知能プログラミング演習 I

第5回: 再帰型ニューラルネットワーク

梅津 佑太

2号館 404A: umezu.yuta@nitech.ac.jp

課題のダウンロード

前回作ったディレクトリに移動して今日の課題のダウンロードと解凍

```
step1: cd ./DLL
```

step2: wget http://www-als.ics.nitech.ac.jp/~umezu/DLL19/Lec5.zip

step3: unzip Lec5.zip

✓ まだ DLL のフォルダを作ってない人は、step1 の前に mkdir -p DLL でフォルダを作成する

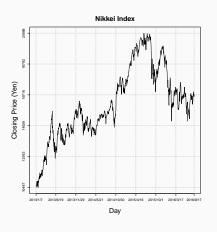
1

今日の講義内容

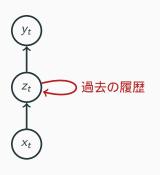
1. 再帰型ニューラルネットワークの学習

空間依存性

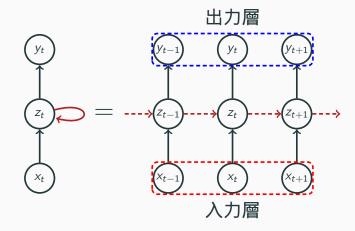
● 時系列, 動画, テキストなどは, "現在の値"が "過去の値"に依存



ネットワーク表現 (モデル)

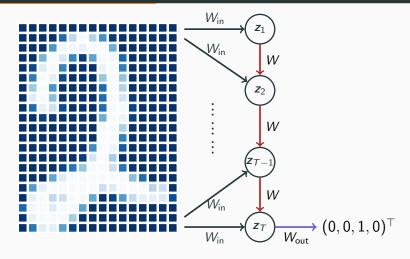


再帰型ニューラルネットワーク



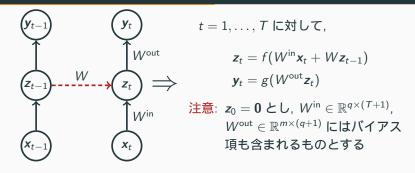
- 入力ごとに出力を計算する再帰型ニューラルネットワークを many-to-many ネットワークと呼ぶ
- 各時点における入出力関係は3層ニューラルネットワーク
- 中間層のネットワークにより、全体として "深い"モデルを構成

画像の分類: many-to-one ネットワーク



- 行ごとにデータをスキャンし、最終的な出力を用いて予測する
- ◆ 分類問題の場合, 誤差関数はクロスエントロピー, 最後の中間層から出力層への活性化関数はソフトマックス関数を用いる

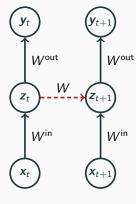
再帰型ニューラルネットワークの関数表現



順伝播は,線形結合 + 活性化関数の繰り返し.簡単のため,誤差関数は引数を省略して書くと $E=\sum_{t=1}^T E_t$.ただし,

(回帰)
$$E_t = \frac{1}{2} (y_t - g(W^{\text{out}} f(W^{\text{in}} \mathbf{x}_t + W \mathbf{z}_{t-1})))^2$$
(分類)
$$E_t = -\sum_{k=1}^m y_{tk} \log g(\mathbf{w}_k^{\text{out}} f(\mathbf{w}_k^{\text{in}} \mathbf{x}_t + \mathbf{w}_k^{\top} \mathbf{z}_{t-1}))$$

誤差逆伝播



 $oldsymbol{u}_t = W^{\mathsf{in}} oldsymbol{x}_t + W oldsymbol{z}_{t-1}, oldsymbol{v}_t = W^{\mathsf{out}} oldsymbol{z}_t$, および以下を定義する:

$$\delta_t = \frac{\partial E}{\partial u_t}, \quad \delta_t^{out} = \frac{\partial E}{\partial v_t} = g(v_t) - y_t$$

このとき,

$$\begin{split} \boldsymbol{W}^{\mathsf{in}} \qquad \boldsymbol{\delta}_t &= \frac{\partial E}{\partial \boldsymbol{u}_{t+1}} \frac{\partial \boldsymbol{u}_{t+1}}{\partial \boldsymbol{u}_t} + \frac{\partial E}{\partial \boldsymbol{v}_t} \frac{\partial \boldsymbol{v}_t}{\partial \boldsymbol{u}_t} \\ &= (\boldsymbol{W}^{\top} \boldsymbol{\delta}_{t+1} + \tilde{\boldsymbol{W}}^{\mathsf{out} \top} \boldsymbol{\delta}_t^{\mathsf{out}}) \odot \nabla f(\boldsymbol{u}_t) \end{split}$$

で誤差を逆伝播する.

ただし, \tilde{W}^{out} は W^{out} の 1 列目を取り除いた行列で, $\delta_{T+1}=\mathbf{0}$ とする. また, $u_t^{\text{in}}=W^{\text{in}}x_t$ としたとき,

$$\frac{\partial E}{\partial \textbf{\textit{u}}_{t}^{\text{in}}} = \frac{\partial E}{\partial \textbf{\textit{u}}_{t}} \frac{\partial \textbf{\textit{u}}_{t}}{\partial \textbf{\textit{u}}_{t}^{\text{in}}} = \frac{\partial E}{\partial \textbf{\textit{u}}_{t}}$$

なので、中間層から入力層への逆伝播にも $\partial E/\partial u_t$ を利用できる.

