

GRADIENTS IN A DEEP NEURAL NETWORK AND THEIR PYTHON IMPLEMENTATIONS

YOUNG HO PARK

ABSTRACT. This is an expository article about the gradients in deep neural network. It is hard to find a place where gradients in a deep neural network are dealt in details in a systematic and mathematical way. We review and compute the gradients and Jacobians to derive formulas for gradients which appear in the backpropagation and implement them in vectorized forms in Python.

1. Introduction

This is an expository article about the gradients appearing in neural networks. A *deep neural network* $\mathbf{n}_{\mathbf{W}, \mathbf{b}}$ with weights \mathbf{W} and biases \mathbf{b} consists of an input layer, one or more hidden layers HL_1, \dots, HL_{L-1} and an output layer.

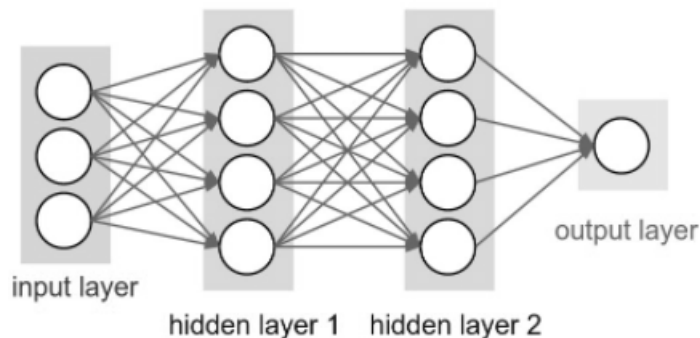


FIGURE 1. Deep neural networks

Mathematically, each hidden layer or a output layer is a function

$$HL_\ell : \mathbb{R}^{n[\ell-1]} \rightarrow \mathbb{R}^{n[\ell]}$$

$$\mathbf{a}[\ell - 1] \mapsto \mathbf{a}[\ell]$$

which is a composition of a linear transformation

$$\mathbf{z}[\ell] = \mathbf{W}[\ell]\mathbf{a}[\ell - 1] + \mathbf{b}[\ell]$$

Received December 5, 2021. Revised February 22, 2022. Accepted March 1, 2022.

2010 Mathematics Subject Classification: 68T07.

Key words and phrases: gradients, deep neural networks, backpropagations, machine learning.

© The Kangwon-Kyungki Mathematical Society, 2022.

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/3.0/>) which permits unrestricted non-commercial use, distribution and reproduction in any medium, provided the original work is properly cited.

and an activation function

$$g[\ell] : \mathbb{R}^{n[\ell-1]} \rightarrow \mathbb{R}^{n[\ell]}, \quad \mathbf{a}[\ell] = g[\ell](\mathbf{z}[\ell]).$$

Here $\mathbf{W}[\ell] \in \text{Mat}_{n[\ell] \times n[\ell-1]}(\mathbb{R})$ is a *weight* matrix and $\mathbf{b}[\ell]$ is a *bias* vector. The most widely used activation function of the hidden layers are relu functions. The structure of the output layer is the same as hidden layers except that the activation function is usually the sigmoid for binary classification or the softmax function for multiclass classification and the identity function for the regression.

Let the training data for the network $\mathbf{n}_{\mathbf{W}, \mathbf{b}}$ be given by

$$\{(\mathcal{X}_k, \mathbf{y}_k) \mid k = 1, \dots, N\},$$

where $\mathcal{X}_k \in \mathbb{R}^n$, and $\mathbf{y}_k \in \mathbb{R}^m$. $\{\mathcal{X}_k\}$ are samples and \mathbf{y}_k is the *target* or *label* of \mathcal{X}_k . Let $\mathbf{o}_k = \mathbf{n}_{\mathbf{W}, \mathbf{b}}(\mathcal{X}_k)$ be the network's output. We will consider the cost function

$$J(\mathbf{W}, \mathbf{b}) = \frac{1}{N} \sum_{k=1}^N \mathcal{L}(\mathbf{o}_k, \mathbf{y}_k)$$

defined by the average of the losses $\mathcal{L}(\mathbf{o}_k, \mathbf{y}_k)$ over the N samples. The loss function \mathcal{L} is usually one of the following:

1. L^2 -loss (regression)

$$\mathcal{L}(\mathbf{o}, \mathbf{y}) = \frac{1}{2} \|\mathbf{o} - \mathbf{y}\|_2^2 = \frac{1}{2} (\mathbf{o} - \mathbf{y}) \cdot (\mathbf{o} - \mathbf{y}),$$

2. binary cross-entropy loss (binary classification, $m = 1$)

$$\mathcal{L}(\mathbf{o}, \mathbf{y}) = -[\mathbf{y} \log \mathbf{o} + (1 - \mathbf{y}) \log(1 - \mathbf{o})],$$

3. categorical cross-entropy loss (multiclass classification)

$$\mathcal{L}(\mathbf{o}, \mathbf{y}) = - \sum_{j=1}^m y_j \log o_j.$$

We want to find a neural network which maps \mathcal{X}_k as close to \mathbf{y}_k as possible and the learning process is to find the parameters \mathbf{W} and \mathbf{b} such that the total cost $J(\mathbf{W}, \mathbf{b})$ is as small as we want. This is basically done by the gradient descent algorithm:

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \frac{\partial J}{\partial \mathbf{W}}, \quad \mathbf{b} \leftarrow \mathbf{b} - \alpha \frac{\partial J}{\partial \mathbf{b}},$$

where α is a constant called the *learning rate*. The exact meaning of these gradients will be given later.

The python packages such as TensorFlow and PyTorch automatically compute the gradients according to the structure of the network, and hence the users of such packages do not have to worry about the actual computation of gradients. See [1] for example. However, we mathematicians care very much about how the formulas are derived. One can find thousands of web pages in the internet, which explain the mathematics of the gradients. However, it is hard to find an article which derive the formulas in a systematic way. Furthermore, a lot of web pages use ugly notations and sometimes the wrong formulas. Even a very famous professor gives partially incorrect formulas for gradients, incorrect in the middle by a constant factor but correct in the final formulas (Compare $d\mathbf{Z}$, $d\mathbf{W}$, $d\mathbf{b}$ in [2] with ours in Python Code 4). This is the motivation of this article presenting the exact mathematical computation of gradients

and their python implementation. We hope that this article will be of some help to those who study mathematics for neural networks.

