

# Développement d'une bibliothèque de machine learning

Soutenance de projet

---

Aurélien DELVAL

IATIC5

17 février 2021



université PARIS-SACLAY

ISTY

Institut des Sciences et Techniques des Yvelines

CAMPUS DE MANTES EN YVELINES

CAMPUS DE SAINT-QUENTIN-EN-YVELINES

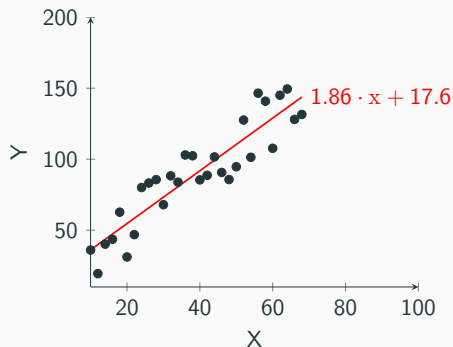
# Table des matières

1. Introduction théorique
2. Différentiation automatique
3. Implémentation des modèles
4. Optimisation
5. Mise en pratique
6. Conclusion

# Introduction théorique

---

## Principe de base : régression



On veut minimiser le **coût** :  $c = \sum_{i=1}^n (a \cdot x_i + b - y_i)^2$

# Introduction aux réseaux de neurones : MNIST

- On va considérer l'exemple de la base de données MNIST

## Introduction aux réseaux de neurones : MNIST

- On va considérer l'exemple de la base de données MNIST
- Reconnaissance de chiffres manuscrits (problème très simple)

# Introduction aux réseaux de neurones : MNIST

- On va considérer l'exemple de la base de données MNIST
- Reconnaissance de chiffres manuscrits (problème très simple)

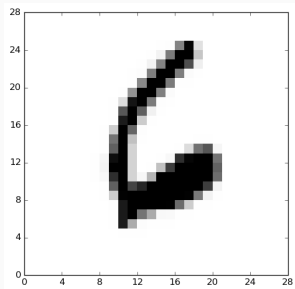
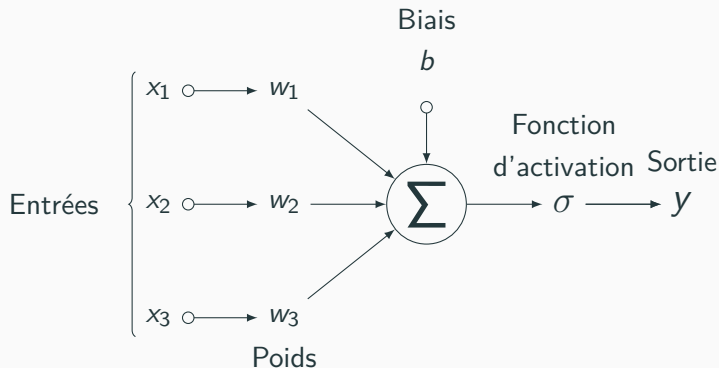


Image d'entrée ( $28 \times 28$  px = 784 entrées)

$$y^{pred} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

Vecteur de sortie attendu

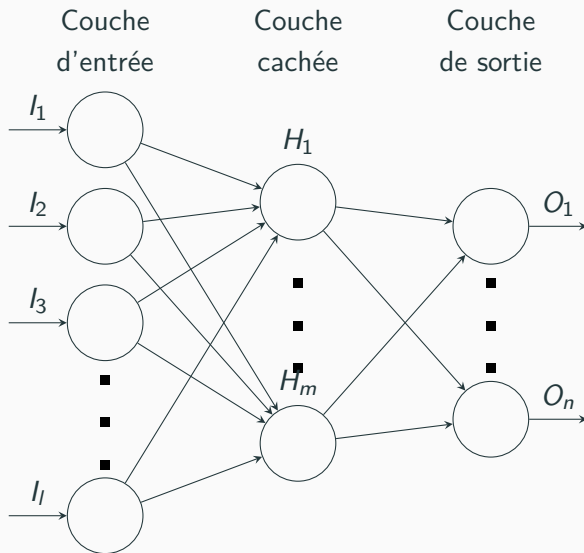
## Introduction aux réseaux de neurones : Neurone artificiel



$$f : \begin{cases} \mathbb{R} \rightarrow [0, 1] \\ x \mapsto f_{act}(\sum_{i=1}^n x_i w_i + b) \end{cases}$$