パラメータの更新ルール!

バイアス項を含む x_t (元データ $X \in \mathbb{R}^{T \times T}$ の行べクトルの 1 列目に 1 を加えた $T \times (T+1)$ 行列) に対して,

$$\frac{\partial E}{\partial W^{\text{in}}} = \sum_{t=1}^{T} \frac{\partial E}{\partial \boldsymbol{u}_{t}} \frac{\partial \boldsymbol{u}_{t}}{\partial W^{\text{in}}} = \sum_{t=1}^{T} \boldsymbol{\delta}_{t} \boldsymbol{x}_{t}^{\top}$$

となる. $D = (\delta_1, \ldots, \delta_T)$ とすれば,

$$\frac{\partial E}{\partial W^{\rm in}} = DX$$

となる. 同様に, $ilde{Z}=(extbf{z}_0, extbf{z}_1,\dots, extbf{z}_{T-1})^ op (extbf{z}_0= extbf{0})$ とすれば,

$$\frac{\partial E}{\partial W} = \sum_{t=1}^{T} \frac{\partial E}{\partial u_t} \frac{\partial u_t}{\partial W} = \sum_{t=1}^{T} \delta_t \mathbf{z}_{t-1}^{\top} = D\tilde{Z}$$

となる.

パラメータの更新ルール

さらに, $Z=(z_1,\ldots,z_T)$ とすれば,

$$\frac{\partial E}{\partial W^{\text{out}}} = \sum_{t=1}^{T} \frac{\partial E}{\partial \mathbf{v}_{t}} \frac{\partial \mathbf{v}_{t}}{\partial W^{\text{out}}} = \sum_{t=1}^{T} (g(\mathbf{v}_{t}) - \mathbf{y}_{t}) \mathbf{z}_{t}^{\top} = (g(V) - Y) Z$$

が得られるので、これらを用いてパラメータを更新する. ただし、

$$g(V) - Y = (g(\mathbf{v}_1) - \mathbf{y}_1, \dots, g(\mathbf{v}_T) - \mathbf{y}_T)$$

• many-to-one ネットワークの場合には, 上記の $\frac{\partial E}{\partial W^{\mathrm{out}}}$ の代わりに,

$$\frac{\partial E}{\partial W^{\text{out}}} = (g(\mathbf{v}_T) - \mathbf{y}_T) \mathbf{z}_T^\top \in \mathbb{R}^{m \times T}$$

を用いる

疑似コード: many-to-many ネットワークの場合

入力: $\{(\mathbf{y}_t, \mathbf{x}_t)\}_{t=1}^T$ に関する n 組のデータ, 各層における活性化関数 (f,g)

- 1. パラメータの初期化: Wⁱⁿ, W^{out}, W
- 2. ランダムにデータを読み込む: $\{(y_t, x_t)\}$
- 3. 順伝播: t = 1, 2, ..., T に対して,

3.1
$$\mathbf{z}_t \leftarrow f(\mathbf{W}^{\mathsf{in}}\mathbf{x}_t + \mathbf{W}\mathbf{z}_{t-1})$$
 $(\mathbf{z}_0 = \mathbf{0})$

3.2
$$\mathbf{y}_t \leftarrow g(W^{\text{out}}\mathbf{z}_t)$$

4. 逆伝播: t = T, ..., 1 に対して,

4.1
$$\delta_t \leftarrow (W^{\top} \delta_{t+1} + \tilde{W}^{\text{out} \top} (g(\mathbf{v}_t) - \mathbf{y}_t)) \odot \nabla f(\mathbf{u}_t) \quad (\delta_{T+1} = \mathbf{0})$$

5. 重みの更新: Adam などを用いてパラメータを更新する. 通常の確率的勾配法なら下記の通り.

5.1
$$W^{\text{in}} \leftarrow W^{\text{in}} - \eta DX$$

5.2
$$W \leftarrow W - \eta D\tilde{Z}$$

5.3
$$W^{\text{out}} \leftarrow W^{\text{out}} - \eta(g(V) - Y)Z$$

6. 収束するまで 2-5 を繰り返す

結果の例Ⅰ

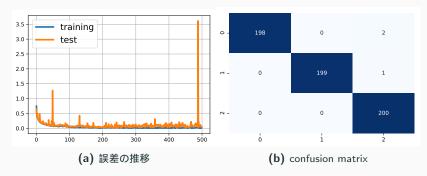
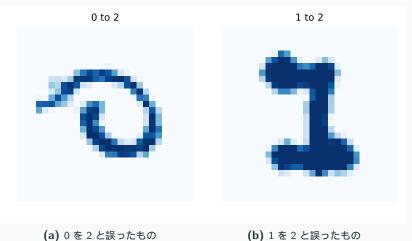


Figure 1: 500 エポックでの誤差関数の推移と confusion matrix. 活性化関数は シグモイド関数を用いた.

結果の例Ⅱ



(b) 1 を 2 と誤ったもの

Figure 2: 誤って分類したデータの図示.