2. Jacobians

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with $\mathbf{y} = f(\mathbf{x})$. The **Jacobian** of f at a point $\mathbf{x} \in \mathbb{R}^n$ is the $m \times n$ matrix of partial derivatives

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \left[\frac{\partial y_i}{\partial x_j} \right] = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial y_m}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \in M_{m \times n}(\mathbb{R}), \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}.$$

Jacobian tells us the relation between the changes in the input and changes in the output

$$f(\mathbf{x} + \Delta \mathbf{x}) \simeq \mathbf{y} + \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Delta \mathbf{x}.$$

Moreover, if we have another function $g : \mathbb{R}^m \rightarrow \mathbb{R}^k$ with $\mathbf{z} = g(\mathbf{y})$, then $g \circ f : \mathbb{R}^n \rightarrow \mathbb{R}^k$ and

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \frac{\partial \mathbf{z}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}}.$$

For a scalar function $g : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{y} = g(\mathbf{x})$, we define its **gradient** as the column vector

$$\nabla_{\mathbf{x}} g = \begin{bmatrix} \frac{\partial g}{\partial x_1} \\ \vdots \\ \frac{\partial g}{\partial x_n} \end{bmatrix}.$$

Notice that the *gradient is the transpose of Jacobian*:

$$\nabla_{\mathbf{x}} g = \left[\frac{\partial g}{\partial \mathbf{x}} \right]^T.$$

See any calculus text, such as [5], for Jacobians and gradients,

3. Activation functions

We will mainly consider the following activation functions.

- **relu** : The function $\text{relu} : \mathbb{R} \rightarrow \mathbb{R}$ is defined by

$$\text{relu}(z) = \max(0, z).$$

We naturally extend relu to matrices as $\text{relu} : \text{Mat}_{m \times n}(\mathbb{R}) \rightarrow \text{Mat}_{m \times n}(\mathbb{R})$ by

$$\text{relu}([x_{ij}]) = [\text{relu}(x_{ij})].$$

- **sigmoid**: The sigmoid function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is defined by

$$\sigma(z) = \frac{1}{1 + e^{-z}}.$$

The sigmoid also extends to matrices.

- **softmax**: The softmax is a function from $\mathbf{s} : \mathbb{R}^m \rightarrow (0, 1]^m$ defined by

$$\mathbf{s}(\mathbf{z}) = \left[\frac{e^{z_i}}{\sum_{k=1}^m e^{z_k}} \right] \in M_{m \times 1}(\mathbb{R}).$$

It is a generalization of sigmoid function.

The relu function itself is not differentiable at $z = 0$. However we use a little trick to modify the definition as follows:

$$\text{relu}_\epsilon(x) = \begin{cases} -\frac{2}{\epsilon^3}z^4 - \frac{3}{\epsilon^2}z^3 + z, & z \in [-\epsilon, 0], \\ \text{relu}(z), & \text{elsewhere.} \end{cases}$$

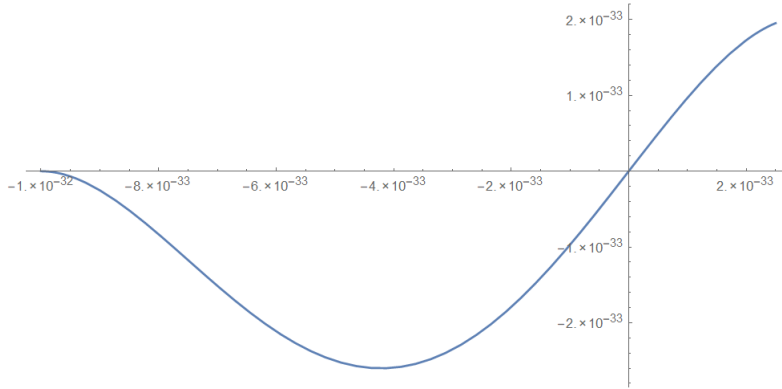


FIGURE 2. Plot of relu_ϵ with $\epsilon = 10^{-32}$

The following is easy to check.

LEMMA 3.1. relu_ϵ is a differentiable function such that

$$\text{relu}'_\epsilon(z) = \begin{cases} 1, & \text{if } z \geq 0 \\ 0, & \text{if } z \leq -\epsilon. \end{cases}$$

We take ϵ to be smaller than the precision of the computer, so that any $\delta < \epsilon$ will be treated as 0 and then relu and relu_ϵ will be the same functions to the computer. Thus we may replace relu with relu_ϵ and simply put

$$(1) \quad \text{relu}'(z) = \begin{cases} 1, & z \geq 0, \\ 0, & z < 0. \end{cases}$$

In a similar manner, we might choose that $\text{relu}'(z)$ is 1 for $z > 0$ and 0 for $z \leq 0$.

On the other hand, the derivative of the sigmoid activation σ is given by

$$(2) \quad \sigma'(z) = \sigma(z)(1 - \sigma(z)).$$

PYTHON CODE 1. Activation functions and their derivatives are defined for numpy arrays. The softmax activation is used only at the output layer and its derivative is not explicitly required but given in the equation (39).

```

import numpy as np

def sigmoid(Z):
    A = 1/(1+np.exp(-Z))
    return A

def d_sigmoid(Z):
    s = sigmoid(Z)
    return s * (1 - s)

def relu(Z):
    A = np.maximum(0,Z)
    return A

def d_relu(Z):
    return (Z>=0)*1.0

def softmax(Z):
    Z_exp = np.exp(Z)
    Z_sum = np.sum(Z_exp, axis=1, keepdims=True)
    A = Z_exp/Z_sum
    return A

```

4. Flow of network functions

In order to simplify notations, we fix a hidden layer or the output layer and omit the index ℓ . Output layer will be treated at a later section in details.

