

1. 完成任务：

- 1) Back Propagation of Neural Network
- 2) Back Propagation of Binary Neural Network

2. Future Work:

1) 代码部分：

- 国庆期间，看完 util.py 中各个函数的具体写法和作用
- 国庆期间，学习 Pytorch 0.4(自己之前用的一直是 0.3)

2) 理论部分：

- 国庆期间，看 W1 和 W2 额外的四篇论文

3. 未解决部分：

- 1) XNOR-NET 论文当中，关于二值化网络的训练部分引用的是 NIPS 2016-Binaryconnect-training-deep-neural-networks-with-binary-weights-during-propagations 这篇论文，其中所给公式为：

$$\frac{\partial C}{\partial W_i} = \frac{\partial C}{\tilde{W}_i} \left(\frac{1}{n} + \frac{\partial \text{sign}}{\partial W_i} \times \alpha \right)$$

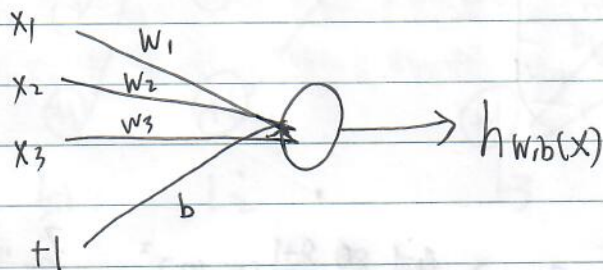
但在 github@jiecaoyu 的代码当中，给出了不一样的解释，核心区别在于：

$$\frac{\partial C}{\partial W_i} = \sum_{j=1}^n \left(\frac{\partial C}{\partial \tilde{W}_j} \cdot \frac{\partial \tilde{W}_j}{\partial W_i} \right)$$

未能理解哪一个究竟是对的，已提至 issue

Back Propagation of neural network

A fixed training set: $\{(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})\}$, 总共 m 个 sample



$$h_{w,b}(x) = f(W^T x) = \sum_{i=1}^3 W_i x_i + b$$

$$\begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

Sigmoid: $f(z) = \frac{1}{1 + \exp(-z)}$

tanh: $f(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

n_L : number of layers

$$f'(z) = f(z)(1 - f(z))$$

$$a_2^{(2)} = f(W_{21}^{(1)} x_1 + W_{22}^{(1)} x_2 + W_{23}^{(1)} x_3 + b_2^{(1)})$$

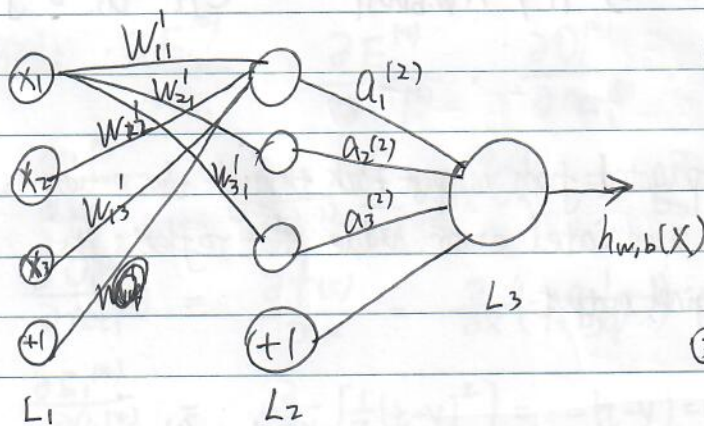
$$a_3^{(2)} = f(W_{31}^{(1)} x_1 + W_{32}^{(1)} x_2 + W_{33}^{(1)} x_3 + b_3^{(1)})$$

$$h_{w,b}(x) = a_1^{(3)} = f(W_{11}^{(2)} a_1^{(2)} + W_{12}^{(2)} a_2^{(2)} + W_{13}^{(2)} a_3^{(2)} + b_1^{(2)})$$

$$\textcircled{1} (W, b) = (W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)})$$

$\textcircled{2} W_{ij}^{(l)}$: weight between unit j in layer l and unit i in layer $l+1$

$\textcircled{3} b_i^{(l)}$: bias associated with unit i in layer $l+1$.



$\textcircled{4} S_L$: the number of nodes in layer L (not counting bias unit)

$\textcircled{5} a_i^{(l)}$: activation of unit i in layer L .

$$a_1^{(2)} = f(W_{11}^{(1)} x_1 + W_{12}^{(1)} x_2 + W_{13}^{(1)} x_3 + b_1^{(1)})$$

$\textcircled{6} z_i^{(l)}$: total weighted sum of inputs to unit i in layer L .

$$W^{(1)} = \begin{bmatrix} W_{11} & W_{12} & W_{13} \\ W_{21} & W_{22} & W_{23} \\ W_{31} & W_{32} & W_{33} \end{bmatrix}$$

