

Neural Network Derivatives

Table of Contents

1	Preface.....	2
2	Common Derivatives.....	3
3	Perceptron.....	5
	Output equations.....	6
	Back propagation.....	6
4	Multi Layer Perceptron.....	9
	Output equations.....	10
	Second layer derivatives.....	10
	First layer derivatives.....	11
5	Autoencoder.....	13
	Output equations.....	14
	Fourth layer derivatives.....	15
	Third layer derivatives.....	15
	Second layer derivatives.....	15
	First layer derivatives.....	16
6	Ensemble learning via negative correlation.....	18
	Output equations.....	19
	Negative correlation learning equations.....	19
	Second layer derivatives.....	20
	First layer derivatives.....	21
7	Ensemble Learning Using Decorrelated Nns.....	22
	Decorrelation equations.....	23
	Output equations.....	24
	Second layer derivatives.....	24
	First layer derivatives.....	25
8	Reference.....	26

1 Preface

This is my personal collection of neural network derivatives. I started this document during the machine learning course at Universidade de São Paulo, when I created the algorithms from scratch.

Derivatives are the cornerstone of neural networks, they're essential part of backpropagation. The algorithm doesn't converge if the any small detail is missing or misplaced. It requires time, patience and practice to master the skills of calculating derivatives.

Well, someone might ask: “why should I bother learning derivatives and building neural networks from scratch as there're plenty of libraries available?”. I would say for a couple of good reasons. While off the shelf libraries usually are optimized for speed, understanding how they work internally is essential when setting up the parameters and selecting the architecture, and becomes crucial for debugging or optimizing the model for speed or handling large datasets.

Besides this, the frequency that commercial libraries are updated is too long after new models are published. The fast development of machine learning algorithms in recent years leaves commercial libraries well behind the scientific research and the academia. Being able to deploy new models soon after a paper is released gives enormous advantage.

I use this document to implement network architectures from scientific publications, and as a reference for my codes. I tested all derivatives in my code, which are available in my GitHub account. If you find any bug or typo, please submit a pull request. I hope you enjoy this document as much as I did writing it.

August, 2020

Victor Ivamoto

Sao Paulo, Brazil

2 Common Derivatives

Sigmoid

$$f(x) = \frac{1}{1+e^{-x}} \rightarrow \frac{d f(x)}{dx} = f(x)(1-f(x)) \quad (2.1)$$

Softmax

$$f(x) = \frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}} \rightarrow \frac{d f(x)}{dx} = f(x_i)(\delta_{ij} - f(x_j)), \delta_{ik} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases} \quad (2.2)$$

Hyperbolic tangent

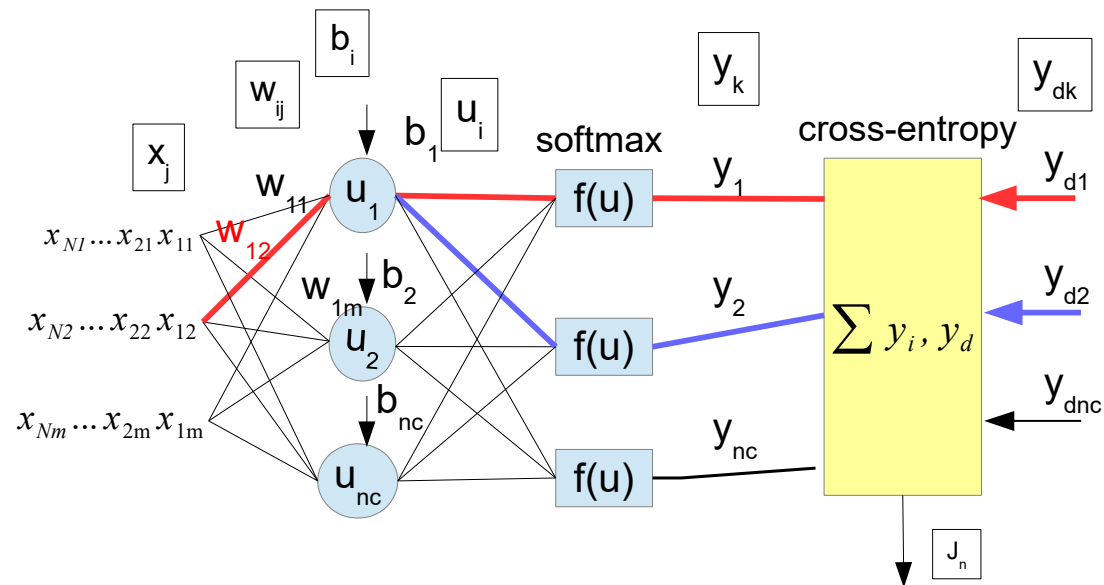
$$f(x) = \tanh(x) \rightarrow \frac{d f(x)}{dx} = 1 - \tanh^2(x) \quad (2.3)$$

Common Functions	Function	Derivative
Constant	c	0
Line	x	1
	ax	a
Square	x^2	$2x$
Square Root	\sqrt{x}	$\frac{1}{2\sqrt{x}}$
Exponential	e^x	e^x
	a^x	$\ln(a)a^x$
Logarithms	$\ln(x)$	$\frac{1}{x}$
	$\log_a(x)$	$\frac{1}{x \ln(a)}$

Common Derivatives

Rules	Function	Derivative
Multiplication by constant	$c f(x)$	$c f'(x)$
Power Rule	x^n	$n x^{n-1}$
Sum Rule	$f(x) + g(x)$	$f'(x) + g'(x)$
Difference Rule	$f(x) - g(x)$	$f'(x) - g'(x)$
Product Rule	$f(x)g(x)$	$f'(x)g(x) + f(x)g'(x)$
Quotient Rule	$\frac{f(x)}{g(x)}$	$\frac{f'(x)g(x) - f(x)g'(x)}{g^2(x)}$
Reciprocal Rule	$\frac{1}{f(x)}$	$-\frac{f'(x)}{f^2(x)}$
Chain Rule		$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx}$

3 Perceptron



Dataset with N elements $D = \{(\vec{x}_1, y_1), \dots, (\vec{x}_N, y_N)\}$, where \vec{x} is a vector and $y = \{0, 1\}$ is one-hot encoded.

Activation function: softmax.

Cost function: cross-entropy.

