

# Suggested Notation for Machine Learning

Beijing Academy of Artificial Intelligence\*

June 15, 2020

## Summary

The field of machine learning is evolving rapidly in recent years. Communication between different researchers and research groups becomes increasingly important. A key challenge for communication arises from inconsistent notation usages among different papers. This proposal suggests a standard for commonly used mathematical notation for machine learning. In this first version, only some notation are mentioned and more notation are left to be done. This proposal will be regularly updated based on the progress of the field. We look forward to more suggestions to improve this proposal in future versions.

## Contents

<b>1 Dataset</b>	<b>2</b>
<b>2 Function</b>	<b>2</b>
<b>3 Loss function</b>	<b>2</b>
<b>4 Activation function</b>	<b>3</b>
<b>5 Two-layer neural network</b>	<b>3</b>
<b>6 General deep neural network</b>	<b>3</b>
<b>7 Complexity</b>	<b>4</b>
<b>8 Training</b>	<b>4</b>
<b>9 Fourier Frequency</b>	<b>4</b>
<b>10 Convolution</b>	<b>4</b>
<b>11 Notation table</b>	<b>5</b>

---

\*This document is published by Beijing Academy of Artificial Intelligence (北京智源人工智能研究院) jointly with Peking University (北京大学) and Shanghai Jiao Tong University (上海交通大学). The first version is drafted by Zhi-Qin John Xu (Corresponding: xuzhiqin@sjtu.edu.cn, Shanghai Jiao Tong University), Tao Luo (Purdue University), Zheng Ma (Purdue University), Yaoyu Zhang (Institute for Advanced Study).

## 1 Dataset

Dataset  $S = \{\mathbf{z}_i\}_{i=1}^n = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  is sampled from a distribution  $\mathcal{D}$  over a domain  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ .

$\mathcal{X}$  is the instance domain (a set),  $\mathcal{Y}$  is the label domain (a set), and  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$  is the example domain (a set).

Usually,  $\mathcal{X}$  is a subset of  $\mathbb{R}^d$  and  $\mathcal{Y}$  is a subset of  $\mathbb{R}^{d_o}$ , where  $d$  is the input dimension,  $d_o$  is the output dimension.

$n = \#S$  is the number of samples. Without specification,  $S$  and  $n$  are for the training set.

## 2 Function

A hypothesis space is denoted by  $\mathcal{H}$ . A hypothesis function is denoted by  $f_{\boldsymbol{\theta}}(\mathbf{x}) \in \mathcal{H}$  or  $f(\mathbf{x}; \boldsymbol{\theta}) \in \mathcal{H}$  with  $f_{\boldsymbol{\theta}} : \mathcal{X} \rightarrow \mathcal{Y}$ .

$\boldsymbol{\theta}$  denotes the set of parameters of  $f_{\boldsymbol{\theta}}$ .

If there exists a target function, it is denoted by  $f^*$  or  $f : \mathcal{X} \rightarrow \mathcal{Y}$  satisfying  $\mathbf{y}_i = f^*(\mathbf{x}_i)$  for  $i = 1, \dots, n$ .

## 3 Loss function

A loss function, denoted by  $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+ := [0, +\infty)$ , measures the difference between a predicted label and a true label, e.g.,  $L^2$  loss:

$$\ell(f_{\boldsymbol{\theta}}, \mathbf{z}) = \frac{1}{2}(f_{\boldsymbol{\theta}}(\mathbf{x}) - \mathbf{y})^2,$$

where  $\mathbf{z} = (\mathbf{x}, \mathbf{y})$ .  $\ell(f_{\boldsymbol{\theta}}, \mathbf{z})$  can also be written as

$$\ell(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y})$$

for convenience. Empirical risk or training loss for a set  $S = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  is denoted by  $L_S(\boldsymbol{\theta})$  or  $L_n(\boldsymbol{\theta})$  or  $R_n(\boldsymbol{\theta})$  or  $R_S(\boldsymbol{\theta})$ ,

$$L_S(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i). \quad (1)$$

The population risk or expected loss is denoted by  $L_{\mathcal{D}}(\boldsymbol{\theta})$  or  $R_{\mathcal{D}}(\boldsymbol{\theta})$

$$L_{\mathcal{D}}(\boldsymbol{\theta}) = \mathbb{E}_{\mathcal{D}} \ell(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}), \quad (2)$$

where  $\mathbf{z} = (\mathbf{x}, \mathbf{y})$  follows the distribution  $\mathcal{D}$ .

