

Nerva Libraries Implementation

Wieger Wesselink

September 12, 2024

Contents

1	Introduction	1
2	Mathematical background	1
3	Matrix operations	3
4	Layer equations	4
4.1	Derivations	8
5	Activation functions	11
5.1	Softmax functions	11
5.2	Derivations	12
6	Loss functions	12
6.1	Derivations	14
7	Weight initialization	14
8	Optimization	15
9	Learning Rate Schedulers	15

1 Introduction

This document contains algorithm specifications and implementation details of the Nerva Libraries <https://github.com/wiegerw/nerva>. The goal of this document is to provide precise mathematical specifications of key parts of the implementation. The code in the repository <https://github.com/wiegerw/nerva-rowwise> closely matches the equations in this document.

2 Mathematical background

In this section we provide notational conventions, we provide definitions for the gradient and the Jacobian, and a comprehensive table with matrix operations that are necessary for the execution of multilayer perceptrons.

In the rest of this document we use the following notation. Vectors are denoted using lowercase symbols x, y, z , and matrices using uppercase symbols X, Y, Z . The columns of a matrix $X \in \mathbb{R}^{m \times n}$ are denoted as x^1, \dots, x^n , while the rows are denoted as x_1, \dots, x_m . Occasionally, the subscript notation x_i is also used to denote the i^{th} element of the vector x , but this will be clear from context. To distinguish between row and

column vectors, we denote their domains as $\mathbb{R}^{1 \times n}$ and $\mathbb{R}^{m \times 1}$, respectively. We use the dot symbol \cdot for matrix multiplication, and not for the dot product. This notation is used primarily to enhance the readability of expressions.

In the context of neural networks, we typically use x, X for inputs, y, Y for outputs, z, Z for intermediate values, and t, T for targets. The number of inputs of a layer is denoted by D and the number of outputs by K . The number of examples in a mini-batch is denoted by N . We follow the convention of the major neural network frameworks where data is stored using a row layout. This means that by default x, y , and z are considered row vectors, and each row of an input matrix X represents an example of a dataset. For example, if $X \in \mathbb{R}^{N \times D}$, then $x_i \in \mathbb{R}^{1 \times D}$ represents the i -th example.

Definition 1 (Gradient). *Let $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ be a function with input X that has elements x_{ij} , with $m, n \in \mathbb{N}^+$. Then the gradient $\nabla_X f$ is defined as*

$$\nabla_X f(X) = \begin{bmatrix} \frac{\partial f(X)}{\partial x_{11}} & \cdots & \frac{\partial f(X)}{\partial x_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f(X)}{\partial x_{m1}} & \cdots & \frac{\partial f(X)}{\partial x_{mn}} \end{bmatrix}. \quad (1)$$

Definition 2 (Gradient of the loss function). *We consider neural networks with output $Y \in \mathbb{R}^{N \times K}$ and target $T \in \mathbb{R}^{N \times K}$. We assume that there is a fixed loss function $\mathcal{L} : \mathbb{R}^{N \times K} \times \mathbb{R}^{N \times K} \rightarrow \mathbb{R}$, such that $\mathcal{L}(Y, T)$ is the loss corresponding to the output Y , given the target T . We use the following shorthand notation for the gradient of the loss:*

$$DY = \nabla_Y \mathcal{L}(Y, T). \quad (2)$$

If Y depends on a parameter Z , i.e. $Y = Y(Z)$, then we define

$$DZ = \nabla_Z \mathcal{L}(Y(Z), T). \quad (3)$$

Definition 3 (Jacobian). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be a function with input $x \in \mathbb{R}^n$. Then the Jacobian $\frac{\partial f}{\partial x}$ is defined as¹*

$$\frac{\partial f}{\partial x}(x) = \begin{bmatrix} \frac{\partial f_1(x)}{\partial x_1} & \cdots & \frac{\partial f_1(x)}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m(x)}{\partial x_1} & \cdots & \frac{\partial f_m(x)}{\partial x_n} \end{bmatrix}. \quad (4)$$

The execution of MLPs depends on a small number of matrix operations. In Table 1 we provide a mathematical notation, a code representation, and a definition for the most important operations. These operations are used in the implementation of activation functions, loss functions, and the feedforward and backpropagation equations of neural network layers. Basic operations like matrix assignment, matrix indexing, and matrix dimensions are omitted, to keep the resulting code familiar to users of the respective Python frameworks. There is some redundancy in the table, to allow for more efficient, or numerically stable implementations. We refrain from using broadcasting notation, as commonly found in frameworks like NumPy, where arrays of different dimensions can be added, and use the standard notation of matrix calculus instead. This approach with explicit dimensions is needed to validate the correctness using SymPy. For example, to calculate the sum of elements in a column vector $x \in \mathbb{R}^{n \times 1}$, we express it as the dot product $1_n^\top \cdot x$, with $1_n = (1, \dots, 1)^\top$ a column vector of ones.

¹We use the convention that the Jacobian does not depend on whether x and $f(x)$ are row vectors or column vectors. This is consistent with SymPy, but other conventions are also in use.

3 Matrix operations

Table 1: An overview of matrix operations that are needed for the implementation of the class of MLPs described in this paper. We assume that $X, Y \in \mathbb{R}^{m \times n}$ and $Z \in \mathbb{R}^{n \times k}$, where $k, m, n \in \mathbb{N}^+$. Wherever possible, we use m to denote the number of rows and n to denote the number of columns of a matrix or vector. The function σ is the sigmoid function as defined in section 5.