Definitions W_{ij} : connects u_i to x_j nc : number of classes m : number of attributes N : number of instances	Definitions X ($N \times m$): Input values w ($nc \times m$): weights b ($nc \times 1$): biases u ($nc \times 1$): input function y ($N \times nc$): output function - softmax y_d ($N \times nc$): desired output
---	---

Output equations

Input function is the linear combination of input values.

$$u_i = \sum_{j=1}^m x_j w_{ij} + b_i \quad (3.1)$$

Activation function: softmax.

$$y_k = f(u) = f(u_1, \dots, u_{nc}) = \frac{e^{u_k}}{\sum_{t=1}^{nc} e^{u_t}} \quad (3.2)$$

Cost function: cross-entropy.

$$J(n) = - \sum_{k=1}^{nc} y_{d_k} \log y_k \quad (3.3)$$

Total cost function:

$$J_T = \frac{1}{N} \sum_{n=1}^N J(n) \quad (3.4)$$

Back propagation

$$\frac{\partial J}{\partial w_{ij}} = \frac{\partial J}{\partial y_k} \frac{\partial y_k}{\partial u_i} \frac{\partial u_k}{\partial w_{ij}} \quad (3.5)$$

$$\frac{\partial J}{\partial y_k} = - \sum_{k=1}^{nc} \frac{y_{dk}}{y_k} \quad (3.6)$$

For $i=k$

$$\frac{\partial y_k}{\partial u_i} = \frac{e^{u_k} \sum_{t=1}^{nc} e^{u_t} - e^{u_k} e^{u_i}}{\left(\sum_{t=1}^{nc} e^{u_t} \right)^2} = \frac{e^{u_k} \sum_{t=1}^{nc} e^{u_t}}{\left(\sum_{t=1}^{nc} e^{u_t} \right)^2} - \frac{e^{u_k} e^{u_i}}{\left(\sum_{t=1}^{nc} e^{u_t} \right)^2} = y_k - y_k y_i = y_k (1 - y_i) \quad (3.7)$$

For $i \neq k$

$$\frac{\partial y_k}{\partial u_i} = \frac{0 - e^{u_k} e^{u_i}}{\left(\sum_{t=1}^{nc} e^{u_t}\right)^2} = \frac{-e^{u_k}}{\left(\sum_{t=1}^{nc} e^{u_t}\right)} \cdot \frac{e^{u_i}}{\left(\sum_{t=1}^{nc} e^{u_t}\right)} = -y_k y_i \quad (3.8)$$

Combining (3.7) and (3.8):

$$\frac{\partial y_k}{\partial u_i} = y_k (\delta_{ik} - y_i), \delta_{ik} = \begin{cases} 1 & \text{if } i=k \\ 0 & \text{if } i \neq k \end{cases} \quad (3.9)$$

$$\frac{\partial u_i}{\partial w_{ij}} = x_j \quad (3.10)$$

$$\frac{\partial u_i}{\partial b_i} = 1 \quad (3.11)$$

Combining (3.6), (3.9) and (3.10) into (3.5):

$$\frac{\partial J}{\partial w_{ij}} = - \sum_{k=1}^{nc} \frac{y_{dk}}{y_k} y_k (\delta_{ik} - y_i) x_j = \sum_{k=1}^{nc} y_{dk} (y_i - \delta_{ik}) x_j = \sum_{k=1}^{nc} y_{dk} y_i x_j - \sum_{k=1}^{nc} y_{dk} \delta_{ik} x_j \quad (3.12)$$

Since $\sum_{k=1}^{nc} y_{dk} = 1$ and replacing δ_{ik} from equation 9 into 12:

$$\frac{\partial J}{\partial w_{ij}} = y_i x_j - y_{di} x_j = (y_i - y_{di}) x_j \quad (3.13)$$

Finally:

$$\frac{\partial J_T}{\partial w_{ij}} = \frac{1}{N} \sum_{n=1}^N \frac{\partial J_n}{\partial w_{ij}} \quad (3.14)$$

$$\frac{\partial J_T}{\partial w_{ij}} = \frac{1}{N} \sum_{n=1}^N (y_i - y_{di}) x_j \quad (3.15)$$

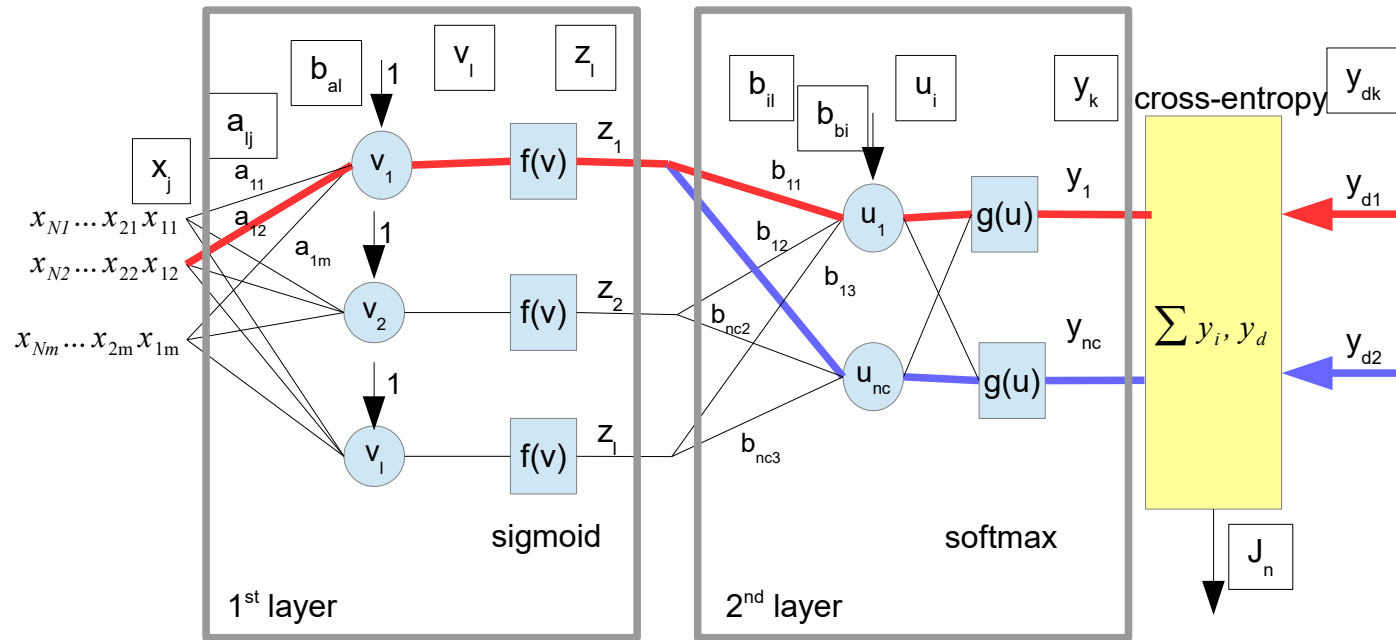
Similarly, we compute the derivative of J w.r.t. bias B :

