

# MathDNN - A deep mathematical understanding of DNNs

James JIANG  
Data Engineer / Scientist  
France

Alex JIANG  
Preparatory class for the Grandes Écoles  
France

iLoveDataJia Github

Version: 0.00

Date: January 26, 2022

## Abstract

Frameworks such as [TensorFlow](#) or [PyTorch](#) make deep learning developments easy. They have made this field wide spread for every enthusiast. Implementations only needs an instinctive understanding of deep learning. The proper math aspect is little by little forgotten. Topology, Normalized vector space, Limit plus continuity, Taylor series expansion, Riemann integral theory, Matrix, Finite dimensional linear algebra and Linear application matrix theories are supposed known. The objective is to do a collection of the important propositions explaining dense neural network (DNN) theories. All the propositions will be mathematically proven as far as possible and under assumptions if necessary. The subject used as reference is a multi-class classification problem with – dense layers, activation layers, Categorical cross-entropy loss and Stochastic gradient descent optimizer with momentum. But all the elements below can be easily re-used or re-defined to cover regressions.

**Keywords:** Dense neural network, Differentiability, Continuous optimization

## 1 Fundamentals

### 1.1 Matrices

*Convention 1.* All sets considered are not empty.

*Notation 1.* Let  $a_{i,j} \in \mathbb{R}$  for  $i \in \llbracket 1, n \rrbracket$  and  $j \in \llbracket 1, m \rrbracket$ . Then a real matrix of dimension  $n * m$  will noted as

$$A = \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,m} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n,1} & a_{n,2} & \cdots & a_{n,m} \end{bmatrix}$$

The following notations are also considered

$$\forall i \in \llbracket 1, n \rrbracket, \forall j \in \llbracket 1, m \rrbracket, A_{i,j} = a_{i,j}$$

$$\forall j \in \llbracket 1, m \rrbracket, A_{:,j} = \begin{bmatrix} a_{1,j} \\ \vdots \\ a_{n,j} \end{bmatrix}$$

$$\forall i \in \llbracket 1, n \rrbracket, j \in \llbracket 1, m \rrbracket, A_{i,j} = \begin{bmatrix} a_{i,1} & \cdots & a_{i,n} \end{bmatrix}$$

The notation  $\mathcal{M}_{n,m}$  means the matrix set of dimension  $n \times m$  with coefficients in  $\mathbb{R}$ .

The notation  $\mathcal{M}_{n,m}(E)$  means the matrix set of dimension  $n \times m$  with coefficients in  $E \subseteq \mathbb{R}$ .

*Convention 2.* Let  $E \subseteq \mathbb{R}$ .

A vector is a matrix with only one row. Thus, the vector set  $E^n$  is equivalent to  $\mathcal{M}_{1,n}(E)$ .

A  $m$ -tuple of vectors is a matrix with  $m$  rows. Thus, the cartesian products of vectors  $(E^n)^m$  is equivalent to  $\mathcal{M}_{m,n}(E)$ .

*Notation 2.* Let  $A \in \mathcal{M}_{n,m}$  and  $B \in \mathcal{M}_{m,p}$ . Let the product noted  $A * B$  be

$$C = A * B$$

where  $C \in \mathcal{M}_{n,p}$  with

$$\forall i \in \llbracket 1, n \rrbracket, \forall j \in \llbracket 1, p \rrbracket, C_{i,j} = \sum_{k=1}^n A_{i,k} * B_{k,j}$$

*Notation 3.* The matrix transpose operation will be noted as  $A^T$ .

*Notation 4.* The notation  $I_n$  means the identity matrix of size  $n \times n$ .

$$I_n = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1 \end{bmatrix}$$

*Notation 5.* Let  $a \in \mathbb{R}^n$ . The euclidean norm on  $\mathbb{R}^n$  will be noted as  $\|a\|_n$ .

$$\|a\|_n = \sqrt{a * a^T}$$

## 1.2 Differential calculus

*Notation 6.* Let  $E \subseteq \mathbb{R}^n$ ,  $F \subseteq \mathbb{R}^m$ ,  $a \in \mathbb{R}^n$  and  $r \in \mathbb{R}^+ *$ .

The notation  $\overset{\circ}{E}$  means the interior of  $E$ .

The notation  $f : E \longrightarrow F$  means the application from  $E$  to  $F$ .

The notation  $\mathcal{F}(E, F)$  means the set of applications from  $E$  to  $F$ .

The notation  $\mathcal{C}(E, F)$  means the set of continuous applications from  $E$  to  $F$ .

The notation  $\mathcal{L}(E, F)$  means the set of linear applications from  $E$  to  $F$ .

The notation  $\mathcal{B}(a, r)$  means the set  $\{x \in \mathbb{R}^n \mid \|x - a\|_n \leq r\}$ .

**Definition 1.1.** Let  $E \subseteq \mathbb{R}^n$  and  $F \subseteq \mathbb{R}^m$ . Then  $f$  differentiable on  $E$  means

$$\forall a \in \mathring{E}, \exists \frac{\partial f}{\partial \cdot}(a) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m), \quad (1)$$

$$\exists \eta \in \mathbb{R}^{+*}, \forall h \in \mathcal{B}(0_{\mathbb{R}^n}, \eta), f(a+h) = f(a) + \frac{\partial f}{\partial h}(a) + o_{h \rightarrow 0}(\|h\|_n)$$

$\frac{\partial f}{\partial \cdot}(a)$  is named differential of  $f$  on  $a$ .

The notation  $\mathcal{D}(E, F)$  means the set of differentiable applications from  $E$  to  $F$ .

**Proposition 1.1.** Let  $E \subseteq \mathbb{R}^n$ ,  $F \subseteq \mathbb{R}^m$ ,  $f \in \mathcal{D}(E, F)$  and  $a \in \mathring{E}$ . Then  $\frac{\partial f}{\partial \cdot}(a)$  is unique and  $\mathcal{D}(E, F) \subset \mathcal{C}(E, F)$ .

**Proof.** Suppose  $\phi_1$  and  $\phi_2$  two differentials of  $f$  on  $a$ .

$$\begin{aligned} & \exists \eta \in \mathbb{R}^{+*}, \forall h \in \mathcal{B}(0_{\mathbb{R}^n}, \eta), \phi_2(h) - \phi_1(h) = o_{h \rightarrow 0}(\|h\|_n) \\ \implies_{def} & \forall \epsilon \in \mathbb{R}^{+*}, \exists \eta \in \mathbb{R}^{+*}, \forall h \in \mathcal{B}(0_{\mathbb{R}^n}, \eta), \|\phi_2(h) - \phi_1(h)\|_m \leq 2 * \|h\|_n * \epsilon \\ \implies_{\phi_2 - \phi_1 \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)} & \forall \epsilon \in \mathbb{R}^{+*}, \forall h \in \mathbb{R}^n, \|\phi_2(h) - \phi_1(h)\|_m \leq 2 * \|h\|_n * \epsilon \\ \implies_{\epsilon \rightarrow 0} & \forall h \in \mathbb{R}^n, \phi_2(h) = \phi_1(h) \end{aligned}$$

Let  $f \in \mathcal{D}(E, F)$ . and  $a \in \mathring{E}$ .

$$\begin{aligned} \frac{\partial f}{\partial \cdot}(a) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m) & \implies \frac{\partial f}{\partial \cdot}(a) \in \mathcal{C}(\mathbb{R}^n, \mathbb{R}^m), \frac{\partial f}{\partial 0_{\mathbb{R}^n}}(a) = 0_{\mathbb{R}^m} \\ & \xRightarrow{(1)} f(a+h) \xrightarrow{h \rightarrow 0} f(a) \end{aligned}$$

□

**Definition 1.2.** Let  $E \subseteq \mathbb{R}^n$ ,  $F \subseteq \mathbb{R}^m$  and  $f = (f_1 \dots f_m) \in \mathcal{D}(E, F)$ . Then  $f_i$  is differentiable on  $E$  for all  $i \in \llbracket 1, m \rrbracket$ . The jacobian is defined as

$$\begin{aligned} \mathcal{J}_f & : \mathring{E} \longrightarrow \mathcal{M}_{m,n} \\ a & \longmapsto \left[ \frac{\partial f}{\partial e_1}(a) \quad \dots \quad \frac{\partial f}{\partial e_n}(a) \right] = \begin{bmatrix} \frac{\partial f_1}{\partial e_1}(a) & \dots & \frac{\partial f_1}{\partial e_n}(a) \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial e_1}(a) & \dots & \frac{\partial f_m}{\partial e_n}(a) \end{bmatrix} \end{aligned} \quad (2)$$

$(e_i)_{i \in \llbracket 1, n \rrbracket}$  means the matrices  $e_i = \begin{bmatrix} 0 & \dots & 1 & \dots & 0 \end{bmatrix}$  at column  $i$  corresponding to  $\mathbb{R}^n$  standard basis.