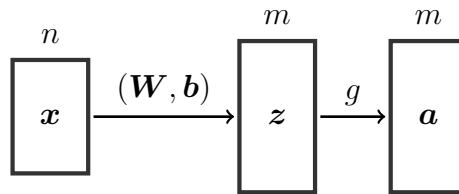


FIGURE 3. Hidden or output layer

Let $\mathbf{x} \in \mathbb{R}^n$ be the input vector to the layer, $\mathbf{W} = [w_{ij}] \in \text{Mat}_{m \times n}(\mathbb{R})$ be a weight matrix, $\mathbf{b} \in \mathbb{R}^m$ be a bias vector and

$$(3) \quad \mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b} \in \mathbb{R}^m$$

be the linear output. Let $g : \mathbb{R} \rightarrow \mathbb{R}$ be the activation function in the layer. We then always extends g to $g : \mathbb{R}^m \rightarrow \mathbb{R}^m$ in an obvious way.

$$(4) \quad \mathbf{a} = g(\mathbf{z}) = \begin{bmatrix} g(z_1) \\ \vdots \\ g(z_m) \end{bmatrix}.$$

The case of the softmax activation function will be treated in later section. Let $\mathbf{a} = g(\mathbf{z})$ be the output of the layer.

Equation (1) can be written componentwise:

$$(5) \quad z_i = \sum_{j=1}^n w_{ij}x_j + b_i.$$

Now we compute various Jacobians. From (5), we get

$$\frac{\partial z_i}{\partial x_j} = w_{ij}, \quad \frac{\partial z_i}{\partial b_j} = \delta_{ij}.$$

From (4), we have that $a_i = g(z_i)$, and hence

$$\frac{\partial a_i}{\partial z_j} = g'(z_i) \frac{\partial z_i}{\partial z_j} = g'(z_i) \delta_{ij}.$$

Let S be any real-valued function of \mathbf{z} . Then $S = S(\mathbf{z}(\mathbf{W}))$ is a function of \mathbf{W} . Notice that $\nabla_{\mathbf{W}} S$ is a $mn \times 1$ matrix. However, it will be inconvenient with this notation to do the gradient descent algorithm

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \nabla_{\mathbf{W}} S.$$

Hence we follow the convention to reshape the gradients:

CONVENTION. The shape of the gradient is the same as the shape of parameters.

Following this convention, we reshape the gradient of S with respect to \mathbf{W} as

$$\nabla_{\mathbf{W}} S = \left[\frac{\partial S}{\partial w_{ij}} \right] \in M_{m \times n}(\mathbb{R}).$$

Note that this convention can also be applied to vectors, without a contradiction. For example, the convention says that

$$\nabla_{\mathbf{x}} g = \begin{bmatrix} \frac{\partial g}{\partial x_1} \\ \vdots \\ \frac{\partial g}{\partial x_n} \end{bmatrix} = \left[\frac{\partial g}{\partial \mathbf{x}} \right]^T,$$

which agrees with the usual definition.

Since S is a function of components z_k 's of \mathbf{z} , we have that

$$(6) \quad \frac{\partial S}{\partial w_{ij}} = \sum_{k=1}^m \frac{\partial S}{\partial z_k} \frac{\partial z_k}{\partial w_{ij}}.$$

Recalling the equation (5)

$$z_k = w_{k1}x_1 + w_{k2}x_2 + \cdots + w_{kn}x_n + b_k,$$

we get $\frac{\partial z_k}{\partial w_{ij}} = \delta_{ik}x_j$. Thus

$$\frac{\partial S}{\partial w_{ij}} = \sum_{k=1}^m \frac{\partial S}{\partial z_k} \delta_{ki}x_j = \frac{\partial S}{\partial z_i} x_j.$$

Therefore,

$$\begin{aligned}\nabla_{\mathbf{w}} S &= \left[\frac{\partial S}{\partial w_{ij}} \right] = \left[\frac{\partial S}{\partial z_i} x_j \right] = \begin{bmatrix} \frac{\partial S}{\partial z_1} x_1 & \frac{\partial S}{\partial z_1} x_2 & \cdots & \frac{\partial S}{\partial z_1} x_n \\ \frac{\partial S}{\partial z_2} x_1 & \frac{\partial S}{\partial z_2} x_2 & \cdots & \frac{\partial S}{\partial z_2} x_n \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial S}{\partial z_m} x_1 & \frac{\partial S}{\partial z_m} x_2 & \cdots & \frac{\partial S}{\partial z_m} x_n \end{bmatrix} \\ &= \left[\frac{\partial S}{\partial \mathbf{z}} \right]^T \mathbf{x}^T = (\nabla_{\mathbf{z}} S) \mathbf{x}^T.\end{aligned}$$

Now we collect various gradients we have computed.

THEOREM 4.1. For $\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$, $\mathbf{a} = g(\mathbf{z})$ and any scalar function $S = S(\mathbf{z})$,

$$(7) \quad \frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \mathbf{W},$$

$$(8) \quad \frac{\partial \mathbf{z}}{\partial \mathbf{b}} = I,$$

$$(9) \quad \frac{\partial \mathbf{a}}{\partial \mathbf{z}} = \text{diag}(g'(\mathbf{z})),$$

$$(10) \quad \nabla_{\mathbf{w}} S = (\nabla_{\mathbf{z}} S) \mathbf{x}^T.$$

Here $\text{diag}(g'(\mathbf{z})) = \text{diag}(g'(z_1), \dots, g'(z_m))$ for $g : \mathbb{R} \rightarrow \mathbb{R}$.

5. Backpropagation through a layer

5.1. Function flow through a network. Consider the function flow of the ℓ -th layer of a network with the input $\mathbf{x}[\ell]$, the weight $\mathbf{W}[\ell]$, the bias $\mathbf{b}[\ell]$ and the activation $g[\ell]$:

$$\mathbf{x}[\ell] = \mathbf{a}[\ell - 1] \in \mathbb{R}^{n[\ell-1]} \quad (\text{input into } \ell\text{-th layer})$$

$$\mathbf{z}[\ell] = \mathbf{W}[\ell]\mathbf{x}[\ell] + \mathbf{b}[\ell] \in \mathbb{R}^{n[\ell]} \quad (\text{linear output})$$

$$\mathbf{a}[\ell] = g[\ell](\mathbf{z}[\ell]) \in \mathbb{R}^{n[\ell]} \quad (\text{activation}).$$

We will simplify the notations a little bit by dropping the indices and writing $\mathbf{a}[\ell]$ as \mathbf{a} and $\mathbf{a}[\ell - 1]$ as $\mathbf{a}_{\triangleleft}$, etc.