$$\frac{\partial J}{\partial b_i} = - \sum_{k=1}^{nc} \frac{y_{dk}}{y_k} y_k (\delta_{ik} - y_i) = y_i - y_{di} \quad (3.16)$$

$$\frac{\partial J_T}{\partial b_i} = \frac{1}{N} \sum_{n=1}^N \frac{\partial J_n}{\partial b_i} \quad (3.17)$$

$$\frac{\partial J_T}{\partial b_i} = \frac{1}{N} \sum_{n=1}^N (y_i - y_{di}) \quad (3.18)$$

4 Multi Layer Perceptron



First layer activation function: sigmoid.

Second layer activation function: softmax

Cost function: cross-entropy.

Dataset with N elements $D = \{(\vec{x}_1, y_1), \dots, (\vec{x}_N, y_N)\}$, where \vec{x} is a vector and $y = \{0, 1\}$ is one-hot encoded.

Definitions a_{ij} : connects v_i to x_j nc : number of classes m : number of attributes N : number of instances L : number of 1 st layer neurons	Definitions X ($N \times m$): Input values a ($L \times m$): 1 st layer weights b_a ($L \times 1$): 1 st layer biases v : ($L \times 1$): 1 st layer input z ($L \times 1$): 1 st layer output b ($nc \times L$): 2 nd layer weights b_b ($nc \times 1$): 2 nd layer biases u ($nc \times 1$): 2 nd layer input y ($N \times nc$): 2 nd layer output y_d ($N \times nc$): desired output – one-hot encoded J_n ($N \times 1$): cost function
--	--

Output equations

First layer input.

$$v_l = \sum_{j=1}^m x_j a_{lj} + b_{al} \quad (4.1)$$

First layer output – sigmoid.

$$z_l = f(v) = \frac{1}{1 + e^{-v_l}} \quad (4.2)$$

Second layer input.

$$u_i = \sum_{l=1}^L z_l b_{il} + b_{bi} \quad (4.3)$$

Second layer output – softmax.

$$y_k = f(u) = f(u_1, \dots, u_{nc}) = \frac{e^{u_k}}{\sum_{t=1}^{nc} e^{u_t}} \quad (4.4)$$

Cost function – cross-entropy.

$$J(n) = - \sum_{k=1}^{nc} y_{d_k} \log y_k \quad (4.5)$$

Total cost function.

$$J_T = \frac{1}{N} \sum_{n=1}^N J(n) \quad (4.6)$$

Second layer derivatives

$$\frac{\partial J(n)}{\partial b_{il}} = \frac{\partial J(n)}{\partial y_k} \frac{\partial y_k}{\partial u_i} \frac{\partial u_i}{\partial b_{il}} \quad (4.7)$$

$$\frac{\partial J(n)}{\partial y_k} = - \sum_{k=1}^{nc} \frac{y_{d_k}}{y_k} \quad (4.8)$$

$$\frac{\partial y_k}{\partial u_i} = \frac{e^{u_k} \sum_{l=1}^{nc} e^{u_l} - e^{u_k} e^{u_i}}{\left(\sum_{l=1}^{nc} e^{u_l} \right)^2} = y_k - y_k y_i \quad (4.9)$$

$$\frac{\partial y_k}{\partial u_i} = \sum_{k=1}^{nc} y_k (\delta_{ik} - y_i), \delta_{ik} = \begin{cases} 1 & \text{if } i=k \\ 0 & \text{if } i \neq k \end{cases} \quad (4.10)$$

$$\frac{\partial u_i}{\partial b_{il}} = z_l \quad (4.11)$$

Combining (4.8), (4.10) and (4.11) into (4.7):

$$\frac{\partial J(n)}{\partial b_{il}} = - \sum_{k=1}^{nc} \frac{y_{d_k}}{y_k} y_k (\delta_{ik} - y_i) z_l = \sum_{k=1}^{nc} y_{d_k} (y_i - \delta_{ik}) z_l \quad (4.12)$$

$$\frac{\partial J(n)}{\partial b_{il}} = \sum_{k=1}^{nc} y_{d_k} y_i z_l - \sum_{k=1}^{nc} y_{d_k} \delta_{ik} z_l \quad (4.13)$$

Since y_d is one-hot encoded, $\sum_{k=1}^{nc} y_{d_k} = 1$. Also, $\delta_{ik} = 1$ if $i=k$ and $\delta_{ik}=0$ otherwise.

$$\frac{\partial J(n)}{\partial b_{il}} = y_i z_l - y_{d_i} z_l = (y_i - y_{d_i}) z_l \quad (4.14)$$

Derivative of J w.r.t. weights:

$$\frac{\partial J_T}{\partial b_{il}} = \frac{1}{N} \sum_{n=1}^N (y_i - y_{d_i}) z_l \quad (4.15)$$

Derivative of J w.r.t. bias:

$$\frac{\partial J_T}{\partial b_{bi}} = \frac{1}{N} \sum_{n=1}^N (y_i - y_{d_i}) \quad (4.16)$$

First layer derivatives

$$\frac{\partial J(n)}{\partial a_{lj}} = \frac{\partial J(n)}{\partial y_k} \frac{\partial y_k}{\partial u_i} \frac{\partial u_i}{\partial v_l} \frac{\partial v_l}{\partial a_{lj}} \quad (4.17)$$

$$\frac{\partial u_i}{\partial v_l} = b_{il} z_l (1 - z_l) \quad (4.18)$$

$$\frac{\partial v_l}{\partial a_{lj}} = x_j \quad (4.19)$$

Combining (4.8), (4.10), (4.18) and (4.19) into (4.17)

$$\frac{\partial J(n)}{\partial a_{lj}} = - \sum_{k=1}^{nc} \frac{y_{d_k}}{y_k} y_k (\delta_{ik} - y_i) b_{il} z_l (1 - z_l) x_j \quad (4.20)$$

$$\frac{\partial J(n)}{\partial a_{lj}} = \sum_{k=1}^{nc} y_{d_k} (y_i - \delta_{ik}) b_{il} z_l (1 - z_l) x_j = \sum_{k=1}^{nc} y_{d_k} y_i b_{il} z_l (1 - z_l) x_j - \sum_{k=1}^{nc} y_{d_k} \delta_{ik} b_{il} z_l (1 - z_l) x_j \quad (4.21)$$