OPERATION	CODE	DEFINITION
0_m	<code>zeros(m)</code>	$m \times 1$ column vector with elements equal to 0
0_{mn}	<code>zeros(m, n)</code>	$m \times n$ matrix with elements equal to 0
1_m	<code>ones(m)</code>	$m \times 1$ column vector with elements equal to 1
1_{mn}	<code>ones(m, n)</code>	$m \times n$ matrix with elements equal to 1
\mathbb{I}_n	<code>identity(n)</code>	$n \times n$ identity matrix
X^\top	<code>X.T</code>	transposition
cX	<code>c * X</code>	scalar multiplication, $c \in \mathbb{R}$
$X + Y$	<code>X + Y</code>	addition
$X - Y$	<code>X - Y</code>	subtraction
$X \cdot Z$	<code>X @ Z</code> or <code>X * Z</code>	matrix multiplication, also denoted as XZ
$x^\top y$ or xy^\top	<code>dot(x, y)</code>	dot product, $x, y \in \mathbb{R}^{m \times 1}$ or $x, y \in \mathbb{R}^{1 \times n}$
$X \odot Y$	<code>hadamard(X, Y)</code>	element-wise product of X and Y
$\text{diag}(X)$	<code>diag(X)</code>	column vector that contains the diagonal of X
$\text{Diag}(x)$	<code>Diag(x)</code>	diagonal matrix with x as diagonal, $x \in \mathbb{R}^{1 \times n}$ or $x \in \mathbb{R}^{m \times 1}$
$1_m^\top \cdot X \cdot 1_n$	<code>elements_sum(X)</code>	sum of the elements of X
$x \cdot 1_n^\top$	<code>column_repeat(x, n)</code>	n copies of column vector $x \in \mathbb{R}^{m \times 1}$
$1_m \cdot x$	<code>row_repeat(x, m)</code>	m copies of row vector $x \in \mathbb{R}^{1 \times n}$
$1_m^\top \cdot X$	<code>columns_sum(X)</code>	$1 \times n$ row vector with sums of the columns of X
$X \cdot 1_n$	<code>rows_sum(X)</code>	$m \times 1$ column vector with sums of the rows of X
$\max(X)_{\text{col}}$	<code>columns_max(X)</code>	$1 \times n$ row vector with maximum values of the columns of X
$\max(X)_{\text{row}}$	<code>rows_max(X)</code>	$m \times 1$ column vector with maximum values of the rows of X
$(1_m^\top \cdot X)/n$	<code>columns_mean(X)</code>	$1 \times n$ row vector with mean values of the columns of X
$(X \cdot 1_n)/m$	<code>rows_mean(X)</code>	$m \times 1$ column vector with mean values of the rows of X
$f(X)$	<code>apply(f, X)</code>	element-wise application of $f : \mathbb{R} \rightarrow \mathbb{R}$ to X
e^X	<code>exp(X)</code>	element-wise application of $f : x \rightarrow e^x$ to X
$\log(X)$	<code>log(X)</code>	element-wise application of the natural logarithm $f : x \rightarrow \ln(x)$ to X
$1/X$	<code>reciprocal(X)</code>	element-wise application of $f : x \rightarrow 1/x$ to X
\sqrt{X}	<code>sqrt(X)</code>	element-wise application of $f : x \rightarrow \sqrt{x}$ to X
$X^{-1/2}$	<code>inv_sqrt(X)</code>	element-wise application of $f : x \rightarrow x^{-1/2}$ to X
$\log(\sigma(X))$	<code>log_sigmoid(X)</code>	element-wise application of $f : x \rightarrow \log(\sigma(x))$ to X ,

In Table 2 we provide implementations in several frameworks of the fundamental matrix-form operations for multilayer perceptrons.

Operation	NumPy + JAX	PyTorch	TensorFlow	Eigen
zeros(m,n)	zeros((m,n))	zeros(m,n)	zeros([m,n])	Zero(m,n)
ones(m,n)	ones((m,n))	ones(m,n)	ones([m,n])	Ones(m,n)
identity(n)	eye(n)	eye(n)	eye(n)	Identity(n,n)
X.T	X.T	X.T	X.T	X.transpose()
c * X	c * X	c * X	c * X	c * X
X + Y	X + Y	X + Y	X + Y	X + Y
X - Y	X - Y	X - Y	X - Y	X - Y
X * Z	X @ Z	X @ Z	X @ Z	X * Z
hadamard(X,Y)	X * Y	X * Y	multiply(X,Y)	X.array() * Y.array()
diag(X)	diag(X)	diag(X)	diag_part(X)	X.diagonal()
Diag(x)	diag(x)	diag(x.flatten())	diag(reshape(x,[-1]))	x.asDiagonal()
elements_sum(X)	sum(X)	sum(X)	reduce_sum(X)	X.sum()
column_repeat(x,n)	tile(x,(1,n))	x.repeat(1,n)	tile(x,[1,n])	x.replicate(1,n)
row_repeat(x,m)	tile(x,(m,1))	x.repeat(m,1)	tile(x,[m,1])	x.replicate(m,1)
columns_sum(X)	sum(X,axis=0)	sum(X,dim=0)	reduce_sum(X,axis=0)	X.colwise().sum()
rows_sum(X)	sum(X,axis=1)	sum(X,dim=1)	reduce_sum(X,axis=1)	X.rowwise().sum()
columns_max(X)	max(X,axis=0)	max(X,dim=0).values	reduce_max(X,axis=0)	X.colwise().maxCoeff()
rows_max(X)	max(X,axis=1)	max(X,dim=1).values	reduce_max(X,axis=1)	X.rowwise().maxCoeff()
columns_mean(X)	mean(X,axis=0)	mean(X,dim=0)	reduce_mean(X,axis=0)	X.colwise().mean()
rows_mean(X)	mean(X,axis=1)	mean(X,dim=1)	reduce_mean(X,axis=1)	X.rowwise().mean()
apply(f,X)	f(X)	f(X)	f(X)	X.unaryExpr(f)
exp(X)	exp(X)	exp(X)	exp(X)	X.array().exp()
log(X)	log(X)	log(X)	log(X)	X.array().log()
reciprocal(X)	1 / X	1 / X	inverse(X)	X.array().inverse()
sqrt(X)	sqrt(X)	sqrt(X)	sqrt(X)	X.array().sqrt()
inv_sqrt(X)	reciprocal(sqrt(X+e))	reciprocal(sqrt(X+e))	reciprocal(sqrt(X+e))	reciprocal(sqrt(X.array()+e))
log_sigmoid(X)	-logaddexp(0,-X)	-softplus(-X)	-softplus(-X)	-log1p(exp(-X.array()))

Table 2: Implementation of matrix operations in NumPy, JAX, PyTorch, TensorFlow and Eigen. We assume that **e** is a given small positive constant that is used to avoid division by zero. Note that some of the operations are located in Python submodules.

4 Layer equations

A major contribution of this paper is the following overview of the feedforward and backpropagation equations of layers in **matrix-form**. For each of the equations, the corresponding Python implementation is given. For several equations a derivation is included, while the correctness of the equations and derivations has been validated using SymPy. All layers are implemented in our Nerva Python packages, and a validation of all equations and derivations can be found in the **tests** directory of the **nerva-sympy** package. We consider input batches X in row layout, i.e. each row represents a single example. We use the following notations:

- $X \in \mathbb{R}^{N \times D}$ is the input batch, where N is the number of examples and D is the input dimension.
- $Y \in \mathbb{R}^{N \times K}$ is the output batch, where K is the output dimension.
- $W \in \mathbb{R}^{K \times D}$ is the weight matrix, which maps the input features to the output features.
- $b \in \mathbb{R}^{1 \times K}$ is the bias vector.
- $Z \in \mathbb{R}^{N \times K}$ is a matrix with intermediate values.
- $\beta, \gamma, \Sigma \in \mathbb{R}^{1 \times K}$ are the parameters of batch normalization.
- $a_l, t_l, a_r, t_r \in \mathbb{R}$ are the coefficients of an SReLU activation function.
- $R \in \mathbb{R}^{K \times D}$ is a dropout matrix.