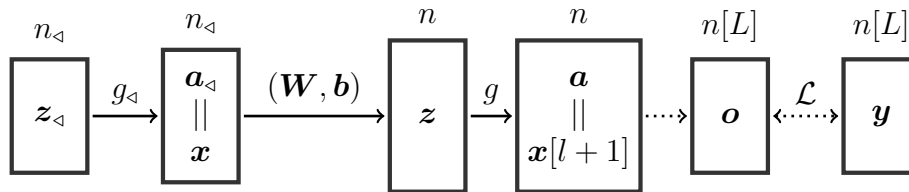


FIGURE 4. $(\ell - 1)$ -th and ℓ -th layers with output layer

Let $\mathcal{L} = \mathcal{L}(\mathbf{x}, \mathbf{W}, \mathbf{b}, \dots)$ be any loss function such as the cross-entropy we previously considered. We assume that we already have computed

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}} = (\nabla_{\mathbf{z}} \mathcal{L})^T \in M_{1 \times m}(\mathbb{R}).$$

We will inductively compute the gradients

$$\nabla_{\mathbf{W}}\mathcal{L}, \quad \nabla_{\mathbf{b}}\mathcal{L}, \quad \nabla_{\mathbf{z}_{\triangleleft}}\mathcal{L}.$$

It is customary to write

$$\boldsymbol{\delta} = \nabla_{\mathbf{z}}\mathcal{L}, \quad \boldsymbol{\delta}_{\triangleleft} = \nabla_{\mathbf{z}_{\triangleleft}}\mathcal{L}$$

excerpted from the name *delta rule* [4]. These are the error signals passed down to \mathbf{z} when we do the backpropagation. We will find the recurrence relation between $\boldsymbol{\delta}$ and $\boldsymbol{\delta}_{\triangleleft}$. Equation (10) gives us

$$(11) \quad \nabla_{\mathbf{W}}\mathcal{L} = (\nabla_{\mathbf{z}}\mathcal{L})\mathbf{x}^T = \boldsymbol{\delta}\mathbf{x}^T = \boldsymbol{\delta}\mathbf{a}_{\triangleleft}^T.$$

The chain rule $\frac{\partial\mathcal{L}}{\partial\mathbf{b}} = \frac{\partial\mathcal{L}}{\partial\mathbf{z}}\frac{\partial\mathbf{z}}{\partial\mathbf{b}}$ together with (8) gives us

$$(12) \quad \nabla_{\mathbf{b}}\mathcal{L} = \boldsymbol{\delta}.$$

From (7) we obtain

$$\frac{\partial\mathcal{L}}{\partial\mathbf{x}} = \frac{\partial\mathcal{L}}{\partial\mathbf{z}}\frac{\partial\mathbf{z}}{\partial\mathbf{x}} = \frac{\partial\mathcal{L}}{\partial\mathbf{z}}\mathbf{W}.$$

Thus

$$(13) \quad \nabla_{\mathbf{x}}\mathcal{L} = \nabla_{\mathbf{a}_{\triangleleft}}\mathcal{L} = \left(\frac{\partial\mathcal{L}}{\partial\mathbf{x}}\right)^T = \mathbf{W}^T\boldsymbol{\delta}.$$

By the chain rule again, we obtain

$$\begin{aligned} \boldsymbol{\delta}_{\triangleleft}^T &= \frac{\partial\mathcal{L}}{\partial\mathbf{z}_{\triangleleft}} = \frac{\partial\mathcal{L}}{\partial\mathbf{a}_{\triangleleft}}\frac{\partial\mathbf{a}_{\triangleleft}}{\partial\mathbf{z}_{\triangleleft}} \\ &= (\boldsymbol{\delta}^T\mathbf{W}) \begin{bmatrix} g'_{\triangleleft}(z_{\triangleleft 1}) & 0 & \cdots & 0 \\ 0 & g'_{\triangleleft}(z_{\triangleleft 2}) & \cdots & 0 \\ \vdots & & & \\ 0 & 0 & \cdots & (g'_{\triangleleft}(z_{\triangleleft n})) \end{bmatrix}, \end{aligned}$$

which implies

$$(14) \quad \boldsymbol{\delta}_{\triangleleft} = \begin{bmatrix} g'_{\triangleleft}(z_{\triangleleft 1}) & 0 & \cdots & 0 \\ 0 & g'_{\triangleleft}(z_{\triangleleft 2}) & \cdots & 0 \\ \vdots & & & \\ 0 & 0 & \cdots & (g'_{\triangleleft}(z_{\triangleleft n})) \end{bmatrix} (\mathbf{W}^T\boldsymbol{\delta}).$$

Now $\mathbf{W}^T\boldsymbol{\delta}$ has shape $(n \times m)(m \times 1) = (n \times 1)$, and thus it is a column vector. By noticing that

$$\text{diag}(d_1, d_2, \dots, d_n) \cdot (c_1, c_2, \dots, c_n)^T = (d_1c_1, d_2c_2, \dots, d_nc_n)^T \in \text{Mat}_{n \times 1}(\mathbb{R})$$

and using the notation $g(\mathbf{z}) = (g(z_1), \dots, g(z_n))^T$ as before, we obtain

$$(15) \quad \boldsymbol{\delta}_{\triangleleft} = g'_{\triangleleft}(\mathbf{z}_{\triangleleft}) * (\mathbf{W}^T\boldsymbol{\delta}) = (\mathbf{W}^T\boldsymbol{\delta}) * g'_{\triangleleft}(\mathbf{z}_{\triangleleft})$$

where $*$ denotes the componentwise multiplication. Now we have proved the following theorem:

