

# A deep mathematical understanding of DNNs

Jiang J.

Data Engineer, Data Scientist

ENSEEIH Computer Science Engineering Degree

INP Toulouse Dual MSc Research Degree in AI, Big Data and Ops

France

## 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, 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. These propositions will be mathematically proven. The subject used as reference is a multi-class classification problem with – dense layers, *ReLU* and *SoftMax* activation layers, Categorical cross-entropy loss and Stochastic gradient descent optimizer. But all the elements below can be easily re-used or re-defined to cover regressions.

## 1 – Fundamentals and Notations

### 1.1 – Matrices

**Convention** – All sets considered are non empty.

**Convention** – A vector is a matrix with only one row. Thus, the real vector set  $\mathbb{R}^m$  is equivalent to  $M_{1,m}$ .

**Notation** – 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  $n \times m$  will be 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

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

$$\forall j \in \llbracket 1, m \rrbracket, A_{:,j} = \begin{bmatrix} a_{1,j} \\ a_{2,j} \\ \vdots \\ a_{n,j} \end{bmatrix}$$

$$\forall i \in \llbracket 1, n \rrbracket, A_{i,:} = [a_{i,1} \quad a_{i,2} \quad \cdots \quad a_{i,n}]$$

The notation  $M_{n,m}$  means the matrix set of dimension  $n \times m$  with coefficients in  $\mathfrak{R}$ .

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

**Notation** – Let  $A \in M_{n,m}$ , and  $B \in M_{m,p}$ , and let the product noted  $A \times B$  or  $AB$  be

$$C = A \times B = AB$$

where  $C$  is in  $M_{n,p}$  with

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

**Notation** – Let  $a \in \mathfrak{R}$  and  $B \in M_{n,m}$ . Let the scalar wise product noted as  $a \times B$  be

$$C = a \times B = B \times a$$

where  $C$  is in  $M_{n,m}$  with

$$\forall i \in \llbracket 1, n \rrbracket, \forall j \in \llbracket 1, m \rrbracket, C_{i,j} = a \times B_{i,j}$$

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

**Notation** – Let  $a \in \mathfrak{R}^n$  and  $b \in \mathfrak{R}^n$ . Let the scalar product on  $\mathfrak{R}^n$  between two vectors noted as  ${}_n(a|b)$  be

$${}_n(a|b) = a \times b^T = b \times a^T = {}_n(b|a)$$

Let  $c \in \mathfrak{R}^n$ . Let the Euclidean norm on  $\mathfrak{R}^n$  noted as  ${}_n\|c\|$  be

$${}_n\|c\| = \sqrt{(c|c)} = \sqrt{c \times c^T}$$

## 1.2 – Differential calculus

**Notation** – Let  $E \subseteq \mathbb{R}^n$  and  $F \subseteq \mathbb{R}^m$  .

The notation  $\overset{\circ}{E}$  means the set with only the interior point of  $E$  .

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

The notation  $\zeta(E, F)$  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  $F$  .

**Definition** – Let  $E \subseteq \mathbb{R}^n$  and  $F \subseteq \mathbb{R}^m$  . Let  $f: E \rightarrow F$  . Then  $f$  differentiable on  $E$  is equivalent to

$$\forall a \in \overset{\circ}{E} , \exists \frac{\partial f}{\partial \cdot}(a): \begin{cases} \mathbb{R}^n \rightarrow \mathbb{R}^m \\ h \mapsto \frac{\partial f}{\partial h}(a) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m) \end{cases} ,$$

$$\forall h \in \mathbb{R}^n , f(a+h) = f(a) + \frac{\partial f}{\partial h}(a) \cdot h + o_{h \rightarrow 0}(\|h\|)$$

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

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

**Proposition** – Let  $E \subseteq \mathbb{R}^n$  ,  $F \subseteq \mathbb{R}^m$  and  $f \in D(E, F)$  . Then the differential  $\frac{\partial f}{\partial \cdot}(a)$  is unique and  $D(E, F) \subset \zeta(E, F)$  .

Proof: Suppose  $\phi_1$  and  $\phi_2$  two differential of  $f$  on  $a$  . Then

$$\forall h \in \mathbb{R}^n , \phi_2(h) - \phi_1(h) = o_{h \rightarrow 0}(2 \times \|h\|)$$

$$\text{so } \forall \epsilon > 0 , \exists \eta > 0 , \forall h \in \mathbb{R}^n , \|h\| \leq \eta \Rightarrow \|\phi_2(h) - \phi_1(h)\| \leq 2 \times \|h\| \times \epsilon$$

$$\phi_2 - \phi_1 \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)$$

$$\text{so } \forall \epsilon > 0 , \forall h \in \mathbb{R}^n , \|\phi_2(h) - \phi_1(h)\| < 2 \times \|h\| \times \epsilon$$

Then Squeeze theorem with  $\epsilon \rightarrow 0$  gives the differential uniqueness

$$\forall h \in \mathbb{R}^n, \phi_2(h) = \phi_1(h)$$

Let  $a \in E$ . Then

$$\frac{\partial f}{\partial \cdot}(a) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)$$

$$\text{so } \frac{\partial f}{\partial \cdot}(a) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m) \text{ and } \frac{\partial f}{\partial 0_{\mathbb{R}^n}}(a) = 0_{\mathbb{R}^m}$$