If $i = k$

$$\frac{\partial J(n)}{\partial a_{lj}} = \sum_{k=1}^{nc} y_{d_k} y_i b_{il} z_l (1 - z_l) x_j - y_{d_i} b_{il} z_l (1 - z_l) x_j = (y_i - y_{d_i}) [b_{il} z_l (1 - z_l) x_j] \quad (4.22)$$

$$\frac{\partial J(n)}{\partial a_{lj}} = (y_i - y_{d_i}) [b_{il} z_l (1 - z_l) x_j] \quad (4.23)$$

Derivative w.r.t. weights:

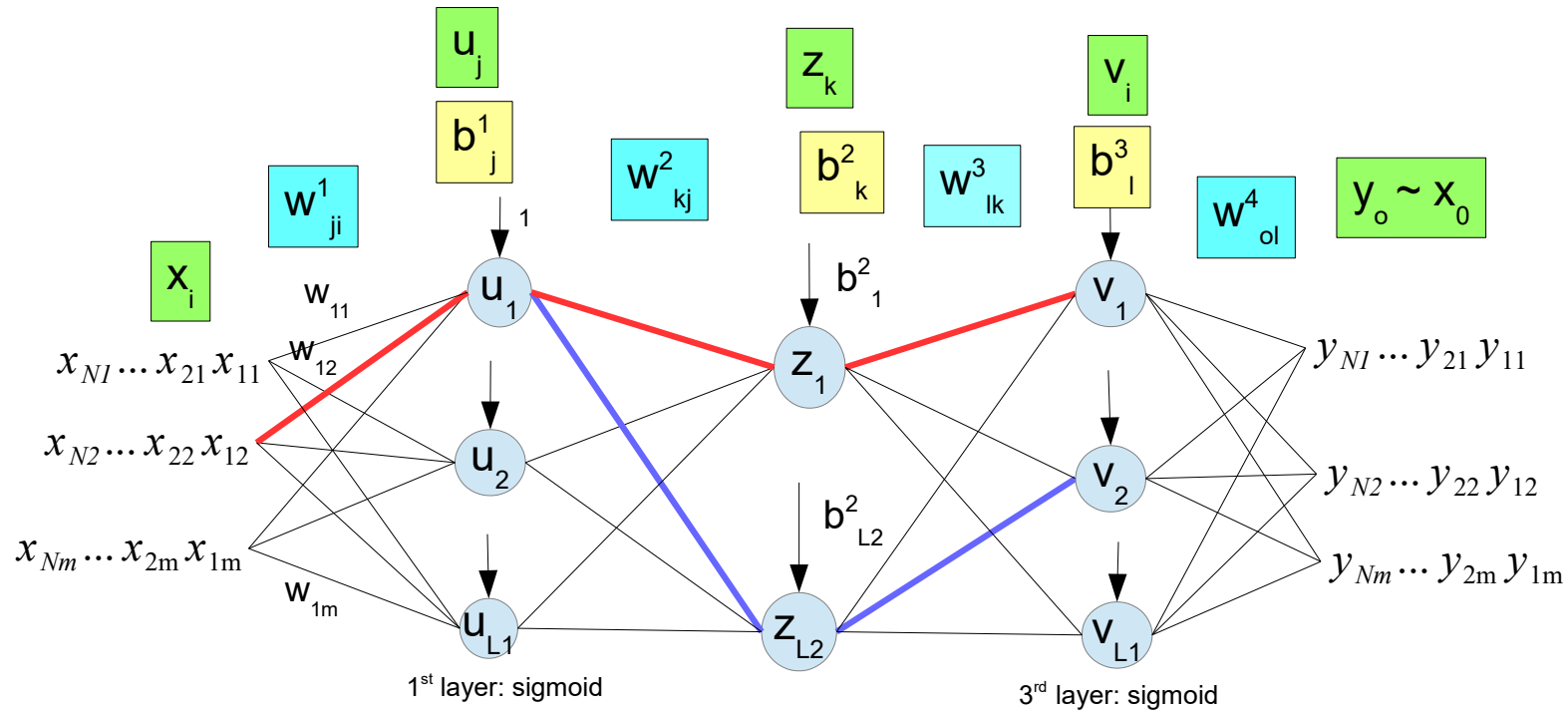
$$\frac{\partial J_T}{\partial a_{lj}} = \frac{1}{N} \sum_{n=1}^N (y_i - y_{d_i}) [b_{il} z_l (1 - z_l) x_j] \quad (4.24)$$

$$\frac{\partial v_l}{\partial b_{al}} = 1 \quad (4.25)$$

Derivative w.r.t. bias:

$$\frac{\partial J_T}{\partial b_{al}} = \frac{1}{N} \sum_{n=1}^N (y_i - y_{d_i}) [b_{il} z_l (1 - z_l)] \quad (4.26)$$

5 Autoencoder



Definitions: L1: 1 st and 3 rd layer neurons L2: 2 nd layer neurons m: number of inputs N: number of instances 1 st and 3 rd layers activation function: sigmoid 2 nd and 4 th layers activation function: none Cost function: MSE	Definitions: X (N x m): input values W ¹ (m x L1): 1 st layer weights W ² (L1 x L2): 2 nd layer weights W ³ (L2 x L1): 3 rd layer weights W ⁴ (L1 x m): 4 th layer weights b ¹ (L1): 1 st layer biases b ² (L2): 2 nd layer biases b ³ (L1): 3 rd layer biases b ⁴ (m): 4 th layer biases U (N x L1): 1 st layer neuron Z (N x L2): encoded values of X. V (N x L1): 3 rd layer neuron Y (N x m): decoded values of Z
---	---

Output equations

First layer input.

$$Uin_j = \sum_{i=1}^m X_i w_{ji}^1 + b_j^1 \quad (5.1)$$

First layer output – sigmoid.

$$U_j = \frac{1}{1 + e^{-Uin_j}} \quad (5.2)$$

The encoded values of X is just a linear combination of 1st layer neurons.