THEOREM 5.1. *For any scalar function $\mathcal{L} = \mathcal{L}(\mathbf{x}, \mathbf{W}, \mathbf{b}, \dots)$ with $\boldsymbol{\delta} = \nabla_{\mathbf{z}}\mathcal{L}$ and $\boldsymbol{\delta}_{\triangleleft} = \nabla_{\mathbf{z}_{\triangleleft}}\mathcal{L}$, we have that*

$$\boldsymbol{\delta}_{\triangleleft} = g'_{\triangleleft}(\mathbf{z}_{\triangleleft}) * (\mathbf{W}^T\boldsymbol{\delta}).$$

Moreover,

$$\begin{aligned}\nabla_{\mathbf{a}_{\triangleleft}} \mathcal{L} &= \mathbf{W}^T \boldsymbol{\delta}, \\ \boldsymbol{\delta} &= g'(\mathbf{z}) * \nabla_{\mathbf{a}} \mathcal{L}, \\ \nabla_{\mathbf{W}} \mathcal{L} &= \boldsymbol{\delta} \mathbf{a}_{\triangleleft}^T, \\ \nabla_{\mathbf{b}} \mathcal{L} &= \boldsymbol{\delta}.\end{aligned}$$

If we use the layer numbers $\ell = 1, \dots, L$ (the 0-th layer is the input layer), the identities become

$$(16) \quad \nabla_{\mathbf{a}[\ell-1]} \mathcal{L} = \mathbf{W}[\ell]^T \boldsymbol{\delta}[\ell],$$

$$(17) \quad \boldsymbol{\delta}[\ell-1] = g[\ell-1]'(\mathbf{z}[\ell-1]) * (\mathbf{W}[\ell]^T \boldsymbol{\delta}[\ell]),$$

$$(18) \quad \nabla_{\mathbf{W}[\ell]} \mathcal{L} = \boldsymbol{\delta}[\ell] \mathbf{a}[\ell-1]^T,$$

$$(19) \quad \nabla_{\mathbf{b}[\ell]} \mathcal{L} = \boldsymbol{\delta}[\ell].$$

Here, $\mathbf{a}[0]$ denotes a sample input data in $\{\mathcal{X}_k\}$.

5.2. Vectorization. Now we vectorize all these results for a fixed ℓ -th layer. Vectorization is a style of computer programming where operations are applied to whole arrays instead of individual samples. The numpy package automatically vectorize many numpy operations. Even though these activation functions are defined for real numbers, numpy automatically vectorize these functions and we can apply relu to numpy arrays. Utilizing this property, we can simplify the network computation and get simple formulas for gradients.

As earlier, the training dataset is given by $(\mathcal{X}_k, \mathbf{y}_k)$ for $k = 1, \dots, N$. Each sample \mathcal{X}_k flows through the each layer and gives \mathbf{z}_k 's, \mathbf{a}_k 's for each layer until it produces the network output \mathbf{o}_k at the output layer. The vectorization means to pack up the input dataset into a big matrix

$$\mathcal{X} = [\mathcal{X}_1 \ \mathcal{X}_2 \ \cdots \ \mathcal{X}_N] \in \text{Mat}_{n \times N}(\mathbb{R})$$

whose columns are samples \mathcal{X}_k , $k = 1, \dots, N$ and process it in a unified fashion.

REMARK 5.2. However, we must note that in this current world of data science, it is customary that data samples are packed into rows of the dataset under the influence of engineers, which makes us to take transposed versions of our formulas. Refer to [1] for example and for generalities of the machine learning.

We use the capital letters $\mathbf{Z} = \mathbf{Z}[\ell]$ and $\mathbf{A} = \mathbf{A}[\ell]$ to denote the packages of linear outputs and activation outputs at the ℓ -th layer (dropping the index ℓ temporarily)

$$\mathbf{Z} = [\mathbf{z}_1 \ \mathbf{z}_2 \ \cdots \ \mathbf{z}_N], \quad \mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_N] \in \text{Mat}_{m \times N}(\mathbb{R})$$

where

$$\mathbf{z}_k = \mathbf{W} \mathbf{a}_{\triangleleft, k} + \mathbf{b}, \quad \mathbf{a}_k = g(\mathbf{z}_k)$$

corresponding to the initial sample \mathcal{X}_k . Notice that the parameters $\mathbf{W} = \mathbf{W}[\ell]$ and $\mathbf{b} = \mathbf{b}[\ell]$ are independent on \mathbf{Z}, \mathbf{A} , but only dependent on the layers. Weights and biases are updated after all inputs \mathcal{X}_k are processed (under the assumption that we use the batch gradient descent algorithm). The total cost function is

$$J(\mathbf{W}, \mathbf{b}, \dots) = \frac{1}{N} \sum_{k=1}^N \mathcal{L}_k(\mathbf{W}, \mathbf{b}, \dots),$$

where $\mathcal{L}_k(\mathbf{W}, \mathbf{b}) = \mathcal{L}(\mathbf{o}_k, \mathbf{y}_k, \mathbf{W}, \mathbf{b})$ is the individual loss from an input \mathcal{X}_k . As before, we write

$$\boldsymbol{\delta}_k = \nabla_{\mathbf{z}_k} \mathcal{L} \in \text{Mat}_{n \times 1}(\mathbb{R})$$

for $k = 1, \dots, N$ and

$$\boldsymbol{\Delta} = \nabla_{\mathbf{Z}} \mathcal{L}.$$

More precisely,

$$\boldsymbol{\Delta} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial z_{11}} & \cdots & \frac{\partial \mathcal{L}}{\partial z_{N1}} \\ \vdots & & \vdots \\ \frac{\partial \mathcal{L}}{\partial z_{1n}} & \cdots & \frac{\partial \mathcal{L}}{\partial z_{Nn}} \end{bmatrix} = [\boldsymbol{\delta}_1 \ \boldsymbol{\delta}_2 \ \cdots \ \boldsymbol{\delta}_N].$$

where

$$\mathbf{Z} = \begin{bmatrix} z_{11} & \cdots & z_{N1} \\ \vdots & & \vdots \\ z_{1n} & \cdots & z_{Nn} \end{bmatrix} \in \text{Mat}_{n \times N}(\mathbb{R})$$