$\frac{\partial f}{\partial e_i}$  is named the partial derivative of  $f$  according the  $i^{th}$  variable.

The jacobian is also named gradient when  $m = 1$  and is noted as  $\nabla_f = \mathcal{J}_f$ .

The jacobian is also named derivative when  $m = 1$  with  $n = 1$  and is noted as  $f' = \nabla_f = \mathcal{J}_f$ .

**Proof.** Let  $i \in \llbracket 1, m \rrbracket$ ,  $a \in \mathring{E}$ .

$$\exists \eta \in \mathbb{R}^{+*}, \forall h \in \mathcal{B}(0_{\mathbb{R}^n}, \eta),$$

$$\begin{aligned}
f_i(a+h) &= f(a+h)_i \\
&= f(a)_i + \frac{\partial f}{\partial h}(a)_i + o_{h \rightarrow 0}(\|h\|_n)_i \\
&= f_i(a) + \frac{\partial f}{\partial h}(a)_i + o_{h \rightarrow 0}(\|h\|_n)_i \\
\frac{\partial f}{\partial \cdot}(a)_i \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}) &\xRightarrow{\text{prop 1.1}} \frac{\partial f_i}{\partial \cdot}(a) = \frac{\partial f}{\partial \cdot}(a)_i
\end{aligned}$$

□

**Corollary.** Let  $E \subseteq \mathbb{R}^n$ ,  $F \subseteq \mathbb{R}^m$  and  $f \in \mathcal{D}(E, F)$ . The jacobian of  $f$  on  $a \in \mathring{E}$  fixed is the canonical associated matrix to the differential of  $f$  on  $a$ .

**Notes:** It means a function differentiability can also be proved by exhibiting its jacobian.

**Proof.** Let  $a \in \mathring{E}$ .  $\frac{\partial f}{\partial \cdot}(a) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)$  and any linear application in finite dimension with values in  $\mathbb{R}$  has an unique associated matrix in the standard basis called canonical associated matrix.

□

**Proposition 1.2.** Let  $E \subseteq \mathbb{R}^n$ ,  $F \subseteq \mathbb{R}^m$ ,  $f \in \mathcal{D}(E, F)$  and  $g \in \mathcal{D}(E, F)$ . Then  $g + f \in \mathcal{D}(E, F)$  and

$$\begin{aligned}
\mathcal{J}_{g+f} &: \mathring{E} \longrightarrow F \\
a &\longmapsto \mathcal{J}_g(a) + \mathcal{J}_f(a)
\end{aligned} \tag{3}$$

**Proof.** Let  $a \in \mathring{E}$ .

$$\exists \eta \in \mathbb{R}^{+*}, \forall h \in \mathcal{B}(0_{\mathbb{R}^n}, \eta),$$

$$\begin{aligned}
(g+f)(a+h) &= g(a+h) + f(a+h) \\
&\stackrel{(1)}{=} g(a) + f(a) + \frac{\partial g}{\partial h}(a) + \frac{\partial f}{\partial h}(a) + 2 * o_{h \rightarrow 0}(\|h\|_n) \\
&= (g+f)(a) + \frac{\partial g}{\partial h}(a) + \frac{\partial f}{\partial h}(a) + o_{h \rightarrow 0}(\|h\|_n) \\
\frac{\partial g}{\partial \cdot}(a) + \frac{\partial f}{\partial \cdot}(a) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m) &\xRightarrow{\text{prop 1.1}} g+f \in \mathcal{D}(E, F), \frac{\partial(g+f)}{\partial \cdot}(a) = \frac{\partial g}{\partial \cdot}(a) + \frac{\partial f}{\partial \cdot}(a) \\
&\xRightarrow{\text{mat}} \mathcal{J}_{g+f}(a) = \mathcal{J}_g(a) + \mathcal{J}_f(a)
\end{aligned}$$

**Note:** *mat* indicates in canonical associated matrix way.

□

**Notation 7.** Let  $f \in \mathcal{F}(E, F)$  and  $g \in \mathcal{F}(F, G)$ . Then the notation  $g \circ f$  means the application

$$\begin{aligned}
g \circ f &: E \longrightarrow G \\
x &\longmapsto g(f(x))
\end{aligned}$$

Let  $f_i \in \mathcal{F}(E_i, E_{i+1})$  for  $i \in \llbracket 1, n \rrbracket$ . Then the notation  $\bigcirc_{i=1}^n f_i$  means the application

$$\begin{aligned}
\bigcirc_{i=1}^n f_i &: E_1 \longrightarrow E_{n+1} \\
x &\longmapsto f_n(\dots f_2(f_1(x)))
\end{aligned}$$

**Theorem 1.3.** Let  $E \subseteq \mathbb{R}^n$ ,  $F \subseteq \mathbb{R}^m$ ,  $G \subseteq \mathbb{R}^p$ ,  $f \in \mathcal{D}(E, F)$  and  $g \in \mathcal{D}(F, G)$ . Then  $g \circ f \in \mathcal{D}(E, G)$  and

$$\begin{aligned} \mathcal{J}_{g \circ f} &: \mathring{E} \longrightarrow G \\ a &\longmapsto \mathcal{J}_g(f(a)) * \mathcal{J}_f(a) \end{aligned} \quad (4)$$

**Note:** This theorem is named the chain rule.

**Proof.** Let  $a \in \mathring{E}$ .

$$\begin{aligned} \exists \eta \in \mathbb{R}^{+*}, \forall h \in \mathcal{B}(0_{\mathbb{R}^n}, \eta), \\ (g \circ f)(a + h) &= g(f(a) + \frac{\partial f}{\partial h}(a) + o_{h \rightarrow 0}(\|h\|_n)) \\ &= g(f(a)) + \frac{\partial g}{\partial(\frac{\partial f}{\partial h}(a) + o_{h \rightarrow 0}(\|h\|_n))}(f(a)) + o_{h \rightarrow 0}(\left\| \frac{\partial f}{\partial h}(a) + o_{h \rightarrow 0}(\|h\|_n) \right\|_n) \\ &\stackrel{\frac{\partial f}{\partial \cdot}(a) \in \mathcal{C}(\mathbb{R}^n, \mathbb{R}^m), \frac{\partial f}{\partial 0_{\mathbb{R}^n}}(a) = 0_{\mathbb{R}^m}}{=} g(f(a)) + \frac{\partial g}{\partial(\frac{\partial f}{\partial h}(a) + o_{h \rightarrow 0}(\|h\|_n))}(f(a)) + o_{h \rightarrow 0}(\|h\|_n) \\ &\stackrel{\frac{\partial g}{\partial \cdot}(a) \in \mathcal{L}(\mathbb{R}^m, \mathbb{R}^p)}{=} g(f(a)) + \frac{\partial g}{\partial(\frac{\partial f}{\partial h}(a))}(f(a)) + \frac{\partial g}{\partial(o_{h \rightarrow 0}(\|h\|_n))}(f(a)) + o_{h \rightarrow 0}(\|h\|_n) \\ &\stackrel{\frac{\partial g}{\partial \cdot}(a) \in \mathcal{C}(\mathbb{R}^m, \mathbb{R}^p), \frac{\partial g}{\partial 0_{\mathbb{R}^m}}(a) = 0_{\mathbb{R}^p}}{=} g(f(a)) + \frac{\partial g}{\partial(\frac{\partial f}{\partial h}(a))}(f(a)) + o_{h \rightarrow 0}(\|h\|_n) \\ \frac{\partial g}{\partial(\frac{\partial f}{\partial \cdot}(a))}(f(a)) &= \frac{\partial g}{\partial \cdot}(f(a)) \circ \frac{\partial f}{\partial \cdot}(a) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^p) \xrightarrow{\text{prop 1.1}} g \circ f \in \mathcal{D}(E, G), \frac{\partial(g \circ f)}{\partial \cdot}(a) = \frac{\partial g}{\partial \cdot}(f(a)) \circ \frac{\partial f}{\partial \cdot}(a) \\ &\xRightarrow{\text{mat}} \mathcal{J}_{g \circ f}(a) = \mathcal{J}_g(f(a)) * \mathcal{J}_f(a) \end{aligned}$$

□

### 1.3 Others

*Notation 8.* Let  $E$  and  $F$  two sets and  $(E_i)_{i \in [1, n]}$   $n$  sets.