# Introduction aux réseaux de neurones : MLP



# Introduction aux réseaux de neurones : Couches denses

Valeur d'un neurone  $i$  :

$$H_i = f_{act}\left(\sum_{j=1}^I I_j w_{j,i} + b_i\right)$$

# Introduction aux réseaux de neurones : Couches denses

Valeur d'un neurone  $i$  :

$$H_i = f_{act}\left(\sum_{j=1}^l I_j w_{j,i} + b_i\right)$$

Plutôt que de considérer chaque neurone individuellement, on peut remarquer que :

$$H = f_{act}\left(\begin{bmatrix} w_{1,1} & w_{1,2} & \cdots & w_{1,l} \\ w_{2,1} & w_{2,2} & \cdots & w_{2,l} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m,1} & w_{m,2} & \cdots & w_{m,l} \end{bmatrix} \cdot \begin{bmatrix} I_1 \\ I_2 \\ \vdots \\ I_l \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}\right)$$

# Introduction aux réseaux de neurones : Optimisation

Comment exprimer le coût du modèle ?

$$\vec{c} = \begin{pmatrix} (\hat{y}_1 - y_1^{pred})^2 \\ (\hat{y}_2 - y_2^{pred})^2 \\ \vdots \\ (\hat{y}_n - y_n^{pred})^2 \end{pmatrix}$$

On exprime la norme de ce vecteur :

$$c = \sum_{i=1}^n c_i = \sum_{i=1}^n (\hat{y}_i - y_i^{pred})^2$$

On obtient la **fonction de coût du réseau**

- On calculera le gradient  $\nabla$  de la fonction de coût.
- en prenant  $-\nabla$ , on s'approchera d'un minimum.
- Cependant, la fonction de coût est très complexe.

- On calculera le gradient  $\nabla$  de la fonction de coût.
- en prenant  $-\nabla$ , on s'approchera d'un minimum.
- Cependant, la fonction de coût est très complexe.

On peut calculer  $\nabla$  grâce à la **différentiation automatique**

# Différentiation automatique

---

## Principe et chain rule

La différentiation automatique ne permet pas d'exprimer, mais seulement de calculer un gradient.

Elle repose sur la **chain rule** :

$$(f \circ g)' = (f' \circ g).g'$$

Que l'on peut aussi écrire :

$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$$



## Exemple en mode forward

Prenons la fonction

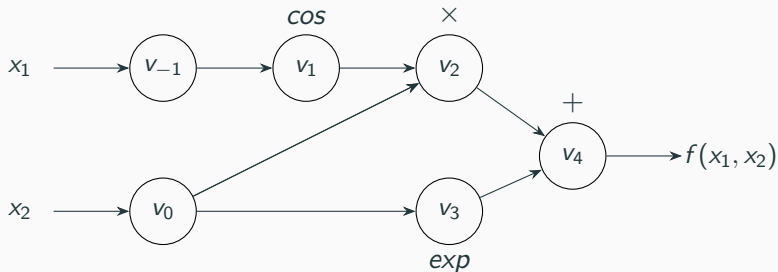
$$f : \begin{cases} \mathbb{R}^2 \rightarrow \mathbb{R} \\ (x_1, x_2) \mapsto x_2 \cdot \cos(x_1) + e^{x_2} \end{cases}$$

## Exemple en mode forward

Prenons la fonction

$$f : \begin{cases} \mathbb{R}^2 \rightarrow \mathbb{R} \\ (x_1, x_2) \mapsto x_2 \cdot \cos(x_1) + e^{x_2} \end{cases}$$

Son **graphe de calcul** est :



## Exemple en mode forward

$$\begin{cases} f(\pi, 1) = 1.\cos(\pi) + e^1 = e - 1 \\ \frac{\partial f}{\partial x_1} = -1.\sin(\pi) = 0 \\ \frac{\partial f}{\partial x_2} = \cos(\pi) + e^1 = e - 1 \end{cases}$$

## Exemple en mode forward

$$\begin{cases} f(\pi, 1) = 1.\cos(\pi) + e^1 = e - 1 \\ \frac{\partial f}{\partial x_1} = -1.\sin(\pi) = 0 \\ \frac{\partial f}{\partial x_2} = \cos(\pi) + e^1 = e - 1 \end{cases}$$