($\boldsymbol{\Delta}$ and \mathbf{Z} have the same shape.) Notice that $\boldsymbol{\Delta} \neq \nabla_{\mathbf{Z}} J$. In fact,

$$\nabla_{\mathbf{z}_k} J = \frac{1}{N} \sum_{k'=1}^N \nabla_{\mathbf{z}_k} \mathcal{L}_{k'} = \frac{1}{N} \nabla_{\mathbf{z}_k} \mathcal{L}_k = \frac{1}{N} \boldsymbol{\delta}_k,$$

so that

$$\nabla_{\mathbf{Z}} J = \frac{1}{N} \nabla_{\mathbf{Z}} \mathcal{L} = \frac{1}{N} \boldsymbol{\Delta}.$$

Similarly,

$$\nabla_{\mathbf{A}_{\triangleleft}} J = [\nabla_{\mathbf{a}_{\triangleleft 1}} J \ \cdots \ \nabla_{\mathbf{a}_{\triangleleft N}} J] = \frac{1}{N} [\mathbf{W}^T \boldsymbol{\delta}_1 \ \cdots \ \mathbf{W}^T \boldsymbol{\delta}_N] = \frac{1}{N} \mathbf{W}^T \boldsymbol{\Delta} = \mathbf{W}^T (\nabla_{\mathbf{Z}} J).$$

Now the gradient of \mathbf{W} is computed as follow:

$$\begin{aligned} \nabla_{\mathbf{W}} J &= \frac{1}{N} \sum_{k=1}^N \nabla_{\mathbf{W}} \mathcal{L}_k = \frac{1}{N} \sum_{k=1}^N \boldsymbol{\delta}_k \mathbf{a}_{\triangleleft k}^T \\ &= \frac{1}{N} [\boldsymbol{\delta}_1 \ \boldsymbol{\delta}_2 \ \cdots \ \boldsymbol{\delta}_N] \begin{bmatrix} \mathbf{a}_{\triangleleft 1}^T \\ \mathbf{a}_{\triangleleft 2}^T \\ \vdots \\ \mathbf{a}_{\triangleleft N}^T \end{bmatrix} = \frac{1}{N} \boldsymbol{\Delta} \mathbf{A}_{\triangleleft}^T = (\nabla_{\mathbf{Z}} J) \mathbf{A}_{\triangleleft}^T. \end{aligned}$$

For the gradient $\nabla_{\mathbf{b}} J$, we use the equation (12) to obtain

$$(20) \quad \nabla_{\mathbf{b}} J = \frac{1}{N} \sum_{k=1}^N \nabla_{\mathbf{b}} \mathcal{L}_k = \frac{1}{N} \sum_{k=1}^N \boldsymbol{\delta}_k = \sum_{k=1}^N \nabla_{\mathbf{z}_k} J,$$

which is the sum of the columns of $\nabla_{\mathbf{Z}} J$. We have finally proved the following formulas about the gradients we need for the backpropagations.

REMARK 5.3. For a function $J = J(\mathbf{T})$, let us use the notation

$$d\mathbf{T} = \nabla_{\mathbf{T}} J.$$

This notation is useful when we carry out the python coding.

THEOREM 5.4. For $\ell = 1, \dots, L$, let $\Delta[\ell] = \nabla_{\mathbf{Z}[\ell]} \mathcal{L}$. Then

$$(21) \quad d\mathbf{Z}[\ell] = \frac{1}{N} \Delta[\ell]$$

and the vectorized gradients at the ℓ -layer for the backpropagation are

$$(22) \quad d\mathbf{A}[\ell - 1] = \mathbf{W}[\ell]^T d\mathbf{Z}[\ell]$$

$$(23) \quad d\mathbf{Z}[\ell - 1] = g'(\mathbf{Z}[\ell - 1]) * d\mathbf{A}[\ell - 1]$$

$$(24) \quad d\mathbf{W}[\ell] = d\mathbf{Z}[\ell] \mathbf{A}[\ell - 1]^T$$

$$(25) \quad d\mathbf{b}[\ell] = \text{sum of columns of } d\mathbf{Z}[\ell]$$

PYTHON CODE 2. (propagation at a relu layer)

We assume that the activation function is given by relu at the given i -th layer ($1 \leq i < L$). The letters $\mathbf{Z}, \mathbf{A}, \mathbf{W}, \mathbf{b}$ are written as Z, A, W, b at the code.

```
Z[i] = W[i] @ A[i-1] + b[i]
A[i] = relu(Z[i])
```

```
dA[i] = W[i+1].T @ dZ[i+1]
dZ[i] = d_relu(Z[i]) * dA[i]
dW[i] = dZ[i] @ A[i-1].T
db[i] = np.sum(dZ[i], axis=1, keepdims=True)
```

6. At the output layer

We consider the function flow of the output layer (L -th layer), where $\mathbf{x} = \mathbf{a}_{L-1}$:

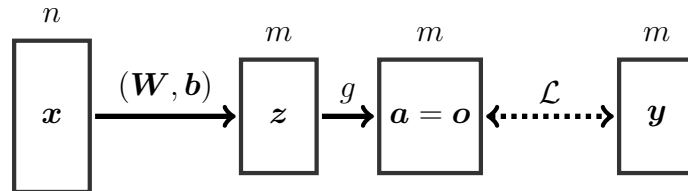


FIGURE 5. Last(L -th) layer

6.1. Linear logistic regression. Suppose we are given a set of data $\{(\mathcal{X}_k, \mathbf{y}_k) \mid k = 1, \dots, N\}$. We want to find a ‘good’ linear function (called hypothesis)

$$h(\mathcal{X}) = \mathbf{W}\mathcal{X} + \mathbf{b}$$

such that $h(\mathcal{X}_k) \approx \mathbf{y}_k$ for all k , where $\mathbf{W} \in \text{Mat}_{1 \times n}(\mathbb{R})$ and $\mathbf{b} \in \mathbb{R}$. Actually, the good function means that the total error

$$J = \frac{1}{N} \sum_{k=1}^N \mathcal{L}(h(\mathcal{X}_k), \mathbf{y}_k)$$

is as small as possible.

$$(26) \quad \mathcal{L}(h(\mathcal{X}), \mathbf{y}) = \frac{1}{2} \|h(\mathcal{X}) - \mathbf{y}\|^2$$

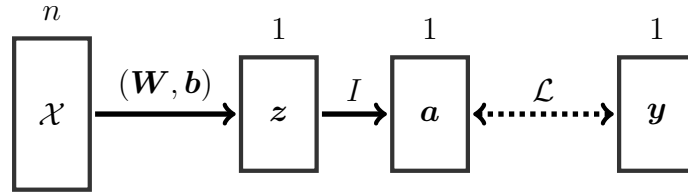


FIGURE 6. Neural network for a logistic regression

is the L^2 -loss (or mean squared error) for the data sample $(\mathcal{X}, \mathbf{y})$. In this way, we can use $h(\mathcal{X})$ to predict the value of \mathbf{y} at a new input \mathcal{X} .