Then  $f$  differentiable definition gives the continuity

$$f(a+h) = f(a) + \frac{\partial f}{\partial h}(a) \cdot h + o(\|h\|) \xrightarrow{h \rightarrow 0} f(a)$$

**Definition & Proposition** – Let  $E \subseteq \mathbb{R}^n$ ,  $F \subseteq \mathbb{R}^m$  and  $f = (f_1 \dots f_m) \in D(E, F)$ . Then  $f_i$  is differentiable on  $E$  for all  $i \in \llbracket 1, m \rrbracket$  and its Jacobian matrix is noted as the application

$$J_f: \begin{cases} E \xrightarrow{o} M_{m,n} \\ a \mapsto \left[ \frac{\partial f}{\partial e_1}(a) \quad \dots \quad \frac{\partial f}{\partial e_n}(a) \right] = \begin{bmatrix} \frac{\partial f_1}{\partial e_1}(a) & \dots & \frac{\partial f_1}{\partial e_n}(a) \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial e_1}(a) & \dots & \frac{\partial f_m}{\partial e_n}(a) \end{bmatrix} \end{cases}$$

$(e_i)_{i \in \llbracket 1, n \rrbracket}$  means the matrices  $e_i = [0 \dots 0 \underset{\text{at index } i}{1} 0 \dots 0] \in \mathbb{R}^n$  corresponding to the  $\mathbb{R}^n$  standard basis.

The Jacobian matrix of  $f$  on  $a$  fixed is also the canonical matrix associated to the differential of  $f$  on  $a$ .

$\frac{\partial f}{\partial e_i}(a)$  is named the partial derivative according the  $i^{\text{th}}$  variable.

The Jacobian is also named Gradient when  $m=1$  and will be noted as  $\nabla_f = J_f$ .

Proof: Let  $f = (f_1 \dots f_m) \in D(E, F)$ . Then

$$\forall i \in \llbracket 1, m \rrbracket, \forall h \in \mathbb{R}^n,$$

$$\begin{aligned} f_i(a+h) &= f(a+h)_i \\ &= f(a)_i + \frac{\partial f}{\partial h}(a)_i + o_n(\|h\|) \\ &= f_i(a) + \frac{\partial f_i}{\partial h}(a) + o_n(\|h\|) \end{aligned}$$

with  $\forall i \in \llbracket 1, m \rrbracket, \frac{\partial f_i}{\partial \cdot}(a) = \frac{\partial f}{\partial \cdot}(a)_i \in \mathcal{L}(\mathbb{R}^n, \mathbb{R})$

**Notation** – Let  $E \subseteq \mathbb{R}^n, F \subseteq \mathbb{R}^m, G \subseteq \mathbb{R}^p, f: E \rightarrow F$  and  $g: F \rightarrow G$ . Then the notation  $g \circ f$  means the application  $g \circ f: \begin{cases} E \rightarrow F \\ x \mapsto g(f(x)) \end{cases}$ .

Let  $(f_i)_{i \in \llbracket 1, n \rrbracket}$  with  $f_i: E_i \rightarrow E_{i+1}$  for  $i \in \llbracket 1, n \rrbracket$ . Then the notation  $\overset{n}{\circ} f_i$  means the application  $\overset{n}{\circ} f_i: \begin{cases} E_1 \rightarrow E_{n+1} \\ x \mapsto f_n(f_{n-1}(\dots f_2(f_1(x)))) \end{cases}$ .

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

$$J_{g \circ f}: \begin{cases} \overset{o}{E} \rightarrow M_{p,n} \\ a \mapsto J_g(f(a)) \times J_f(a) \end{cases}$$

Proof: Let  $E \subseteq \mathbb{R}^n, F \subseteq \mathbb{R}^m, G \subseteq \mathbb{R}^p, f \in D(E, F)$  and  $g \in D(F, G)$ . Let  $a \in \overset{o}{E}$ . Then

$$\forall h \in \mathbb{R}^n,$$

$$\begin{aligned} g \circ f(a+h) &= g(f(a) + \frac{\partial f}{\partial h}(a) + o_n(\|h\|)) \\ &= g(f(a)) + \frac{\partial g}{\partial (\frac{\partial f}{\partial h}(a) + o_n(\|h\|))}(f(a)) + o_n(\|\frac{\partial f}{\partial h}(a) + o_n(\|h\|)\|) \end{aligned}$$

Then  $\frac{\partial g}{\partial \cdot}(a) \in \mathcal{L}(\mathbb{R}^m, \mathbb{R}^p), \frac{\partial g}{\partial \cdot}(a) \in \xi(\mathbb{R}^m, \mathbb{R}^p), \frac{\partial g}{\partial 0_{\mathbb{R}^m}}(a) = 0_{\mathbb{R}^p}, \frac{\partial f}{\partial \cdot}(a) \in \xi(\mathbb{R}^n, \mathbb{R}^m)$  and

$$\frac{\partial f}{\partial 0_{\mathbb{R}^n}}(a) = 0_{\mathbb{R}^m} \text{ gives}$$

$$\forall h \in \mathbb{R}^n ,$$

$$g \circ f(a+h) = g(f(a)) + \frac{\partial g}{\partial \left(\frac{\partial f}{\partial h}(a)\right)}(f(a)) + o_{h \rightarrow 0}(\|h\|)$$

$$\text{It means } g \circ f \in D(E, G) \text{ and } \frac{\partial g \circ f}{\partial \cdot}(a) = \frac{\partial g}{\partial \cdot}(f(a)) \circ \frac{\partial f}{\partial \cdot}(a) .$$

