

Feed Forward Neural Network

Masahiro Ogawa

January 20, 2016

1 Notation

- \mathbf{x} : vector
- \mathbf{A} : matrix
- ${}^t\mathbf{A}$: transpose of \mathbf{A}
- $f.(\mathbf{A}) = \begin{pmatrix} f(a_{1,1}) & f(a_{1,2}) & \cdots & f(a_{1,n}) \\ f(a_{2,1}) & f(a_{2,2}) & \cdots & f(a_{2,n}) \\ \vdots & \vdots & \ddots & \vdots \\ f(a_{m,1}) & f(a_{m,2}) & \cdots & f(a_{m,n}) \end{pmatrix}$
- $(\mathbf{A} * \mathbf{B})_{i,j} = a_{i,j}b_{i,j}$
- $\mathbf{1}_{m \times n} = \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 1 \end{pmatrix} : m \times n \text{ dim}$
- $\tilde{\mathbf{y}} = \begin{pmatrix} 1 \\ \mathbf{y} \end{pmatrix} : \text{homogeneous vector}$
- n : data number
- N : total number of data
- l : layer number
- L : output layer number. thus, the total number of layer is $L + 1$
- k : learning iteration number
- K : final number of the learning iteration
- J : cost
- $f(x)$: activation function

- $\mathbf{X}_0 = (\mathbf{x}_{0,1}, \dots, \mathbf{x}_{0,N}) = \begin{pmatrix} x_{0,1,1} & \cdots & x_{0,1,N} \\ \vdots & \ddots & \vdots \\ x_{0,d_0,1} & \cdots & x_{0,d_0,N} \end{pmatrix} : d_0 \times N \text{ dim} : \text{input vectors}$
- $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_N) = \begin{pmatrix} b_{1,1} & \cdots & b_{1,N} \\ \vdots & \ddots & \vdots \\ b_{d_L,1} & \cdots & b_{d_L,N} \end{pmatrix} : d_L \times N \text{ dim} : \text{instruction signals}$
- $\tilde{\mathbf{W}}_l(k) = (\tilde{\mathbf{w}}_{l,1}(k), \dots, \tilde{\mathbf{w}}_{l,d_{l+1}}(k)) = \begin{pmatrix} \tilde{w}_{l,0,1}(k) & \cdots & \tilde{w}_{l,0,d_{l+1}}(k) \\ \vdots & \ddots & \vdots \\ \tilde{w}_{l,d_l,1}(k) & \cdots & \tilde{w}_{l,d_l,d_{l+1}}(k) \end{pmatrix} :$
weights : $(d_l + 1) \times d_{l+1} \text{ dim}$

2 Input

parameters

- L : output layer number.
- $\{d_l\}_{l=0}^L$: neuron numbers in each layer l
- $\{\tilde{\mathbf{W}}_l(0)\}_{l=0}^L \mid \tilde{\mathbf{W}}_l(0) = (\tilde{\mathbf{w}}_{l,1}(0), \dots, \tilde{\mathbf{w}}_{l,d_{l+1}}(0)) : (d_l + 1) \times d_{l+1} \text{ dim} :$
initial weights
- ρ : learning rate parameter
- T_e : error threshold
- T_K : threshold of the maximum number of iteration

data

- $\mathbf{X}_0 = (\mathbf{x}_{0,1}, \dots, \mathbf{x}_{0,N}) : d_0 \times N \text{ dim} : \text{input vectors}$
- $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_N) : d_L \times N \text{ dim} : \text{instruction signals}$

3 Output

- $\{\tilde{\mathbf{W}}_l(K)\}_{l=0}^L \mid \tilde{\mathbf{W}}_l(K) = (\tilde{\mathbf{w}}_{l,1}(K), \dots, \tilde{\mathbf{w}}_{l,d_{l+1}}(K)) : (d_l + 1) \times d_{l+1} \text{ dim} :$
final (after K th iteration) learned weights
- $\mathbf{Y}_L = (\mathbf{y}_{L,1}, \dots, \mathbf{y}_{L,N}) : d_L \times N \text{ dim} : \text{output vectors}$
- $J = \sum_{n=1}^N J_n : \text{cost, final error}$

