \frac{\partial E}{\partial c^t} \frac{\partial}{\partial c^{t-1}} \{ a^t - c^{t-1} \} = \boxed{\frac{\partial E}{\partial c^t}}$$

勾配を1にするために導入されたのが CEC

(1) CEC の課題

入力データについて、時間依存度に関係なく重みが一律である。

——>ニューラルネットワークの学習特性が無いということ。

入力層->隠れ層への重み：入力重み衝突

隠れ層->出力層への重み：出力重み衝突

CEC は覚えることしかしないので、周りに学習機能を配置する。

2. 3 入力ゲートと出力ゲート

CEC の取り巻き

入力ゲート：CEC に覚え方を支持する。

出力ゲート：CEC のデータをどのように使用するか。

入力・出力ゲートの役割とは？

入力・出力ゲートを追加することで、
それぞれのゲートへの入力値の重みを、
重み行列W,Uで可変可能とする。



CECの課題を解決。

2. 4 忘却ゲート

(1) LSTM の現状

CEC は過去の情報が全て保管されている。

(2) 課題

過去の情報がいらなくなった場合、削除することはできず保管され続ける。

(3) 解決策

過去の情報がいらなくなった場合、そのタイミングで情報を忘却する機能が必要。

——>忘却ゲートの誕生

いらなくなった情報を CEC から消去する。

(4) 演習問題

「映画面白かったね。ところで、とてもお腹が空いたから何か・・・。」

この文章を LSTM に入力し・・・に当てはまる単語を予測したいとする。

文中の「とても」という単語は空欄の予測においてなくなっても影響を及ぼさないと考えられる。

このような場合、どのゲートが作用するか——>忘却ゲート

2. 5 覗き穴結合

(1) 課題

CEC の保存されている過去の情報を、任意のタイミングで他のノードに伝播させたり、あるいは任意のタイミングで忘却させたい。

CEC 自身の値は、ゲート制御に影響を与えていない。

(2) 覗き穴結合とは

CEC 自身の値に、重み行列を介して伝播可能にした構造。

(3) 結果

あまり大きな改善は見られなかった。

3. section3_GRU

(1) LSTM の課題

LSTM ではパラメータ数が多く、計算負荷が高くなる問題があった。

(2) GRU とは

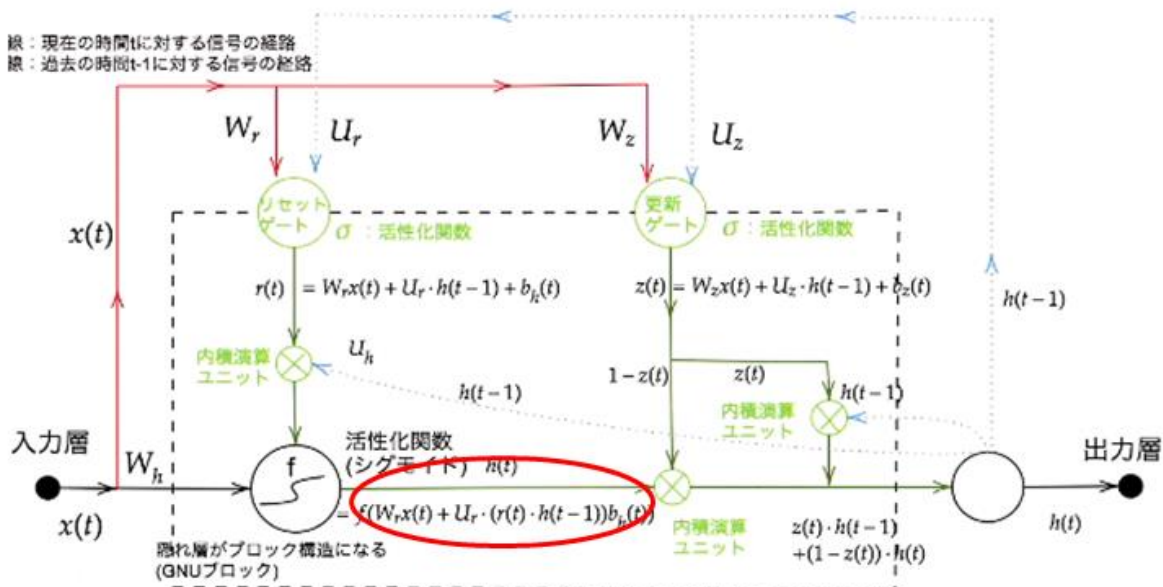
従来の LSTM では、パラメータが多数存在していたため、計算負荷が大きかった。