Evaluation de $f$	Evaluation de $\frac{\partial f}{\partial x_1}$	Evaluation de $\frac{\partial f}{\partial x_2}$
$v_{-1} = x_1 = \pi$ $v_0 = x_2 = 1$	$\dot{v}_{-1} = \dot{x}_1 = 1$ $\dot{v}_0 = \dot{x}_2 = 0$	$\dot{v}_{-1} = \dot{x}_1 = 0$ $\dot{v}_0 = \dot{x}_2 = 1$
$v_1 = \cos(v_{-1}) = -1$ $v_2 = v_0.v_1 = -1$ $v_3 = e^{v_0} = e$	$\dot{v}_1 = -\dot{v}_{-1}.\sin(v_{-1}) = 0$ $\dot{v}_2 = \dot{v}_0.v_1 + v_0.\dot{v}_1 = 0$ $\dot{v}_3 = \dot{v}_0.e^{v_0} = 0$	$\dot{v}_1 = -\dot{v}_{-1}.\sin(v_{-1}) = 0$ $\dot{v}_2 = \dot{v}_0.v_1 + v_0.\dot{v}_1 = -1$ $\dot{v}_3 = \dot{v}_0.e^{v_0} = e$
$v_4 = v_2 + v_3 = e - 1$	$\dot{v}_4 = \dot{v}_2 + \dot{v}_3 = 0$	$\dot{v}_4 = \dot{v}_2 + \dot{v}_3 = e - 1$

## Exemple en mode backward

Le principe reste le même, mais en réécrivant la chain rule :

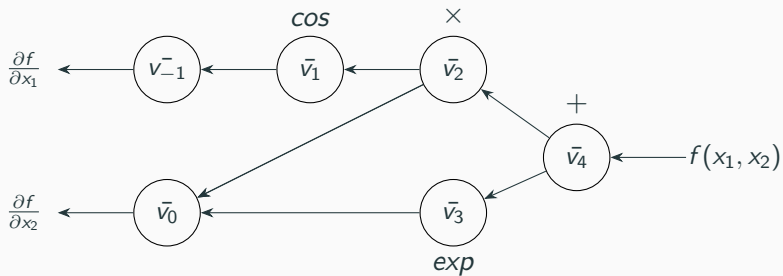
$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$$

Devient :

$$\frac{dx}{dz} = \frac{dy}{dz} \cdot \frac{dx}{dy}$$

## Exemple en mode backward

On va parcourir le graphe de calcul depuis la fin :



## Exemple en mode backward

Evaluation de $f$	Evaluation des $\frac{\partial f}{\partial x_i}$
$v_{-1} = x_1 = \pi$	$v_{-1} = \bar{v}_1 \cdot \frac{\partial v_1}{\partial v_{-1}} = 0$
$v_0 = x_2 = 1$	$\bar{v}_0 = \bar{v}_3 \cdot \frac{\partial v_3}{\partial v_0} + \bar{v}_2 \cdot \frac{\partial v_2}{\partial v_0} = e - 1$
$v_1 = \cos(v_{-1}) = -1$	$\bar{v}_1 = \bar{v}_2 \cdot \frac{\partial v_2}{\partial v_1} = -1$
$v_2 = v_0 \cdot v_1 = -1$	$\bar{v}_2 = \bar{v}_4 \cdot \frac{\partial v_4}{\partial v_2} = -1$
$v_3 = e^{v_0} = e$	$\bar{v}_3 = \bar{v}_4 \cdot \frac{\partial v_4}{\partial v_3} = e$
$v_4 = v_2 + v_3 = e - 1$	$\bar{v}_4 = 1$

On n'a fait qu'un seul passage !

Dans le cas général, si on a  $m$  entrées et  $n$  sorties :



Dans le cas général, si on a  $m$  entrées et  $n$  sorties :

- $n > m$ , il est plus facile de calculer  $\nabla$  avec le mode forward ( $m$  passages)

Dans le cas général, si on a  $m$  entrées et  $n$  sorties :

- $n > m$ , il est plus facile de calculer  $\nabla$  avec le mode forward ( $m$  passages)
- $n < m$ , il est plus facile de calculer  $\nabla$  avec le mode backward ( $n$  passages)

# Différentiation automatique dans le contexte du machine learning

Dans le cas général, si on a  $m$  entrées et  $n$  sorties :

- $n > m$ , il est plus facile de calculer  $\nabla$  avec le mode forward ( $m$  passages)
- $n < m$ , il est plus facile de calculer  $\nabla$  avec le mode backward ( $n$  passages)

Dans notre contexte, seul le **mode backward** est intéressant.