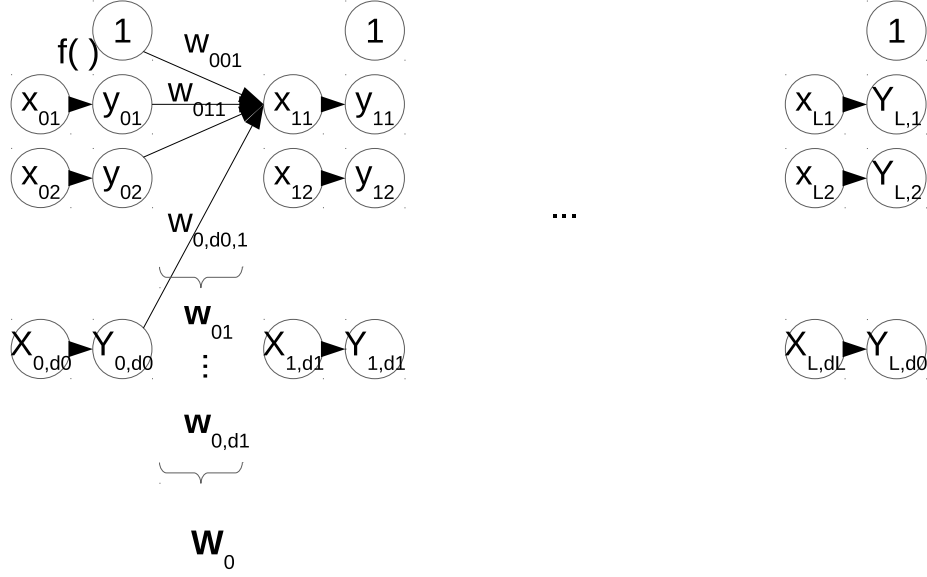


Figure 1: Neural network architecture

4 Algorithm(Iterative version, Stochastic Gradient descent)

4.1 Iteration

Iterate feed forward and back propagation algorithm for each input datum until the error become smaller than the threshold or the iteration number exceeds the maximum iteration number. That is,

$$J < T_e, \text{ or } k < T_K. \quad (1)$$

The iteration means update of the weights as,

$$\tilde{w}_{lij}(k+1) = \tilde{w}_{lij}(k) - \rho \frac{\partial J_n}{\partial \tilde{w}_{lij}(k)} \quad (2)$$

$$\Longleftrightarrow \tilde{\mathbf{W}}_l(k+1) = \tilde{\mathbf{W}}_l(k) - \rho \frac{\partial J_n}{\partial \tilde{\mathbf{W}}_l(k)} \quad (3)$$

From this equation, once we could compute $\frac{\partial J_n}{\partial \tilde{\mathbf{W}}_l(k)}$, we can update the weights. Therefore we concentrate on computing $\frac{\partial J_n}{\partial \tilde{\mathbf{W}}_l(k)}$ below.

4.2 Feed Forward

The below are true for all data and all iteration. So we omit the data number n and iteration number k as e.g. $\mathbf{x}_{l,n} = \mathbf{x}_l$, $\tilde{\mathbf{W}}_l(k) = \tilde{\mathbf{W}}_l$.

$$\mathbf{x}_l = {}^t\tilde{\mathbf{W}}_{l-1}\tilde{\mathbf{y}}_{l-1} \quad (4)$$

$$\mathbf{y}_l = f.(\mathbf{x}_l) \quad (5)$$

$$(6)$$

4.2.1 activation function for hidden layers

$$f(x) = \begin{cases} (1 + e^{-x})^{-1} & : \text{logistic function} \\ \tanh(x) & \end{cases} \quad (7)$$

4.2.2 activation function for an output layer

To fit output value as desired, we have to choose appropriate activation function for an output layer.

$$f(x) = \begin{cases} x & : \text{identity function for regression} \\ (1 + e^{-x})^{-1} & : \text{logistic function for } [0,1] \text{ output binary classification} \\ \frac{\exp(x)}{\sum_j \exp(x_{Lj})} & : \text{soft max function for multiclass classification} \\ \tanh(x) & : \text{for } [-1,1] \text{ output} \end{cases} \quad (8)$$

Note that we divide binary and multiclass case, because if we use soft max function for binary class classification, $f(x)$ always become 1.