GRU では、そのパラメータを大幅に削減し、精度は同等又はそれ以上が望めるようになった構造をしている。

(3) メリット

計算負荷が低い

3. 1 GRU の全体像



全体的にシンプル。C E C がなくなった

リセットゲート・更新ゲート追加

$h(t-1)$: 前回の記憶

更新ゲート : 前回の記憶と今回の記憶をどれくらいの割合で混ぜ合わせるか

(1) 問題

LTM と CEC が抱える課題について、それぞれ簡潔に述べよ

パラメータが多いので計算時間がかかる。

(1) L S T M の課題

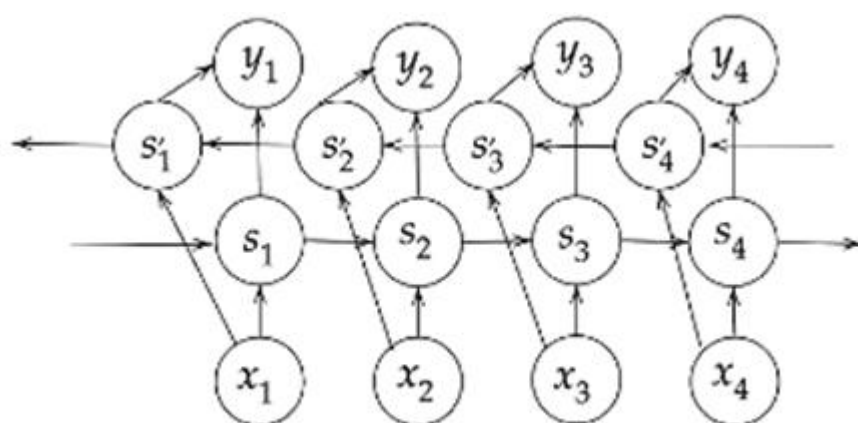
入力ゲート、出力ゲート、忘却ゲート、C E C はパラメータが多い

(2) C E C の課題

勾配が 1 なので学習機能が無い

4. section4_双方向 RNN

過去の情報だけでなく、未来の情報を加味することで、精度を向上させるためのモデル



未来から過去

過去から未来

文章：過去の情報と未来の情報

実用例

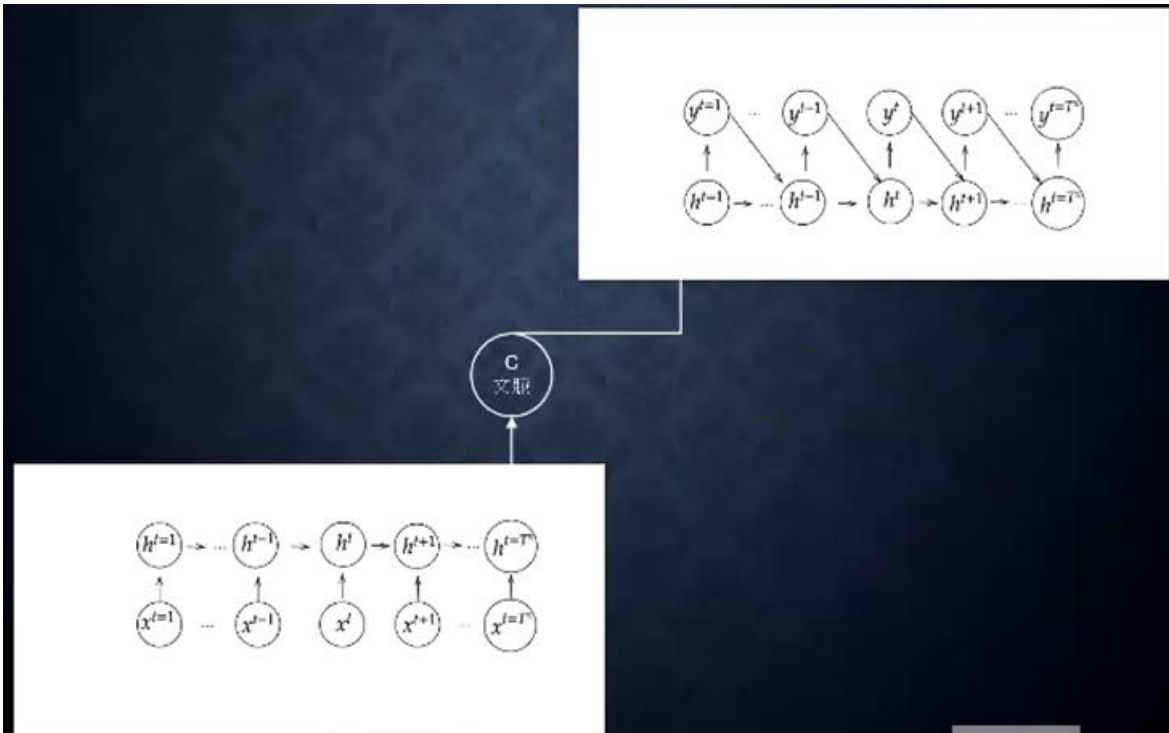
文章の推敲や、機械翻訳等

5. section5_Seq2Seq

Encoder-Decoder モデルの一種を指す。

具体的な用途は、機械対話ら、機械翻訳などに使用される。

5. 1 全体図



2つのニューラルネットワーク

1つ目のニューラルネットワークには、単語が入っていく。