# Implémentation

On va construire le graphe de calcul au fur et à mesure en utilisant une **liste de Wengert** :



Chaque noeud contient :

- les dérivées partielles de l'opération associée
- les indices des noeuds dont il dépend

# Implémentation

Pour calculer le gradient, l'algorithme est le suivant :

---

**Algorithme 1** : Calcul de gradient depuis une liste de Wengert

---

**Entrée** : liste de Wengert  $W$ , variable  $v$  associée au noeud  $W_k$

**Résultat** : gradient  $\nabla$

```
1  $\nabla \leftarrow \{0\}_0^{|W|}$ 
2  $\nabla_k \leftarrow 1$ 
   /* Parcours depuis la fin                                     */
3 pour  $i \leftarrow k$  à 1 faire
4    $C \leftarrow \text{enfants}(W_k)$ 
5   pour  $j \leftarrow 0$  à  $|C|$  faire
6     /* Indice dans  $C$  != Indice dans  $W$  ou  $\nabla$                        */
7      $j' \leftarrow \text{indice}_{\nabla}(C_j)$ 
8      $\nabla_{j'} \leftarrow \nabla_{j'} + \frac{\partial C_j}{\partial W_k} \cdot \frac{\partial C_j}{\partial C_i}$ 
9   fin
10 fin
```

---

Dans la pratique, on construira la liste au fur et à mesure des calculs avec de la **surcharge d'opérateur**.

Cette approche a plusieurs avantages :

Dans la pratique, on construira la liste au fur et à mesure des calculs avec de la **surcharge d'opérateur**.

Cette approche a plusieurs avantages :

- Quasiment transparente

Dans la pratique, on construira la liste au fur et à mesure des calculs avec de la **surcharge d'opérateur**.

Cette approche a plusieurs avantages :

- Quasiment transparente
- Calcul du gradient en un appel de fonction



Dans la pratique, on construira la liste au fur et à mesure des calculs avec de la **surcharge d'opérateur**.

Cette approche a plusieurs avantages :

- Quasiment transparente
- Calcul du gradient en un appel de fonction
- Rapide à calculer

- Pour un réseau de neurones “réaliste”, la liste de Wengert sera énorme.

- Pour un réseau de neurones “réaliste”, la liste de Wengert sera énorme.
- On peut encapsuler non plus des scalaires mais des tenseurs.

- Pour un réseau de neurones “réaliste”, la liste de Wengert sera énorme.
- On peut encapsuler non plus des scalaires mais des tenseurs.
- Tout au long du développement, on utilisera la bibliothèque Eigen.

- Pour un réseau de neurones “réaliste”, la liste de Wengert sera énorme.
- On peut encapsuler non plus des scalaires mais des tenseurs.
- Tout au long du développement, on utilisera la bibliothèque Eigen.
- On encapsulera le type `Eigen::Array` dans `ts::Tensor`.

# Implémentation des modèles

---

On veut pouvoir créer des modèles :

On veut pouvoir créer des modèles :

- de manière simple (un utilisateur pourrait le faire)



On veut pouvoir créer des modèles :

- de manière simple (un utilisateur pourrait le faire)
- utilisables avec tous les algorithmes d'optimisation

On veut pouvoir créer des modèles :

- de manière simple (un utilisateur pourrait le faire)
- utilisables avec tous les algorithmes d'optimisation

On va créer une classe virtuelle `ts::Model` avec une méthode `compute` qui va encapsuler tous les calculs.

On veut pouvoir créer des modèles :

- de manière simple (un utilisateur pourrait le faire)
- utilisables avec tous les algorithmes d'optimisation

On va créer une classe virtuelle `ts::Model` avec une méthode `compute` qui va encapsuler tous les calculs.

Il suffira de définir un constructeur et cette méthode de calcul pour créer un modèle.