It also gives with canonical associated matrices

$$\begin{aligned} \forall h \in \mathbb{R}^n , \quad J_{g \circ f}(a) \times h &= J_g(f(a)) \times (J_f(a) \times h) = (J_g(f(a)) \times J_f(a)) \times h \\ \text{so } J_{g \circ f}(a) &= J_g(f(a)) \times J_f(a) \end{aligned}$$

### 1.3 – Function convexity and smoothness

**Definition** – Let  $E \subseteq \mathbb{R}^n$  convex,  $F \subseteq \mathbb{R}$  and  $f \in \zeta(E, F)$ . Then  $f$  convex is equivalent to

$$\begin{aligned} \forall (x, y) \in E^2 , \quad \forall t \in [0, 1] , \\ f(t \times x + (1-t) \times y) \leq t \times f(x) + (1-t) \times f(y) \end{aligned}$$

**Proposition** – Let  $E \subseteq \mathbb{R}^n$  convex,  $F \subseteq \mathbb{R}$  and  $f \in D(E, F)$  convex. Then  $f$  convex is equivalent to

$$\begin{aligned} \forall (x, y) \in E^2 , \\ f(x) + \nabla_f(x) \times (y - x)^T \leq f(y) \end{aligned}$$

Proof: TO DO.

**Proposition** – Let  $\Omega \neq \emptyset$  convex of  $\mathbb{R}^m$  and  $f: x \mapsto f(x)$  such as  $f \in D(\Omega, \mathbb{R})$  and convex. Then

$$\exists X^* \subset \Omega \setminus \{\emptyset\} , \quad \forall x^* \in X^* , \quad f: x \mapsto f(x^*) \leq f(x) \quad \text{<TODO existing and that's all>}$$

**Definition** – Let  $\Omega \neq \emptyset$  subset of  $\mathbb{R}^m$ ,  $\|\cdot\|_m$  a norm on  $\mathbb{R}^m$ , and  $f: x \mapsto f(x)$  such as  $f \in D(\Omega, \mathbb{R})$ . Let  $L > 0$ . Then  $f$   $L$ -smooth on  $\Omega$  is equivalent to

$$\forall (y, z) \in \Omega^2, \quad \left\| \frac{df}{dx}(y) - \frac{df}{dx}(z) \right\| \leq L \times_m \|y - z\|$$

**Proposition** – Let  $\Omega \neq \emptyset$  convex of  $\mathbb{R}^m$  and  $f: x \mapsto f(x)$  such as  $f \in D(\Omega, \mathbb{R})$ . Then  $f$  first order integral form Taylor expansion is

$$\forall (y, z) \in \Omega^2, \\ f(y) = f(z) + \int_0^1 \frac{df}{dx}(z + \tau(y - z))(y - z)^T d\tau$$

Proof: TO DO.

**Proposition** – Let  $\Omega \neq \emptyset$  convex of  $\mathbb{R}^m$ ,  $L > 0$ , and  $f: x \mapsto f(x)$  such as  $f \in D(\Omega, \mathbb{R})$  plus  $L$ -smooth on  $\Omega$ . Then

$$\forall (y, z) \in \Omega^2, \\ f(z) \leq f(y) + \frac{df}{dx}(y) \times (z - y)^T + \frac{L}{2} \times_m \|z - y\|^2$$

$$\forall x \in \Omega, \\ f\left(x - \frac{1}{L} \times \frac{df}{dx}(x)\right) - f(x) \leq -\frac{1}{2L} \times_m \left\| \frac{df}{dx}(x) \right\|^2$$

Proof: TO DO.

**Proposition** – Let  $\Omega \neq \emptyset$  convex of  $\mathbb{R}^m$ ,  $\|\cdot\|_m$  a norm on  $\mathbb{R}^m$ , and  $f: x \mapsto f(x)$  such as  $f \in D(\Omega, \mathbb{R})$ , convex and  $L$ -smooth. Then  $f$  is co-coercive

$$\forall (y, z) \in \Omega^2, \\ \frac{1}{L} \times_m \left\| \frac{df}{dx}(y) - \frac{df}{dx}(z) \right\|^2 \leq \left( \frac{df}{dx}(y) - \frac{df}{dx}(z) \right) \times (x - y)^T$$

Proof: TO DO.