In this case, the network has no hidden layer, and the activation function is the identity so that $\mathbf{z} = \mathbf{a} = \mathbf{o} = h(\mathcal{X})$. Thus

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}} = (\mathbf{a} - \mathbf{y})^T$$

from (26) and hence

$$(27) \quad \boldsymbol{\delta} = \nabla_{\mathbf{z}} \mathcal{L} = \mathbf{a} - \mathbf{y}$$

$$(28) \quad \nabla_{\mathbf{W}} \mathcal{L} = \boldsymbol{\delta} \mathbf{x}^T$$

$$(29) \quad \nabla_{\mathbf{b}} \mathcal{L} = \boldsymbol{\delta}$$

In vectorized notations, we have that

$$(30) \quad d\mathbf{Z} = \frac{1}{N}(\mathbf{A} - \mathbf{Y})$$

$$(31) \quad d\mathbf{W} = d\mathbf{Z} \mathbf{X}^T$$

$$(32) \quad d\mathbf{b} = \text{sum of columns of } d\mathbf{Z}$$

where $\mathbf{Y} = [\mathbf{y}_1 \cdots \mathbf{y}_N]$.

PYTHON CODE 3. The training data is given by (X, y) with $X \in \text{Mat}_{n \times N}(\mathbb{R})$, $y \in \text{Mat}_{1 \times N}(\mathbb{R})$.

```
# Linear Regression
```

```
# After reading in the data (X,y):
```

```
# X = matrix of shape (n, N) whose columns are samples
```

```
# y = array of shape (1,N)
```

```
# initialization
```

```
n, N = X.shape
```

```
W = np.random.rand(1,n)*0.01
```

```
b = 0
```

```
def forward_propagation(X, W, b):
```

```
    Z = W @ X + b
```

```
    A = Z # not necessary in practice
```

```
    return A
```

```

def back_propagation(X, y, W, b)
    A = forward_propagation(X, W, b)
    N = X.shape[1]
    dZ = 1/N * (A - y)
    dW = dZ @ X.T
    db = np.sum(dZ, axis=1, keepdims=True)
    return dW, db

def update_params(W, b, dW, db, learning_rate=0.1):
    W -= learning_rate * dW
    b -= learning_rate * db
    return W, b

def cost(A,y):
    return np.sum(((A-y)**2)) / (2*N)

# Gradient Descent Algorithm
# hyperparameters
n_iter = 300
learning_rate = 0.01

cost_history = []
for i in range(n_iter):
    A = forward_propagation(X, W, b)
    dW, db = backward_propagation(X, y, W, b)
    W, b = update_params(W, b, dW, db, learning_rate)
    cost_history.append(cost(A,y))

print(f'Final cost : {cost_history[-1]}')
print(f'W = {W}')
print(f'b = {b}')

```

6.2. Binary classification. In this case, notice that $z, a, y \in \mathbb{R}$.

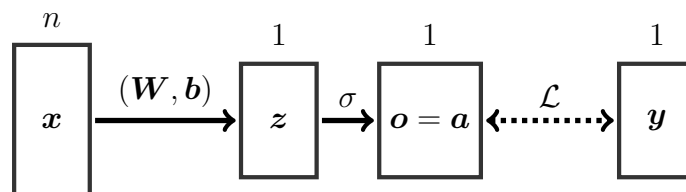


FIGURE 7. Last layer with sigmoid activation

We have the sigmoid activation function $g = \sigma$ and the binary cross-entropy loss

$$\mathcal{L}(a, y) = -\left[y \log a + (1 - y) \log(1 - a) \right].$$

From $\mathbf{a} = \sigma(\mathbf{z})$, we know that $\mathbf{a}'(\mathbf{z}) = \mathbf{a}(1 - \mathbf{a})$. Thus

$$\frac{d\mathcal{L}}{d\mathbf{z}} = -\mathbf{y} \frac{\mathbf{a}(1 - \mathbf{a})}{\mathbf{a}} - (1 - \mathbf{y}) \frac{-\mathbf{a}(1 - \mathbf{a})}{1 - \mathbf{a}} = \mathbf{a} - \mathbf{y}.$$

In a vectorized notation, we therefore have that

$$(33) \quad \boldsymbol{\delta} = \nabla_{\mathbf{z}} \mathcal{L} = \mathbf{a} - \mathbf{y},$$

$$(34) \quad \nabla_{\mathbf{W}} \mathcal{L} = \Delta \mathbf{x}^T,$$

$$(35) \quad \nabla_b \mathcal{L} = \Delta.$$

Notice that they are exactly the same as the linear logistic regression case. The vectorized gradients at this last L -th layer are

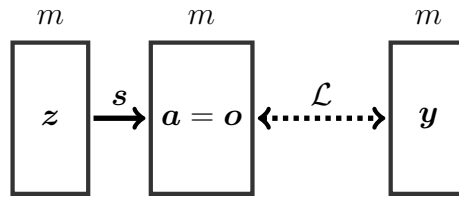
$$(36) \quad d\mathbf{Z}[L] = \frac{1}{N}(\mathbf{A}[L] - \mathbf{Y}),$$

$$(37) \quad d\mathbf{W}[L] = d\mathbf{Z}[L] \mathbf{A}[L-1]^T,$$

$$(38) \quad d\mathbf{b}[L] = \text{sum of columns of } d\mathbf{Z}[L].$$

6.3. Softmax output. Suppose that the softmax activation function $\mathbf{o} = \mathbf{s}(\mathbf{z})$ is used at the last L -th layer. It is defined to be the function from $\mathbb{R}^m \rightarrow (0, 1]^m$ for classifying the inputs into m objects:

$$\mathbf{o} = \mathbf{s}(\mathbf{z}) = \left[\frac{e^{z_i}}{\sum_{k=1}^m e^{z_k}} \right] \in M_{m \times 1}(\mathbb{R}).$$



In this case, we assume that the labels $\{\mathbf{y}_k\}$ are given by the *one-hot vectors*. For example, the class 0 is represented by $\mathbf{y}_k = (1, 0, \dots, 0)^T$.