The notation  $E \times F$  means the cartesian product between  $E$  and  $F$ .

The notation  $\bigcirc_{i=1}^n E_i$  means the cartesian product  $E_n \times \dots \times E_1$ .

*Notation 9.* The notation  $\delta_{\cdot, \cdot}$  means the kronecker delta application

$$\begin{aligned} \delta_{\cdot, \cdot} &: \mathbb{Z} \times \mathbb{Z} \longrightarrow \{0, 1\} \\ (i, j) &\longmapsto \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases} \end{aligned}$$

*Notation 10.* Let  $E \subseteq \mathbb{R}^n$ . The notation  $\mathbb{1}_E$  means the indicator function of  $E$  on  $\mathbb{R}^n$ .

$$\begin{aligned} \mathbb{1}_E &: E \longrightarrow \{0, 1\}^n \\ x &\longmapsto \begin{cases} 1 & x \in E \\ 0 & x \notin E \end{cases} \end{aligned}$$

*Notation 11.* The notation  $\max(0, \cdot)$  means the application

$$\begin{aligned} \max(0, \cdot) &: \mathbb{R} \longrightarrow \mathbb{R}^+ \\ x &\longmapsto \begin{cases} x & x > 0 \\ 0 & x \leq 0 \end{cases} \end{aligned}$$

*Assumption 1.*  $\max(0, \cdot) \in \mathcal{D}(\mathbb{R}, \mathbb{R}^+)$  with

$$\begin{aligned} \max(0, \cdot)' &: \mathbb{R} \longrightarrow \mathbb{R}^+ \\ x &\longmapsto \mathbb{1}_{\mathbb{R}^{++}}(x) \end{aligned}$$

**Note:**  $\max(0, \cdot)$  is actually not *differentiable* on 0 and the notation  $\mathbb{R}^*$  means  $\mathbb{R}_{\setminus \{0\}}$ .

*Notation 12.* Let  $f$  an application with  $n$  inputs and  $m$  outputs.

$$\begin{aligned} f &: E_1 \times \dots \times E_n \longrightarrow F_1 \times \dots \times F_m \\ (x_1, \dots, x_n) &\longmapsto f(x_1, \dots, x_n) \end{aligned}$$

Let  $k \in \llbracket 1, n \rrbracket$ . The notation  $f(x_1, \dots, x_{k-1}, \cdot, x_{k+1}, \dots, x_n)$  means

$$\begin{aligned} f(x_1, \dots, x_{k-1}, \cdot, x_{k+1}, \dots, x_n) &: E_k \longrightarrow F_1 \times \dots \times F_m \\ x_k &\longmapsto f(x_1, \dots, x_{k-1}, x_k, x_{k+1}, \dots, x_n) \end{aligned}$$

## 2 Activation functions

*Notation 13.* Let  $E \subseteq \mathbb{R}^m \times (\mathbb{R}^*)^p$  ( $p$  parameter vectors of any sizes) and  $F \subseteq \mathbb{R}^m$ .

The notation  $\mathcal{F}_{act}(E, F)$  means the set of activation functions from  $E$  to  $F$ .

**Note:** An activation function is an application defined in this section.

**Definition 2.1.** Let the activation function *ReLU* noted as  $\mathcal{R}$  be

$$\begin{aligned} \mathcal{R} &: \mathbb{R}^m \longrightarrow \mathbb{R}^m \\ z &\longmapsto \begin{bmatrix} \max(0, z_1) \\ \vdots \\ \max(0, z_m) \end{bmatrix} \end{aligned}$$

**Proposition 2.1.**  $\mathcal{R} = (\mathcal{R}_1 \dots \mathcal{R}_m) \in \mathcal{D}(\mathbb{R}^m, \mathbb{R}^m)$  and its jacobian is

$$\begin{aligned} \mathcal{J}_{\mathcal{R}} &: \mathbb{R}^m \longrightarrow \mathcal{M}_{m,m} \\ z &\longmapsto \begin{bmatrix} \mathbb{1}_{\mathbb{R}^{++}}(z_1) & 0 & \dots & 0 \\ 0 & \mathbb{1}_{\mathbb{R}^{++}}(z_2) & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \mathbb{1}_{\mathbb{R}^{++}}(z_m) \end{bmatrix} \end{aligned} \tag{5}$$

**Proof.** Let  $i \in \llbracket 1, m \rrbracket$ ,  $j \in \llbracket 1, m \rrbracket$  and  $z \in \mathbb{R}^m$ .

$$\mathcal{R}_i(z) = \max(0, z_i) \xRightarrow{\text{assump 1}} \frac{\partial \mathcal{R}_i}{\partial e_j}(z) = \begin{cases} \mathbb{1}_{\mathbb{R}^{++}}(z_i) & i = j \\ 0 & i \neq j \end{cases}$$

□

**Definition 2.2.** Let the activation function *Softmax* noted as  $\mathcal{S}$  be

$$\begin{aligned} \mathcal{S} &: \mathbb{R}^m \longrightarrow ]0, 1[^m \\ z &\longmapsto \begin{bmatrix} \frac{e^{z_1}}{\sum_{k=1}^m e^{z_k}} \\ \vdots \\ \frac{e^{z_m}}{\sum_{k=1}^m e^{z_k}} \end{bmatrix} \end{aligned}$$

**Proposition 2.2.**  $\mathcal{S} = (\mathcal{S}_1 \dots \mathcal{S}_m) \in \mathcal{D}(\mathbb{R}^m, ]0, 1[^m)$  and its jacobian is

$$\begin{aligned} \mathcal{J}_{\mathcal{S}} &: \mathbb{R}^m \longrightarrow \mathcal{M}_{m,m} \\ z &\longmapsto \begin{bmatrix} \mathcal{S}_1 * (1 - \mathcal{S}_1) & -\mathcal{S}_1 * \mathcal{S}_2 & \cdots & -\mathcal{S}_1 * \mathcal{S}_m \\ -\mathcal{S}_2 * \mathcal{S}_1 & \mathcal{S}_2 * (1 - \mathcal{S}_2) & \ddots & \vdots \\ \vdots & \ddots & \ddots & -\mathcal{S}_{m-1} * \mathcal{S}_m \\ -\mathcal{S}_m * \mathcal{S}_1 & \cdots & -\mathcal{S}_m * \mathcal{S}_{m-1} & \mathcal{S}_m * (1 - \mathcal{S}_m) \end{bmatrix} (z) \end{aligned} \quad (6)$$

**Proof.** Let  $i \in \llbracket 1, m \rrbracket$ ,  $j \in \llbracket 1, m \rrbracket$  and  $z \in \mathbb{R}^m$ .

$$\begin{aligned} \mathcal{S}_i(z) &= \frac{e^{z_i}}{\sum_{k=1}^m e^{z_k}} \\ \implies \frac{\partial \mathcal{S}_i}{\partial e_j}(z) &= \frac{(\delta_{i,j} * e^{z_i}) * \sum_{k=1}^m e^{z_k} - e^{z_j} * e^{z_i}}{(\sum_{k=1}^m e^{z_k})^2} \\ &= \delta_{i,j} * \mathcal{S}_i(z) - \mathcal{S}_j(z) * \mathcal{S}_i(z) \\ &= \mathcal{S}_i(z) * (\delta_{i,j} - \mathcal{S}_j(z)) \end{aligned}$$

□

### 3 Loss

*Notation 14.* Let  $E \subseteq \mathbb{R}^m \times \mathbb{R}^m$ ,  $F \subseteq \mathbb{R}$ .

The notation  $\mathcal{F}_{loss}(E, F)$  means the set of loss functions from  $E$  to  $F$ .