Compact notation: **organizing in matrix**

$$z^{(2)} = W^{(1)}x + b^{(1)} \quad (a^{(1)} = x)$$

$$a^{(2)} = f(z^{(2)})$$

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$$

$$h_{w,b}(x) = a^{(3)} = f(z^{(3)})$$

forward-propagation

$$S-E \quad J(w,b) = \left[\frac{1}{m} \sum_{i=1}^m J(w,b; x^{(i)}, y^{(i)}) \right] + \frac{\lambda}{2} \sum_{l=1}^{n-1} \sum_{j=1}^{s_l} \sum_{i=1}^{s_{l+1}} (w_{ij}^{(l)})^2$$

$$\left\{ \begin{aligned} w_{ij}^{(l)} &= w_{ij}^{(l)} - \alpha \frac{\partial}{\partial w_{ij}^{(l)}} J(w,b) \\ b_i^{(l)} &= b_i^{(l)} - \alpha \frac{\partial}{\partial b_i^{(l)}} J(w,b) \end{aligned} \right.$$

$$\begin{matrix} w_{11} \\ w_{21} \end{matrix}$$

如何计算以上2个偏导数? BP 算法:

o $\delta_i^{(l)}$: error term for each node i in layer l

o Set $\delta_i^{(n)} = \frac{\partial}{\partial z_i^{(n)}} \frac{1}{2} \|y - h_{w,b}(x)\|^2 = -(y_i - a_i^{(n)}) \cdot f'(z_i^{(n)})$
output layer unit i

The error values are propagated from output back through the network, until each neuron has an associated error value that reflects its contribution to the original output.

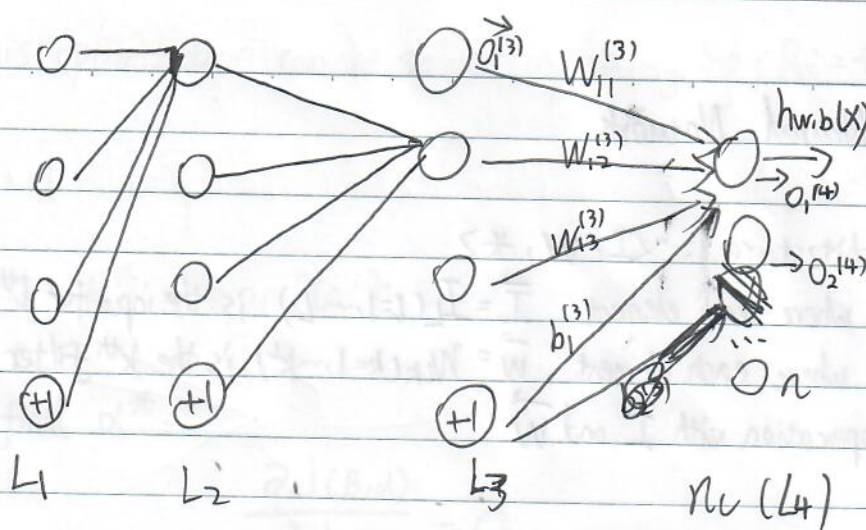
o For $l = n-1, n-2, \dots, 2$

For each node i in layer l , set $\delta_i^{(l)} = \left(\sum_{j=1}^{s_{l+1}} w_{ji}^{(l)} \delta_j^{(l+1)} \right) f'(z_i^{(l)})$

o Desired Partial Derivatives

$$\frac{\partial}{\partial w_{ij}^{(l)}} J(w,b; x, y) = a_j^{(l)} \delta_i^{(l+1)}$$

$$\frac{\partial}{\partial b_i^{(l)}} J(w,b; x, y) = \delta_i^{(l+1)}$$



假设 fully connected,
论文中未考虑 bias

Step 1: Forward Propagation, 产生了一个输出值 $hwb(x)$ (predicted)

Step 2: Backward Propagation:

Suppose M-S-E: $E = \frac{1}{2} (hwb(x) - y)^2$

$$\begin{aligned} \begin{cases} i\text{th neuron in } L_3 \\ i\text{th neuron in } L_4 \end{cases} & \frac{\partial E}{\partial w_{ij}^{(3)}} = \frac{\partial E}{\partial a_j^{(4)}} \cdot \frac{\partial a_j^{(4)}}{\partial z_j^{(3)}} \cdot \frac{\partial z_j^{(3)}}{\partial w_{ij}^{(3)}} \\ \text{(chain Rule)} & = \frac{\partial E}{\partial o_j^{(4)}} \cdot \frac{\partial o_j^{(4)}}{\partial net_j^{(3)}} \cdot \frac{\partial net_j^{(3)}}{\partial w_{ij}^{(3)}} \end{aligned} \quad \left| \begin{aligned} a^{(4)} &= f(z^{(3)}) \\ f(x) &= \frac{1}{1+\exp^{-x}} \\ z^{(3)} &= w_{11}^{(3)} \cdot x_1^{(3)} + w_{12}^{(3)} \cdot x_2^{(3)} + w_{13}^{(3)} \cdot x_3^{(3)} \end{aligned} \right.$$

$$(1) \frac{\partial net_j^{(3)}}{\partial w_{ij}^{(3)}} = \frac{\partial}{\partial w_{ij}^{(3)}} \left(\sum_{k=1}^n w_{kj}^{(3)} \cdot x_k^{(3)} \right) \xrightarrow{\text{only one sum}} \frac{\partial}{\partial w_{ij}^{(3)}} w_{ij}^{(3)} \cdot x_i^{(3)}$$