$$Zin_k = \sum_{j=1}^{L1} U_j w_{kj}^2 + b_k^2 \quad (5.3)$$

Second layer input.

$$Vin_l = \sum_{k=1}^{L2} Z_k w_{lk}^3 + b_l^3 \quad (5.4)$$

Second layer output – sigmoid.

$$V_l = \frac{1}{1 + e^{-Vin_l}} \quad (5.5)$$

Decoded values of Z is the linear combination of second layer neurons.

$$Y_o = \sum_{l=1}^{L1} V_l w_{ol}^4 + b_o^4 \quad (5.6)$$

Cost function: mean squared error.

$$J = \frac{(Y_o - X_o)^2}{2} \quad (5.7)$$

Total cost function is the mean of all cost functions.

$$J_T = \frac{1}{N} \sum_{n=1}^N J(n) \quad (5.8)$$

Fourth layer derivatives

$$\frac{\partial J_T}{\partial w_{ol}^4} = \frac{1}{N} \sum_{n=1}^N \frac{\partial J}{\partial y_o} \frac{\partial y_o}{\partial w_{ol}^4} = \frac{1}{N} \sum_{n=1}^N (y_o - y_o) v_l \quad (5.9)$$

$$\frac{\partial J_T}{\partial b_o^4} = \frac{1}{N} \sum_{n=1}^N \frac{\partial J}{\partial y_o} \frac{\partial y_o}{\partial b_o^4} = \frac{1}{N} \sum_{n=1}^N (y_o - y_o) \quad (5.10)$$

Third layer derivatives

$$\frac{\partial J_T}{\partial w_{lk}^3} = \frac{1}{N} \sum_{n=1}^N \frac{\partial J}{\partial y} \frac{\partial y}{\partial v_l} \frac{\partial v_l}{\partial w_{lk}^3} \quad (5.11)$$

$$\frac{\partial J}{\partial y} = \sum_{o=1}^m (y_o - x_o) \quad (5.12)$$

$$\frac{\partial y_o}{\partial v_l} = w_{ol}^4 v_l (1 - v_l) \quad (5.13)$$

$$\frac{\partial v_l}{\partial w_{lk}^3} = z_k \quad (5.14)$$

$$\frac{\partial J_T}{\partial w_{lk}^3} = \frac{1}{N} \sum_{n=1}^N \sum_{o=1}^m (y_o - x_o) w_{ol}^4 v_l (1 - v_l) z_k \quad (5.15)$$

$$\frac{\partial J_T}{\partial b_l^3} = \frac{1}{N} \sum_{n=1}^N \sum_{o=1}^m (y_o - x_o) w_{ol}^4 v_l (1 - v_l) \quad (5.16)$$

Second layer derivatives

$$\frac{\partial J_T}{\partial w_{kj}^2} = \frac{1}{N} \sum_{n=1}^N \frac{\partial J}{\partial y} \frac{\partial y}{\partial v} \frac{\partial v}{\partial z_k} \frac{\partial z_k}{\partial w_{kj}^2} \quad (5.17)$$

$$\frac{\partial y}{\partial v} = \sum_{l=1}^{LI} w_{ol}^4 v_l (1 - v_l) \quad (5.18)$$

$$\frac{\partial v}{\partial z_k} = w_{lk}^3 \quad (5.19)$$

$$\frac{\partial z_k}{\partial w_{kj}^2} = u_j \quad (5.20)$$

$$\frac{\partial J_T}{\partial w_{kj}^2} = \frac{1}{N} \sum_{n=1}^N \sum_{o=1}^m \sum_{l=1}^{LI} (y_o - x_o) w_{ol}^4 v_l (1 - v_l) w_{lk}^3 u_j \quad (5.21)$$

$$\frac{\partial J_T}{\partial b_k^2} = \frac{1}{N} \sum_{n=1}^N \frac{\partial J}{\partial y} \frac{\partial y}{\partial v} \frac{\partial v}{\partial z_k} \frac{\partial z_k}{\partial b_k^2} \quad (5.22)$$

$$\frac{\partial J_T}{\partial b_k^2} = \frac{1}{N} \sum_{n=1}^N \sum_{o=1}^m \sum_{l=1}^{LI} (y_o - x_o) w_{ol}^4 v_l (1 - v_l) w_{lk}^3 \quad (5.23)$$

First layer derivatives

$$\frac{\partial J_T}{\partial w_{ji}^1} = \frac{1}{N} \sum_{n=1}^N \frac{\partial J}{\partial y} \frac{\partial y}{\partial v} \frac{\partial v}{\partial z} \frac{\partial z}{\partial u_j} \frac{\partial u_j}{\partial w_{ji}^1} \quad (5.24)$$

$$\frac{\partial y}{\partial v} = \sum_{l=1}^{LI} w_{ol}^4 v_l (1 - v_l) \quad (5.25)$$