4.3 Back Propagation

In this section, we compute $\frac{\partial J_n}{\partial \tilde{\mathbf{W}}_{l-1}(k)}$ instead of $\frac{\partial J_n}{\partial \tilde{\mathbf{W}}_l(k)}$ for just index convenience. And we omit n as $J_n = J$ in this section.

$$\frac{\partial J}{\partial \tilde{w}_{l-1,ij}} = \frac{\partial J}{\partial x_{lj}} \frac{\partial x_{lj}}{\partial \tilde{w}_{l-1,ij}} \quad (9)$$

$$= \epsilon_{lj} \frac{\partial (\sum_m \tilde{w}_{l-1,mj} \tilde{y}_{l-1,m})}{\partial \tilde{w}_{l-1,ij}} \quad | \quad \epsilon_{lj} := \frac{\partial J}{\partial x_{lj}} \quad (10)$$

$$= \epsilon_{lj} \tilde{y}_{l-1,i} \quad (11)$$

$$\iff \frac{\partial J}{\partial \tilde{\mathbf{W}}_{l-1}} = \tilde{\mathbf{y}}_{l-1}^t \boldsymbol{\epsilon}_l \quad (12)$$

$$\Rightarrow \tilde{\mathbf{W}}_{l-1}(k+1) = \tilde{\mathbf{W}}_{l-1}(k) - \rho \tilde{\mathbf{y}}_{l-1}^t \boldsymbol{\epsilon}_l \quad (13)$$

$$\epsilon_{lj} = \frac{\partial J}{\partial x_{lj}} \quad (14)$$

$$= \frac{\partial J}{\partial y_{lj}} \frac{\partial y_{lj}}{\partial x_{lj}} \quad (15)$$

$$= \frac{\partial J}{\partial y_{lj}} \frac{\partial f(x_{lj})}{\partial x_{lj}} \quad (16)$$

$$(17)$$

4.3.1 back propagation for hidden layers

$$\epsilon_{lj} = \frac{\partial J}{\partial y_{lj}} \frac{\partial f(x_{lj})}{\partial x_{lj}} \quad (18)$$

$$= \sum_h \frac{\partial J}{\partial x_{l+1,h}} \frac{\partial x_{l+1,h}}{\partial y_{lj}} f'(x_{lj}) \quad (19)$$

$$= \sum_h \epsilon_{l+1,h} \frac{(\sum_m w_{lmh} y_{lm} + w_{l0h})}{\partial y_{lj}} f'(x_{lj}) \quad (20)$$

$$= \sum_h \epsilon_{l+1,h} w_{ljh} f'(x_{lj}) \quad (21)$$

$$\Rightarrow \boldsymbol{\epsilon}_l = (\mathbf{W}_l \boldsymbol{\epsilon}_{l+1}) \cdot * f'(\mathbf{x}_l) \quad (22)$$

$$= \begin{cases} (\mathbf{W}_l \boldsymbol{\epsilon}_{l+1}) \cdot * (\mathbf{y}_l - \mathbf{y}_l \cdot * \mathbf{y}_l) & \text{if } f(x) = (1 + e^{-x})^{-1} \\ (\mathbf{W}_l \boldsymbol{\epsilon}_{l+1}) \cdot * (\mathbf{1}_{dl} - \mathbf{y}_l \cdot * \mathbf{y}_l) & \text{if } f(x) = \tanh(x) \end{cases} \quad (23)$$

4.3.2 back propagation for an output layer

To compute back propagation for an output layer, we need cost function form. Depend on that, we have to choose appropriate cost function for each problem as below.

$$J = \begin{cases} \frac{1}{2} \|\mathbf{y}_L - \mathbf{b}\|^2 & \text{regression} \\ -(b \ln y_L + (1-b) \ln(1-y_L)) & \text{binary classification} \\ -\sum \mathbf{b} \ln \mathbf{y}_L & \text{multiclass classification} \end{cases} \quad (24)$$

Using cost(24) and activation function(8), ϵ_L can be computed as below. Note that $f(x_h)$ is related to x_j , so we have to think about the term $\frac{\partial f(x_{Lh})}{\partial x_{Lj}}$.

$$\epsilon_{Lj} = \sum_h \frac{\partial J}{\partial y_{Lh}} \frac{\partial f(x_{Lh})}{\partial x_{Lj}} \quad (25)$$

$$= \begin{cases} \sum_h \left(\frac{\partial(1/2\|\mathbf{y}_L - \mathbf{b}\|^2)}{\partial y_{Lh}} \right) \frac{\partial x_{Lh}}{\partial x_{Lj}} & \text{regression} \\ -\left(\frac{b}{y_L} - \frac{1-b}{1-y_L} \right) \frac{\partial}{\partial x_L} (1 + e^{-x_L})^{-1} & \text{binary classification} \\ -\sum_h \frac{b_h}{y_{L,h}} \frac{\partial}{\partial x_{Lj}} \left(\frac{\exp(x_{Lh})}{\sum_i \exp(x_{Li})} \right) & \text{multiclass classification} \end{cases} \quad (26)$$