Each of these matrices has a gradient with the same dimensions, denoted using the same symbol preceded by D, e.g. DX is the gradient corresponding to X . The implementation uses the same names. The input X , and layer parameters like W , b and their gradients DX , DW , Db are stored in the attributes of a layer. The only exception is the output Y and its gradient DY , which are stored in the X and DX attributes of the next layer.

linear layer

FEEDFORWARD EQUATIONS

$$Y = XW^\top + 1_N \cdot b$$

$$Y = X * W.T + \text{row_repeat}(b, N)$$

BACKPROPAGATION EQUATIONS

$$DW = DY^\top \cdot X$$

$$Db = 1_N^\top \cdot DY$$

$$DX = DY \cdot W$$

$$DW = DY.T * X$$

$$Db = \text{columns_sum}(DY)$$

$$DX = DY * W$$

activation layer

Let $\text{act} : \mathbb{R} \rightarrow \mathbb{R}$ be an activation function, for example relu .

FEEDFORWARD EQUATIONS

$$Z = XW^\top + 1_N \cdot b$$

$$Y = \text{act}(Z)$$

$$Z = X * W.T + \text{row_repeat}(b, N)$$

$$Y = \text{act}(Z)$$

BACKPROPAGATION EQUATIONS

$$DZ = DY \odot \text{act}'(Z)$$

$$DW = DZ^\top \cdot X$$

$$Db = 1_N^\top \cdot DZ$$

$$DX = DZ \cdot W$$

$$DZ = \text{hadamard}(DY, \text{act.gradient}(Z))$$

$$DW = DZ.T * X$$

$$Db = \text{columns_sum}(DZ)$$

$$DX = DZ * W$$

srelu-layer

FEEDFORWARD EQUATIONS

$$Z = XW^\top + 1_N \cdot b$$

$$Y = \text{srelu}(Z)$$

$$Z = X * W.T + \text{repeat_row}(b, N)$$

$$Y = \text{apply}(\text{act}, Z)$$

BACKPROPAGATION EQUATIONS

$$DZ = DY \odot \text{srelu}'(Z)$$

$$DW = DZ^\top \cdot X$$

$$Db = 1_N^\top \cdot DZ$$

$$DX = DZ \cdot W$$

$$Da_l = 1_N^\top \cdot (DY \odot A^l) \cdot 1_K$$

$$Da_r = 1_N^\top \cdot (DY \odot A^r) \cdot 1_K$$

$$Dt_l = 1_N^\top \cdot (DY \odot T^l) \cdot 1_K$$

$$Dt_r = 1_N^\top \cdot (DY \odot T^r) \cdot 1_K$$

$$DZ = \text{hadamard}(DY, \text{apply}(\text{act_prime}, Z))$$

$$DW = DZ.T * X$$

$$Db = \text{sum_columns}(DZ)$$

$$DX = DZ * W$$

$$A_l = \text{apply}(a_l, Z)$$

$$A_r = \text{apply}(a_r, Z)$$

$$T_l = \text{apply}(t_l, Z)$$

$$T_r = \text{apply}(t_r, Z)$$

$$Dal = \text{sum_elements}(\text{hadamard}(DY, A_l))$$

$$Dar = \text{sum_elements}(\text{hadamard}(DY, A_r))$$

$$Dtl = \text{sum_elements}(\text{hadamard}(DY, T_l))$$

$$Dtr = \text{sum_elements}(\text{hadamard}(DY, T_r))$$

where

$$A_{ij}^l = \begin{cases} Z_{ij} - t_l & \text{if } Z_{ij} \leq t_l \\ 0 & \text{otherwise} \end{cases}$$

$$A_{ij}^r = \begin{cases} 0 & \text{if } Z_{ij} \leq t_l \vee Z_{ij} < t_r \\ Z_{ij} - t_r & \text{otherwise} \end{cases}$$

$$T_{ij}^l = \begin{cases} 1 - a_l & \text{if } Z_{ij} \leq t_l \\ 0 & \text{otherwise} \end{cases}$$

$$T_{ij}^r = \begin{cases} 1 - a_r & \text{if } Z_{ij} \geq t_r \\ 0 & \text{otherwise} \end{cases}$$

softmax layer

FEEDFORWARD EQUATIONS

$$Z = XW^\top + 1_N \cdot b$$

$$Y = \text{softmax}(Z)$$

$$Z = X * W.T + \text{row_repeat}(b, N)$$

$$Y = \text{softmax}(Z)$$

BACKPROPAGATION EQUATIONS

$$DZ = Y \odot (DY - \text{diag}(DY \cdot Y^\top) \cdot 1_K^\top)$$

$$DW = DZ^\top \cdot X$$

$$Db = 1_N^\top \cdot DZ$$

$$DX = DZ \cdot W$$

$$DZ = \text{hadamard}(Y, \text{DY} - \text{column_repeat}(\text{diag}(DY * Y.T), K))$$

$$DW = DZ.T * X$$

$$Db = \text{columns_sum}(DZ)$$

$$DX = DZ * W$$

log-softmax layer

FEEDFORWARD EQUATIONS

$$Z = XW^T + 1_N \cdot b$$

$$Y = \text{logsoftmax}(Z)$$

$$\begin{aligned} Z &= X * W.T + \text{row_repeat}(b, N) \\ Y &= \text{log_softmax}(Z) \end{aligned}$$

batch normalization layer

FEEDFORWARD EQUATIONS

$$R = X - \frac{1_N \cdot 1_N^T}{N} \cdot X$$

$$\Sigma = \frac{1}{N} \cdot \text{diag}(R^T R)^T$$

$$Z = (1_N \cdot \Sigma^{-\frac{1}{2}}) \odot R$$

$$Y = (1_N \cdot \gamma) \odot Z + 1_N \cdot \beta$$

$$\begin{aligned} R &= X - \text{row_repeat}(\text{columns_mean}(X), N) \\ \text{Sigma} &= \text{diag}(R.T * R).T / N \\ \text{inv_sqrt_Sigma} &= \text{inv_sqrt}(\text{Sigma}) \\ Z &= \text{hadamard}(\text{row_repeat}(\text{inv_sqrt_Sigma}, N), R) \\ Y &= \text{hadamard}(\text{row_repeat}(\gamma, N), Z) + \text{row_repeat}(\beta, N) \end{aligned}$$

linear dropout layer

FEEDFORWARD EQUATIONS

$$Y = X(W^T \odot R) + 1_N \cdot b$$

$$Y = X * \text{hadamard}(W.T, R) + \text{row_repeat}(b, N)$$

BACKPROPAGATION EQUATIONS

$$DZ = DY - \text{softmax}(Z) \odot (DY \cdot 1_K \cdot 1_K^T)$$

$$DW = DZ^T \cdot X$$

$$Db = 1_N^T \cdot DZ$$

$$DX = DZ \cdot W$$

$$\begin{aligned} DZ &= DY - \text{hadamard}(\text{softmax}(Z), \\ &\quad \text{column_repeat}(\text{rows_sum}(DY), K)) \end{aligned}$$