**Note:** A loss function is an application defined in this section.

**Definition 3.1.** Let the loss function *Categorical cross-entropy* noted as  $\xi$  be

$$\begin{aligned} \xi &: ]0, 1[^m \times \{0, 1\}^m \longrightarrow \mathbb{R} \\ (y, y^*) &\longmapsto -\sum_{k=1}^m y_k^* * \log(y_k) \end{aligned}$$

**Proposition 3.1.** Let  $y^* \in \{0, 1\}^m$ .  $\xi(\cdot, y^*) \in \mathcal{D}(]0, 1[^m, \mathbb{R})$  and its gradient is

$$\begin{aligned} \nabla_{\xi(\cdot, y^*)} &: ]0, 1[^m \longrightarrow \mathbb{R} \\ y &\longmapsto -\begin{bmatrix} \frac{y_1^*}{y_1} & \cdots & \frac{y_m^*}{y_m} \end{bmatrix} \end{aligned} \quad (7)$$

**Proof.** Let  $j \in \llbracket 1, m \rrbracket$  and  $y \in ]0, 1[^m$ .

$$\xi(y, y^*) = -\sum_{k=1}^m y_k^* * \log(y_k) \implies \frac{\partial \xi(\cdot, y^*)}{\partial e_j}(y) = -\frac{y_j^*}{y_j}$$

□

## 4 Layers

*Notation 15.* Let  $E \subseteq \mathbb{R}^n \times (\mathbb{R}^p)^p$  ( $p$  parameter vectors of any sizes) and  $F \subseteq \mathbb{R}^m$ .

The notation  $\mathcal{F}_{\text{layer}}(E, F)$  means the set of layer functions from  $E$  to  $F$ .

**Note:** A layer function is an application defined in this section.

**Definition 4.1.** Let the layer function *Dense layer* noted as  $\mathbb{L}$  be

$$\begin{aligned} \mathbb{L} &: \mathbb{R}^n \times \mathcal{M}_{m,n} \times \mathbb{R}^m \longrightarrow \mathbb{R}^m \\ (y, W, b) &\longmapsto y * W^T + b \end{aligned}$$

**Note:**  $\mathbb{L} : \mathbb{R}^n \times \mathcal{M}_{m,n} \times \mathbb{R}^m \longrightarrow \mathbb{R}^m$  is equivalent to  $\mathbb{L} : \mathbb{R}^n \times (\mathbb{R}^n)^m \times \mathbb{R}^m \longrightarrow \mathbb{R}^m$ .

**Proposition 4.1.** Let  $W \in \mathcal{M}_{m,n}$  and  $b \in \mathbb{R}^m$ .  $\mathbb{L}(\cdot, W, b) = (\mathbb{L}_1(\cdot, W, b) \dots \mathbb{L}_m(\cdot, W, b)) \in \mathcal{D}(\mathbb{R}^n, \mathbb{R}^m)$  and its gradient is

$$\begin{aligned} \mathcal{J}_{\mathbb{L}(\cdot, W, b)} &: \mathbb{R}^n \longrightarrow \mathcal{M}_{m,n} \\ y &\longmapsto W \end{aligned} \tag{8}$$

**Proof.** Let  $i \in \llbracket 1, m \rrbracket$ ,  $j \in \llbracket 1, n \rrbracket$  and  $y \in \mathbb{R}^n$ .

$$\begin{aligned} \mathbb{L}(y, W, b) = y * W^T + b &\implies \mathbb{L}_i(y, W, b) = y * W_{i,:}^T + b_i \\ &\implies \frac{\partial \mathbb{L}_i(\cdot, W, b)}{\partial e_j}(y) = W_{i,j} \end{aligned}$$

□

**Proposition 4.2.** Let  $y \in \mathbb{R}^n$ ,  $(w^{(k)})_{k \in \llbracket 1, m-1 \rrbracket} \in (\mathbb{R}^n)^{m-1}$ ,  $b \in \mathbb{R}^m$ .

$\forall i^* \in \llbracket 1, m \rrbracket$ ,  $\mathbb{L}(y, w^{(1)}, \dots, w^{(i^*-1)}, \cdot, w^{(i^*)}, \dots, w^{(m-1)}, b) \in \mathcal{D}(\mathbb{R}^n, \mathbb{R}^m)$  and jacobians are

$$\forall i^* \in \llbracket 1, m \rrbracket,$$

$$\begin{aligned} \mathcal{J}_{\mathbb{L}(y, w^{(1)}, \dots, w^{(i^*-1)}, \cdot, w^{(i^*)}, \dots, w^{(m-1)}, b)} &: \mathbb{R}^n \longrightarrow \mathcal{M}^{m,n} \\ w &\longmapsto \begin{bmatrix} (0) \\ y_1 & \cdots & y_n \\ (0) \end{bmatrix} \text{ at row } i^* \end{aligned} \tag{9}$$

**Note:** For  $i^* = 1$  and  $i^* = m$ , the applications  $\mathbb{L}(y, \cdot, w^{(1)}, \dots, w^{(m-1)}, b)$  and  $\mathbb{L}(y, w^{(1)}, \dots, w^{(m-1)}, \cdot, b)$  are meant respectively.

**Proof.** Let  $i^* \in \llbracket 1, m \rrbracket$ ,  $i \in \llbracket 1, m \rrbracket$ ,  $j \in \llbracket 1, n \rrbracket$  and  $w \in \mathbb{R}^n$ .

$$\begin{aligned} &\mathbb{L}(y, w^{(1)}, \dots, w^{(i^*-1)}, \cdot, w^{(i^*)}, \dots, w^{(m-1)}, b)(w) \\ &= \begin{bmatrix} y * w^{(1)T} & \cdots & y * w^{(i^*-1)T} & y * w^T & y * w^{(i^*)T} & \cdots & y * w^{(m-1)T} \end{bmatrix} + b \\ &\implies \frac{\partial \mathbb{L}_i(y, w^{(1)}, \dots, w^{(i^*-1)}, \cdot, w^{(i^*)}, \dots, w^{(m-1)}, b)}{\partial e_j}(w) = \begin{matrix} y_j & i = i^* \\ 0 & i \neq i^* \end{matrix} \end{aligned}$$

□



**Proposition 4.3.** Let  $y \in \mathbb{R}^n$  and  $W \in \mathcal{M}_{m,n}$ .  $\mathbb{L}(y, W, \cdot) = (\mathbb{L}_1(y, W, \cdot) \dots \mathbb{L}_m(y, W, \cdot)) \in \mathcal{D}(\mathbb{R}^m, \mathbb{R}^m)$  and its gradient is

$$\begin{aligned} \mathcal{J}_{\mathbb{L}(y, W, \cdot)} &: \mathbb{R}^m \longrightarrow \mathcal{M}_{m,m} \\ b &\longmapsto I_m \end{aligned} \quad (10)$$

**Proof.** Let  $i \in \llbracket 1, m \rrbracket$ ,  $j \in \llbracket 1, n \rrbracket$  and  $b \in \mathbb{R}^m$ .

$$\begin{aligned} \mathbb{L}(y, W, b) = y * W^T + b &\implies \mathbb{L}_i(y, W, b) = y * W_{i,:}^T + b_i \\ &\implies \frac{\partial \mathbb{L}_i(y, W, \cdot)}{\partial e_j}(b) = \delta_{i,j} \end{aligned}$$

□

## 5 Neural network

### 5.1 Simplified jacobian matrices

**Proposition 5.1.** Let  $\mathcal{F}^{(upstream)} \in \mathcal{D}(\mathbb{R}^m, \mathbb{R})$  and  $\mathcal{R} \in \mathcal{F}_{act}(\mathbb{R}^m, \mathbb{R}^m)$ .  $\mathcal{F}^{(upstream)} \circ \mathcal{R} \in \mathcal{D}(\mathbb{R}^m, \mathbb{R})$  and its gradient is

$$\begin{aligned} \nabla_{\mathcal{F}^{(upstream)} \circ \mathcal{R}} &: \mathbb{R}^m \longrightarrow \mathbb{R}^m \\ z &\longmapsto \left[ \nabla_{\mathcal{F}^{(upstream)}}(y)_1 * \mathbb{1}_{\mathbb{R}^{++}}(z_1) \quad \dots \quad \nabla_{\mathcal{F}^{(upstream)}}(y)_m * \mathbb{1}_{\mathbb{R}^{++}}(z_m) \right] \end{aligned} \quad (11)$$

where

$$y = \mathcal{R}(z)$$

**Note:** It means such a gradient can be implemented without matrix multiplication.

**Proof.** Let  $j \in \llbracket 1, m \rrbracket$  and  $z \in \mathbb{R}^m$ .

$$\nabla_{\mathcal{F}^{(upstream)} \circ \mathcal{R}} \stackrel{(4)}{=} \nabla_{\mathcal{F}^{(upstream)}}(\mathcal{R}(z)) * \mathcal{J}_{\mathcal{R}}(z) \stackrel{(5)}{\implies} \frac{\partial \mathcal{F}^{(upstream)} \circ \mathcal{R}}{\partial e_j}(z) = \nabla_{\mathcal{F}^{(upstream)}}(\mathcal{R}(z))_j * \mathbb{1}_{\mathbb{R}^{++}}(z_j)$$

□

**Proposition 5.2.** Let  $\mathcal{F}^{(upstream)} \in \mathcal{D}(\mathbb{R}^m, \mathbb{R})$ ,  $\mathbb{L} \in \mathcal{F}_{layer}(\mathbb{R}^n \times (\mathbb{R}^n)^m \times \mathbb{R}^m, \mathbb{R}^m)$ ,  $y \in \mathbb{R}^n$ ,  $(w^{(k)})_{k \in \llbracket 1, m-1 \rrbracket} \in (\mathbb{R}^n)^{m-1}$  and  $b \in \mathbb{R}^m$ .

