

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.

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 – Prerequisites and Notations

1.1 – Matrices

Definition – Let $a_{i,j} \in \mathbb{R}$ for $i \in \llbracket 1, n \rrbracket$ and $j \in \llbracket 1, m \rrbracket$. Then the ordered rectangular array

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

is a real matrix of dimension $n \times m$.

The following notations are considered

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

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

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

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 $E \subseteq \mathfrak{R}$.

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

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

Definition – 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 a $n \times p$ matrix with

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

Definition – 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 – Let $a \in \mathfrak{R}^n$ and $b \in \mathfrak{R}^n$. Then the scalar product between two vectors is noted as

$$(a|b) = a \times b^T = b \times a^T = (b|a)$$

1.2 – Function regularity

Definition – Let Ω a non-empty open subset of \mathfrak{R}^m , $\|\cdot\|_m$ a norm on \mathfrak{R}^m , $\|\cdot\|_{m'}$ a norm on $\mathfrak{R}^{m'}$, and $f: \Omega \rightarrow \mathfrak{R}^{m'}$. Then f continuous function is equivalent to

$$\forall a \in \Omega,$$

$$\forall \epsilon > 0, \exists \eta > 0, \forall x \in \Omega, {}_m\|x - a\| \leq \eta \Rightarrow {}_m\|f(x) - f(a)\| \leq \epsilon$$

The notation $\xi(\Omega, \mathfrak{R}^{m'})$ means the set of continuous functions from Ω to $\mathfrak{R}^{m'}$.

The notation $\xi(\mathfrak{R}^m)$ means the set of continuous functions from \mathfrak{R}^m to \mathfrak{R}^m .

Definition – Let Ω a non-empty open subset of \mathfrak{R}^m , ${}_m\|\cdot\|$ a norm on \mathfrak{R}^m , ${}_{m'}\|\cdot\|$ a norm on $\mathfrak{R}^{m'}$, and $f: \Omega \rightarrow \mathfrak{R}^{m'}$. Then f derivable is equivalent to

$$\forall a \in \Omega, \exists f'(a) \in \mathfrak{R}^{m'},$$

$$\forall \epsilon > 0, \exists \eta > 0, \forall x \in \Omega, {}_m\|x - a\| \leq \eta \Rightarrow {}_{m'}\left\|\frac{f(x) - f(a)}{x - a} - f'(a)\right\| \leq \epsilon$$

The notation $D(\Omega, \mathfrak{R}^{m'})$ means the set of derivable functions from Ω to $\mathfrak{R}^{m'}$.

The notation $D(\mathfrak{R}^m)$ means the set of derivable functions from \mathfrak{R}^m to \mathfrak{R}^m .

Definition – Let Ω a non-empty open subset of \mathfrak{R}^m , ${}_m\|\cdot\|$ a norm on \mathfrak{R}^m , ${}_{m'}\|\cdot\|$ a norm on $\mathfrak{R}^{m'}$, and $f: \Omega \rightarrow \mathfrak{R}^{m'}$. Then f piece-wise derivable is equivalent to

$\forall K \subset \Omega$ such as K compact and bounded,

$$\exists (K_i)_{i \in \llbracket 0, n \rrbracket} \text{ non-empty open subsets such as } \bigcup_{i=0}^n \overline{K_i} = K \text{ and } \forall (i, i') \in \llbracket 0, n \rrbracket^2, \\ i \neq i' \Rightarrow K_i \cap K_{i'} = \emptyset,$$

$$\forall i \in \llbracket 0, n \rrbracket, \exists f_i \in D^0(\overline{K_i}, \mathfrak{R}^{m'}), \forall x \in K_i, f_i(x) = f(x)$$

The notation $D_{pw}(\Omega, \mathfrak{R}^{m'})$ means the set of piece-wise derivable functions from Ω to $\mathfrak{R}^{m'}$.

The notation $D_{pw}(\mathfrak{R}^m)$ means the set of piece-wise derivable functions from \mathfrak{R}^m to \mathfrak{R}^m .