$$DW = DZ.T * X$$

$$Db = \text{columns_sum}(DZ)$$

$$DX = DZ * W$$

BACKPROPAGATION EQUATIONS

$$DZ = (1_N \cdot \gamma) \odot DY$$

$$D\beta = 1_N^T \cdot DY$$

$$D\gamma = 1_N^T \cdot (Z \odot DY)$$

$$\begin{aligned} DX &= \left(\frac{1}{N} \cdot 1_N \cdot \Sigma^{-\frac{1}{2}} \right) \odot \\ &\quad ((N \cdot \mathbb{I}_N - 1_N \cdot 1_N^T) \cdot DZ - Z \odot (1_N \cdot \text{diag}(Z^T \cdot DZ)^T)) \end{aligned}$$

$$DZ = \text{hadamard}(\text{row_repeat}(\gamma, N), DY)$$

$$Dbeta = \text{columns_sum}(DY)$$

$$Dgamma = \text{columns_sum}(\text{hadamard}(Z, DY))$$

$$\begin{aligned} DX &= \text{hadamard}(\text{row_repeat}(\text{inv_sqrt_Sigma} / N, N), \\ &\quad (N * \text{identity}(N) - \text{ones}(N, N)) * DZ - \\ &\quad \text{hadamard}(Z, \text{row_repeat}(\text{diag}(Z.T * DZ).T, N))) \end{aligned}$$

BACKPROPAGATION EQUATIONS

$$DW = (DY^T \cdot X) \odot R^T$$

$$Db = 1_N^T \cdot DY$$

$$DX = DY(W \odot R^T)$$

$$DW = \text{hadamard}(DY.T * X, R.T)$$

$$Db = \text{columns_sum}(DY)$$

$$DX = DY * \text{hadamard}(W, R.T)$$

activation dropout layer

FEEDFORWARD EQUATIONS

$$Z = X(W^\top \odot R) + 1_N \cdot b$$

$$Y = \text{act}(Z)$$

$$\begin{aligned} Z &= X * \text{hadamard}(W.T, R) + \text{row_repeat}(b, N) \\ Y &= \text{act}(Z) \end{aligned}$$

BACKPROPAGATION EQUATIONS

$$DZ = DY \odot \text{act}'(Z)$$

$$DW = (DZ^\top \cdot X) \odot R^\top$$

$$Db = 1_N^\top \cdot DZ$$

$$DX = DZ(W \odot R^\top)$$

$$DZ = \text{hadamard}(DY, \text{act.gradient}(Z))$$

$$DW = \text{hadamard}(DZ.T * X, R.T)$$

$$Db = \text{columns_sum}(DZ)$$

$$DX = DZ * \text{hadamard}(W, R.T)$$

4.1 Derivations

In this section we give some derivations of the backpropagation equations. We give some applications of the product rule and chain rule for vector functions, and show how some properties on the rows and columns of a matrix can be generalized into matrix-form.

Lemma 1 (Product Rule for vector functions). *The product rule for scalar functions u and v is given by*

$$(u \cdot v)' = u' \cdot v + u \cdot v'.$$

It can be generalized to vector functions, but the result is sensitive to the orientation of the operands. Below we give four concrete applications of the product rule for vector functions. Let $x \in \mathbb{R}^p$, $A \in \mathbb{R}^{m \times n}$, and $h(x) = f(x)g(x)$ for $m, n, p \in \mathbb{N}^+$.

$$\text{Let } f(x) \in \mathbb{R}^{n \times 1}, \text{ and let } g(x) \in \mathbb{R}, \text{ then } \frac{\partial h}{\partial x} = \frac{\partial f}{\partial x} g + f \frac{\partial g}{\partial x}. \quad (5)$$

$$\text{Let } f(x) \in \mathbb{R}^{1 \times n}, \text{ and let } g(x) \in \mathbb{R}, \text{ then } \frac{\partial h}{\partial x} = \frac{\partial f}{\partial x} g + f^\top \frac{\partial g}{\partial x}. \quad (6)$$

$$\text{Let } f(x) = A, \text{ and let } g(x) \in \mathbb{R}^{n \times 1}, \text{ then } \frac{\partial h}{\partial x} = A \frac{\partial g}{\partial x}. \quad (7)$$

$$\text{Let } f(x) \in \mathbb{R}^{1 \times m}, \text{ and let } g(x) = A, \text{ then } \frac{\partial h}{\partial x} = A^\top \frac{\partial f}{\partial x}. \quad (8)$$

Lemma 2 (Chain rule for vector functions). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, let $g : \mathbb{R}^m \rightarrow \mathbb{R}^p$, and let $h(x) = g(y)$ with $y = f(x)$. Then we have*

$$\frac{\partial h}{\partial x} = \frac{\partial g}{\partial y} \cdot \frac{\partial f}{\partial x} \quad (9)$$

Note that this equations holds irrespective of whether $f(x)$ is a row or a column vector.

Property 1 (Matrix properties). *Let $X, Y, Z \in \mathbb{R}^{m \times n}$. We denote the i^{th} row of matrix A as a_i and the j^{th} column of matrix A as a^j . The element at position (i, j) of A is denoted as a_{ij} . Below we give a number of properties on the rows and columns of matrices, and the generalization of these properties into matrix-form.*

$$\text{If } z^j = x^j \cdot (x^j)^\top \cdot y^j \quad (1 \leq j \leq n), \text{ then } Z = X \odot (1_m \cdot \text{diag}(X^\top \cdot Y)^\top). \quad (10)$$

$$\text{If } z^j = 1_m \cdot (x^j)^\top \cdot y^j \quad (1 \leq j \leq n), \text{ then } Z = 1_m \cdot \text{diag}(X^\top Y)^\top. \quad (11)$$

$$\text{If } z^j = x^j \cdot 1_m^\top \cdot y^j \quad (1 \leq j \leq n), \text{ then } Z = X \odot (1_m \cdot 1_m^\top \cdot Y) \quad (12)$$

$$\text{If } z_i = x_i \cdot y_i^\top \cdot y_i \quad (1 \leq i \leq m), \text{ then } Z = (\text{diag}(X \cdot Y^\top) \cdot 1_n^\top) \odot Y \quad (13)$$

$$\text{If } z_i = x_i \cdot y_i^\top \cdot 1_n^\top \quad (1 \leq i \leq m), \text{ then } Z = \text{diag}(X \cdot Y^\top) \cdot 1_n^\top \quad (14)$$

$$\text{If } z_i = x_i \cdot 1_n \cdot y_i \quad (1 \leq i \leq m), \text{ then } Z = (X \cdot 1_n \cdot 1_n^\top) \odot Y \quad (15)$$