$\forall i^* \in \llbracket 1, m \rrbracket$ ,  $\mathcal{F}^{(upstream)} \circ \mathbb{L}(y, w^{(1)}, \dots, w^{(i^*-1)}, \underset{\text{at index } i^*}{\cdot}, w^{(i^*)}, \dots, w^{(m-1)}, b) \in \mathcal{D}(\mathbb{R}^n, \mathbb{R})$  and gradients are

$$\forall i^* \in \llbracket 1, m \rrbracket,$$

$$\begin{aligned} \nabla_{\mathcal{F}^{(upstream)} \circ \mathbb{L}(y, w^{(1)}, \dots, w^{(i^*-1)}, \cdot, w^{(i^*)}, \dots, w^{(m-1)}, b)} &: \mathbb{R}^n \longrightarrow \mathbb{R}^n \\ w &\longmapsto \nabla_{\mathcal{F}^{(upstream)}}(z)_{i^*} * y \end{aligned} \quad (12)$$

where

$$z = \mathbb{L}(y, w^{(1)}, \dots, w^{(i^*-1)}, w, w^{(i^*)}, \dots, w^{(m-1)}, b)$$

**Note:** It means these gradients for  $i^* \in \llbracket 1, m \rrbracket$  can be implemented with  $\nabla_{\mathcal{F}^{(upstream)}}(z)^T * y$ .

**Proof.** Let  $i^* \in \llbracket 1, m \rrbracket$ ,  $j \in \llbracket 1, n \rrbracket$  and  $w \in \mathbb{R}^n$ . Let  $z = \mathbb{L}(y, w^{(1)}, \dots, w^{(i^*-1)}, w, w^{(i^*)}, \dots, w^{(m-1)}, b)$ .

$$\begin{aligned} \nabla_{\mathcal{F}^{(upstream)} \circ \mathbb{L}(y, w^{(1)}, \dots, w^{(i^*-1)}, \cdot, w^{(i^*)}, \dots, w^{(m-1)}, b)}(w) &= \nabla_{\mathcal{F}^{(upstream)}}(z) * \mathcal{J}_{\mathbb{L}(y, w^{(1)}, \dots, w^{(i^*-1)}, \cdot, w^{(i^*)}, \dots, w^{(m-1)}, b)}(w) \\ &\stackrel{(9)}{\implies} \frac{\partial \nabla_{\mathcal{F}^{(upstream)} \circ \mathbb{L}(y, w^{(1)}, \dots, w^{(i^*-1)}, \cdot, w^{(i^*)}, \dots, w^{(m-1)}, b)}}{\partial e_j}(w) = \nabla_{\mathcal{F}^{(upstream)}}(z)_{i^*} * y_j \end{aligned}$$

□

**Proposition 5.3.** Let  $\mathcal{S} \in \mathcal{F}_{act}(\mathbb{R}^m, ]0, 1[^m)$  and  $\xi \in \mathcal{F}_{loss}(]0, 1[^m \times \{0, 1\}^m, \mathbb{R})$ .

Let  $y^* \in \{0, 1\}^m$  with  $\|y^*\|_m = 1$ .  $\xi(\cdot, y^*) \circ \mathcal{S} \in \mathcal{D}(\mathbb{R}^m, \mathbb{R})$  and its gradient is

$$\begin{aligned} \nabla_{\xi(\cdot, y^*) \circ \mathcal{S}} &: \mathbb{R}^m \longrightarrow \mathbb{R} \\ z &\longmapsto \mathcal{S}(z) - y^* \end{aligned} \tag{13}$$

**Proof.** Let  $j \in \llbracket 1, m \rrbracket$  and  $z \in \mathbb{R}^m$ .

$$\begin{aligned} \nabla_{\xi(\cdot, y^*) \circ \mathcal{S}}(z) &= \nabla_{\xi(\cdot, y^*)}(\mathcal{S}(z)) * \mathcal{J}_{\mathcal{S}}(z) \\ \implies \frac{\partial \xi(\cdot, y^*) \circ \mathcal{S}}{\partial e_j}(z) &\stackrel{(6),(7)}{=} -y_j^* + \sum_{k=1}^m y_k^* * \mathcal{S}_j(z) \\ &\stackrel{y^* \in \{0, 1\}^m, \|y\|_m = 1}{=} \mathcal{S}_j(z) - y_j^* \end{aligned}$$

□

## 5.2 Definitions

*Notation 16.* Let  $E \subseteq \mathbb{R}^n \times (\mathbb{R})^p$  ( $p$  parameter vectors of any sizes) and  $F \subseteq \mathbb{R}^m$ .

The notation  $\mathcal{F}_{net}(E, F)$  means the set of neural network functions from  $E$  to  $F$ .

**Note:** A neural network is an application defined in this section.

**Definition 5.1.** Let  $(m_k)_{k \in [0, p]} \in (\mathbb{N}^*)^p$ . Let the neural network *Multi-class dense neural network* noted as  $\mathcal{N}_{c^+}$  be

$$\begin{aligned} \mathcal{N}_{c^+} &: \mathbb{R}^{m_0} \times \left( \bigtimes_{k=1}^p \mathcal{M}_{m_k, m_{k-1}} \right) \times \left( \bigtimes_{k=1}^p \mathbb{R}^{m_k} \right) \longrightarrow \mathbb{R}^{m_p} \\ (x, (W^{(k)})_{k \in [1, p]}, (b^{(k)})_{k \in [1, p]}) &\longmapsto (\mathcal{S} \circ \mathbb{L}^{(p)}(\cdot, W^{(p)}, b^{(p)})) \circ \left( \bigcirc_{k=1}^{p-1} \mathcal{R}^{(k)} \circ \mathbb{L}^{(k)}(\cdot, W^{(k)}, b^{(k)}) \right)(x) \end{aligned}$$

where

$$\begin{aligned} (\mathbb{L}^{(k)})_{k \in [1, p]} &\in \bigtimes_{k=1}^p \mathcal{F}_{layer}(\mathbb{R}^{m_{k-1}} \times \mathcal{M}_{m_k, m_{k-1}} \times \mathbb{R}^{m_k}, \mathbb{R}^{m_k}) \\ (\mathcal{R}^{(k)})_{k \in [1, p-1]} &\in \bigtimes_{k=1}^{p-1} \mathcal{F}_{act}(\mathbb{R}^{m_k}, \mathbb{R}^{m_k}) \\ \mathcal{S} &\in \mathcal{F}_{act}(\mathbb{R}^{m_p}, ]0, 1[^{m_p}) \end{aligned}$$

**Note:**  $\mathcal{N}_{c^+} : \mathbb{R}^{m_0} \times \left( \bigtimes_{k=1}^p \mathcal{M}_{m_k, m_{k-1}} \right) \times \left( \bigtimes_{k=1}^p \mathbb{R}^{m_k} \right) \longrightarrow \mathbb{R}^{m_p}$  is equivalent to  $\mathcal{N}_{c^+} : \mathbb{R}^{m_0} \times \left( \bigtimes_{k=1}^p (\mathbb{R}^{m_{k-1}})^{m_k} \right) \times \left( \bigtimes_{k=1}^p \mathbb{R}^{m_k} \right) \longrightarrow \mathbb{R}^{m_p}$ .