## Exemple du MLP

On veut stocker :

## Exemple du MLP

On veut stocker :

- les matrices de poids

## Exemple du MLP

On veut stocker :

- les matrices de poids
- les vecteurs de biais

## Exemple du MLP

On veut stocker :

- les matrices de poids
- les vecteurs de biais

L'algorithme de calcul est :

---

**Algorithme 2** : Calcul de la sortie d'un MLP

---

**Entrée** : vecteur d'entrée  $v$ , matrices de poids  $\{w\}_{i=0}^n$ , vecteurs de biais  $\{b\}_{i=0}^n$ ,  
fonction  $f_{act}$

**Résultat** : vecteur  $H_n$

```
1  $H_0 \leftarrow v$ 
2 pour  $i \leftarrow 1$  à  $n$  faire
3    $H_i \leftarrow f_{act}(w_i \cdot H_{i-1} + b_i)$ 
4 fin
```

---

On est déjà capables

- de calculer la sortie de réseaux de neurones simples
- d'en implémenter de plus complexes
- d'obtenir très facilement leurs gradients (un appel de fonction)

Il reste à utiliser les calculs de gradient pour ajuster les paramètres.

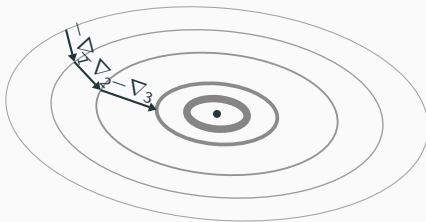


# Optimisation

---

# Principe de la descente de gradient

Le vecteur  $-\nabla$  nous indique la direction d'un minimum de la fonction de coût :



En ajustant petit à petit les paramètres du modèle, nos prédictions deviendront meilleures.

L'algorithme SGD (*stochastic gradient descent*) de base est la version la plus simple :

L'algorithme SGD (*stochastic gradient descent*) de base est la version la plus simple :

- On définit un taux d'apprentissage scalaire  $\eta$  (très faible)

L'algorithme SGD (*stochastic gradient descent*) de base est la version la plus simple :

- On définit un taux d'apprentissage scalaire  $\eta$  (très faible)
- On calcule  $\nabla$  sur des données d'entraînement

L'algorithme SGD (*stochastic gradient descent*) de base est la version la plus simple :

- On définit un taux d'apprentissage scalaire  $\eta$  (très faible)
- On calcule  $\nabla$  sur des données d'entraînement
- On incrémente le modèle de  $-\eta\nabla$

L'algorithme SGD (*stochastic gradient descent*) de base est la version la plus simple :

- On définit un taux d'apprentissage scalaire  $\eta$  (très faible)
- On calcule  $\nabla$  sur des données d'entraînement
- On incrémente le modèle de  $-\eta\nabla$

Après des centaines ou milliers d'itérations, la fonction de coût converge.

# Algorithme SGD

Dans la pratique, on utilise souvent des **batches** et plusieurs **epochs**.

---

**Algorithme 3** : Descente stochastique de gradient

---

**Entrée** : fonction  $f$ , paramètres  $\theta_0$ , ensemble de batches  $B$ , nombre d'epochs  $e$ ,  
taux d'apprentissage  $\eta$

**Résultat** : paramètres  $\theta_t$

```
1 pour  $i \leftarrow 1$  à  $e$  faire
2   pour  $t \leftarrow 1$  à  $|B|$  faire
3      $b \leftarrow B_t$ 
4      $\nabla_{sum} \leftarrow \{0\}_{k=1}^{|\theta_0|}$ 
5     pour  $j \leftarrow 1$  à  $|b|$  faire
6        $\nabla_{sum} \leftarrow \nabla_{sum} + \nabla f_t(\theta_{t-1}, b_j)$ 
7     fin
8      $\theta_t \leftarrow \theta_{t-1} - \eta \frac{\nabla_{sum}}{|b|}$ 
9   fin
10 fin
```

---



## Mise en pratique

---

# Résolution de MNIST

On a maintenant tous les éléments pour résoudre MNIST.

Voici comment se fait l'utilisation de la bibliothèque sur ce problème :

```
1  ts::MultiLayerPerceptron<float> model(  
2      EXPECTED_IMAGE_SIZE ,  
3      {512, 128, 10}  
4  );  
5  
6  ts::AdamOptimizer<float> optimizer;  
7  
8  std::vector<std::vector<std::vector< float >>> losses =  
9  optimizer.run(model, trainingData);
```

**Listing 1** – Définition et utilisation du MLP pour MNIST

# Résolution de MNIST

Phase d'entraînement : **10** epochs et **1000** batches de taille **5**.