All properties have been validated with SymPy. Naturally there is a lot of symmetry between the row and column properties. By means of example we will prove equation (13). From $z_i = x_i \cdot y_i^\top \cdot y_i$ we derive that $z_{ij} = (x_i \cdot y_i^\top) y_{ij}$ for $1 \leq j \leq n$. Hence we can write $Z = R \odot Y$, where R is defined using $r_{ij} = (x_i \cdot y_i^\top)$. We observe that $\text{diag}(X \cdot Y^\top) = (x_1 \cdot y_1^\top, \dots, x_m \cdot y_m^\top)^\top$. From the definition of R we can see that it consists of n copies of the column vector $\text{diag}(X \cdot Y^\top)$. Hence we have $R = \text{diag}(X \cdot Y^\top) \cdot 1_n^\top$.

linear layer

We have $Y = X \cdot W^\top + 1_N \cdot b$. Let x_i, y_i, b_i be the i^{th} row of $X, Y, 1_N \cdot b$ for $1 \leq i \leq N$, hence $y_i = x_i W^\top + b_i$, where $b_i = b$. Furthermore, let w_j be the j^{th} row of W , and let e_i be the i^{th} row of the unit matrix \mathbb{I}_N . Let $\mathcal{L}(Y) = \sum_{i=1}^N \mathcal{L}(y_i)$ be the corresponding loss. We calculate

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial x_i} &\stackrel{(9)}{=} \sum_{k=1}^N \frac{\partial \mathcal{L}}{\partial y_k} \frac{\partial y_k}{\partial x_i} = \frac{\partial \mathcal{L}}{\partial y_i} \frac{\partial y_i}{\partial x_i} \stackrel{(8)}{=} \frac{\partial \mathcal{L}}{\partial y_i} \cdot W, \text{ hence } D x_i = D y_i \cdot W \\ \frac{\partial \mathcal{L}}{\partial b} &\stackrel{(9)}{=} \sum_{k=1}^N \frac{\partial \mathcal{L}}{\partial y_k} \frac{\partial y_k}{\partial b} \stackrel{(8)}{=} \sum_{k=1}^N \frac{\partial \mathcal{L}}{\partial y_k} \cdot \mathbb{I}_K = \sum_{k=1}^N \frac{\partial \mathcal{L}}{\partial y_k}, \text{ hence } D b = 1_N^\top \cdot D Y \\ \text{We have } y_k &= x_k W^\top + b \text{ and } \frac{\partial y_k}{\partial w_{ij}} = x_{kj} e_i, \text{ hence} \end{aligned} \quad (16)$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_{ij}} &= \sum_{k=1}^N \frac{\partial \mathcal{L}}{\partial y_k} \frac{\partial y_k}{\partial w_{ij}} = \sum_{k=1}^N \frac{\partial \mathcal{L}}{\partial y_k} x_{kj} e_i = \sum_{k=1}^N \sum_{l=1}^K \frac{\partial \mathcal{L}}{\partial y_{kl}} x_{kj} (e_i)_l \\ &= \sum_{k=1}^N (e_i)_l \sum_{l=1}^K \frac{\partial \mathcal{L}}{\partial y_{kl}} x_{kj} = \sum_{k=1}^N (e_i)_l (D Y^\top \cdot X)_{lj} = (D Y^\top \cdot X)_{ij} \end{aligned}$$

This can be generalized to matrix equations: $D X = D Y \cdot W$, $D b = 1_N^\top \cdot D Y$, and $D W = D Y^\top \cdot X$.

log-softmax layer

We have $Y = \text{log-softmax}(Z)$. Let y_i, z_i be the i^{th} row of Y, Z for $1 \leq i \leq N$, hence $y_i = \text{log-softmax}(z_i)$. We calculate

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial z_i} &\stackrel{(9)}{=} \frac{\partial \mathcal{L}}{\partial y_i} \frac{\partial}{\partial z_i} \text{log-softmax}(z_i) \stackrel{(28)}{=} \frac{\partial \mathcal{L}}{\partial y_i} (\mathbb{I}_K - 1_K \cdot \text{softmax}(z_i)) \\ &\iff D z_i = D y_i \cdot (\mathbb{I}_K - 1_K \cdot \text{softmax}(z_i)) = D y_i - D y_i \cdot 1_K \cdot \text{softmax}(z_i). \end{aligned}$$

This can be generalized to matrices using property 15:

$$D Z = D Y - \text{softmax}(Z) \odot (D Y \cdot 1_K \cdot 1_K^\top) \quad (17)$$

batch normalization layer

We will derive the equation for $\mathbf{D}X$, which is the most complicated one. Following the approach of [Yeh17], we first derive the equations for a single column. Let $x^j, r^j, z^j \in \mathbb{R}^{N \times 1}$ be the j^{th} column of X , R and Z , and let $\sigma \in \mathbb{R}$ be the j^{th} element of Σ , with $1 \leq j \leq D$. Then we obtain the following equations:

$$r^j = x^j - \frac{1_N \cdot 1_N^\top}{N} \cdot x^j = (\mathbb{I}_N - \frac{1_N \cdot 1_N^\top}{N}) \cdot x^j \quad (18)$$

$$\sigma = \frac{(r^j)^\top r^j}{N} \quad (19)$$

$$z^j = r^j \cdot \sigma^{-\frac{1}{2}} \quad (20)$$

We calculate

$$\frac{\partial z^j}{\partial r^j} \stackrel{(5)}{=} \frac{\partial r^j}{\partial r^j} \cdot \sigma^{-\frac{1}{2}} + r^j \cdot \frac{\partial \sigma^{-\frac{1}{2}}}{\partial r^j} \stackrel{(9)}{=} \sigma^{-\frac{1}{2}} - r^j \cdot \frac{\sigma^{-\frac{3}{2}}}{2} \cdot \frac{\partial \sigma}{\partial r^j} = \sigma^{-\frac{1}{2}} - r^j \cdot \frac{\sigma^{-\frac{3}{2}}}{2} \cdot \left(\frac{2(r^j)^\top}{N} \right) \quad (21)$$

$$= \frac{\sigma^{-\frac{1}{2}}}{N} (N \cdot \mathbb{I}_N - \sigma^{-1} r^j (r^j)^\top) = \frac{\sigma^{-\frac{1}{2}}}{N} \cdot (N \cdot \mathbb{I}_N - z^j (z^j)^\top). \quad (22)$$