具体举例: $\frac{\partial E}{\partial w_{11}^{(3)}} = \frac{\partial E}{\partial o_1^{(4)}} \cdot \frac{\partial o_1^{(4)}}{\partial net_1^{(3)}} \cdot \frac{\partial net_1^{(3)}}{\partial w_{11}^{(3)}} \quad (\text{输出层 } L \rightarrow L-1 \text{ 层})$

$$(1): \frac{\partial net_1^{(3)}}{\partial w_{11}^{(3)}} = \frac{\partial}{\partial w_{11}^{(3)}} \sum_{k=1}^n w_{1k}^{(3)} o_k^{(3)} = \frac{\partial}{\partial w_{11}^{(3)}} (w_{11}^{(3)} o_1^{(3)}) = o_1^{(3)}$$

$$(2) \frac{\partial o_1^{(4)}}{\partial net_1^{(3)}} = \frac{\partial f(x)}{\partial x} = \frac{\partial}{\partial x} \left(\frac{1}{1+\exp^{-x}} \right) = o_1^{(4)} (1 - o_1^{(4)}) \quad f(x) \rightarrow \text{sigmoid}$$

$$(3) \frac{\partial E}{\partial o_1^{(4)}} = \frac{\partial}{\partial y} \left[\frac{1}{2} (t - y)^2 \right] = -(t - y) = y - t = o_1^{(4)} - t_1^{(4)}$$

So $\frac{\partial E}{\partial w_{11}^{(3)}} = [o_1^{(4)} - t_1^{(4)}] \cdot o_1^{(4)} \cdot (1 - o_1^{(4)}) \cdot o_1^{(3)} = \delta_1^{(3)} \cdot o_1^{(3)}$ error term

如果是在内层: $\delta_i^{(l)} = \left(\sum_{j=1}^{L+1} w_{ij}^{(l)} \delta_j^{(l+1)} \right) f'(z_i^{(l)})$

Paper:

① Binary Convolutional Neural Network

• an L-layer CNN Architecture: $\langle I, W, * \rangle$

I : set of tensors, where each element $\vec{I} = I_L (L=1, \dots, L)$ is the input for L^{th} layer

W : set of tensors, where each element $\vec{W} = W_{lk} (k=1, \dots, k')$ is the k^{th} filter in L

$*$: A convolutional operation with \vec{I} and \vec{W}

BWN: elements of W are binary tensors

XNOR-Network: elements of ~~XNOR-Network~~ both I and W are binary tensors.

3) Binary Filter $B \in \{+1, -1\}^{c \times w \times h}$, scaling factor α such that $W \approx \alpha B$

$$I * W \approx (I \oplus B) \alpha \quad \oplus \text{ indicates a conv without 乘法}$$

BWN: $\langle I, B, A, \oplus \rangle$ $\begin{cases} A: \text{a set of positive scalars} \\ B: \text{a set of binary tensors} \\ (W_{lk} = A_{lk} B_{lk}) \end{cases}$

• Estimating binary weights:

optimization $J(B, \alpha) = \|W - \alpha B\|^2$
 $= \alpha^2 B^T B - 2\alpha W^T B + W^T W$

$$\alpha^*, B^* = \underset{\alpha, B}{\operatorname{argmin}} J(B, \alpha) \rightarrow \text{函数 } J() \text{ 取得最小值的 } \alpha, B \text{ 的集合}$$

$B^T B$ and $W^T W$ is constant.

$$\text{So } J(B, \alpha) = \alpha^2 \eta - 2\alpha W^T B + C$$

optimization becomes:

$$\Rightarrow B^* = \underset{B}{\operatorname{argmax}} \{W^T B\} \text{ s.t. } B \in \{+1, -1\}^n$$

This optimization can be solved by assigning $B_i = +1$ if $W_i \geq 0$
 $B_i = -1$ if $W_i < 0$

$$\Rightarrow B^* = \text{sign}(W)$$

To find α^* :

$$\frac{\partial J(B, \alpha)}{\partial \alpha} = 0$$

$$2\alpha\eta - 2W^T B = 0$$

$$\alpha\eta = W^T B$$

$$\alpha = \frac{W^T B}{\eta}$$

$$\text{So } \alpha^* = \frac{W^T B^*}{\eta} = \frac{W^T \text{sign}(W)}{\eta} = \frac{\sum |W_i|}{\eta}$$

Training:

weight update methods: SGD or Adam