# Résolution de MNIST

Phase d'entraînement : **10** epochs et **1000** batches de taille **5**.

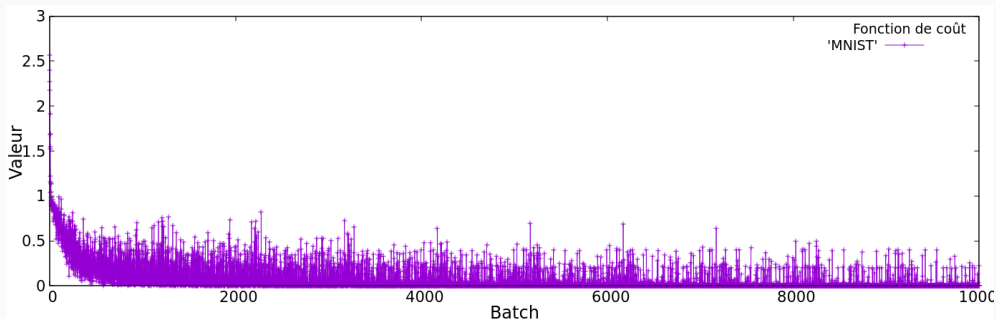
Après quelques minutes d'entraînement, on obtient une précision de **96%**.

# Résolution de MNIST

Phase d'entraînement : **10** epochs et **1000** batches de taille **5**.

Après quelques minutes d'entraînement, on obtient une précision de **96%**.

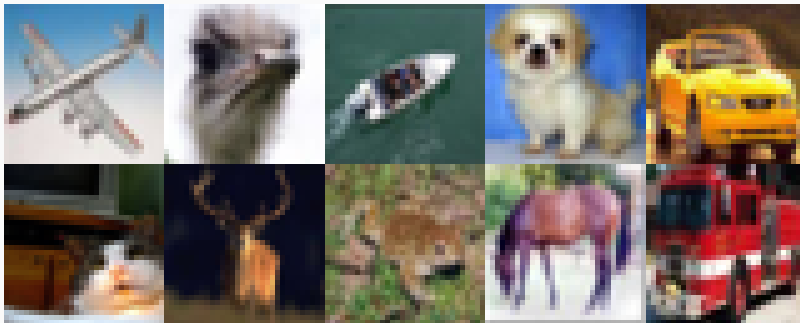
On observe aussi que la fonction de coût a convergé comme prévu :



On va maintenant essayer de résoudre CIFAR, un problème similaire à MNIST, mais plus complexe :

# CIFAR et CNN

On va maintenant essayer de résoudre CIFAR, un problème similaire à MNIST, mais plus complexe :



Phase d'entraînement : **15** epochs et **1000** batches de taille **5**.



Phase d'entraînement : **15** epochs et **1000** batches de taille **5**.

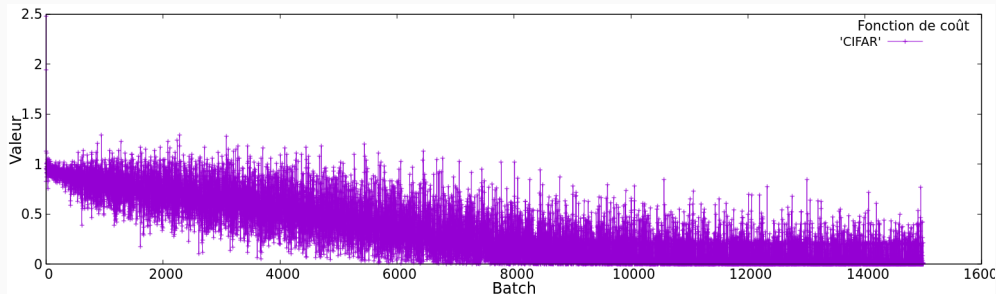
On obtient une précision de **49%**.

# CIFAR et CNN

Phase d'entraînement : **15** epochs et **1000** batches de taille **5**.

On obtient une précision de **49%**.