### 1.3 – Others

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

$$1_E: \begin{cases} E \rightarrow \mathfrak{R}^n \\ x \mapsto \begin{cases} 1 & x \in E \\ 0 & x \notin E \end{cases} \end{cases}$$

**Notation** – Let  $f: \begin{cases} E_1 \times \dots \times E_n \rightarrow F_1 \times \dots \times F_m \\ (x_1, \dots, x_n) \mapsto f(x_1, \dots, x_n) \end{cases}$  an application with  $n$  parameters and  $m$  outputs. Then for  $k \in \llbracket 1, n \rrbracket$  the notation  $f(x_1, \dots, x_{k-1}, \cdot, x_{k+1}, \dots, x_n)$  means the application

$$f(x_1, \dots, x_{k-1}, \cdot, x_{k+1}, \dots, x_n): \begin{cases} \Omega_k \rightarrow \Omega'_1 \times \dots \times \Omega'_m \\ x_k \mapsto f(x_1, \dots, x_{k-1}, x_k, x_{k+1}, \dots, x_n) \end{cases}.$$

## 2 – Activation functions

**Definition** – Let  $F_{act} \in D(\mathfrak{R}^m)$ . Then the vector wise application

$$F_{act}: \begin{cases} \mathfrak{R}^m \rightarrow \mathfrak{R}^m \\ z \mapsto f(z) \end{cases}$$

is an activation function.

**Definition** –  $ReLU$  is the following vector wise application

$$ReLU: \begin{cases} \mathfrak{R}^m \rightarrow \mathfrak{R}^m \\ z \mapsto \max(0, z) \end{cases}$$

with  $\max$  the element-wise maximum operation between two vectors.

**Hypothesis** – The notation  $ReLU_j$  means the application corresponding to the coefficient  $j$  of the function  $ReLU$ . Let  $z \in \mathfrak{R}^m$  then

$$\forall j \in \llbracket 1, m \rrbracket, \quad ReLU_j(z_j) = \max(0, z_j) = ReLU(z)_j$$

$ReLU$  is supposed derivable on every coefficients at 0



$$\forall j \in \llbracket 1, m \rrbracket, \text{ReLU}_j'(0) = 0$$

**Proposition** –  $\text{ReLU}$  is an activation function. Its Jacobian matrix is

$$\frac{d \text{ReLU}}{dz} : \left\{ \begin{array}{l} \mathfrak{R}^m \rightarrow M_{m,m} \\ z \mapsto \begin{bmatrix} 1_{\mathfrak{R}_{\setminus \{0\}}^+(z_1)} & 0 & \cdots & 0 \\ 0 & 1_{\mathfrak{R}_{\setminus \{0\}}^+(z_2)} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1_{\mathfrak{R}_{\setminus \{0\}}^+(z_m)} \end{bmatrix} \end{array} \right.$$

Proof: TO DO.

**Proposition** – The following vector wise application is an activation function

$$\text{SoftMax} : \left\{ \begin{array}{l} \mathfrak{R}^m \rightarrow ]0, 1[^m \\ z \mapsto \frac{e^{z_j}}{\sum_{j'=1}^m e^{z_{j'}}} \end{array} \right.$$

with  $e$  the element-wise exponential operation.

The  $\text{SoftMax}$  function will be denoted as  $S$  for simplicity.

Its Jacobian matrix is

$$\frac{dS}{dz} : \left\{ \begin{array}{l} \mathfrak{R}^m \rightarrow M_{m,m} \\ z \mapsto \frac{dS}{dz}(z) \end{array} \right.$$

where  $\forall z \in \mathfrak{R}^m, \forall (j, j') \in \{1, 2, \dots, m\}^2$ ,

$$\frac{dS}{dz}(z)_{j,j'} = S(z)_j \times (\delta_{j,j'} - S(z)_{j'})$$

with  $\delta_{j,j'}$  the Kronecker delta.

Proof: TO DO.

### 3 – Loss

**Definition** – Let  $\hat{\Omega} \in M_{n,m}$  and  $\Omega \subseteq M_{n,m}$  non empty subsets. Let  $\hat{y} \in \hat{\Omega}$  and  $F_{loss}^{\hat{y}} \in D(\Omega, \mathbb{R})$ . Then  $F_{loss}^{\hat{y}}$  is a loss function is equivalent to the application

$$F_{loss}^{\hat{y}} \circ g: \begin{cases} E \rightarrow \mathbb{R} \\ \epsilon \mapsto (F_{loss}^{\hat{y}} \circ g)(\epsilon) = F_{loss}^{\hat{y}}(\hat{y} + \epsilon) \end{cases}$$

is an increasing function according each coefficient with  $E \subseteq M_{n,m}$  such as  $F_{loss}^{\hat{y}} \circ g$  is always defined.

The  $\hat{y}$  matrix is named the ground truth.

**Proposition** – Let  $\hat{y} \in \{0,1\}^m$  a ground truth matrix. Then the application

$$\xi^{\hat{y}}: \begin{cases} ]0,1[^m \rightarrow \mathbb{R} \\ y \mapsto -\sum_{j=1}^m \hat{y}_j \log(y_j) \end{cases}$$

is a loss function. The application is named Categorical cross-entropy loss.

Its Gradient matrix is

$$\frac{d\xi^{\hat{y}}}{dy}: \begin{cases} ]0,1[^m \rightarrow \mathbb{R}^m \\ y \mapsto -\begin{bmatrix} \frac{\hat{y}_1}{y_1} & \dots & \frac{\hat{y}_m}{y_m} \end{bmatrix} \end{cases}$$

Proof: TO DO.

