

A Proposal on Standard Mathematical Notations of Machine Learning

January 16, 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 over different papers. This note suggests a standard for commonly used mathematical notations of machine learning. This standard will be regularly updated based on the progress of the field. We look forward to more suggestions to improve this note in future versions.

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.

Function

Hypothesis space is denoted by \mathcal{H} . 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 $i = 1, \dots, n$.

Loss function

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:

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

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 $\mathcal{R}_n(\boldsymbol{\theta})$,

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

Without ambiguity, L is also used for L_S .

The population risk or expected loss is denoted by

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

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

Activation function 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. $\sigma(x) = \cos x, \sin x$.

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 σ is the activation function, \mathbf{w}_j is the input weight, a_j is the output weight, b_j is the bias term.

General deep neural network

The counting of the layer number excludes the input layer. The 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 $\mathbf{W}^{[l]} \in \mathbb{R}^{m_{l+1} \times m_l}$, $\mathbf{b}^{[l]} \in \mathbb{R}^{m_{l+1}}$, $m_0 = d_{\text{in}} = d$, $m_L = d_{\text{out}} = d$, σ 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.

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

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

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

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

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 η .

Gram matrix

The Gram matrix is denoted by K_n .

Fourier Frequency

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

Convolution

The convolution operation is denoted by $*$.

More notations

More notations about GAN, RNN, CNN, Resnet etc. are left to be done.

Notation table

symbol	meaning	L ^A T _E 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}	$= \mathcal{X} \times \mathcal{Y}$ examples domain (a set)	<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 with parameters $\boldsymbol{\theta}$	<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}	a distribution over some set (usually \mathcal{Z})	<code>\mathcal{D}</code>	<code>\fD</code>
$S = \{\mathbf{z}_i\}_{i=1}^n$	$= \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ is sampled from \mathcal{D} over \mathcal{Z}		
$L_S(\boldsymbol{\theta}), L_n(\boldsymbol{\theta}), \mathcal{R}(\boldsymbol{\theta})$	$= \sum_{i=1}^n \ell(f_{\boldsymbol{\theta}}, \mathbf{z}_i)$ empirical risk or training loss		
$L_{\mathcal{D}}(\boldsymbol{\theta})$	$= \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} \ell(f_{\boldsymbol{\theta}}, \mathbf{z})$ 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 $f(\mathbf{x}; \boldsymbol{\theta})$	$= \sum_{j=1}^m a_j \sigma(\mathbf{w}_j \cdot \mathbf{x} + b_j)$ two-layer neural network	<code>\f_{\bm{\theta}}</code>	<code>\f_{\vtheta}</code>
$\text{VCdim}(\mathcal{H})$	the VC-dimension of a hypothesis space \mathcal{H}		
$\text{Rad}(\mathcal{H} \circ S)$	the Rademacher complexity of \mathcal{H} on S		
GD	gradient descent		
SGD	stochastic gradient descent		
B	a batch set	<code>B</code>	
b	batch size	<code>b</code>	
η	learning rate	<code>\eta</code>	
K_n	Gram matrix	<code>K_n</code>	
\mathbf{k}	discretized frequency	<code>\bm{k}</code>	<code>\vk</code>
ξ	continuous frequency	<code>\bm{\xi}</code>	<code>\vxi</code>
$*$	convolution operation	<code>*</code>	

L -layer neural network

symbol	meaning	L ^A T _E 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>	
$\boldsymbol{\theta}$	$= (\mathbf{W}^{[0]}, \mathbf{W}^{[1]}, \dots, \mathbf{W}^{[L-1]}, \mathbf{b}^{[0]}, \mathbf{b}^{[1]}, \dots, \mathbf{b}^{[L-1]})$, parameters	<code>\bm{\theta}</code>	<code>\vtheta</code>
$f_{\boldsymbol{\theta}}^{[0]}(\mathbf{x})$	$= \mathbf{x}$		
$f_{\boldsymbol{\theta}}^{[l]}(\mathbf{x})$	$= \sigma \circ (\mathbf{W}^{[l-1]} f_{\boldsymbol{\theta}}^{[l-1]}(\mathbf{x}) + \mathbf{b}^{[l-1]})$, l -th layer output		
$f_{\boldsymbol{\theta}}(\mathbf{x})$	$= f_{\boldsymbol{\theta}}^{[L]}(\mathbf{x}) = \mathbf{W}^{[L-1]} f_{\boldsymbol{\theta}}^{[L-1]}(\mathbf{x}) + \mathbf{b}^{[L-1]}$, L -layer NN		