La fonction de coût converge, mais plus lentement :



## Détection d'objets avec un CNN

Dans MNIST, on se contentait de classifier des images : on va maintenant essayer de faire de la **détection d'objets**. On va :

# Détection d'objets avec un CNN

Dans MNIST, on se contentait de classifier des images : on va maintenant essayer de faire de la **détection d'objets**. On va :

- entraîner un CNN à reconnaître des véhicules grâce à CIFAR

## Détection d'objets avec un CNN

Dans MNIST, on se contentait de classifier des images : on va maintenant essayer de faire de la **détection d'objets**. On va :

- entraîner un CNN à reconnaître des véhicules grâce à CIFAR
- utiliser un autre jeu de données (TrafficNet) pour les prédictions

# Détection d'objets avec un CNN

Dans MNIST, on se contentait de classifier des images : on va maintenant essayer de faire de la **détection d'objets**. On va :

- entraîner un CNN à reconnaître des véhicules grâce à CIFAR
- utiliser un autre jeu de données (TrafficNet) pour les prédictions
- isoler plusieurs zones des images de test pour détecter les véhicules

# Détection d'objets avec un CNN

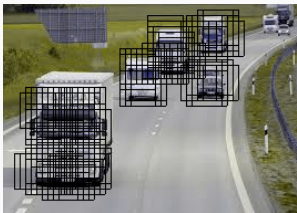
Dans MNIST, on se contentait de classifier des images : on va maintenant essayer de faire de la **détection d'objets**. On va :

- entraîner un CNN à reconnaître des véhicules grâce à CIFAR
- utiliser un autre jeu de données (TrafficNet) pour les prédictions
- isoler plusieurs zones des images de test pour détecter les véhicules

**Remarque :** La mise en place sera une version très simplifiée d'algorithmes tels que R-CNN ou YOLO.

# Détection d'objets avec un CNN

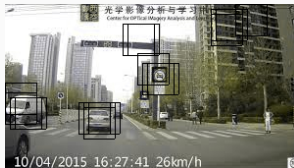
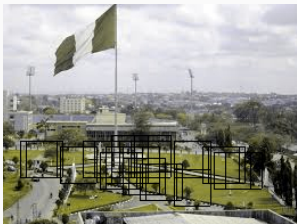
Voici quelques exemples de détection satisfaisants :





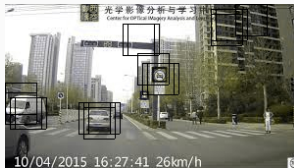
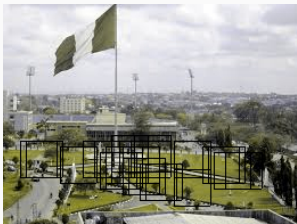
# Détection d'objets avec un CNN

Et d'autres comportant des erreurs



# Détection d'objets avec un CNN

Et d'autres comportant des erreurs



Il est normal que les résultats ne soient pas parfaits car :

- Les jeux de données utilisés n'étaient pas prévus à cet effet.
- Le CNN utilisé n'est pas tout à fait optimal (87% de précision sur la classification binaire).

## Conclusion

---

Utilisation de plusieurs bibliothèques :

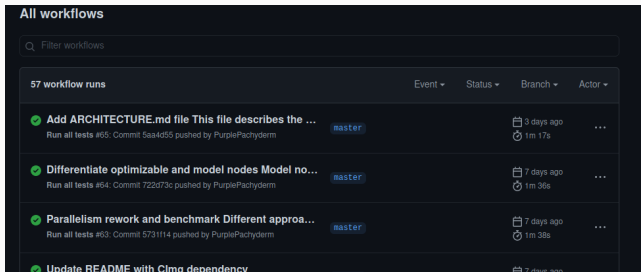
- Eigen
- Google Benchmarks
- Google Tests

# Processus de développement

Utilisation de plusieurs bibliothèques :

- Eigen
- Google Benchmarks
- Google Tests

Intégration continue des tests unitaires avec les Github Workflows :



- Son approche est très générale.

## Conclusion sur la bibliothèque

- Son approche est très générale.
- Elle n'est évidemment pas complète...

## Conclusion sur la bibliothèque

- Son approche est très générale.
- Elle n'est évidemment pas complète...
- ... mais fournit des outils pour de nombreuses utilisations potentielles.