**Corollary.**  $\mathcal{N}_{c^+} \in \mathcal{D}(\mathbb{R}^{m_0} \times \left( \bigtimes_{k=1}^p (\mathbb{R}^{m_{k-1}})^{m_k} \right) \times \left( \bigtimes_{k=1}^p \mathbb{R}^{m_k} \right), \mathbb{R}^{m_p})$  and the total number of parameter is

$$\sum_{k=1}^p m_k * (m_{k-1} + 1)$$

**Proof.**  $\mathcal{N}_{c^+}$  is a composition of *differentiable* applications so it is *differentiable* by the **chain rule** theorem. Let  $a \in (\prod_{k=1}^p (\mathbb{R}^{m_{k-1}})^{m_k}) \times (\prod_{k=1}^p \mathbb{R}^{m_k})$  then  $a$  has  $\sum_{k=1}^p m_k * m_{k-1} + \sum_{k=1}^p m_k$  coefficients. □

**Definition 5.2.** Let  $(m_k)_{k \in [0, p]}$ ,  $\mathcal{N}_{c^+} \in \mathcal{F}_{net}(\mathbb{R}^{m_0} \times (\prod_{k=1}^p \mathcal{M}_{m_k, m_{k-1}}) \times (\prod_{k=1}^p \mathbb{R}^{m_k}), \mathbb{R}^{m_p})$ ,  $X = (x^{(i)})_{i \in [1, n]} \in (\mathbb{R}^{m_0})^n$  and  $Y^* = (y^{*(i)})_{i \in [1, n]} \in (\{0, 1\}^{m_p})^n$  with  $\forall i \in [1, n]$ ,  $\|y^{*(i)}\|_{m_p} = 1$ .

Let the *Multi-class optimization problem* noted as  $(\mathcal{P}_{c^+})$  be

$$(\mathcal{P}_{c^+}) : \min_{(W^{(k)})_{k \in [1, p]}, (b^{(k)})_{k \in [1, p]}} \sum_{i=1}^n \xi(y^{(i)}, y^{*(i)}) \quad (14)$$

where

$$\begin{aligned} \xi &\in \mathcal{F}_{loss}([0, 1]^{m_p \times \{0, 1\}^{m_p}}, \mathbb{R}) \\ y^{(i)} &= \mathcal{N}_{c^+}(x^{(i)}, (W^{(k)})_{k \in [1, p]}, (b^{(k)})_{k \in [1, p]}) \end{aligned}$$

$\sum_{i=1}^n \xi(\cdot, y^{*(i)}) \circ \mathcal{N}_{c^+}(x^{(i)}, \cdot, \cdot)$  is named the objective function and will be noted as  $\mathcal{O}_{c^+}(X, Y^*, \cdot, \cdot)$ .

**Corollary.**  $\mathcal{O}_{c^+}(X, Y^*, \cdot, \cdot) \in \mathcal{D}((\prod_{k=1}^p \mathcal{M}_{m_k, m_{k-1}}) \times (\prod_{k=1}^p \mathbb{R}^{m_k}), \mathbb{R})$ .

**Note:** Its gradients for each variable can be computed recursively through each composition using (1.2), (4), (8), (12), (10), (11) and (13).

**Proof.**  $\mathcal{O}_{c^+}(X, Y^*, \cdot, \cdot)$  is a sum and composition of *differentiable* applications so it is *differentiable* by the proposition 1.2 and **chain rule** theorem. □

## 6 Gradient descent

### 6.1 Optimization fundamentals

**Definition 6.1.** Let  $f \in \mathcal{D}(\mathbb{R}^m, \mathbb{R})$ .  $f$  *convex* means

$$\begin{aligned} \forall x \in \mathbb{R}^m, \forall y \in \mathbb{R}^m, \\ \forall \tau \in [0, 1], f(\tau * x + (1 - \tau) * y) \leq \tau * f(x) + (1 - \tau) * f(y) \end{aligned} \quad (15)$$

**Proposition 6.1.** Let  $f \in \mathcal{D}(\mathbb{R}^m, \mathbb{R})$ .  $f$  *convex* is equivalent to

$$\begin{aligned} \forall x \in \mathbb{R}^m, \forall y \in \mathbb{R}^m, \\ f(x) + \nabla f(x) * (y - x)^T \leq f(y) \end{aligned} \quad (16)$$

**Proof.** Suppose  $f$  *convex* (15). Let  $x \in \mathbb{R}^m$  and  $y \in \mathbb{R}^m$ .

$$\begin{aligned} \exists \eta \in \mathbb{R}^{+*}, \forall \tau \in [-\eta, \eta], \\ f(x + \tau * (y - x)) \underset{\|\cdot\|_m \in \mathcal{C}(\mathbb{R}^m, \mathbb{R}), (1)}{=} f(x) + \frac{\partial f}{\partial(\tau * (y - x))}(x) + \underset{\tau \rightarrow 0}{o}(\|\tau * (y - x)\|_m) \end{aligned}$$

$$\begin{aligned}
& \exists \eta \in \mathbb{R}^{+*}, \forall \tau \in [-\eta, \eta], \\
& \xRightarrow{(15)} f(x) + \frac{\partial f}{\partial(\tau * (y-x))}(x) + \underset{\tau \rightarrow 0}{o}(\|\tau * (y-x)\|_m) \leq \tau * f(y) + (1-\tau) * f(x) \\
& \xRightarrow{\frac{\partial f}{\partial} \in \mathcal{L}(\mathbb{R}^m, \mathbb{R})} \exists \eta \in \mathbb{R}^{+*}, \forall \tau \in [-\eta, \eta], \frac{\partial f}{\partial(y-x)}(x) + \underset{\tau \rightarrow 0}{o}(\|y-x\|_m) \leq f(y) - f(x) \\
& \xRightarrow{\tau \rightarrow 0} \frac{\partial f}{\partial(y-x)}(x) \leq f(y) - f(x) \\
& \xRightarrow{mat} f(x) + \nabla_f(x) * (y-x)^T \leq f(y)
\end{aligned}$$

Suppose (16). Let  $x \in \mathbb{R}^m$ ,  $y \in \mathbb{R}^m$  and  $\tau \in [0, 1]$ . Let  $z = \tau * x + (1-\tau) * y$ .

$$\begin{aligned}
(a): \quad f(z) - (1-\tau) * \nabla_f(z) * (y-x)^T & \underset{(16)}{=} f(z) + \nabla_f(z) * (x-z)^T \leq f(x) \\
(b): \quad f(y) + \tau * \nabla_f(z) * (y-x)^T & \underset{(16)}{=} f(z) + \nabla_f(z) * (y-z)^T \leq f(y) \\
& \xRightarrow{\tau * (a) + (1-\tau) * (b)} f(z) \leq \tau * f(x) + (1-\tau) * f(y)
\end{aligned}$$

□

**Definition 6.2.** Let  $f \in \mathcal{D}(\mathbb{R}^m, \mathbb{R})$  and  $L \in \mathbb{R}^{+*}$ .  $f$   $L$ -smooth means

$$\begin{aligned}
& \forall x \in \mathbb{R}^m, \forall y \in \mathbb{R}^m, \\
& \|\nabla_f(x) - \nabla_f(y)\|_m \leq L * \|x - y\|_m
\end{aligned} \tag{17}$$

**Proposition 6.2.** Let  $f \in \mathcal{D}(\mathbb{R}^m, \mathbb{R})$  and  $L \in \mathbb{R}^{+*}$ . If  $f$   $L$ -smooth then

$$\begin{aligned}
& \forall x \in \mathbb{R}^m, \forall y \in \mathbb{R}^m, \\
& f(y) \leq f(x) + \nabla_f(x) * (y-x)^T + \frac{L}{2} * \|y-x\|_m^2
\end{aligned} \tag{18}$$