## 4 Activation function

An activation function is denoted by  $\sigma(x)$ .

**Example 1.** *Some commonly used activation functions are*

1.  $\sigma(x) = \text{ReLU}(x) = \max(0, x)$ ;
2.  $\sigma(x) = \text{sigmoid}(x) = \frac{1}{1+e^{-x}}$ ;
3.  $\sigma(x) = \tanh x$ ;
4.  $\sigma(x) = \cos x, \sin x$ .

## 5 Two-layer neural network

The neuron number of the hidden layer is denoted by  $m$ . The two-layer neural network is

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{j=1}^m a_j \sigma(\mathbf{w}_j \cdot \mathbf{x} + b_j), \quad (3)$$

where  $\sigma$  is the activation function,  $\mathbf{w}_j$  is the input weight,  $a_j$  is the output weight,  $b_j$  is the bias term. We denote the set of parameters by

$$\boldsymbol{\theta} = (a_1, \dots, a_m, \mathbf{w}_1, \dots, \mathbf{w}_m, b_1, \dots, b_m).$$

## 6 General deep neural network

The counting of the layer number excludes the input layer. An  $L$ -layer neural network is denoted by

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbf{W}^{[L-1]} \sigma \circ (\mathbf{W}^{[L-2]} \sigma \circ (\dots (\mathbf{W}^{[1]} \sigma \circ (\mathbf{W}^{[0]} \mathbf{x} + \mathbf{b}^{[0]}) + \mathbf{b}^{[1]}) \dots) + \mathbf{b}^{[L-2]}) + \mathbf{b}^{[L-1]}, \quad (4)$$

where  $\mathbf{W}^{[l]} \in \mathbb{R}^{m_{l+1} \times m_l}$ ,  $\mathbf{b}^{[l]} \in \mathbb{R}^{m_{l+1}}$ ,  $m_0 = d_{\text{in}} = d$ ,  $m_L = d_{\text{out}}$ ,  $\sigma$  is a scalar function and “ $\circ$ ” means entry-wise operation. We denote the set of parameters by

$$\boldsymbol{\theta} = (\mathbf{W}^{[0]}, \mathbf{W}^{[1]}, \dots, \mathbf{W}^{[L-1]}, \mathbf{b}^{[0]}, \mathbf{b}^{[1]}, \dots, \mathbf{b}^{[L-1]}),$$

and an entry of  $\mathbf{W}^{[l]}$  by  $\mathbf{W}_{ij}^{[l]}$ . This can also be defined recursively.

$$f_{\boldsymbol{\theta}}^{[0]}(\mathbf{x}) = \mathbf{x}, \quad (5)$$

$$f_{\boldsymbol{\theta}}^{[l]}(\mathbf{x}) = \sigma \circ (\mathbf{W}^{[l-1]} f_{\boldsymbol{\theta}}^{[l-1]}(\mathbf{x}) + \mathbf{b}^{[l-1]}) \quad 1 \leq l \leq L-1, \quad (6)$$

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = f_{\boldsymbol{\theta}}^{[L]}(\mathbf{x}) = \mathbf{W}^{[L-1]} f_{\boldsymbol{\theta}}^{[L-1]}(\mathbf{x}) + \mathbf{b}^{[L-1]}. \quad (7)$$

## 7 Complexity

The VC-dimension of a hypothesis class  $\mathcal{H}$  is denoted  $\text{VCdim}(\mathcal{H})$ .