蓄積されたデータ：中間層には入力されたものの記憶が溜まる。

意味がベクトル表現で保持される。

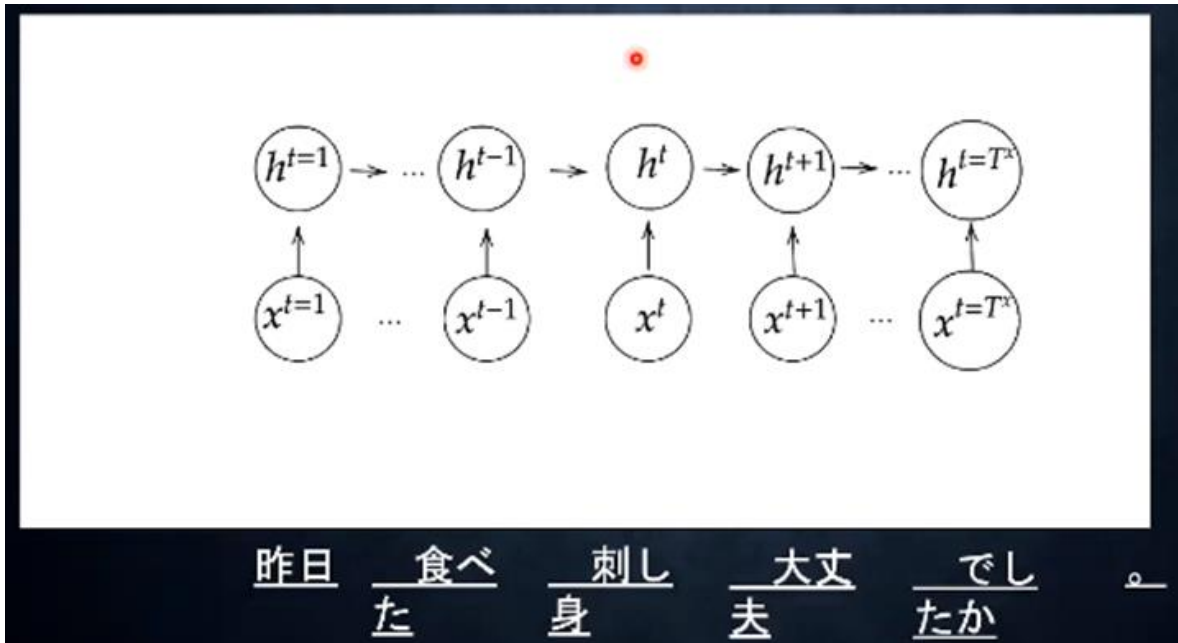
これをもう1つのニューラルネットワークに渡す。

翻訳の形式

- (1) エンコーダ
- (2) デコーダ

5. 2 EncoderRNN

ユーザーがインプットしたテキストを、単語等のトークンに区切って渡す構造。



文の意味を集約する。

全ての単語が集まって1つの文章を形成する。

(1) Taking : 文章を単語等のトークン毎に分割し、トークン毎の ID に分割する。

(2) Embedding : ID から、そのトークンを表す分散表現ベクトルに変換する。

(3) EncoderRNN : ベクトルを順番に RNN に入力していく。

自然言語のベクトル化

単語が無限にある。どのように登録するか

単語に番号を付ける。

One Hot Vector

5. 2. 1 EncoderRNN 処理手順

(1) vec1 を RNN に入力し、hidden state を出力。

この hidden state と次の入力 vec2 をまた RNN に入力してきた hidden state を出力という流れを繰り返す。

(2) 最後の vec を入れた時の hidden state を final state として取っておく。この final state が thought vector と呼ばれ、入力した文の意味を表すベクトルとなる。