**Proof.** Let  $x \in \mathbb{R}^m$  and  $y \in \mathbb{R}^m$ . Let

$$\begin{aligned}
g & : [0, 1] \longrightarrow \mathbb{R}^m \\
\tau & \longmapsto x + \tau * (y-x)
\end{aligned}$$

$$\begin{aligned}
\forall \tau \in [0, 1], (f \circ g)(\tau) = f(x + \tau * (y-x)) & \xRightarrow{(4)} \forall \tau \in [0, 1], (f \circ g)'(\tau) = \nabla_f(g(\tau)) * (y-x)^T \\
& \xRightarrow{f} f(y) - f(x) = \int_0^1 \nabla_f(g(\tau)) * (y-x)^T d\tau
\end{aligned}$$

$$\begin{aligned}
f(y) &= f(x) + \int_0^1 \nabla_f(g(\tau)) * (y-x)^T d\tau \\
&= f(x) + \nabla_f(x) * (y-x)^T + \int_0^1 (\nabla_f(g(\tau)) - \nabla_f(x)) * (y-x)^T d\tau \\
&\stackrel{\text{Cauchy-Schwarz}}{\leq} f(x) + \nabla_f(x) * (y-x)^T + \int_0^1 \|\nabla_f(g(\tau)) - \nabla_f(x)\|_m * \|y-x\|_m d\tau \\
&\stackrel{(17)}{\leq} f(x) + \nabla_f(x) * (y-x)^T + L * \|y-x\|_m^2 * \int_0^1 \tau d\tau
\end{aligned}$$

□

**Proposition 6.3.** Let  $L \in \mathbb{R}^{+*}$  and  $f \in \mathcal{D}(\mathbb{R}^m, \mathbb{R})$ . If  $f$  convex and  $L$ -smooth then

$$\begin{aligned}
& \forall x \in \mathbb{R}^m, \forall y \in \mathbb{R}^m, \\
& \frac{1}{L} * \|\nabla_f(y) - \nabla_f(x)\|_m^2 \leq (\nabla_f(y) - \nabla_f(x)) * (y-x)^T
\end{aligned} \tag{19}$$

**Notes:** This proposition is named the gradient co-coercivity.

**Proof.** Let  $x \in \mathbb{R}^m$ ,  $y \in \mathbb{R}^m$  and  $z = x - \frac{1}{L}(\nabla_f(y) - \nabla_f(x))$ .

$$\begin{aligned} f(y) - f(x) &= f(y) - f(z) + f(z) - f(x) \\ &\stackrel{(16),(18)}{\leq} \nabla_f(y) * (y - z)^T + \nabla_f(x) * (z - x) + \frac{L}{2} * \|z - x\|_m^2 \\ &\leq \nabla_f(y) * (y - x)^T - \frac{1}{2L} * \|\nabla_f(x) - \nabla_f(y)\|_m^2 \end{aligned}$$

The inequality is true for all  $x \in \mathbb{R}^m$  and  $y \in \mathbb{R}^m$ .

Let  $x \in \mathbb{R}^m$  and  $y \in \mathbb{R}^m$ . The previous inequality gives

$$\begin{aligned} (a): \quad f(y) - f(x) &\leq \nabla_f(y) * (y - x)^T - \frac{1}{2L} * \|\nabla_f(x) - \nabla_f(y)\|_m^2 \\ (b): \quad f(x) - f(y) &\leq \nabla_f(x) * (x - y)^T - \frac{1}{2L} * \|\nabla_f(y) - \nabla_f(x)\|_m^2 \\ \stackrel{(a)+(b)}{\implies} 0 &\leq (\nabla_f(y) - \nabla_f(x)) * (y - x)^T - \frac{1}{L} * \|\nabla_f(y) - \nabla_f(x)\|_m^2 \end{aligned}$$

□

**Definition 6.3.** Let  $f \in \mathcal{F}(\mathbb{R}^m, \mathbb{R})$  and  $x^* \in \mathbb{R}^m$ . *a global minimum* of  $f$  means

$$\forall x \in \mathbb{R}^m, f(x^*) \leq f(x) \quad (20)$$

**Proposition 6.4.** Let  $f \in \mathcal{D}(\mathbb{R}^m, \mathbb{R})$  and  $x^* \in \mathbb{R}^m$ . If  $x^*$  *global minimum* of  $f$  then

$$\nabla_f(x^*) = 0_{\mathbb{R}^m} \quad (21)$$

**Proof.** Let  $x^*$  *global minimum* of  $f$ ,  $v \in \mathbb{R}^m$  and

$$\begin{aligned} g &: \mathbb{R} \longrightarrow \mathbb{R}^m \\ \tau &\longmapsto x^* + \tau * v \end{aligned}$$

$$\begin{aligned} \forall \tau \in \mathbb{R}, (f \circ g)(\tau) = f(x^* + \tau * v) &\implies \begin{aligned} (f \circ g)'(0) &\stackrel{(4)}{=} \nabla_f(x^*) * v^T \\ \forall \tau \in \mathbb{R}, (f \circ g)(0) &\stackrel{(20)}{\leq} (f \circ g)(\tau) \end{aligned} \end{aligned}$$

$$\forall \tau \in \mathbb{R}, (f \circ g)(0) \leq (f \circ g)(\tau)$$

$$\begin{aligned} \implies \exists \eta \in \mathbb{R}^{+*}, \forall \tau \in ]0, \eta], &\quad \begin{aligned} 0 &\leq (f \circ g)'(0) * \tau + \underset{\tau \rightarrow 0}{o}(\tau) \\ 0 &\leq (f \circ g)'(0) * (-\tau) + \underset{\tau \rightarrow 0}{o}(\tau) \end{aligned} \\ \implies \exists \eta \in \mathbb{R}^{+*}, \forall \tau \in ]0, \eta], &\quad \underset{\tau \rightarrow 0}{o}(\tau) \leq (f \circ g)'(0) \leq \underset{\tau \rightarrow 0}{o}(\tau) \\ \implies \underset{\tau \rightarrow 0}{\nabla_f(x^*) * v^T} = (f \circ g)'(0) &= 0 \end{aligned}$$

The equality is true for all  $v \in \mathbb{R}^m$  in particular for the vectors  $(e_i)_{i \in [1, n]}$  corresponding to  $\mathbb{R}^m$  standard basis thus

$$\nabla_f(x^*) = 0_{\mathbb{R}^m}$$

□

## 6.2 Algorithms

*Notation 17.* The notation  $\mathbb{R}^{\mathbb{N}}$  means the set of numerical sequences.

The notation  $\mathcal{G}_{descent}$  means the set of gradient descent sequences.

**Note:**  $\mathcal{G}_{descent} \subset \mathbb{R}^{\mathbb{N}}$  and a gradient descent sequence is a numerical sequence defined in this section.

**Definition 6.4.** Let  $f \in \mathcal{D}(\mathbb{R}^m, \mathbb{R})$  and  $\alpha \in \mathbb{R}^{+*}$ . Let the *Gradient descent* be the numerical sequence

$$(x_v^{(n)})_{n \in \mathbb{N}} = \begin{cases} x_v^{(0)} \in \mathbb{R}^m & n = 0 \\ x_v^{(n+1)} = x_v^{(n)} - \alpha \nabla f(x_v^{(n)}) & n \in \mathbb{N}^* \end{cases} \quad (22)$$

$\alpha$  is named the learning rate.

$n \in \mathbb{N}$  fixed is named an epoch.

**Note:** For the *Multi-class classification problem*  $(\mathcal{P}_{c^+})$ ,  $\mathcal{O}_{c^+}(X, Y^*, \cdot, \cdot)$  with all parameters fixed but one  $W_{i^*, \cdot}^{(k)}$  or  $b^{(k)}$  will be  $f$  and  $W_{i^*, \cdot}^{(k)}$  or  $b^{(k)}$  will be the numerical sequence  $(x_v^{(n)})_{n \in \mathbb{N}}$ . This is only for one parameter  $W^{(k)}$  or  $b^{(k)}$  but all the parameters are actually optimized simultaneously. It means for each iteration  $n$  all  $W^{(k)}$  or  $b^{(k)}$  receive an update.

**Note:** The *Stochastic gradient descent* numerical sequence is similar to the *Gradient descent* sequence but instead of having  $\nabla f(x_v^{(n)})$  it uses  $\nabla f(x_v^{(n)} | \mathcal{S})$  an estimate of the actual gradient. For the *Multi-class classification problem*  $(\mathcal{P}_{c^+})$  with  $\mathcal{O}_{c^+}(X, Y^*, \cdot, \cdot)$ , the computations can be intensives with large matrices  $X$  and  $Y^*$ . For this reason, several sub-samples (of size  $b = 32$  generally) of  $X$  and  $Y^*$  called *batches* are used when computing the gradient estimate. More precisely, a *batch* results from a sampling without replacement by couple of rows  $X$  and  $Y^*$ . The last batch can be of size inferior or equal to  $b$ .  $\nabla f(x_v^{(n)} | \mathcal{S})$  the estimate is equal to the mean of the gradients of each batch.