**Proposition** – Let  $\hat{y} \in \{0,1\}^m$  a ground truth matrix. Let  $S: \mathbb{R}^m \rightarrow ]0,1[^m$  and  $\xi^{\hat{y}}: ]0,1[^m \rightarrow \mathbb{R}$  the *SoftMax* activation and Categorical cross-entropy loss functions. Then  $\xi^{\hat{y}} \circ S: \mathbb{R}^m \rightarrow \mathbb{R}$  is derivable on  $\mathbb{R}^m$  and its Gradient matrix is

$$\frac{d(\xi^{\hat{y}} \circ S)}{dz}: \begin{cases} \mathbb{R}^m \rightarrow \mathbb{R}^m \\ z \mapsto S(z) - \hat{y} \end{cases}$$

Proof: TO DO.

## 4 – Dense layers

**Definition** – The application

$$L: \begin{cases} \mathbb{R}^m \times M_{m',m} \times \mathbb{R}^{m'} \rightarrow \mathbb{R}^{m'} \\ (y, W, b) \mapsto y \times W^T + b \end{cases}$$

defines a dense layer with  $y$  named the input vector,  $W$  named the weight matrix and  $b$  named the bias matrix.

The notation  $L_j$  means the application  $L_j: \mathbb{R}^m \times \mathbb{R}^m \times \mathbb{R} \rightarrow \mathbb{R}$  corresponding to the row  $j$  of the second matrix component of the dense layer  $L$ . Let  $y \in \mathbb{R}^m$  an input vector,  $W \in M_{m',m}$  a weight matrix and  $b \in \mathbb{R}^{m'}$  a bias matrix then

$$\forall j \in \llbracket 1, m' \rrbracket, \quad L_j(y, W_{j,:}, b_j) = y \times (W_{j,:})^T + b_j = L(y, W, b)_j$$

**Proposition** – Let  $L: \mathbb{R}^m \times M_{m',m} \times \mathbb{R}^{m'} \rightarrow \mathbb{R}^{m'}$  a dense layer function. Then  $L$  is derivable according the first and third variables on  $\mathbb{R}^m$  and  $\mathbb{R}^{m'}$  respectively.

Let  $y \in \mathbb{R}^m$  an input vector and  $b \in \mathbb{R}^{m'}$  a bias matrix. Then  $L_{j,:}: \mathbb{R}^m \times \mathbb{R}^m \times \mathbb{R} \rightarrow \mathbb{R}$  is also derivable according the second variable for all  $j \in \llbracket 1, m' \rrbracket$ .

Its Gradient or Jacobian matrices are

$$\begin{aligned} \frac{\partial L}{\partial y} &: \begin{cases} \mathbb{R}^m \rightarrow M_{m',m} \\ y \mapsto W \end{cases} \\ \forall j \in \llbracket 1, m' \rrbracket, \quad \frac{\partial L_j}{\partial w} &: \begin{cases} \mathbb{R}^m \rightarrow \mathbb{R}^m \\ w \mapsto y \end{cases} \\ \frac{\partial L}{\partial b} &: \begin{cases} \mathbb{R}^{m'} \rightarrow M_{m',m'} \\ b \mapsto I_{m'} \end{cases} \end{aligned}$$

with  $I_{m'}$  the identity matrix of size  $m' \times m'$ .

Proof: TO DO.

**Proposition** – Let  $L: \mathfrak{R}^m \times M_{m',m} \times \mathfrak{R}^{m'} \rightarrow \mathfrak{R}^{m'}$  and  $ReLU: \mathfrak{R}^{m'} \rightarrow \mathfrak{R}^{m'}$  the dense layer and

$ReLU$  activation functions. Let  $F^{upstream}: \begin{cases} \mathfrak{R}^{m'} \rightarrow \mathfrak{R} \\ y' \mapsto F^{upstream}(y') \end{cases}$  such as  $F^{upstream} \in D(\mathfrak{R}^{m'}, \mathfrak{R})$ .

Then  $F^{upstream} \circ ReLU \circ L(\cdot, W, b): \mathfrak{R}^m \times M_{m',m} \times \mathfrak{R}^{m'} \rightarrow \mathfrak{R}$  is derivable according the first and third variables on  $\mathfrak{R}^m$  and  $\mathfrak{R}^{m'}$  respectively.

The notation  $F_{j'}^{upstream}$  means the application corresponding to the coefficient  $j'$  of  $F^{upstream}$ .  
Let  $y' \in \mathfrak{R}^{m'}$  then

$$\forall j' \in \llbracket 1, m' \rrbracket, F_{j'}^{upstream}(y'_j) = F^{upstream}(y')_{j'}$$

Let  $y \in \mathfrak{R}^m$  an input vector and  $b \in \mathfrak{R}^{m'}$  a bias matrix. Then

$F_{j'}^{upstream} \circ ReLU_{j'} \circ L_{j'}(\cdot, w, b_{j'}) : \mathfrak{R} \rightarrow \mathfrak{R}$  is also derivable for all  $j' \in \llbracket 1, m' \rrbracket$ .

