

Denoising Autoencoders における確率的勾配降下法

Stochastic Gradient Descent for Denoising Autoencoders

Yusuke Sugomori

<http://yusugomori.com>

February 2013

概要

Denoising Autoencoder を実装する際に必要となる数式の導出過程をまとめた .

Denoising Autoencoder (DA) [Vincent 2008] では , 入力データを x , 入力データの一部を損傷させたデータを $\tilde{x} \sim q_{\mathcal{D}}(\tilde{x}|x)$ としたとき , \tilde{x} から x を復元するように訓練を行う . x が 2 値のベクトル , すなわち $x \in (0, 1)^d$ であるとき , DA における各写像は

$$\text{Encode: } \mathbf{y} = f_{\theta}(\tilde{\mathbf{x}}) = s(W\tilde{\mathbf{x}} + \mathbf{b}) \quad (1)$$

$$\text{Decode: } \mathbf{z} = g_{\theta'}(\mathbf{y}) = s(W'\mathbf{y} + \mathbf{b}') \quad (2)$$

で与えられる . ただし , $s(\cdot)$ はシグモイド関数である . また , $W' = W^T$ として考える . 式 (2) における z が , \tilde{x} から復元されたデータに相当する .

DA の実装には , 確率的勾配降下法 (stochastic gradient descent, SGD) が用いられる . 誤差関数は , 交差エントロピー誤差関数

$$L_H(x, z) = -x \log z - (1 - x) \log(1 - z) \quad (3)$$

で与えられるため , $\theta = (W, \mathbf{b}, \mathbf{b}')$ に対しての L_H の勾配を求めればよい . 以下 , 符号の煩わしさを防ぐため , $L := -L_H = x \log z + (1 - x) \log(1 - z)$

とし , L の勾配を導出する . 導出に際し , 以下の式で表される $\mathbf{h}_1 = \mathbf{h}_1(\tilde{\mathbf{x}})$, $\mathbf{h}_2 = \mathbf{h}_2(\mathbf{y})$ を定義する .

$$\mathbf{h}_1(\tilde{\mathbf{x}}) = W\tilde{\mathbf{x}} + \mathbf{b} \quad (4)$$

$$\mathbf{h}_2(\mathbf{y}) = W^T\mathbf{y} + \mathbf{b}' \quad (5)$$

$\mathbf{h}_1, \mathbf{h}_2$ を用いると , $\mathbf{y} = s(\mathbf{h}_1)$, $\mathbf{z} = s(\mathbf{h}_2)$ と表せる .

$W, \mathbf{b}, \mathbf{b}'$ それぞれによる微分は

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial \mathbf{h}_1} \frac{\partial \mathbf{h}_1}{\partial W} + \frac{\partial L}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial W} \quad (6)$$

$$\frac{\partial L}{\partial \mathbf{b}} = \frac{\partial L}{\partial \mathbf{h}_1} \frac{\partial \mathbf{h}_1}{\partial \mathbf{b}} = \frac{\partial L}{\partial \mathbf{h}_1} \quad (7)$$

$$\frac{\partial L}{\partial \mathbf{b}'} = \frac{\partial L}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{b}'} = \frac{\partial L}{\partial \mathbf{h}_2} \quad (8)$$

であり , シグモイド関数の微分 $s'(x) = s(x)(1 - s(x))$ を利用すると ,

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{h}_1} &= \frac{\partial L}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{h}_1} \\ &= \frac{\partial L}{\partial \mathbf{y}} * \mathbf{y} * (\mathbf{1} - \mathbf{y})\end{aligned}\quad (9)$$

$$\frac{\partial L}{\partial \mathbf{h}_2} = \mathbf{x} - \mathbf{z} \quad (10)$$

が得られるため，式 (9),(10) を式 (6),(7),(8) に代入すればよい．ただし，* はベクトルの要素積を表している．これは，一般に実装の際に用いられる記号に合わせて記述した．また，

$$\frac{\partial L}{\partial \mathbf{y}} = W(\mathbf{x} - \mathbf{z}) \quad (11)$$

となるため，最終的に，

$$\begin{aligned}\frac{\partial L}{\partial W} &= \left(W(\mathbf{x} - \mathbf{z}) * \mathbf{y} * (\mathbf{1} - \mathbf{y}) \right) \tilde{\mathbf{x}}^T \\ &\quad + \left((\mathbf{x} - \mathbf{z}) \mathbf{y}^T \right)^T\end{aligned}\quad (12)$$

$$\frac{\partial L}{\partial \mathbf{b}} = W(\mathbf{x} - \mathbf{z}) * \mathbf{y} * (\mathbf{1} - \mathbf{y}) \quad (13)$$

$$\frac{\partial L}{\partial \mathbf{b}'} = \mathbf{x} - \mathbf{z} \quad (14)$$

なる勾配が得られる．

N 個の入力データ群 $X = \{ \mathbf{x}_1, \dots, \mathbf{x}_N \}$ が与えられたとすると， $L = -L_H$ であるから，符号に注意すると，DA における確率的勾配降下法によるパラメータの更新は以下の式で与えられる．ただし， ϵ は学習率である．

$$W^{new} = W^{old} + \frac{\epsilon}{N} \sum_{n=1}^N \frac{\partial L}{\partial W} \quad (15)$$

$$\mathbf{b}^{new} = \mathbf{b}^{old} + \frac{\epsilon}{N} \sum_{n=1}^N \frac{\partial L}{\partial \mathbf{b}} \quad (16)$$

$$\mathbf{b}'^{new} = \mathbf{b}'^{old} + \frac{\epsilon}{N} \sum_{n=1}^N \frac{\partial L}{\partial \mathbf{b}'} \quad (17)$$

参考文献

- [Vincent 2008] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol: Extracting and Composing Robust Features with Denoising Autoencoders, *Proceedings of the 25th International Conference on Machine Learning*, (2008), pp.1096-1103