$$\frac{\partial v}{\partial z} = \sum_{k=1}^{L2} w_{lk}^3 \quad (5.26)$$

$$\frac{\partial z}{\partial u_j} = w_{kj}^2 u_j (1 - u_j) \quad (5.27)$$

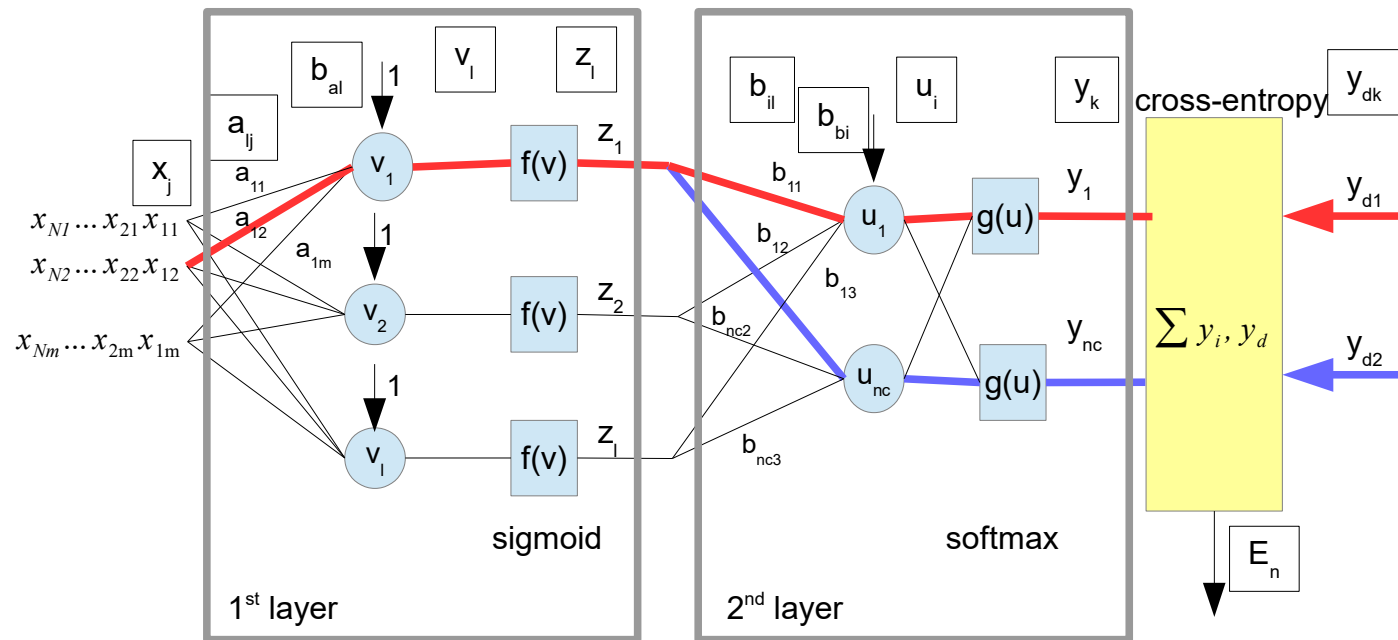
$$\frac{\partial u_j}{\partial w_{ji}^1} = x_i \quad (5.28)$$

$$\frac{\partial J_T}{\partial w_{ji}^1} = \frac{1}{N} \sum_{n=1}^N \sum_{o=1}^m \sum_{l=1}^{LI} \sum_{k=1}^{L2} (y_o - x_o) w_{ol}^4 v_l (1 - v_l) w_{lk}^3 w_{kj}^2 u_j (1 - u_j) x_i \quad (5.29)$$

$$\frac{\partial J_T}{\partial b_j^1} = \frac{1}{N} \sum_{n=1}^N \frac{\partial J}{\partial y} \frac{\partial y}{\partial v} \frac{\partial v}{\partial z} \frac{\partial z}{\partial u_j} \frac{\partial u_j}{\partial b_j^1} \quad (5.30)$$

$$\frac{\partial J_T}{\partial b_j^1} = \frac{1}{N} \sum_{n=1}^N \sum_{o=1}^m \sum_{l=1}^{L1} \sum_{k=1}^{L2} (y_o - x_o) w_{ol}^4 v_l (1 - v_l) w_{lk}^3 w_{kj}^2 u_j (1 - u_j) \quad (5.31)$$

6 Ensemble learning via negative correlation



Definitions a_{ij} : connects v_i to x_j nc : number of classes m : number of attributes N : number of instances L : number of 1 st layer neurons	Definitions X ($N \times m$): input vector, normalized values a ($L \times m$): 1 st layer weights b_a ($N \times L$): 1 st layer biases v ($N \times L$): 1 st layer input z ($N \times L$): 1 st layer output b ($nc \times L$): 2 nd layer weights b_b ($N \times nc$): 2 nd layer biases u ($N \times nc$): 2 nd layer input y ($N \times nc$): 2 nd layer output y_d ($N \times nc$): desired output E_n ($N \times nc$): error function
--	---

Output equations

First layer input equation.

$$v_i = \sum_{j=1}^m x_j a_{ij} + b_{ai} \quad (6.1)$$

First layer output equation – sigmoid.

$$z_i = f(v) = \frac{1}{1 + e^{-v_i}} \quad (6.2)$$

Second layer input equation.

$$u_i = \sum_{l=1}^L z_l b_{il} + b_{bi} \quad (6.3)$$

Second layer output – softmax.

$$y_k = f(u) = f(u_1, \dots, u_{nc}) = \frac{e^{u_k}}{\sum_{t=1}^{nc} e^{u_t}} \quad (6.4)$$

Negative correlation learning equations

Dataset:

$$D = \{(x(1), d(1)), \dots, (x(N), d(N))\} \quad (6.5)$$

Where $x \in \mathbb{R}^p$, d is a scalar, and N is the size of the training set. The assumption that the output d is a scalar has been made merely to simplify exposition of ideas without loss of generality.

Error function for network i (Equation 2):

$$E_i = \frac{1}{N} \sum_{n=1}^N E_i(n) = \frac{1}{N} \sum_{n=1}^N \frac{1}{2} (F_i(n) - d(n))^2 + \frac{1}{N} \sum_{n=1}^N \lambda p_i(n) \quad (6.6)$$

where $E_i(n)$ is the value of the error function of network i at presentation of the n th training pattern. The parameter $0 \leq \lambda \leq 1$ is used to adjust the strength of the penalty.

Correlation penalty function p_i (Equation 3):

$$p_i(n) = (F_i(n) - F(n)) \sum_{j \neq i} (F_j(n) - F(n)) \quad (6.7)$$

Partial derivative of $E_i(n)$ with respect to the output of network i on the n th training pattern (Equation 4):

$$\frac{\partial E_p(n)}{\partial F_p(n)} = (1 - \lambda)(F_p(n) - y_d(n)) + \lambda(F(n) - y_d(n)) \quad (6.8)$$

$$F_p = y, d = y_d \quad (6.9)$$

Second layer derivatives

$$\frac{\partial E_p(n)}{\partial b_{il}(n)} = \frac{\partial E_p}{\partial F_p} \frac{\partial F_p}{\partial u_i} \frac{\partial u_i}{\partial b_{il}} \quad (6.10)$$

$$\frac{\partial F_p}{\partial u_i} = \frac{\partial y_k}{\partial u_i} = y_k(\delta_{ik} - y_i), \delta_{ik} = \begin{cases} 1 & \text{if } i = k \\ 0 & \text{if } i \neq k \end{cases} \quad (6.11)$$

$$\frac{\partial u_i}{\partial b_{il}} = z_l \quad (6.12)$$

$$\frac{\partial E_p(n)}{\partial b_{il}(n)} = [(1 - \lambda)(F_p(n) - y_d(n)) + \lambda(F(n) - y_d(n))] y_k(\delta_{ik} - y_i) z_l \quad (6.13)$$

