

A mathematical understanding of deep learning

Jiang J.

Data Engineer, Data Scientist

ENSEEIH Computer Science Engineering Degree, INP Toulouse Dual MSc Research Degree in AI / Big Data / Ops
France

1 – Introduction

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 summary of the important propositions. These propositions will be mathematically proven. The subject tackled is a multi-class classification problem with – dense layers, *ReLU* and *SoftMax* activation layers, Categorical cross-entropy loss, Stochastic gradient descent optimizer. All the elements below are defined for classifications but can be re-used or easily re-defined to cover regressions.

2 – Notation and Nomenclature

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

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

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

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

The notation $\mathcal{C}(\mathbb{R}^m)$ means the set of continuous functions from \mathbb{R}^m to $\mathbb{R}^{m'}$.

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

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

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

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 , $\|\cdot\|_m$ a norm on \mathfrak{R}^m , $\|\cdot\|_{m'}$ a norm on $\mathfrak{R}^{m'}$, and $f: \Omega \rightarrow \mathfrak{R}^{m'}$. Then f 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 .

Definition – Let $(a_{i,j})_{i \in \llbracket 1, n \rrbracket, j \in \llbracket 1, m \rrbracket} \in \mathfrak{R}^{n \times m}$. 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 $m \times p$ 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 $a \in \mathfrak{R}^m$ and $b \in \mathfrak{R}^m$. Let the element wise product noted as $(a|b)$ be

$$c = (a|b) = (b|a)$$

where c is in \mathfrak{R}^m with

$$\forall j \in \llbracket 1, m \rrbracket, c_j = a_j \times b_j$$

3 – Activation functions

Definition – Let $f \in \zeta(\mathfrak{R}^m) \cap D(\mathfrak{R}^m)$. Then the vector wise application

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

is an activation function.

Definition – $ReLU$ is the following vector wise application

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

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

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

$ReLU$ is supposed derivable on every coefficients at 0

$$\forall j \in \llbracket 1, m \rrbracket, \quad \frac{d ReLU_j}{dz_j}(0) = 0$$

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

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

with $1_{\mathfrak{R}_{\setminus\{0\}}^+}$ the $\mathfrak{R}_{\setminus\{0\}}^+$ indicator function on \mathfrak{R}^m .

Proof: TO DO.

Proposition – The following vector wise application is an activation function

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

with e the element-wise exponential operation.

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

Its Jacobian matrix is

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

where $\forall z \in \mathfrak{R}^m, \quad \forall (j, j') \in \{1, 2, \dots, m\}^2, \quad \frac{d S}{dz}(z)_{j,j'} = S(z)_{i,j} \times (\delta_{j,j'} - S(z)_{i,j'})$

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

Proof: TO DO.

3 – Loss

Definition – Let $\hat{\Omega} \in \mathcal{R}^m$ and $\Omega \subseteq \mathcal{R}^m$ non empty subsets. Let $\hat{y} \in \hat{\Omega}$ and $f \in \xi(\Omega, \mathcal{R}) \cap D(\Omega, \mathcal{R})$. Then f is a loss function is equivalent to the application

$$f \circ g: \begin{cases} E \rightarrow \mathcal{R} \\ \epsilon \mapsto (f \circ g)(\epsilon) = f(\hat{y} + \epsilon) \end{cases}$$

is an increasing function according each coefficient with $E \subseteq \mathcal{R}^m$ such as $f \circ g$ is always defined.

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

Proposition – Let $\hat{y} \in \{0,1\}^m$. Then the application

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

is a loss function.

The application is commonly named as Categorical cross-entropy loss.

Its Gradient matrix is

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

Proof: TO DO.

Proposition – Let $S: \mathcal{R}^m \mapsto]0,1[^m$ and $\xi:]0,1[^m \mapsto \mathcal{R}$ the activation and loss functions respectively. Then $S \circ \xi: \mathcal{R}^m \mapsto \mathcal{R}$ is derivable

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 corresponding to the coefficient j of the dense layer L . Let $y \in \mathbb{R}^m$, $W \in M_{m',m}$ and $b \in \mathbb{R}^{m'}$ then

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

Proposition – Its Jacobian matrices are

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

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

Proof: TO DO.

Proposition – Let $y \in M_{n,m}$ the input vector, $W \in M_{m,m'}$ the weight matrix and $b \in M_{1,m'}$ the bias matrix. Then the sequential operations

$$\begin{aligned} \forall i \in \llbracket 1, n \rrbracket, \forall j' \in \llbracket 1, m' \rrbracket, z'_{i,j'} &= y_{i,:} \times (W^T)_{:,j'} + b \\ \forall i \in \llbracket 1, n \rrbracket, y'_{i,:} &= S(z'_{i,:}) \end{aligned}$$

defines a dense layer.

The gradient matrices are the following

$$\forall i \in [1, n] ,$$

5 – Neural Network

Definition – Suppose a data set with n samples. Each sample have m features and a corresponding one-hot encoded label among l possible labels.

Let $X \in M_{n,m}$, and $Y \in M_{n,l}$ the matrices defining the features and the labels respectively for each sample.

Suppose a neural network with k layers.

6 – References