We need to compute $\frac{\partial \mathbf{o}}{\partial \mathbf{z}} = \left[\frac{\partial o_i}{\partial z_j} \right]$. Since

$$o_i = \frac{e^{z_i}}{e^{z_1} + \dots + e^{z_i} + \dots + e^{z_m}}$$

we have that

$$\frac{\partial o_i}{\partial z_j} = \frac{\delta_{ij} e^{z_i} (\sum e^{z_k}) - e^{z_i} e^{z_j}}{(\sum e^{z_k})^2} = \begin{cases} o_i(1 - o_j), & i = j \\ -o_i o_j, & i \neq j \end{cases} = o_i(\delta_{ij} - o_j).$$

That is,

$$(39) \quad \frac{\partial \mathbf{o}}{\partial \mathbf{z}} = [o_i(\delta_{ij} - o_j)] = \begin{bmatrix} o_1(1 - o_1) & -o_1 o_2 & \cdots & -o_1 o_m \\ -o_2 o_1 & o_2(1 - o_2) & \cdots & -o_2 o_m \\ \vdots & \vdots & \ddots & \vdots \\ -o_m o_1 & -o_m o_2 & \cdots & o_m(1 - o_m) \end{bmatrix}.$$

The categorical cross-entropy function for the softmax activation is defined to be

$$\mathcal{L}(\mathbf{o}, \mathbf{y}) = - \sum_{i=1}^m y_i \log(o_i), \quad \mathbf{o} = \mathbf{s}(\mathbf{z})$$

Recalling that $\sum_i o_i = 1$, we get

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial z_j} &= - \sum_i y_i \frac{\partial \log(o_i)}{\partial z_j} = - \sum_i o_i \frac{1}{o_i} \frac{\partial o_i}{\partial z_j} \\ &= - \sum_i o_i \frac{1}{o_i} \cdot o_i (\delta_{ij} - o_j) = - \sum_i o_i (\delta_{ij} - o_j) \\ &= \sum_{i \neq j} y_i o_j + y_j (-1 + o_j) = \left(\sum_i y_i \right) o_j - y_j = o_j - y_j\end{aligned}$$

Hence the final answer is again

$$(40) \quad \boldsymbol{\delta} = \mathbf{o} - \mathbf{y}.$$

THEOREM 6.1. *For each of the following output layer of*

1. *the linear logistic regression (identity activation, L^2 -loss),*
2. *the binary classification (sigmoid activation, cross-entropy loss)*
3. *the multi-class classification (softmax activation, cross-entropy loss),*

we have the same gradient

$$\boldsymbol{\delta} = \nabla_{\mathbf{z}} \mathcal{L} = \mathbf{o} - \mathbf{y}.$$

at the last layer for single data input.

We consider the deep neural network for the classification into m -classes ($m \geq 3$). The training data is given by (X, Y) with $X \in \text{Mat}_{n \times N}(\mathbb{R})$, $Y \in \text{Mat}_{m \times N}(\mathbb{R})$. We assume that the activation at the hidden layers are given by relu function and activation at the output layer is given by the softmax. We only give the forward and backward propagations using Numpy.

PYTHON CODE 4.

```
# multiclass classification (softmax output)

# After reading in the data (X,y):
# X = numpy array of shape (n, N) whose columns are samples.
# Y = numpy array of shape (m,N) whose columns are one-hot vectors.
# W, b are list of weights and biases at the layers.
# dW, db are list of gradients of the cost with respect to weights
#   and biases.
# Z, A are list of linear outputs and activation outputs at the layers.
# dZ, dA are lists of gradients of the cost with respect to Z and A
# Set A[0] = X

def forward_propagation(X, W, b):
    for i in range(1, L):
        Z[i] = W[i] @ A[i-1] + b[i]
        A[i] = relu(Z[i])

    Z[L] = W[L] @ A[L-1] + b[L]
    A[L] = softmax(Z[L])
    return Z, A

def back_propagation(X, Y, W, b)
    Z, A = forward_propagation(X, W, b)
```

```

N = X.shape[1]

dZ[L] = 1/N * (A[L] - Y)
dW[L] = dZ[L] @ A[L-1].T
db[L] = np.sum(dZ[L], axis=1, keepdims=True)

for i in range(L-1,0,-1):
    dA[i] = Ws[i+1].T @ dZ[i+1]
    dZ[i] = d_relu(Z[i]) * dA[i]
    dW[i] = dZ[i] @ A[i-1].T
    db[i] = np.sum(dZ[i], axis=1, keepdims=True)
return dW, db

```

Finally, we carry out an experiment to calculate the gradients using our formulas given in Theorem 5.4 and check if they agree with the gradients computed by autograd in PyTorch. The result can be found in [3].

References

- [1] A. Géron, Hands-on Machine Learning with Scikit-Learn & TensorFlow (핸즈온 머신러닝), Hanbit Media, 2018
- [2] Andrew Ng, Note on neural network and deep learning, <https://github.com/ashishpatel26/Andrew-NG-Notes>, accessed February 23, 2022
- [3] Y.H. Park, Verifying gradient formulas by PyTorch, https://deepmath.kangwon.ac.kr/~yhpark/verify_gradients.pdf, accessed February 23, 2022
- [4] T. Rashid, Make your own neural network (신경망 첫걸음), Hanbit Media, 2017
- [5] J. Stewart, Calculus, Books-Hill, 2021

Young Ho Park

Department of Mathematics, Kangwon National University,
Chuncheon 24341, Korea.

E-mail: yhpark@kangwon.ac.kr