Its Gradient matrices are

$$\frac{\partial(F^{upstream} \circ ReLU \circ L(\cdot, W, b))}{\partial y} : \begin{cases} \mathfrak{R}^m \rightarrow \mathfrak{R}^m \\ y \mapsto \frac{\partial(F^{upstream} \circ ReLU \circ L(\cdot, W, b))}{\partial y}(y) \end{cases}$$

where  $\forall y \in \mathfrak{R}^m, \forall j \in \llbracket 1, m \rrbracket$ ,

$$\frac{\partial(F^{upstream} \circ ReLU \circ L(\cdot, W, b))}{\partial y}(y)_j = \sum_{j'=1}^{m'} \frac{dF^{upstream}}{dy'}(ReLU(L(y, W, b)))_{j'} \times 1_{\mathfrak{R}_{\setminus \{0\}}^+}(L(y, W, b)_{j'}) \times W_{j', j}$$

with  $W \in M_{m',m}$  a weight matrix,  $b \in \mathfrak{R}^{m'}$  a bias matrix.

$$\forall j' \in \llbracket 1, m' \rrbracket, \frac{\partial(F_{j'}^{upstream} \circ ReLU_{j'} \circ L_{j'}(\cdot, w, b_{j'}))}{\partial w} : \begin{cases} \mathfrak{R}^m \rightarrow \mathfrak{R}^m \\ w \mapsto \frac{\partial(F_{j'}^{upstream} \circ ReLU_{j'} \circ L_{j'}(\cdot, w, b_{j'}))}{\partial w}(w) \end{cases}$$

where  $\forall w \in \mathfrak{R}^m, \forall j \in \llbracket 1, m \rrbracket$ ,

$$\frac{d(F_{j'}^{upstream} \circ ReLU_{j'} \circ L_{j'}(\cdot, w, b_{j'}))}{dw}(w)_j = F_{j'}^{upstream'}(ReLU_{j'}(L_{j'}(y, w, b_{j'}))) \times 1_{\mathfrak{R}_{\setminus \{0\}}^+}(L_{j'}(y, w, b_{j'})) \times y_j$$

with  $y \in \mathfrak{R}^m$  an input matrix,  $b \in \mathfrak{R}^{m'}$  a bias matrix.

$$\frac{\partial (F^{upstream} \circ \text{ReLU} \circ L(\cdot, W, b))}{\partial b} : \begin{cases} \mathfrak{R}^{m'} \rightarrow \mathfrak{R}^{m'} \\ b \mapsto \frac{\partial (F^{upstream} \circ \text{ReLU} \circ L(\cdot, W, b))}{\partial b}(b) \end{cases}$$

where  $\forall b \in \mathfrak{R}^{m'}$ ,  $\forall j' \in \llbracket 1, m' \rrbracket$ ,

$$\frac{\partial (F^{upstream} \circ \text{ReLU} \circ L(\cdot, W, b))}{\partial b}(b)_{j'} = \frac{dF^{upstream}}{dy'}(\text{ReLU}(L(y, W, b)))_{j'} \times 1_{\mathfrak{R}_{\setminus \{0\}}^+}(L(y, W, b)_{j'})$$

with  $y \in \mathfrak{R}^m$  a weight matrix and  $W \in M_{m', m}$  a weight matrix.

Proof: TO DO.

## 5 – Neural Networks

**Definition** – A training data set is defined as couples of vectors  $(X^i, \hat{Y}^i) \in \mathfrak{R}^m \times \mathfrak{R}^l$  for  $i \in \llbracket 1, n \rrbracket$ . The  $X^i$  are named input or feature matrices and the  $\hat{Y}^i$  target or label matrices.

**Definition** – Let  $p$  dense layers with activation functions  $F_{act}^k \circ L^k(\cdot, W^k, b^k) : \mathfrak{R}^{m_k} \rightarrow \mathfrak{R}^{m_{k+1}}$  for  $k \in \llbracket 1, p \rrbracket$  with  $W^k \in M_{m_{k+1}, m_k}$  and  $b^k \in \mathfrak{R}^{m_{k+1}}$  the  $L^k$  weight and bias matrices. Let a training data set  $(X^i, \hat{Y}^i) \in \mathfrak{R}^{m_1} \times \mathfrak{R}^{m_{p+1}}$  for  $i \in \llbracket 1, n \rrbracket$ . Let  $F_{loss}^{\hat{Y}^i} : \mathfrak{R}^{m_{p+1}} \rightarrow \mathfrak{R}$  loss functions for  $i \in \llbracket 1, n \rrbracket$  with  $(\hat{Y}^i)_{i \in \llbracket 1, n \rrbracket}$  as ground truth matrices respectively.

Then a neural network is defined as the application  $N : \begin{cases} \mathfrak{R}^{m_1} \rightarrow \mathfrak{R}^{m_{p+1}} \\ y \mapsto \bigcirc_{k=1}^n (F_{act}^k \circ L^k(\cdot, W^k, b^k))(y) \end{cases}$ .

The optimization problem is  $\min_{(W_{1,:}^k, \dots, W_{m_{k+1},:}^k, b^k)_{k \in \llbracket 1, p \rrbracket}} \sum_{i=1}^n F_{loss}^{\hat{Y}^i}(N(X^i))$  and

$F_{loss}^{global} : ((W_{1,:}^k, \dots, W_{m_{k+1},:}^k, b^k)_{k \in \llbracket 1, p \rrbracket}) \mapsto \sum_{i=1}^n F_{loss}^{\hat{Y}^i}(N(X^i))$  is named the objective function or global

loss.