Using the chain rule we find

$$\frac{\partial \mathcal{L}}{\partial r^j} \stackrel{(9)}{=} \frac{\partial \mathcal{L}}{\partial z^j} \frac{\partial z^j}{\partial r^j}, \text{ hence } \mathbf{D}r^j = \frac{\sigma^{-\frac{1}{2}}}{N} (N \cdot \mathbb{I}_N - z^j (z^j)^\top) \cdot \mathbf{D}z^j \quad (23)$$

$$\frac{\partial \mathcal{L}}{\partial x^j} \stackrel{(9)}{=} \frac{\partial \mathcal{L}}{\partial r^j} \frac{\partial r^j}{\partial x^j}, \text{ hence } \mathbf{D}x^j = (\mathbb{I}_N - \frac{1_N \cdot 1_N^\top}{N}) \cdot \mathbf{D}r^j, \quad (24)$$

where one needs to take into account that for a column vector x we have $\mathbf{D}x = (\frac{\partial \mathcal{L}}{\partial x})^\top$. Hence

$$\begin{aligned} \mathbf{D}x^j &= \frac{\sigma^{-\frac{1}{2}}}{N} \cdot (\mathbb{I}_N - \frac{1_N \cdot 1_N^\top}{N}) \cdot (N \cdot \mathbb{I}_N - z^j \cdot (z^j)^\top) \cdot \mathbf{D}z^j \\ &= \frac{\sigma^{-\frac{1}{2}}}{N} \cdot (N \cdot \mathbb{I}_N - z^j \cdot (z^j)^\top - 1_N \cdot 1_N^\top + \frac{1_N \cdot 1_N^\top \cdot z^j \cdot (z^j)^\top}{N}) \cdot \mathbf{D}z^j \\ &\quad \{ \text{In batch normalization the column sum } 1_N^\top \cdot z^j \text{ evaluates to zero. } \} \\ &= \frac{\sigma^{-\frac{1}{2}}}{N} \cdot ((N \cdot \mathbb{I}_N - 1_N \cdot 1_N^\top) \cdot \mathbf{D}z^j - z^j \cdot (z^j)^\top \cdot \mathbf{D}z^j) \end{aligned}$$

Using property (10) we generalize this into matrix-form:

$$\mathbf{D}X = \left(\frac{1}{N} \cdot 1_N \cdot \Sigma^{-\frac{1}{2}} \right) \odot \left((N \cdot \mathbb{I}_N - 1_N \cdot 1_N^\top) \cdot \mathbf{D}Z - Z \odot (1_N \cdot \text{diag}(Z^\top \cdot \mathbf{D}Z)^\top) \right). \quad (25)$$

For the remaining equations, let $y_i, z_i, \beta_i, \gamma_i$ be the i^{th} row of Y , Z , $1_N \cdot \beta$, and $1_N \cdot \gamma$ for $1 \leq i \leq N$, where $\beta_i = \beta$, and $\gamma_i = \gamma$. Hence $y_i = \gamma_i \odot z_i + \beta_i$. From this it follows

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \beta_i} &= \frac{\partial \mathcal{L}}{\partial y_i} \\ \frac{\partial \mathcal{L}}{\partial \gamma_i} &= \frac{\partial \mathcal{L}}{\partial y_i} \odot z_i, \end{aligned}$$

which we can generalize into matrix-form: $\mathbf{D}\beta = 1_N^\top \cdot \mathbf{D}Y$ and $\mathbf{D}\gamma = 1_N^\top \cdot (\mathbf{D}Y \odot Z)$.

5 Activation functions

In this section we give an overview of some commonly used univariate activation functions. These activation functions are typically applied element-wise to the output matrix Y of a neural network.

NAME	FUNCTION	DERIVATIVE
ReLU [Fuk75]	$\text{relu}(x) = \max(0, x)$	$\text{relu}'(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{otherwise} \end{cases}$
Leaky ReLU [Maa13]	$\text{leaky-relu}(x) = \max(\alpha x, x)$	$\text{leaky-relu}'(x) = \begin{cases} \alpha & \text{if } x < \alpha x \\ 1 & \text{otherwise} \end{cases}$
All-ReLU ([CMP21])	$\text{all-relu}(x) = \begin{cases} \alpha x & \text{if } x < 0 \\ x & \text{otherwise} \end{cases}$	$\text{all-relu}'(x) = \begin{cases} \alpha & \text{if } x < 0 \\ 1 & \text{otherwise} \end{cases}$
SReLU ([JXF ⁺ 16])	$\text{srelu}(x) = \begin{cases} t_l + a_l(x - t_l) & \text{if } x \leq t_l \\ x & \text{if } t_l < x < t_r \\ t_r + a_r(x - t_r) & \text{if } x \geq t_r \end{cases}$	$\text{srelu}'(x) = \begin{cases} a_l & \text{if } x \leq t_l \\ 1 & \text{if } t_l < x < t_r \\ a_r & \text{if } x \geq t_r \end{cases}$
Hyperbolic tangent	$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$	$\tanh'(x) = 1 - \tanh(x)^2$
Sigmoid / logistic function [HDY ⁺ 12]	$\sigma(x) = \frac{1}{1 + e^{-x}}$	$\sigma'(x) = \sigma(x)(1 - \sigma(x))$

Note that Leaky ReLU and All-ReLU depend on a parameter $\alpha \in \mathbb{R}$ and SReLU depends on four parameters $a_l, t_l, a_r, t_r \in \mathbb{R}$. N.B. The three cases of the SReLU function overlap if $t_l \geq t_r$. We use the convention that in such a case the first possible alternative must be chosen. For All-ReLU and SReLU layers the sign of the parameter α should be alternated. In a network with layers $1, 2, 3, \dots$ the All-ReLU function should be applied to the layers $2, 3, 4, \dots$ with parameters $-\alpha, \alpha, -\alpha, \dots$.

The SReLU parameters a_l, t_l, a_r, t_r are learnable parameters, i.e. they should be updated periodically based on the values of their gradients with respect to the loss function. As initial values we take $a_l, t_l, a_r, t_r = 0, 0, 0, 1$. This corresponds to the ReLU function.

5.1 Softmax functions

In this section we give an overview of softmax functions and their derivatives, both for single inputs and for batches. These equations are present in the activation function of softmax layers and in some of the loss functions in section 6. We use the same notation as before, where $x \in \mathbb{R}^{1 \times D}$ is a single input with dimension D , and $X \in \mathbb{R}^{N \times D}$ is an input batch, where N is the number of samples in the batch. We denote the rows of X as x_1, \dots, x_N .