特徴量抽出

機械学習で同じ意味の単語をおなじような embedding 表現にする。

MLM - Masked Language Model

私は 昨日 ラーメン を 食べ ました
[...] [...] [...] [...] [...] [...] [...]

モデルを学習する時に 1 か所を見えなくする。

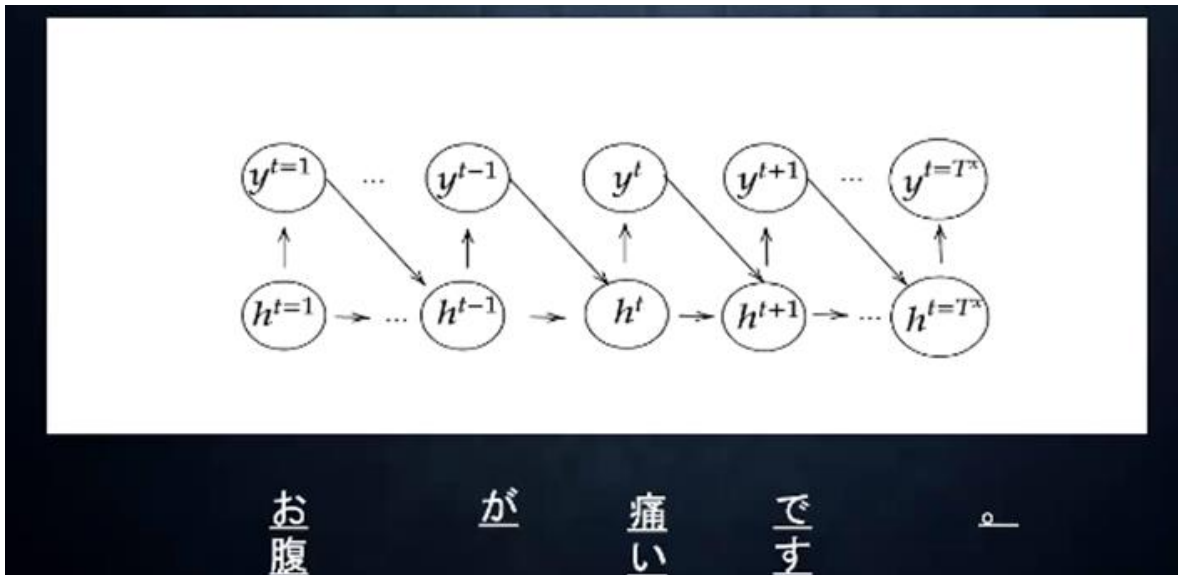
前後の文にや区から予想する。

MLM - Masked Language Model

私は 昨日 ラーメン を 食べ ました
[...] [...] [...] [...] [...] [...] [...]

5. 3 DecoderRNN

システムがアウトプットデータを単語等のトークンごとに生成する構造。



- (1) DecoderRNN : EncoderRNN の final state(thought vector)から、各 token の生成確率を出力していき
まず final state を DecoderRNN の initial state として設定し、Embedding を入力する。
- (2) Sampling : 生成確率に基づいて token をランダムに選ぶ。
- (3) Embedding : 2 で選ばれた token を Embedding して DecoderRNN への次の入力とする。
- (4) Detokenize : 1 - 3 を繰り返し、2 で得られた token を文字列に直す。

単語	ID	one-hot	embedding
私	1	$[1, 0, 0 \dots 0]$	$[0.2 \ 0.4 \ 0.6 \dots 0.1]$
は	2	$[0, 1, 0 \dots 0]$	$[\dots]$
刺身	3	$[0, 0, 1 \dots 0]$	$[\dots]$
昨日	4	$[\dots]$	$[\dots]$
⋮	⋮	⋮	⋮
xxx	10000	$[0, 0, 0 \dots 1]$	$[\dots]$