If $i = k, \delta_{ii} = 1$:

$$\frac{\partial E_p(n)}{\partial b_{il}(n)} = [(1 - \lambda)(F_p(n) - y_d(n)) + \lambda(F(n) - y_d(n))] y_i(1 - y_i) z_l \quad (6.14)$$

If $i \neq k, \delta_{ik} = 0$:

$$\frac{\partial E_p(n)}{\partial b_{il}(n)} = -[(1 - \lambda)(F_p(n) - y_d(n)) + \lambda(F(n) - y_d(n))] y_k y_i z_l \quad (6.15)$$

First layer derivatives

$$\frac{\partial E_p(n)}{\partial a_{lj}(n)} = \frac{\partial E_p}{\partial F_p} \frac{\partial F_p}{\partial u_i} \frac{\partial u_i}{\partial v_l} \frac{\partial v_l}{\partial a_{lj}} \quad (6.16)$$

$$\frac{\partial F_p}{\partial u_i} = \frac{\partial y_k}{\partial u_i} = \sum_{k=1}^{nc} y_k (\delta_{ik} - y_i), \delta_{ik} = \begin{cases} 1 & \text{if } i=k \\ 0 & \text{if } i \neq k \end{cases} \quad (6.17)$$

$$\frac{\partial u_i(n)}{\partial v_l(n)} = b_{il} z_l (1 - z_l) \quad (6.18)$$

$$\frac{\partial v_l(n)}{\partial a_{lj}(n)} = x_j \quad (6.19)$$

$$\frac{\partial E_p(n)}{\partial a_{lj}(n)} = [(1-\lambda)(F_p(n) - y_d(n)) + \lambda(F(n) - y_d(n))] y_k (\delta_{ik} - y_i) b_{il} z_l (1 - z_l) x_j \quad (6.20)$$

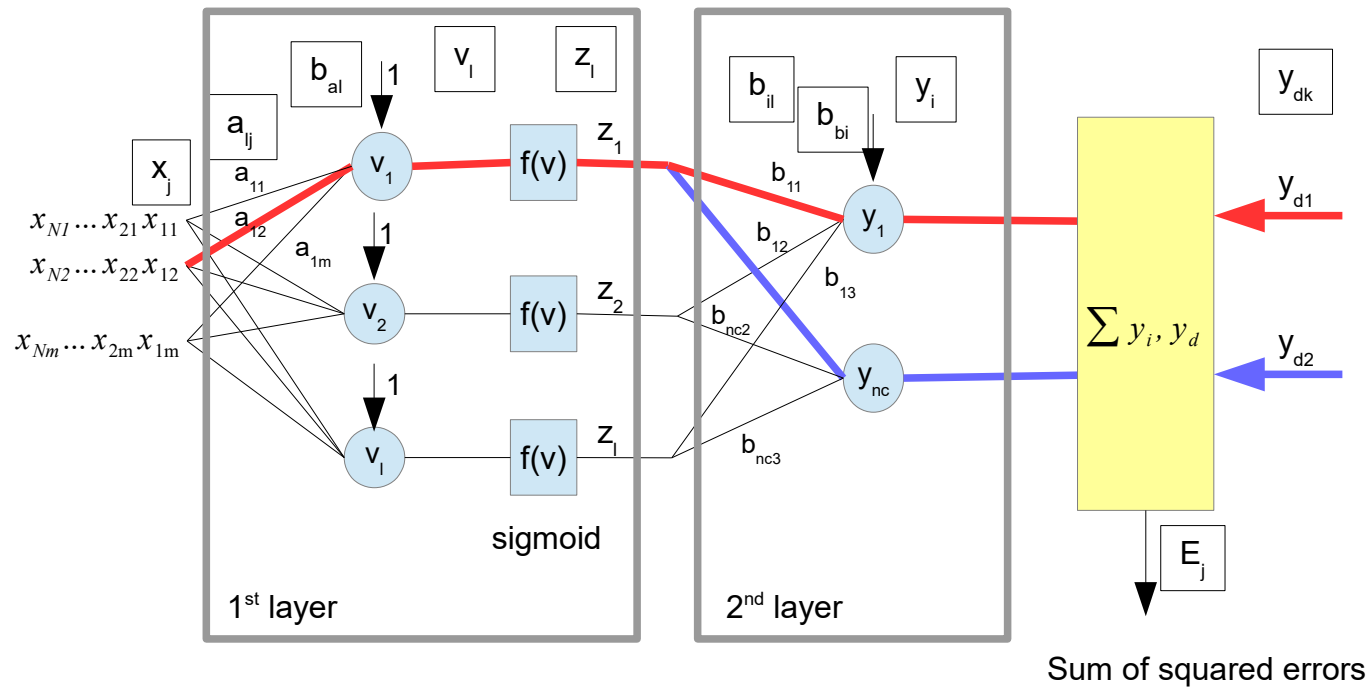
If $i=k, \delta_{ii}=1$:

$$\frac{\partial E_p(n)}{\partial a_{lj}(n)} = [(1-\lambda)(F_p(n) - y_d(n)) + \lambda(F(n) - y_d(n))] y_i (1 - y_i) b_{il} z_l (1 - z_l) x_j \quad (6.21)$$

If $i \neq k, \delta_{ik}=0$:

$$\frac{\partial E_p(n)}{\partial a_{lj}(n)} = -[(1-\lambda)(F_p(n) - y_d(n)) + \lambda(F(n) - y_d(n))] y_k y_i b_{il} z_l (1 - z_l) x_j \quad (6.22)$$

7 Ensemble Learning Using Decorrelated Nns



First layer activation function: sigmoid.

Second layer activation function: none.

Cost function: sum of squared errors.

Definitions a_{ij} : connects v_i to x_j nc : number of classes m : number of attributes N : number of instances L : number of 1 st layer neurons	Definitions X ($N \times m$): input vector, normalized values a ($L \times m$): 1 st layer weights b_a ($N \times L$): 1 st layer biases v ($N \times L$): 1 st layer input z ($N \times L$): 1 st layer output b ($nc \times L$): 2 nd layer weights b_b ($N \times nc$): 2 nd layer biases y ($N \times nc$): 2 nd layer output y_d ($N \times nc$): desired output E_n ($N \times nc$): error function
--	---

Decorrelation equations

The training set is a set of N patterns $\{(\vec{x}_1, y_1), \dots, (\vec{x}_N, y_N)\}$ with the p th pattern defined by some unknown relationship:

$$y_p = g(\vec{x}_p) + \varepsilon \quad (7.1)$$

where g is a regression function, and ε is some mean zero additive noise with finite variance σ^2 .