Below an overview of softmax functions and their derivatives is given. Equations in matrix-form are only given if they are needed later on. Note that the function **log-softmax** is defined using

$$\text{log-softmax}(x) = \log(\text{softmax}(x)).$$

VECTOR FUNCTIONS

$$\text{softmax}(x) = \frac{e^x}{e^x \cdot \mathbf{1}_D}$$

$$\text{stable-softmax}(x) = \frac{e^z}{e^z \cdot \mathbf{1}_D}$$

$$\text{log-softmax}(x) = x - \log(e^x \cdot \mathbf{1}_D) \cdot \mathbf{1}_D^\top$$

$$\text{stable-log-softmax}(x) = z - \log(e^z \cdot \mathbf{1}_D) \cdot \mathbf{1}_D^\top$$

MATRIX FUNCTIONS

$$\text{softmax}(X) = e^X \odot \left(\frac{1}{e^X \cdot \mathbf{1}_D} \cdot \mathbf{1}_D^\top \right)$$

$$\text{stable-softmax}(X) = e^Z \odot \left(\frac{1}{e^Z \cdot \mathbf{1}_D} \cdot \mathbf{1}_D^\top \right)$$

$$\text{log-softmax}(X) = X - \log(e^X \cdot \mathbf{1}_D) \cdot \mathbf{1}_D^\top$$

$$\text{stable-log-softmax}(X) = Z - \log(e^Z \cdot \mathbf{1}_D) \cdot \mathbf{1}_D^\top,$$

where

$$\begin{cases} z = x - \max(x) \cdot \mathbf{1}_D^\top \\ Z = X - \max(X)_{\text{row}} \cdot \mathbf{1}_D^\top. \end{cases}$$

The derivatives of the softmax functions are given by

$$\frac{\partial}{\partial x} \text{softmax}(x) = \frac{\partial}{\partial x} \text{stable-softmax}(x) = \text{Diag}(\text{softmax}(x)) - \text{softmax}(x)^\top \cdot \text{softmax}(x), \quad (26)$$

$$\frac{\partial}{\partial x} \text{log-softmax}(x) = \frac{\partial}{\partial x} \text{stable-log-softmax}(x) = \mathbb{I}_D - \mathbf{1}_D \cdot \text{softmax}(x) \quad (27)$$

5.2 Derivations

log-softmax

Let $y(x) = \text{softmax}(x)$, then we have

$$\begin{aligned} \frac{\partial}{\partial x} \log(\text{softmax}(x)) &\stackrel{(9)}{=} \frac{\partial}{\partial y} \log(y) \frac{\partial}{\partial x} \text{softmax}(x) \\ &\stackrel{(\text{??})}{=} \text{Diag}\left(\frac{1}{y}\right) \cdot (\text{Diag}(y) - y^\top y) \\ &= \mathbb{I}_D - \text{Diag}(y) y^\top y \\ &= \mathbb{I}_D - \mathbf{1}_D \cdot y \\ &= \mathbb{I}_D - \mathbf{1}_D \cdot \text{softmax}(x). \end{aligned} \quad (28)$$

6 Loss functions

In this section we give an overview of some loss functions and their gradients, both for single inputs and for batches. We use the following notations:

1. $y \in \mathbb{R}^{1 \times K}$ is a single output, where K is the output dimension.
2. $t \in \mathbb{R}^{1 \times K}$ is a target for a single output y .
3. $Z \in \mathbb{R}^{N \times K}$ is an output batch, where N is the number of examples in the batch.
4. $T \in \mathbb{R}^{N \times K}$ contains the targets for an output batch Y .

squared error loss

$$\begin{aligned}
\mathcal{L}_{\text{SE}}(y, t) &= (y - t)(y - t)^\top \\
\nabla_y \mathcal{L}_{\text{SE}}(y, t) &= 2(y - t) \\
\mathcal{L}_{\text{SE}}(Y, T) &= \mathbf{1}_N^\top \cdot ((Y - T) \odot (Y - T)) \cdot \mathbf{1}_K \\
\nabla_Y \mathcal{L}_{\text{SE}}(Y, T) &= 2(Y - T)
\end{aligned} \tag{29}$$

The mean squared error loss is a scaled version of the squared error loss.

$$\mathcal{L}_{\text{MSE}}(Y, T) = \frac{\mathcal{L}_{\text{SE}}(Y, T)}{K \cdot N}$$

cross-entropy loss

$$\begin{aligned}
\mathcal{L}_{\text{CE}}(y, t) &= -t \cdot \log(y)^\top \\
\nabla_y \mathcal{L}_{\text{CE}}(y, t) &= -t \odot \frac{1}{y} \\
\mathcal{L}_{\text{CE}}(Y, T) &= -\mathbf{1}_N^\top \cdot (T \odot \log(Y)) \cdot \mathbf{1}_K \\
\nabla_Y \mathcal{L}_{\text{CE}}(Y, T) &= -T \odot \frac{1}{Y}
\end{aligned} \tag{30}$$

softmax cross-entropy loss

$$\begin{aligned}
\mathcal{L}_{\text{SCE}}(y, t) &= -t \cdot \log\text{-softmax}(y)^\top \\
\nabla_y \mathcal{L}_{\text{SCE}}(y, t) &= t \cdot \mathbf{1}_K \cdot \text{softmax}(y) - t \\
\mathcal{L}_{\text{SCE}}(Y, T) &= -\mathbf{1}_N^\top \cdot (T \odot \log\text{-softmax}(Y)) \cdot \mathbf{1}_K \\
\nabla_Y \mathcal{L}_{\text{SCE}}(Y, T) &= \text{softmax}(Y) \odot (T \cdot \mathbf{1}_K \cdot \mathbf{1}_K^\top) - T
\end{aligned} \tag{31}$$

Note that if t is a one-hot encoded target (e.g. in case of a classification problem), it is a vector consisting of one value 1 and all others values 0. In other words we have $t \cdot \mathbf{1}_K = 1$, hence in this case the gradients can be simplified to

$$\begin{aligned}
\nabla_y \mathcal{L}_{\text{SCE}}(y, t) &= \text{softmax}(y) - t \\
\nabla_Y \mathcal{L}_{\text{SCE}}(Y, T) &= \text{softmax}(Y) - T
\end{aligned}$$

logistic cross-entropy loss

$$\begin{aligned}
\mathcal{L}_{\text{LCE}}(y, t) &= -t \cdot \log(\sigma(y))^\top \\
\nabla_y \mathcal{L}_{\text{LCE}}(y, t) &= t \odot \sigma(y) - t \\
\mathcal{L}_{\text{LCE}}(Y, T) &= -\mathbf{1}_N^\top \cdot (T \odot (\log(\sigma(Y)))) \cdot \mathbf{1}_K \\
\nabla_Y \mathcal{L}_{\text{LCE}}(Y, T) &= T \odot \sigma(Y) - T
\end{aligned} \tag{32}$$

negative log-likelihood loss

$$\begin{aligned}
\mathcal{L}_{\text{NLL}}(y, t) &= -\log(y \cdot t^\top) \\
\nabla_y \mathcal{L}_{\text{NLL}}(y, t) &= -\frac{1}{y \cdot t^\top} \cdot t \\
\mathcal{L}_{\text{NLL}}(Y, T) &= -1_N^\top \cdot (\log((Y \odot T) \cdot 1_K)) \cdot 1_K \\
\nabla_Y \mathcal{L}_{\text{NLL}}(Y, T) &= -\left(\frac{1}{(Y \odot T) \cdot 1_K} \cdot 1_K^\top\right) \odot T
\end{aligned} \tag{33}$$