エンコーダーとは逆になる。

5. 4 HRED

(1) Seq2Seq の課題

一問一答しかできない——>問に対して文脈も何もなく、ただ応答が行われ続ける。

この問題から HRED が生まれる。

5. 4. 1 HRED とは

過去 $n-1$ 個の発話から次の発話を生成する。

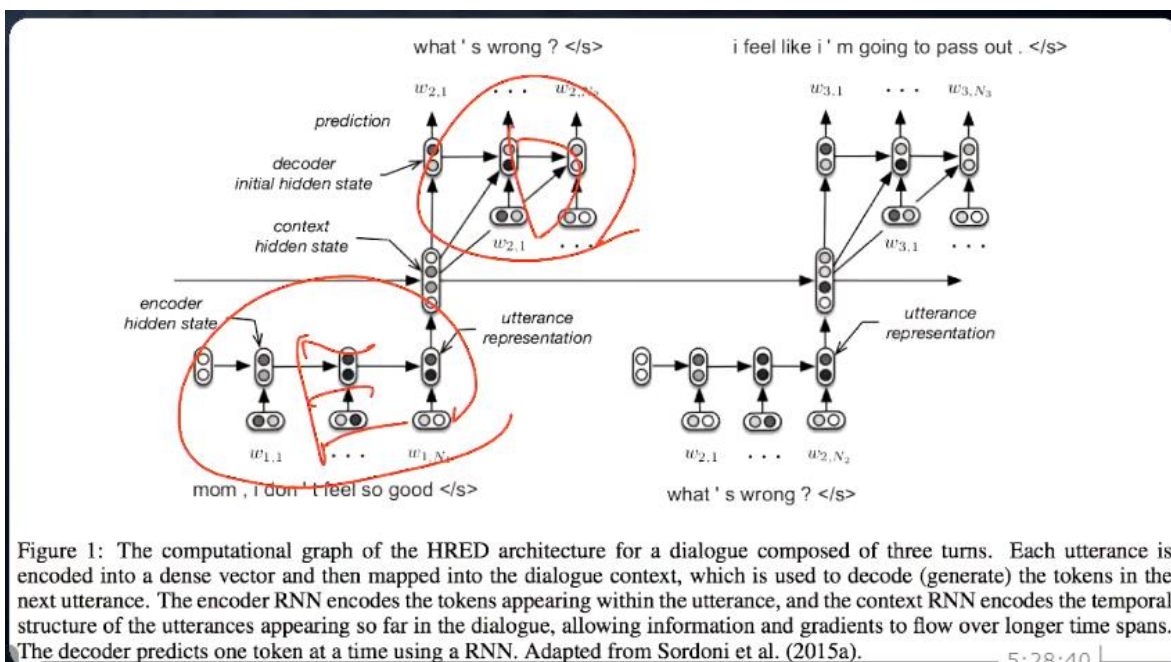
Seq2Seq では、会話の文脈無視で応答がなされたが、HRED では、前の単語の流れに即して応答されるため、より人間らしい文章が生成される。

(例)