$$= \begin{cases} \sum_h (y_{Lh} - b_h) \delta_{h,j} \\ -\left(\frac{b}{y_L} - \frac{1-b}{1-y_L} \right) (y_L (1 - y_L)) \\ -\sum_h \frac{b_h}{y_{L,h}} (\delta_{h,j} y_{Lj} - y_{Lh} y_{Lj}) \end{cases} \quad (27)$$

$$= y_{Lj} - b_j \quad (28)$$

We use the softmax property $\sum_h b_h = 1$ in the last equation of multi-class classification case.

In both case the result becomes $y - b$, so there is no gradient diminishing problem. That's why we use these combination of the cost and final layers activation function.

4.3.3 final result

Finally, we get below results.

$$\epsilon_l = \begin{cases} \mathbf{y}_L - \mathbf{b} & \text{if } l = L \\ (\mathbf{W}_l \epsilon_{l+1}) \cdot * f'(\mathbf{x}_l) & \text{otherwise} \end{cases} \quad (29)$$

$$= \begin{cases} \mathbf{y}_L - \mathbf{b} \\ (\mathbf{W}_l \epsilon_{l+1}) \cdot * (\mathbf{1}_{dl} - \mathbf{y}_l \cdot * \mathbf{y}_l) & \text{if } f(x) = \tanh(x) \end{cases} \quad (30)$$

$$\tilde{\mathbf{W}}_{l-1}(k+1) = \tilde{\mathbf{W}}_{l-1}(k) - \rho \tilde{\mathbf{y}}_{l-1}^t \epsilon_l \quad (31)$$

5 Algorithm(Matrix version)

5.1 Iteration

Iterate feed forward algorithm for all input data and back propagation for all data until the error become smaller than the threshold or the iteration number exceeds the maximum iteration number. That is,

$$J < T_e, \text{ or } k < T_K. \quad (32)$$

5.2 Feed Forward

The below are true for all iteration. So we omit the iteration number as e.g. $\tilde{\mathbf{W}}_l(i) = \tilde{\mathbf{W}}_l$.

$$\mathbf{X}_l = {}^t\tilde{\mathbf{W}}_{l-1}\tilde{\mathbf{Y}}_{l-1} \quad (33)$$

$$\mathbf{Y}_l = f.(\mathbf{X}_l) \quad (34)$$

$f(x)$ is defined by (7) and (8).

5.3 Back Propagation

$$J = \begin{cases} \frac{1}{2}\text{Tr}({}^t(\mathbf{Y}_L - \mathbf{B})(\mathbf{Y}_L - \mathbf{B})) & \text{regression} \\ -\text{Tr}\left\{\begin{pmatrix} {}^t\mathbf{B} \\ \mathbf{1}_{1 \times N} - {}^t\mathbf{B} \end{pmatrix} \ln.(\mathbf{Y}_L \mathbf{1}_{N \times 1} - \mathbf{Y}_L)\right\} & \text{binary classification} \\ -\text{Tr}({}^t\mathbf{B} \ln.(\mathbf{Y}_L)) & \text{multiclass classification} \end{cases} \quad (35)$$

$$\mathbf{E}_l = (\epsilon_{l,1}, \dots, \epsilon_{l,N}) = \begin{cases} \mathbf{Y}_L - \mathbf{B} & \text{if } l = L \\ (\mathbf{W}_l \mathbf{E}_{l+1}). * f'.(\mathbf{X}_l) & \text{otherwise} \end{cases} \quad (36)$$

$$= \begin{cases} \mathbf{Y}_L - \mathbf{B} \\ (\mathbf{W}_l \mathbf{E}_{l+1}). * (\mathbf{1}_{dl \times N} - \mathbf{Y}_l. * \mathbf{Y}_l) & \text{if } f(x) = \tanh(x) \end{cases} \quad (37)$$

$$\tilde{\mathbf{W}}_{l-1}(k+1) = \tilde{\mathbf{W}}_{l-1}(k) - \rho \sum_{n=1}^N \tilde{\mathbf{y}}_{l-1,n} {}^t\epsilon_{l,n} \quad (38)$$

$$= \tilde{\mathbf{W}}_{l-1}(k) - \rho \tilde{\mathbf{Y}}_{l-1} {}^t\mathbf{E}_l \quad (39)$$