Proposition – Let Ω a non-empty open subset of \mathfrak{R}^m . Then

$$D(\Omega, \mathfrak{R}^m) \subset D_{pw}(\Omega, \mathfrak{R}^m) \subset \xi(\Omega, \mathfrak{R}^m)$$

Proof: TO DO.

Notation – Let $f: \Omega \rightarrow \Omega'$ and $g: \Omega' \rightarrow \Omega''$. Then the notation $g \circ f$ means the application

$$g \circ f: \begin{cases} \Omega \rightarrow \Omega'' \\ x \mapsto g(f(x)) \end{cases}.$$

Let $(f_i)_{i \in [1, n]}$ with $f_i: \Omega_i \rightarrow \Omega_{i+1}$ for $i \in [1, n]$. Then the notation $\bigcirc_{i=1}^n f_i$ means the application

$$\bigcirc_{i=1}^n f_i: \begin{cases} \Omega_1 \rightarrow \Omega_{n+1} \\ x \mapsto f_n(f_{n-1}(\dots f_2(f_1(x)))) \end{cases}.$$

Theorem – Let U and V non-empty open subsets of \mathbb{R}^n and \mathbb{R}^m . Let $f: \begin{cases} U \rightarrow V \\ x \mapsto f(x) \end{cases}$ and

$g: \begin{cases} V \rightarrow \mathbb{R}^p \\ y \mapsto g(y) \end{cases}$ such as $f \in D(U, V)$ and $g \in D(V, \mathbb{R}^p)$. Then $g \circ f: \begin{cases} U \rightarrow \mathbb{R}^p \\ x \mapsto g(f(x)) \end{cases}$ is in $D(U, \mathbb{R}^p)$ and its Jacobian is

$$\frac{d(g \circ f)}{dx}: \begin{cases} U \rightarrow M_{p, n} \\ x \mapsto \frac{dg}{dy}(f(x)) \times \frac{df}{dx}(x) \end{cases}$$

Proof: TO DO.

1.3 – Function convexity and smoothness

Definition – Let Ω a non-empty convex open subset of \mathbb{R}^m , $\|\cdot\|_m$ a norm on \mathbb{R}^m , and $f: x \mapsto f(x)$ such as $f \in \zeta(\Omega, \mathbb{R})$. Then f convex on Ω is equivalent to

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

Proposition – Let Ω a non-empty convex open subset of \mathbb{R}^m , $\|\cdot\|_m$ a norm on \mathbb{R}^m , and $f \in D(\Omega, \mathbb{R})$ convex on Ω . Then f convex on Ω is equivalent to

$$\begin{aligned} & \forall (y, z) \in \Omega^2, \\ & f(y) + \frac{df}{dx}(y) \times (z - y)^T \leq f(z) \end{aligned}$$

Proof: TO DO.

Definition – Let Ω a non-empty open subset of \mathbb{R}^m , $\|\cdot\|_m$ a norm on \mathbb{R}^m , and $f: \mathcal{X} \mapsto \mathcal{Y}$ such as $f \in D(\Omega, \mathbb{R})$. Let $L > 0$. Then f L -smooth on Ω 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 Ω a non-empty convex open subset of \mathbb{R}^m , $\|\cdot\|_m$ a norm on \mathbb{R}^m , and $f: \mathcal{X} \mapsto \mathcal{Y}$ such as $f \in D(\Omega, \mathbb{R})$. Then

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

Proof: TO DO.

<Co-coercivity with L-smooth & convex TODO>

1.3 – Others

Notation – Let $f: \begin{cases} \Omega_1 \times \dots \times \Omega_n \rightarrow \Omega'_1 \times \dots \times \Omega'_m \\ (x_1, \dots, x_n) \mapsto f(x_1, \dots, x_n) \end{cases}$ 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 $f(x_1, \dots, x_{k-1}, \cdot, x_{k+1}, \dots, x_n): \begin{cases} \Omega_k \rightarrow \Omega'_1 \times \dots \times \Omega'_m \\ x_k \mapsto f(x_1, \dots, x_{k-1}, x_k, x_{k+1}, \dots, x_n) \end{cases}$.