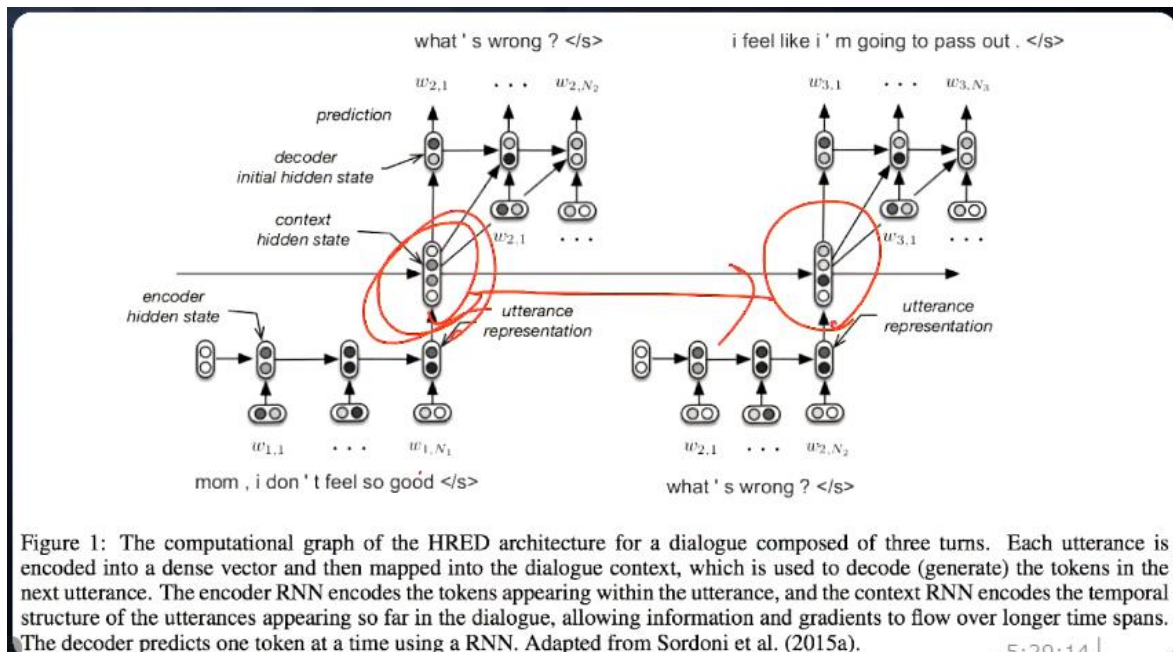
システム：インコかわいいよね。

ユーザー：うん

システム：インコかわいいのわかる。



5:28:40



5. 4. 2 HRED の構造

HRED とは Seq2Seq+Context RNN

Context RNN : Encoder のまとめた各文章の系列をまとめて、これまでの会話コンテキスト全体を表すベクトルに変換する構造。

→過去の発話の履歴を加味した返答をできる。

5. 4. 3 HRED の課題

HRED は確率的な多様性が字面にしかなく、会話の「流れ」のような多様性が無い。

→同じコンテキスト（発話リスト）を与えられても、答の内容が毎回会話の流れとしては同じものしか出せない。

HRED は短く情報量に乏しい答えをしがちである。

→短いよくある答えを学ぶ傾向がある。

(例)「うん」「そうだね」「・・・」など。

5. 5 VRED

5. 5. 1 VRED とは

HRED に、VAE の潜在変数の概念を追加したもの

→HRED の課題を VAE の潜在変数の概念を追加することで解決した構造。

5. 6 VAE

5. 6. 1 オートエンコーダ

(1) オートエンコーダとは

教師なし学習の一つ。

そのため学習時の入力データは訓練データのみで、教師データは利用しない。

(2) オートエンコーダの具体例

MNIST の場合、 28×28 の数字の画像を入れて、同じ画像を出力するニューラルネットワークとなる。

(3) オートエンコーダの構造

入力データから潜在変数 Z に変換するニューラルネットワークを **Encoder**、逆に潜在変数 Z をインプットして元画像を復元するニューラルネットワークを **Decoder** と呼ぶ。