6.1 Derivations

cross-entropy loss

$$\nabla_y \mathcal{L}_{\text{CE}}(y, t) = -\frac{\partial}{\partial y} (t \cdot \log(y)^\top) \stackrel{(7)}{=} -t \cdot \text{Diag}\left(\frac{1}{y}\right) = -t \odot \frac{1}{y} \tag{34}$$

softmax cross-entropy loss

$$\begin{aligned}
\nabla_y \mathcal{L}_{\text{SCE}}(y, t) &= -\frac{\partial}{\partial y} (t \cdot \log\text{-softmax}(y)^\top) \stackrel{(7)}{=} -t \cdot \frac{\partial}{\partial y} \log\text{-softmax}(y) \\
&\stackrel{(28)}{=} -t \cdot (\mathbb{I}_K - 1_K \cdot \text{softmax}(y)) = t \cdot 1_K \cdot \text{softmax}(y) - t
\end{aligned} \tag{35}$$

If the target t is a one-hot encoded vector, we have $t \cdot 1_K = 1$. In that case the gradient simplifies to

$$\nabla_y \mathcal{L}_{\text{SCE-one-hot}}(y, t) = \text{softmax}(y) - t.$$

Using property (15) we can generalize equation (35) to

$$\nabla_Y \mathcal{L}_{\text{SCE}}(Y, T) = \text{softmax}(Y) \odot (T \cdot 1_K \cdot 1_K^\top) - T.$$

logistic cross-entropy loss

$$\begin{aligned}
\nabla_y \mathcal{L}_{\text{LCE}}(y, t) &= \frac{\partial}{\partial y} (-t \cdot \log(\sigma(y))^\top) \stackrel{(7)}{=} -t \cdot \frac{\partial}{\partial y} \log(\sigma(y)) \\
&= -t \cdot \text{Diag}(1_K^\top - \sigma(y)) = t \odot \sigma(y) - t,
\end{aligned}$$

where we used the well known fact that in the univariate case $\frac{d}{dt} \log(\sigma(t)) = 1 - \log(\sigma(t))$.

7 Weight initialization

The initial values of the weights in linear layers must be carefully chosen, as they may have a large impact on the performance of a neural network [NBS22]. Typically, these values are randomly generated on the basis of specific probability distributions. In this section, we give a few commonly used distributions.

NAME	DISTRIBUTION
Uniform	$U(a, b)$
Xavier [GB10]	$U(-\frac{1}{\sqrt{D}}, \frac{1}{\sqrt{D}})$
Normalized Xavier [GB10]	$U(-\frac{\sqrt{6}}{\sqrt{D+K}}, \frac{\sqrt{6}}{\sqrt{D+K}})$
He [HZRS15]	$\mathcal{N}(0, \sqrt{\frac{2}{D}})$

where D is the number of inputs and K the number of outputs of the layer to which the weight matrix belongs. Furthermore, $U(a, b)$ is the uniform distribution in a given interval $[a, b]$, and $\mathcal{N}(\mu, \sigma)$ is the normal distribution with mean μ and standard deviation σ .

Unlike weights, the initial values of bias vectors are typically set to a small constant or zero rather than drawn from probability distributions. However, the Xavier weight distribution can also be used.

8 Optimization

In the optimization step, the parameters θ of a layer are updated based on the value of their gradient $D\theta$ with respect to a given loss function. The goal of this step is to decrease the value of the loss. In this section, we give three common choices for optimization methods: gradient descent, momentum, and Nesterov. All take a learning rate parameter η as input, which is used to control the size of the optimization step. Our equations are equivalent to the ones in Keras [C⁺15], but presented in matrix form. We use a prime symbol to denote updated values.

GRADIENT DESCENT	MOMENTUM	NESTEROV
$\theta' = \theta - \eta \cdot D\theta$	$\Delta'_\theta = \mu \cdot \Delta_\theta - \eta \cdot D\theta$ $\theta' = \theta + \Delta'_\theta$	$\Delta'_\theta = \mu \cdot \Delta_\theta - \eta \cdot D\theta$ $\theta' = \theta + \mu \cdot \Delta'_\theta - \eta \cdot D\theta$

Both momentum and Nesterov depend on a parameter $0 < \mu < 1$. Furthermore, they store an additional parameter Δ_θ with the same dimensions as θ . This parameter is updated in each optimization call and initially contains only zeros.

9 Learning Rate Schedulers

We define a learning rate scheduler as a function $\eta : \mathbb{N} \rightarrow \mathbb{R}$ that returns the learning rate at optimization step i . We assume that η_0 is a given initial learning rate.

FORMULA	DESCRIPTION
$\eta_i = \eta_0$	A constant scheduler with initial value η_0 .
$\eta_{i+1} = \frac{\eta_i}{1 + d \cdot i}$	A time-based scheduler with decay parameter d .
$\eta_i = \eta_0 \cdot d^{\lfloor \frac{1+i}{r} \rfloor}$	A step-based scheduler with change rate d and drop rate r .
$\eta_i = \eta_0 e^{-d \cdot i}$	An exponential scheduler with decay parameter d .
$\eta_i = \eta_0 \Gamma^{\sum_{j=1}^k \lfloor \frac{i}{m_j} \rfloor}$	A multi-step scheduler with decay parameter Γ and milestones $\{m_1, \dots, m_k\}$.

References

- [C⁺15] François Chollet et al. Keras. <https://keras.io>, 2015.
- [CMP21] Selima Curci, Decebal Constantin Mocanu, and Mykola Pechenizkiy. Truly sparse neural networks at scale. *CoRR*, abs/2102.01732, 2021.
- [Fuk75] K. Fukushima. Cognitron: a self-organizing multilayered neural network. *Biological Cybernetics*, 20:121–136, 1975.
- [GB10] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In Yee Whye Teh and Mike Titterton, editors, *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 249–256, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR.
- [HDY⁺12] Geoffrey Hinton, Li Deng, Dong Yu, George E. Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N. Sainath, and Brian Kingsbury. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Processing Magazine*, 29(6):82–97, 2012.
- [HZRS15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1026–1034, 2015.
- [JXF⁺16] Xiaojie Jin, Chunyan Xu, Jiashi Feng, Yunchao Wei, Junjun Xiong, and Shuicheng Yan. Deep learning with s-shaped rectified linear activation units. In Dale Schuurmans and Michael P. Wellman, editors, *AAAI*, pages 1737–1743. AAAI Press, 2016.
- [Maa13] Andrew L. Maas. Rectifier nonlinearities improve neural network acoustic models. 2013.
- [NBS22] Meenal V. Narkhede, Prashant P. Bartakke, and Mukul S. Sutaone. A review on weight initialization strategies for neural networks. *Artif. Intell. Rev.*, 55(1):291–322, 2022.
- [Yeh17] Chris Yeh. Deriving Batch-Norm Backprop Equations. <https://chrisyeh96.github.io/2017/08/28/deriving-batchnorm-backprop.html>, 2017.