**Theorem** – Let  $p$  dense layers with activation functions –  $ReLU^k \circ L^k(\cdot, W^k, b^k): \mathfrak{R}^{m_k} \rightarrow \mathfrak{R}^{m_{k+1}}$  for  $k \in \llbracket 1, p-1 \rrbracket$  and  $S \circ L^p(\cdot, W^p, b^p): \mathfrak{R}^{m_p} \rightarrow \mathfrak{R}^{m_{p+1}}$ .  $W^k \in M_{m_{k+1}, m_k}$  and  $b^k \in \mathfrak{R}^{m_{k+1}}$  are defined as the  $L^k$  weight and bias matrices. Let a training data set  $(X^i, \hat{Y}^i) \in \mathfrak{R}^{m_1} \times \mathfrak{R}^{m_{p+1}}$  for  $i \in \llbracket 1, n \rrbracket$ . Let  $\xi^{\hat{Y}^i}: \mathfrak{R}^{m_{p+1}} \rightarrow \mathfrak{R}$  Categorical cross-entropy losses for  $i \in \llbracket 1, n \rrbracket$  with  $(\hat{Y}^i)_{i \in \llbracket 1, n \rrbracket}$  as ground truth matrices respectively.

Then the following application  $N: \mathfrak{R}^{m_1} \rightarrow \mathfrak{R}^{m_{p+1}}$  with

$$N = S \circ L^p(\cdot, W^p, b^p) \circ \bigcirc_{k=1}^{p-1} (ReLU^k \circ L^k(\cdot, W^k, b^k))$$

is a neural network and its objective function is

$$\xi_{loss}^{global}: ((W_{1,:}^k, \dots, W_{m_{k+1},:}^k, b^k)_{k \in \llbracket 1, p \rrbracket}) \mapsto \sum_{i=1}^n \xi_{loss}^{\hat{Y}^i}(N(X^i))$$

Let  $k \in \llbracket 1, p \rrbracket$ . For all  $i \in \llbracket 1, n \rrbracket$ , let

$$y^{downstream(k), X^i} = \begin{cases} \bigcirc_{l=1}^{k-1} (ReLU^l \circ L^l(\cdot, W^l, b^l))(X^i) & k \geq 2 \\ X^i & k = 1 \end{cases}$$

$$F^{upstream(k), \hat{Y}^i}: \begin{cases} \mathfrak{R}^{m_{k+1}} \rightarrow \mathfrak{R} \\ y \mapsto \begin{cases} \xi^{\hat{Y}^i} \circ S(y) & k = p \\ \xi^{\hat{Y}^i} \circ S \circ L^p(\cdot, W^p, b^p)(y) & k = p-1 \\ \xi^{\hat{Y}^i} \circ S \circ L^p(\cdot, W^p, b^p) \circ \bigcirc_{l=k+1}^{p-1} (ReLU^l \circ L^l(\cdot, W^l, b^l))(y) & k \leq p-2 \end{cases} \end{cases}$$

$$\text{then } \frac{d F^{upstream(k), \hat{Y}^i}}{d y}: \begin{cases} \mathfrak{R}^{m_{k+1}} \rightarrow \mathfrak{R} \\ y \mapsto \begin{cases} S(y) - \hat{Y}^i & k = p \\ (S(y) - \hat{Y}^i) \times W^p & k = p-1 \text{ where} \\ \frac{\partial (F^{upstream(k+1), \hat{Y}^i} \circ ReLU^{k+1} \circ L^{k+1}(\cdot, W, b))}{\partial y}(y) & k \leq p-2 \end{cases} \end{cases}$$

$$\forall k \in \llbracket 1, p-2 \rrbracket, \forall y \in \mathfrak{R}^{m_{k+1}}, \forall j \in \llbracket 1, m \rrbracket,$$

$$\frac{\partial (F^{upstream(k+1), \hat{Y}^i} \circ ReLU^{k+1} \circ L^{k+1}(\cdot, W^{k+1}, b^{k+1}))}{\partial y}(y)_j$$

$$= \sum_{j'=1}^{m'} \frac{d F^{upstream(k+1), \hat{Y}^i}}{d y'} (ReLU^{k+1}(L^{k+1}(y, W^{k+1}, b^{k+1})))_{j'} \times 1_{\mathfrak{R}_{\setminus \{0\}}^+}(L^{k+1}(y, W^{k+1}, b^{k+1}))_{j'} \times W_{j', j}^{k+1}$$

with  $W^{k+1} \in M_{m', m}$  a weight matrix,  $b^{k+1} \in \mathfrak{R}^{m'}$  a bias matrix.