(4) メリット

次元削減が行える。

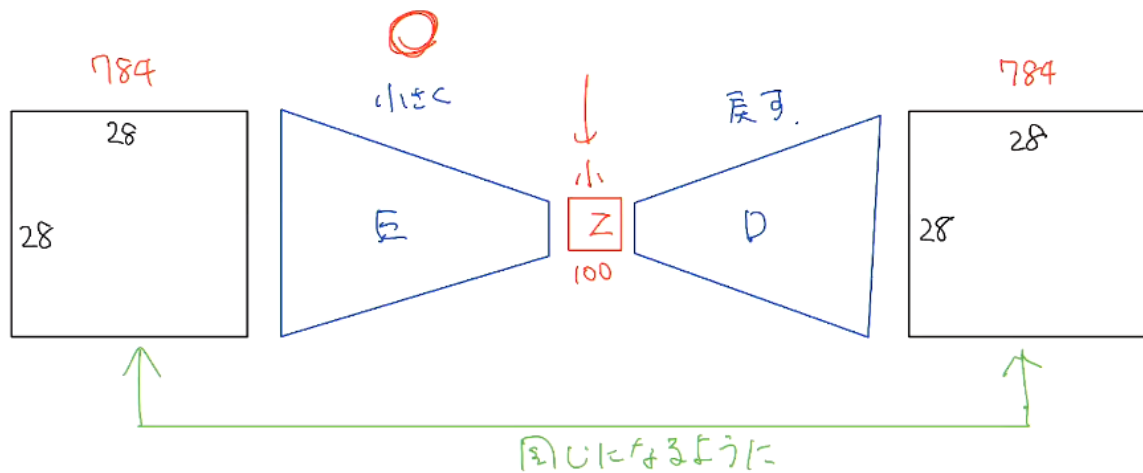
(注) Z の次元が入力データより小さい場合、次元削減とみなすことができる。

潜在変数 Z

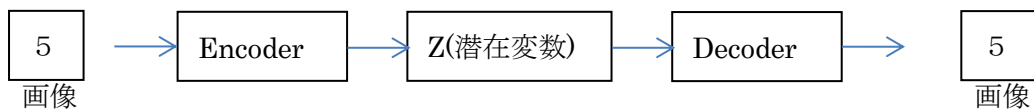
MNIST データ

入力と出力が同じになるように学習する

784 → 100 : 次元削減



(5) オートエンコーダの構造図

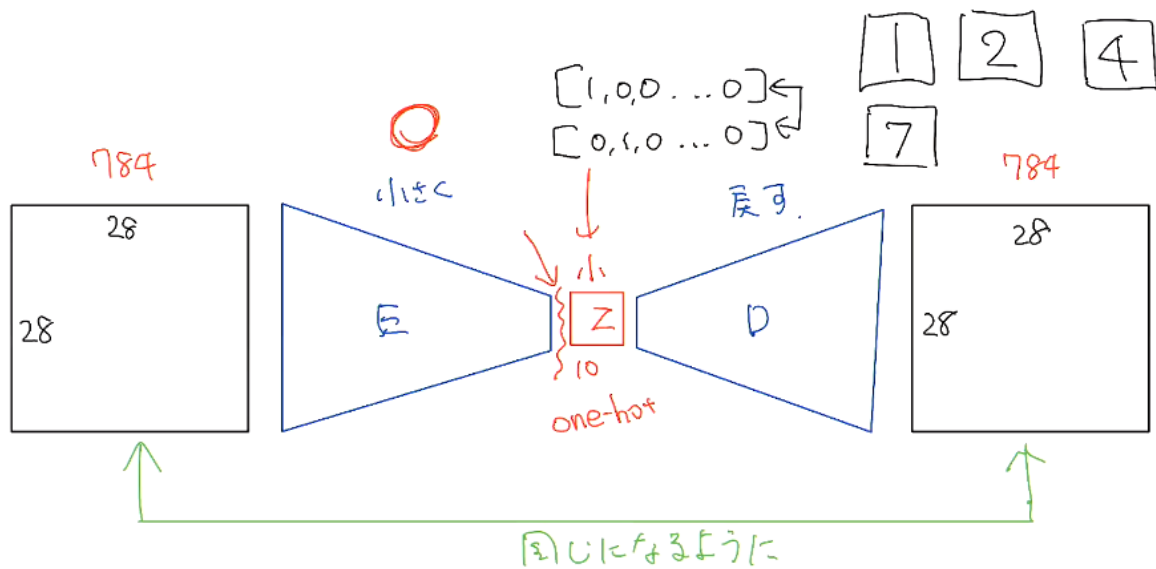


5. 6. 2 VAE

通常のオートエンコーダの場合、何かしら潜在変数 Z にデータを押し込めているものの、その構造がどのようなものかわからない。→VAE はこの潜在変数 Z に確率分布 $Z \sim N(0,1)$ を仮定したもの。

→VAE は、データを潜在変数 Z の確率分布という構造に押し込めることを可能にする。
Encoder の出力に対して、Decoder の入力にランダム性を持たせる。→ノイズを加える。

10 この数字を One Hot Vector を学習→元のデータの近さは無視される。
近いと同じようなベクトルになってほしい



ノイズを加える。

デコーダーは、より汎用的になる。

(1) 問題

VAE とは自己符号化器の潜在変数に確率分布を導入したもの。

6. section6_Word2vec

課題：RNN では、単語のような可変長の文字列を RNN に与えることはできない。

——>固定長形式で単語を表す必要がある。

学習データからボキャブラリを作成

注) わかりやすく 7 語のボキャブラリを作成したら、本来は辞書の単語数だけ出来上がる。

(例) I want to eat a pples.I like apples.