**Proposition 6.5.** Let  $\alpha \in \mathbb{R}^{+*}$ ,  $L \in \mathbb{R}^{+*}$ ,  $f \in \mathcal{D}(\mathbb{R}^m, \mathbb{R})$  and  $(x_v^{(n)})_{n \in \mathbb{N}} \in \mathcal{G}_{descent}$ . If  $f$  convex,  $L$ -smooth, admits  $x_v^*$  as a *global minimum* and  $\alpha < \frac{2}{L}$  then

$$\forall n \in \mathbb{N}^*, f(x_v^{(n)}) - f(x_v^*) \leq \frac{\|x_v^{(0)} - x_v^*\|_m^2}{nC}$$

where

$$C = \alpha - \frac{L\alpha^2}{2}$$

**Note:** It means  $(x_v^{(n)})_{n \in \mathbb{N}}$  converge to the global minimum  $x_v^*$  with a rate of  $\frac{o}{n \rightarrow 0}(n^{-1})$ . In case of  $\mathcal{O}_{c^+}(X, Y^*, \cdot, \cdot)$ , *Stochastic gradient descent* convergence is still an active research subject nowadays [5].

**Proof.** Let  $f$  convex,  $L$ -smooth, admits a *global minimum*  $x_v^*$  and  $\alpha < \frac{2}{L}$  and  $n \in \mathbb{N}$ .

$$\begin{aligned} \|x_v^{(n+1)} - x_v^*\|_m^2 &= \|x_v^{(n)} - x_v^* - \alpha \nabla f(x_v^{(n)})\|_m^2 \\ &= \|x_v^{(n)} - x_v^*\|_m^2 - 2\alpha(x_v^{(n)} - x_v^*)^T \nabla f(x_v^{(n)}) + \alpha^2 \|\nabla f(x_v^{(n)})\|_m^2 \\ &\stackrel{(19)}{\leq} \|x_v^{(n)} - x_v^*\|_m^2 - \underbrace{\left(\frac{2\alpha}{L} - \alpha^2\right)}_{\in \mathbb{R}^{+*}} \|\nabla f(x_v^{(n)})\|_m^2 \\ &\leq_{rec} \|x_v^{(0)} - x_v^*\|_m^2 \end{aligned}$$

It also means  $(\|x_v^{(n)} - x_v^*\|_m)_{n \in \mathbb{N}}$  is a *decreasing numerical sequence*.

**Note:** *rec* means by applying recursively.

$$\begin{aligned}
f(x_v^{(n+1)}) &\stackrel{(22),(18)}{\leq} f(x_v^{(n)}) + \nabla f(x_v^{(n)}) * (-\alpha \nabla f(x_v^{(n)}))^T + \frac{L}{2} \left\| -\alpha \nabla f(x_v^{(n)}) \right\|_m^2 \\
&\leq f(x_v^{(n)}) - \underbrace{\left( \alpha - \frac{L\alpha^2}{2} \right)}_{\in \mathbb{R}^{+*}} \left\| \nabla f(x_v^{(n)}) \right\|^2 \\
\Rightarrow (a): \quad 0 &\stackrel{(20)}{\leq} f(x_v^{(n+1)}) - f(x_v^*) \leq f(x_v^{(n)}) - f(x_v^*) - C \left\| \nabla f(x_v^{(n)}) \right\|^2
\end{aligned}$$

Let  $\forall n \in \mathbb{N}$ ,  $\delta^{(n)} = f(x_v^{(n)}) - f(x_v^*)$ . It also means  $(\delta^{(n)})_{n \in \mathbb{N}}$  is a *decreasing numerical sequence* with a lower bound of 0.

If  $\forall k \in \llbracket 0, n+1 \rrbracket$ ,  $\delta^{(k)} \neq 0$  and  $\|x_v^{(0)} - x_v^*\|_m \neq 0$  then

$$\begin{aligned}
f(x_v^{(n)}) - f(x_v^*) &\stackrel{(16)}{\leq} \nabla f(x_v^{(n)}) * (x_v^{(n)} - x_v^*)^T \\
&\stackrel{\text{Cauchy-Schwarz}}{\leq} \left\| \nabla f(x_v^{(n)}) \right\|_m \left\| x_v^{(n)} - x_v^* \right\|_m \\
&\leq \left\| \nabla f(x_v^{(n)}) \right\|_m \left\| x_v^{(0)} - x_v^* \right\|_m \\
\Rightarrow_{(a)} \delta^{(n+1)} &\leq \delta^{(n)} - \frac{C}{\|x_v^{(0)} - x_v^*\|_m^2} * \delta^{(n)^2} \\
\Rightarrow \frac{C}{\|x_v^{(0)} - x_v^*\|_m^2} &\leq \frac{C}{\|x_v^{(0)} - x_v^*\|_m^2} \frac{\delta^{(n)}}{\delta^{(n+1)}} \leq \frac{1}{\delta^{(n+1)}} - \frac{1}{\delta^{(n)}} \\
\Rightarrow_{\Sigma_0^n, tel} \frac{(n+1)C}{\|x_v^{(0)} - x_v^*\|_m^2} &\leq \frac{1}{\delta^{(n+1)}} - \frac{1}{\delta^{(0)}} \stackrel{(20)}{\leq} \frac{1}{\delta^{(n+1)}}
\end{aligned}$$

**Note:**  $\Sigma_0^n$  means a sum from 0 to  $n$  and *tel* means telescopic cancellation.

Else  $\exists k \in \llbracket 0, n+1 \rrbracket$ ,  $\delta^{(k)} = 0$  or  $\|x_v^{(0)} - x_v^*\|_m = 0$ . Let

$$r = \begin{cases} \min\{k \in \llbracket 0, n+1 \rrbracket \mid \delta^{(k)} = 0\} & \left\| x_v^{(0)} - x_v^* \right\|_m \neq 0 \\ 0 & \left\| x_v^{(0)} - x_v^* \right\|_m = 0 \end{cases}$$

If  $r \neq 0$  the exact same reasoning can be done on  $\llbracket 0, r-1 \rrbracket$  to obtain the inequality. For the rest from  $r$  to  $n+1$  the inequality is also true because  $\forall k \in \llbracket r, n+1 \rrbracket, \delta^{(k)} = 0$

□

**Definition 6.5.** Let  $f \in \mathcal{D}(\mathbb{R}^m, \mathbb{R})$ ,  $\alpha \in \mathbb{R}^{+*}$  and  $\beta \in [0, 1]$ . Let the *Gradient descent with momentum* be the numerical sequence

$$(x_m^{(n)})_{n \in \mathbb{N}} = \begin{cases} x_m^{(0)} \in \mathbb{R}^m & n = 0 \\ x_m^{(n+1)} = x_m^{(n)} - \alpha m^{(n)} & n \in \mathbb{N}^* \end{cases} \quad (23)$$

where

$$(m^{(n)})_{n \in \mathbb{N}} = \begin{cases} 0_{\mathbb{R}^m} & n = 0 \\ m^{(n+1)} = \beta m^{(n)} + (1 - \beta) \nabla f(x_m^{(n)}) & n \in \mathbb{N}^* \end{cases}$$

$(m^{(n)})_{n \in \mathbb{N}}$  is named the momentum.

**Note:** In case of  $\mathcal{O}_{c^+}(X, Y^*, \cdot, \cdot)$ , an estimate  $\nabla_f(x_v^{(n)}|\mathbb{S})$  of the actual gradient is mostly used. This defines the *Stochastic gradient descent with momentum* numerical sequence. Its convergence is still an active research subject nowadays [4].

## References

- [1] Goh, G. (2017). Why momentum really works. *Distill*.
- [2] Gower, R. M. (2019). Convergence theorems for gradient descent.
- [3] Kinsley, H. and Kukiela, D. (2020). Neural networks from scratch in python.
- [4] Liu, Y., Gao, Y., and Yin, W. (2020). An improved analysis of stochastic gradient descent with momentum.
- [5] Nguyen, L. M., Nguyen, P. H., Richtárik, P., Scheinberg, K., Takáč, M., and van Dijk, M. (2019). New convergence aspects of stochastic gradient algorithms.