The Rademacher complexity of a hypothesis space  $\mathcal{H}$  on a sample set  $S$  is denoted by  $\text{Rad}(\mathcal{H} \circ S)$  or  $\text{Rad}_S(\mathcal{H})$ . The complexity  $\text{Rad}_S(\mathcal{H})$  is random because of the randomness of  $S$ . The expectation of the empirical Rademacher complexity over all samples of size  $n$  is denoted by

$$\text{Rad}_n(\mathcal{H}) = \mathbb{E}_S \text{Rad}_S(\mathcal{H}).$$

## 8 Training

The Gradient Descent is often denoted by GD. The Stochastic Gradient Descent is often denoted by SGD.

A batch set is denoted by  $B$  and the batch size is denoted by  $|B|$ .

The learning rate is denoted by  $\eta$ .

## 9 Fourier Frequency

The discretized frequency is denoted by  $\mathbf{k}$ , and the continuous frequency is denoted by  $\boldsymbol{\xi}$ .

## 10 Convolution

The convolution operation is denoted by  $*$ .

## 11 Notation table

symbol	meaning	L <sup>A</sup> T <sub>E</sub> X	simplified
$\mathbf{x}$	input	<code>\bm{x}</code>	<code>\vx</code>
$\mathbf{y}$	output, label	<code>\bm{y}</code>	<code>\vy</code>
$d$	input dimension	<code>d</code>	
$d_o$	output dimension	<code>d_{\rm o}</code>	
$n$	number of samples	<code>n</code>	
$\mathcal{X}$	instances domain (a set)	<code>\mathcal{X}</code>	<code>\fX</code>
$\mathcal{Y}$	labels domain (a set)	<code>\mathcal{Y}</code>	<code>\fY</code>
$\mathcal{Z}$	$= \mathcal{X} \times \mathcal{Y}$ example domain	<code>\mathcal{Z}</code>	<code>\fZ</code>
$\mathcal{H}$	hypothesis space (a set)	<code>\mathcal{H}</code>	<code>\fH</code>
$\boldsymbol{\theta}$	a set of parameters	<code>\bm{\theta}</code>	<code>\vtheta</code>
$f_{\boldsymbol{\theta}} : \mathcal{X} \rightarrow \mathcal{Y}$	hypothesis function	<code>f_{\bm{\theta}}</code>	<code>f_{\vtheta}</code>
$f$ or $f^* : \mathcal{X} \rightarrow \mathcal{Y}$	target function	<code>f, f^*</code>	
$\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}^+$	loss function	<code>\ell</code>	
$\mathcal{D}$	distribution of $\mathcal{Z}$	<code>\mathcal{D}</code>	<code>\fD</code>
$S = \{\mathbf{z}_i\}_{i=1}^n$	$= \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ sample set		
$L_S(\boldsymbol{\theta}), L_n(\boldsymbol{\theta}),$ $R_n(\boldsymbol{\theta}), R_S(\boldsymbol{\theta})$	empirical risk or training loss		
$L_{\mathcal{D}}(\boldsymbol{\theta}), R_{\mathcal{D}}(\boldsymbol{\theta})$	population risk or expected loss		
$\sigma : \mathbb{R} \rightarrow \mathbb{R}^+$	activation function	<code>\sigma</code>	
$\mathbf{w}_j$	input weight	<code>\bm{w}_j</code>	<code>\vw_j</code>
$a_j$	output weight	<code>a_j</code>	
$b_j$	bias term	<code>b_j</code>	
$f_{\boldsymbol{\theta}}(\mathbf{x})$ or $f(\mathbf{x}; \boldsymbol{\theta})$	neural network	<code>f_{\bm{\theta}}</code>	<code>f_{\vtheta}</code>
$\sum_{j=1}^m a_j \sigma(\mathbf{w}_j \cdot \mathbf{x} + b_j)$	two-layer neural network		
$\text{VCdim}(\mathcal{H})$	VC-dimension of $\mathcal{H}$		
$\text{Rad}(\mathcal{H} \circ S), \text{Rad}_S(\mathcal{H})$	Rademacher complexity of $\mathcal{H}$ on $S$		
$\text{Rad}_n(\mathcal{H})$	Rademacher complexity over samples of size $n$		
GD	gradient descent		
SGD	stochastic gradient descent		
$B$	a batch set	<code>B</code>	
$ B $	batch size	<code>b</code>	
$\eta$	learning rate	<code>\eta</code>	
$\mathbf{k}$	discretized frequency	<code>\bm{k}</code>	<code>\vk</code>
$\boldsymbol{\xi}$	continuous frequency	<code>\bm{\xi}</code>	<code>\vxi</code>
*	convolution operation	<code>*</code>	