{aples,eat,I,like,to,want}

(1) メリット

大規模データの分散表現の学習が、現実的な計算速度とメモリ量で実現可能にした。

->X: ボキャブラリ×ボキャブラリだけの重み行列が誕生。

○: ボキャブラリ×任意の単語ベクトル次元で重み行列が誕生。

The diagram illustrates the process of generating a word vector from a vocabulary. On the left, a one-hot vector is shown as $[0, 1, 0, 0, 0, 0, 0, 0, 0, 0, \dots]$. This vector is multiplied by a weight matrix, which is represented as a box containing the following structure:

$$\begin{bmatrix} w_{11}^{(l)} & \dots & w_{1l}^{(l)} \\ \vdots & \ddots & \vdots \\ w_{j1}^{(l)} & \dots & w_{jl}^{(l)} \end{bmatrix}$$

A blue arrow points from the weight matrix to the result on the right, which is the word vector $= [1, 13, 15]$.

ボキャブラリ数×単語ベクトルの次元数の重み行列

7. section7_Attention Mechanism

課題：Seq2Seq の問題は、長い文章への対応が難しい。

Seq2Seq では、2 単語でも、1 0 0 単語でも固定次元ベクトルの中に入力しなければならない。

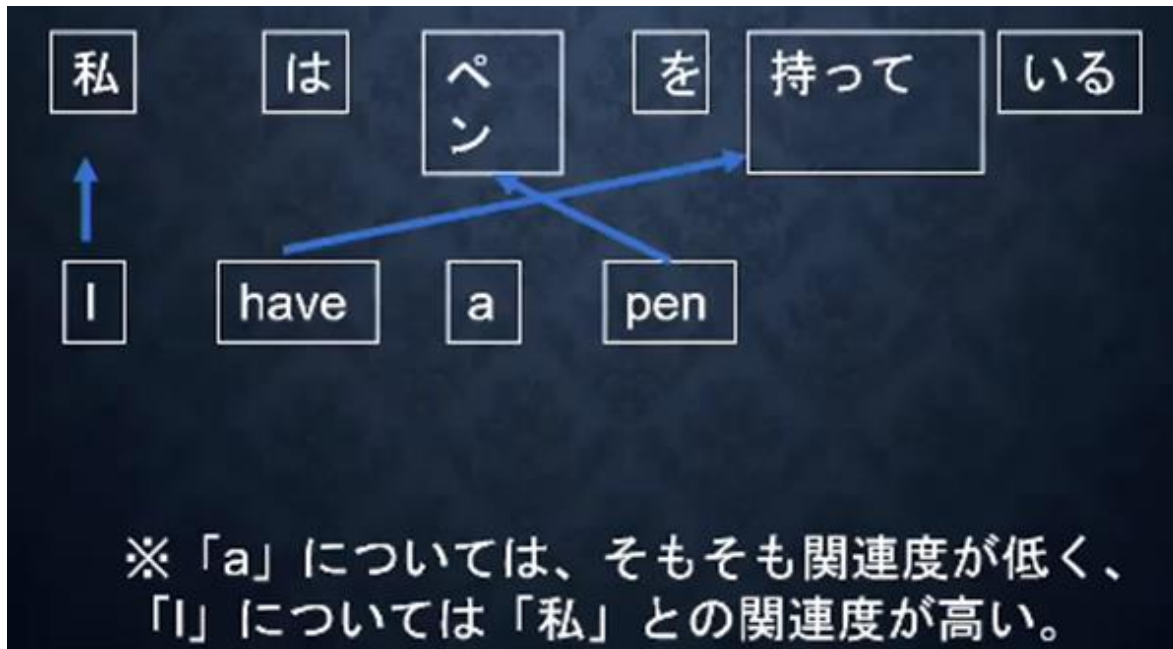
解決策：文章が長くなるほどそのシーケンスの内部表現の次元も大きくなっていく仕組みが必要になる。

→Attention Mechanism



「入力と出力のどの単語が関連しているのか」の関連度を学習する仕組み。

(1) 具体例



(2) 問題 1

RNN と word2vec、seq2seq と Attention の違い

RNN は時系列データを処理するのに適したもの

word2vec :

Attention

それぞれの関連系に重みを付ける手法

(3) 問題 2

Seq2seq と HRED、HRED と VHRED の違い

seq2seq : 1 文の 1 問 1 答に対して処理ができる

ある時系列データからある時系列データを作り出すネットワークのこと。

HRED :

Seq2Seq に文脈の解釈を付け加えたもの

文脈の意味をくみ取った変換

VHRED :

HRED は当たり障りのない回答しかできなかったのを、当たり障りのない単語以上の出力をできるようにしたもの。