Error function for an individual network i (equation 7):

$$E_j = \sum_{p=1}^N \left[(y_p - f_j(\vec{x}_p))^2 + \sum_{i=1}^{j-1} \lambda(t) d(i, j) P(\vec{x}_p, y_p, f_i, f_j) \right] \quad (7.2)$$

Where $f_j(\vec{x})$ is the output of the j th network.

Correlation penalty function P (equation 8):

$$P(\vec{x}, y, f_i, f_j) = (y - f_i(\vec{x}))(y - f_j(\vec{x})) \quad (7.3)$$

The indicator function d specifies which individual networks are to be decorrelated. For example, to penalize an individual network for being correlated with the previously trained network, the indicator function is (eq. 9):

$$d(i, j) = \begin{cases} 1, & \text{if } i = j - 1 \\ 0, & \text{otherwise} \end{cases} \quad (7.4)$$

To allow alternate networks to be trained independently of one another yet decorrelate pairs of networks, the indicator function can be defined as (equation 10)

$$d(i, j) = \begin{cases} 1, & \text{if } i = j - 1 \text{ and } i \text{ is even} \\ 0, & \text{otherwise} \end{cases} \quad (7.5)$$

The scaling function $\lambda(t)$ is either constant, or is time dependent.

$$E_j = \sum_{p=1}^N \left[(y_p - f_j(\vec{x}_p))^2 + \sum_{i=1}^{j-1} \lambda(t) d(i, j) (y - f_i(\vec{x}))(y - f_j(\vec{x})) \right] \quad (7.6)$$

Replacing with our variables:

$$y = y_d, f_j(\vec{x}) = y \quad (7.7)$$

$$E_j = \sum_{p=1}^N \left[(y_{d_p} - y_j)^2 + \sum_{i=1}^{j-1} \lambda(t) d(i, j) (y_d - y_i) (y_d - y_j) \right] \quad (7.8)$$

Output equations

First layer input.

$$v_l = \sum_{j=1}^m x_j a_{lj} + b_{al} \quad (7.9)$$

Fist layer output: sigmoid.

$$z_l = f(v) = \frac{1}{1 + e^{-v_l}} \quad (7.10)$$

MLP output is the linear combination of first layer input.

$$f_j(\vec{x}) = y_i = \sum_{l=1}^L z_l b_{il} + b_{bl} \quad (7.11)$$

Second layer derivatives

$$\frac{\partial E_j}{\partial b_{il}} = \frac{\partial E_j}{\partial f_j(\vec{x}_p)} \frac{\partial f_j(\vec{x}_p)}{\partial b_{il}} \quad (7.12)$$

$$\frac{\partial E_j}{\partial f_j(\vec{x}_p)} = \sum_{p=1}^N \left[-2(y_p - f_j(\vec{x}_p)) f'_j(\vec{x}_p) + \sum_{i=1}^{j-1} \lambda(t) d(i, j) (y - f_i(\vec{x})) (-f'_j(\vec{x})) \right] \quad (7.13)$$

Rearranging (7.13):

$$\frac{\partial E_j}{\partial f_j(\vec{x}_p)} = f'_j(\vec{x}) \sum_{p=1}^N \left[2(f_j(\vec{x}_p) - y_p) + \sum_{i=1}^{j-1} \lambda(t) d(i, j) (f_i(\vec{x}) - y) \right] \quad (7.14)$$

$$\frac{\partial f_j(\vec{x}_p)}{\partial b_{il}} = z_l \quad (7.15)$$

Combining (7.14) and (7.15) into (7.12) (Derivative w.r.t. weights):

$$\frac{\partial E_j}{\partial b_{il}} = \sum_{p=1}^N \left[2(f_j(\vec{x}_p) - y_p) + \sum_{i=1}^{j-1} \lambda(t) d(i, j) (f_i(\vec{x}) - y) \right] z_l \quad (7.16)$$

Derivative w.r.t. bias:

$$\frac{\partial E_j}{\partial b_{il}} = \sum_{p=1}^N \left[2(f_j(\vec{x}_p) - y_p) + \sum_{i=1}^{j-1} \lambda(t) d(i, j) (f_i(\vec{x}) - y) \right] \quad (7.17)$$

First layer derivatives

$$\frac{\partial E_j}{\partial a_{lj}} = \frac{\partial E_j}{\partial f_j(\vec{x}_p)} \frac{\partial f_j(\vec{x}_p)}{\partial v_l} \frac{\partial v_l}{\partial a_{lj}} \quad (7.18)$$

$$\frac{\partial f_j(\vec{x}_p)}{\partial v_l(n)} = b_{il} z_l (1 - z_l) \quad (7.19)$$

$$\frac{\partial v_l(n)}{\partial a_{lj}(n)} = x_j \quad (7.20)$$

Derivative w.r.t. weights:

$$\frac{\partial E_j}{\partial a_{lj}} = \sum_{p=1}^N \left[2(f_j(\vec{x}_p) - y_p) + \sum_{i=1}^{j-1} \lambda(t) d(i, j) (f_i(\vec{x}) - y) \right] b_{il} z_l (1 - z_l) x_j \quad (7.21)$$

Derivative w.r.t. bias:

$$\frac{\partial v_l(n)}{\partial b_{il}(n)} = 1 \quad (7.22)$$

$$\frac{\partial E_j}{\partial b_{il}} = \sum_{p=1}^N \left[2(f_j(\vec{x}_p) - y_p) + \sum_{i=1}^{j-1} \lambda(t) d(i, j) (f_i(\vec{x}) - y) \right] b_{il} z_l (1 - z_l) \quad (7.23)$$

8 Reference

Math Is Fun website - <https://www.mathsisfun.com/calculus/derivatives-rules.html>

Eli Bendersky's website The Softmax function and its derivative
<https://eli.thegreenplace.net/2016/the-softmax-function-and-its-derivative/>

Liu, Yong & Yao, Xin. (2000). Ensemble learning via negative correlation. Neural networks : the official journal of the International Neural Network Society. 12. 1399-1404. 10.1016/S0893-6080(99)00073-8.

Rosen, B. (1996). Ensemble Learning Using Decorrelated Neural Networks. Connect. Sci., 8, 373-384. DOI: 10.1080/095400996116820 <https://doi.org/10.1080/095400996116820>