## 12 $L$ -layer neural network

symbol	meaning	L <sup>A</sup> T <sub>E</sub> X	simplified
$d$	input dimension	<code>d</code>	
$d_o$	output dimension	<code>d_{\rm o}</code>	
$m_l$	the number of $l$ th layer neuron, $m_0 = d$ , $m_L = d_o$	<code>m_l</code>	
$\mathbf{W}^{[l]}$	the $l$ th layer weight	<code>\bm{W}^{[1]}</code>	<code>\mW^{[1]}</code>
$\mathbf{b}^{[l]}$	the $l$ th layer bias term	<code>\bm{b}^{[1]}</code>	<code>\vb^{[1]}</code>
$\circ$	entry-wise operation	<code>\circ</code>	
$\sigma : \mathbb{R} \rightarrow \mathbb{R}^+$	activation function	<code>\sigma</code>	
$\theta$	$= (\mathbf{W}^{[0]}, \dots, \mathbf{W}^{[L-1]}, \mathbf{b}^{[0]}, \dots, \mathbf{b}^{[L-1]})$ , parameters	<code>\bm{\theta}</code>	<code>\vtheta</code>
$f_{\theta}^{[0]}(\mathbf{x})$	$= \mathbf{x}$		
$f_{\theta}^{[l]}(\mathbf{x})$	$= \sigma \circ (\mathbf{W}^{[l-1]} f_{\theta}^{[l-1]}(\mathbf{x}) + \mathbf{b}^{[l-1]})$ , $l$ -th layer output		
$f_{\theta}(\mathbf{x})$	$= f_{\theta}^{[L]}(\mathbf{x}) = \mathbf{W}^{[L-1]} f_{\theta}^{[L-1]}(\mathbf{x}) + \mathbf{b}^{[L-1]}$ , $L$ -layer NN		

## 13 Acknowledgements

Chenglong Bao (Tsinghua), Zhengdao Chen (NYU), Bin Dong (Peking), Weinan E (Princeton), Quanquan Gu (UCLA), Kaizhu Huang (XJTU), Shi Jin (SJTU), Jian Li (Tsinghua), Lei Li (SJTU), Tiejun Li (Peking), Zhenguo Li (Huawei), Zheming Li (NUDT), Shaobo Lin (XJTU), Ziqi Liu (CSRC), Zichao Long (Peking), Chao Ma (Princeton), Chao Ma (SJTU), Yuheng Ma (WHU), Dengyu Meng (XJTU), Wang Miao (Peking), Pingbing Ming (CAS), Zuoqiang Shi (Tsinghua), Jihong Wang (CSRC), Liwei Wang (Peking), Bican Xia (Peking), Zhouwang Yang (USTC), Haijun Yu (CAS), Yang Yuan (Tsinghua), Cheng Zhang (Peking), Lulu Zhang (SJTU), Jiwei Zhang (WHU), Pingwen Zhang (Peking), Xiaoqun Zhang (SJTU), Chengchao Zhao (CSRC), Zhanxing Zhu (Peking), Chuan Zhou (CAS), Xiang Zhou (cityU).