Let  $k = p$ . Then  $\xi_{loss}^{global}$  Gradient matrices are

$$\forall j' \in \llbracket 1, m_{k+1} \rrbracket, \quad \frac{\partial \xi_{loss}^{global}}{\partial W_{j', :}^k} : \left\{ \begin{array}{c} \mathfrak{R}^{m_k} \rightarrow \mathfrak{R}^{m_k} \\ w \mapsto \sum_{i=1}^n \frac{\partial (F_{j'}^{upstream(k), \hat{Y}^i} \circ L_{j'}^k(\cdot, w, b_{j'}^k))}{\partial w}(w) \end{array} \right.$$

where  $\forall i \in \llbracket 1, n \rrbracket$ ,  $\forall w \in \mathfrak{R}^{m_k}$ ,  $\forall j \in \llbracket 1, m_k \rrbracket$ ,

$$\frac{\partial (F_{j'}^{upstream(k), \hat{Y}^i} \circ L_{j'}^k(\cdot, w, b_{j'}^k))}{\partial w}(w)_j$$

$$= F_{j'}^{upstream(k), \hat{Y}^i} (L_{j'}^k(y^{downstream(k), X^i}, w, b_{j'}^k)) \times y_j^{downstream(k), X^i} \quad \text{with } b^k \in \mathfrak{R}^{m_{k+1}} \text{ a bias matrix.}$$

$$\frac{\partial \xi_{loss}^{global}}{\partial b^k} : \left\{ \begin{array}{c} \mathfrak{R}^{m_{k+1}} \rightarrow \mathfrak{R}^{m_{k+1}} \\ b^k \mapsto \sum_{i=1}^n \frac{\partial (F^{upstream(k), \hat{Y}^i} \circ L^k(\cdot, W^k, b^k))}{\partial b^k}(b^k) \end{array} \right.$$

where  $\forall i \in \llbracket 1, n \rrbracket$ ,  $\forall b^k \in \mathfrak{R}^{m_k}$ ,  $\forall j' \in \llbracket 1, m_{k+1} \rrbracket$ ,

$$\frac{\partial (F^{upstream(k), \hat{Y}^i} \circ L^k(\cdot, W^k, b^k))}{\partial b^k}(b^k)_j,$$

$$= \frac{d F^{upstream(k), \hat{Y}^i}}{d y} (L^k(y^{downstream(k), X^i}, W^k, b^k))_j \times 1_{\mathfrak{R}_{\setminus \{0\}}^+}(L^k(y^{downstream(k), X^i}, W^k, b^k))_j,$$

with  $W^k \in M_{m_{k+1}, m_k}$  a weight matrix.

Let  $k \in \llbracket 1, p-1 \rrbracket$ . Then  $\xi_{loss}^{global}$  Gradient matrices are

$$\forall j' \in \llbracket 1, m_{k+1} \rrbracket, \quad \frac{\partial \xi_{loss}^{global}}{\partial W_{j',:}^k} : \left\{ \begin{array}{c} \mathfrak{R}^{m_k} \rightarrow \mathfrak{R}^{m_k} \\ w \mapsto \sum_{i=1}^n \frac{\partial (F_{j'}^{upstream(k), \hat{Y}^i} \circ ReLU_{j'}^k \circ L_{j'}^k(\cdot, w, b_{j'}^k))}{\partial w} (w) \end{array} \right.$$

where  $\forall i \in \llbracket 1, n \rrbracket, \forall w \in \mathfrak{R}^{m_k}, \forall j \in \llbracket 1, m_k \rrbracket,$

$$\frac{\partial (F_{j'}^{upstream(k), \hat{Y}^i} \circ ReLU_{j'}^k \circ L_{j'}^k(\cdot, w, b_{j'}^k))}{\partial w} (w)_j$$

$$= F_{j'}^{upstream(k), \hat{Y}^i} (ReLU_{j'}^k (L_{j'}^k (y^{downstream(k), X^i}, w, b_{j'}^k))) \times 1_{\mathfrak{R}_{\setminus \{0\}}^+} (L_{j'}^k (y^{downstream(k), X^i}, w, b_{j'}^k)) \times y_j^{downstream(k), X^i}$$

with  $b^k \in \mathfrak{R}^{m_{k+1}}$  a bias matrix.

$$\frac{\partial \xi_{loss}^{global}}{\partial b^k} : \left\{ \begin{array}{c} \mathfrak{R}^{m_{k+1}} \rightarrow \mathfrak{R}^{m_{k+1}} \\ b^k \mapsto \sum_{i=1}^n \frac{\partial (F_{j'}^{upstream(k), \hat{Y}^i} \circ ReLU_{j'}^k \circ L_{j'}^k(\cdot, W^k, b^k))}{\partial b^k} (b^k) \end{array} \right.$$

where  $\forall i \in \llbracket 1, n \rrbracket, \forall b^k \in \mathfrak{R}^{m_k}, \forall j' \in \llbracket 1, m_{k+1} \rrbracket,$

$$\frac{\partial (F_{j'}^{upstream(k), \hat{Y}^i} \circ ReLU_{j'}^k \circ L_{j'}^k(\cdot, W^k, b^k))}{\partial b^k} (b^k)_j,$$

$$= \frac{d F_{j'}^{upstream(k), \hat{Y}^i}}{d y} (ReLU_{j'}^k (L_{j'}^k (y^{downstream(k), X^i}, W^k, b^k)))_j \times 1_{\mathfrak{R}_{\setminus \{0\}}^+} (L_{j'}^k (y^{downstream(k), X^i}, W^k, b^k)_j)$$

with  $W^k \in M_{m_{k+1}, m_k}$  a weight matrix.

Proof: TO DO.

## 6 – Optimizations

**Definition** – Let  $f \in D(\Omega, \mathfrak{R})$  with  $\Omega \in \mathfrak{R}^m$ . <TODO>

## 7 – References