Notation – The notation $1_{\mathbb{R}_{\setminus \{0\}}^+}$ means the $\mathbb{R}_{\setminus \{0\}}^+$ indicator function on \mathbb{R} .

$$1_{\mathbb{R}_{\setminus \{0\}}^+}: \begin{cases} \mathbb{R}^m \rightarrow \mathbb{R}^m \\ z \mapsto f(z) \end{cases}$$

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

$$\forall j \in \llbracket 1, m \rrbracket, \quad ReLU_j'(0) = 0$$

Proposition – $ReLU$ is an activation function. Its Jacobian matrix is

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

Proof: TO DO.

Proposition – The following vector wise application is an activation function

$$\text{SoftMax} : \begin{cases} \mathbb{R}^m \rightarrow]0,1[^m \\ z \mapsto \frac{e^{z_j}}{\sum_{j'=1}^m e^{z_{j'}}} \end{cases}$$

with e the element-wise exponential operation.

The *SoftMax* function will be denoted as S for simplicity.

Its Jacobian matrix is

$$\frac{dS}{dz} : \begin{cases} \mathbb{R}^m \rightarrow M_{m,m} \\ z \mapsto \frac{dS}{dz}(z) \end{cases}$$

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

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

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 $F_{\text{loss}}^{\hat{y}} \in D(\Omega, \mathbb{R})$. Then $F_{\text{loss}}^{\hat{y}}$ is a loss function is equivalent to the application

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

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

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

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 $L: \mathbb{R}^m \times M_{m',m} \times \mathbb{R}^{m'} \rightarrow \mathbb{R}^{m'}$ and $ReLU: \mathbb{R}^{m'} \rightarrow \mathbb{R}^{m'}$ the dense layer and

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

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

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

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

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

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

Its Gradient matrices are

$$\frac{\partial (F^{upstream} \circ ReLU \circ L(\cdot, W, b))}{\partial y} : \left\{ y \mapsto \frac{\partial (F^{upstream} \circ ReLU \circ L(\cdot, W, b))}{\partial y}(y) \right.$$

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

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

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

$$\frac{\partial (F_{j'}^{upstream} \circ ReLU_{j'} \circ L_{j'}(\cdot, w, b_{j'}))}{\partial w} : \left\{ w \mapsto \frac{\partial (F_{j'}^{upstream} \circ ReLU_{j'} \circ L_{j'}(\cdot, w, b_{j'}))}{\partial w}(w) \right.$$

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

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

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

$$\frac{\partial (F^{upstream} \circ ReLU \circ L(\cdot, W, b))}{\partial b} : \left\{ b \mapsto \frac{\partial (F^{upstream} \circ ReLU \circ L(\cdot, W, b))}{\partial b}(b) \right.$$

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

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

with $y \in \mathbb{R}^m$ 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

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

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

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

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

$$\begin{aligned} & \frac{\partial (F^{upstream(k+1), \hat{Y}^i} \circ ReLU^{k+1} \circ L^{k+1}(\cdot, W^{k+1}, b^{k+1}))}{\partial y}(y)_j \\ &= \sum_{j'=1}^{m'} \frac{dF^{upstream(k+1), \hat{Y}^i}}{dy'} (ReLU^{k+1}(L^{k+1}(y, W^{k+1}, b^{k+1})))_{j'} \times 1_{\mathcal{R}_{\setminus \{0\}}^+}(L^{k+1}(y, W^{k+1}, b^{k+1}))_{j'} \times W_{j', j}^{k+1} \end{aligned}$$

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

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

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

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

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

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

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

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

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

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

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

Let $k \in \llbracket 1, p-1 \rrbracket$. Then ξ_{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.

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

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

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

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

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

Proof: TO DO.

6 – Optimizations

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

7 – References