

# Résistance aux attaques adversaires

Maximilien de Dinechin

Juin 2018

## Introduction

Les exemples adversaires sont l'une des principales vulnérabilités des réseaux de neurones, en particulier des réseaux classificateurs : une perturbation imperceptible d'une entrée peut conduire à une erreur de classification. Une telle faiblesse est préoccupante, par exemple dans le cadre de la conduite autonome pour l'identification des panneaux de circulation, ou dans tout autre système où l'on souhaite confier un rôle décisionnel à des réseaux de neurones. Une simple modification d'un pixel peut parfois suffire à induire en erreur un réseau [1], et il n'y a pas pour l'instant de consensus pour résoudre ce problème : des publications proposent régulièrement des solutions, mais sont souvent mises en échec peu de temps après par des exemples adversaires plus sophistiqués.

Ce travail, qui tente de répondre à ce problème, est motivé par l'observation d'un phénomène intéressant dans le fonctionnement d'un algorithme particulier d'attaque adversaire, mettant en évidence en premier lieu une corrélation entre la difficulté à mener cette attaque sur une image donnée et la justesse de la prédiction du réseau ; puis ensuite une corrélation plus forte encore entre cette difficulté et le fait que l'image initiale soit un exemple adversaire ou non.

## 1 Les attaques adversaires

### 1.1 Les exemples adversaires

Les réseaux de neurones sont hautement vulnérables aux *exemples adversaires* [2, 3]. Un exemple adversaire est une entrée légèrement perturbée dans le but d'induire en erreur un réseau classificateur. Le plus souvent, ces perturbations sont imperceptibles pour l'œil humain.

Plus concrètement : posons  $\text{Pred}$  la fonction qui à une image associe la catégorie prédite par réseau ; et considérons une image  $img \in [0, 1]^n$  (c'est-à-dire à  $n$  pixels en noir et blanc), on cherche une perturbation  $r \in [-1, 1]^n$ , de norme (le plus souvent euclidienne) minimale, telle que :

- (i)  $img + r \in [0, 1]^n$
- (ii)  $\text{Pred}(img + r) \neq \text{Pred}(img)$

*Remarque* : Dans toute la suite, on utilisera la norme euclidienne, notée  $\|\cdot\|$ . Enfin, la minimalité de  $\|r\|$  n'est pas toujours nécessaire : on verra par la suite que certaines perturbations restent visuellement imperceptibles alors qu'elles ont une norme élevée.

## 1.2 Les attaques adversaires

Une *attaque adversaire* est un algorithme qui détermine un exemple adversaire à partir d'une image donnée.

Proposons l'attaque suivante, dite "par descente de gradient itérée". Soit une image  $img$ , prédite de catégorie  $c$  par le réseau. Introduisons  $Conf_c$  la fonction qui quantifie la probabilité selon le réseau que l'image soit de catégorie  $c$  (les réseaux classificateurs considérés retournent pour chaque entrée une distribution de probabilités sur les différentes classes possibles). On cherche alors à minimiser par descente de gradient, sur  $r$  initialisé à  $0^n$ , la fonction :

$$Loss_1(r) = \begin{cases} \|r\| & \text{si } Conf_c(img + r) \leq 0.2 \\ Conf_c(img + r) + \|r\| & \text{sinon.} \end{cases}$$

Cette attaque échoue presque toujours : la perturbation  $r$  reste "bloquée" en 0. On pourrait corriger ce problème en initialisant la perturbation à une faible valeur aléatoire, mais le plus simple est d'"inciter" la perturbation à grossir en norme en ajoutant un troisième cas de figure quand  $Conf_c(img + r) > 0.9$  :

$$Loss_2(r) = \begin{cases} \|r\| & \text{si } Conf_c(img + r) \leq 0.2 \\ Conf_c(img + r) + \|r\| & \text{si } 0.2 < Conf_c(img + r) \leq 0.9 \\ Conf_c(img + r) - \|r\| & \text{sinon.} \end{cases}$$

*Remarque* : les seuils à 0.2 et 0.9 ont été choisis car efficaces en pratique.

Même si elle échoue encore souvent, cette deuxième fonction se révélera suffisante pour la suite. On appellera alors  $Pert_N$  la fonction qui à une image associe la perturbation obtenue après  $N$  étapes de descente de gradient de la fonction  $Loss_2$  avec  $r$  initialisé à 0 (algorithme **Adam** [4], avec un taux d'apprentissage  $\eta = 10^{-3}$ ). En **Annexe A** sont présentés quelques résultats d'attaques adversaires sur des images de MNIST.

## 1.3 Réseaux classificateurs et bases de données utilisées

On réalisera toute cette étude sur deux réseaux de type **AlexNet** (CNN avec Dropout) [5], appliqués respectivement aux problèmes de la classification des images de **MNIST** [6] et de **FashionMNIST** [7] (images de chiffres manuscrits et d'habits). Ces réseaux et leur entraînement sont décrits plus précisément dans l'**Annexe B**.

On travaillera sur les sous-bases de **test** (constituées de 10000 images) de ces bases de données, afin d'étudier des images qui n'ont pas été utilisées lors de l'apprentissage. Les deux réseaux ont alors les performances suivantes :

- 62 erreurs sur les 10000 images de **test** de **MNIST**,
- 876 erreurs sur les 10000 images de **test** de **FashionMNIST**.

Ces deux bases de données présentent ainsi l'avantage d'être très proche structurellement (même format et nombre d'images), mais de difficultés très différentes dans le problème de leur classification : étudier les deux en même temps permettra de mieux analyser les résultats obtenus.

## 2 Résistance à une attaque

### 2.1 Images "faciles" et "difficiles" à attaquer

On s'intéresse aux valeurs prises par  $\|r\|$  et  $\text{Conf}_c$  au cours de l'attaque adverse décrite en 1.2. La Figure 1 a été obtenue en attaquant deux images différentes de MNIST.

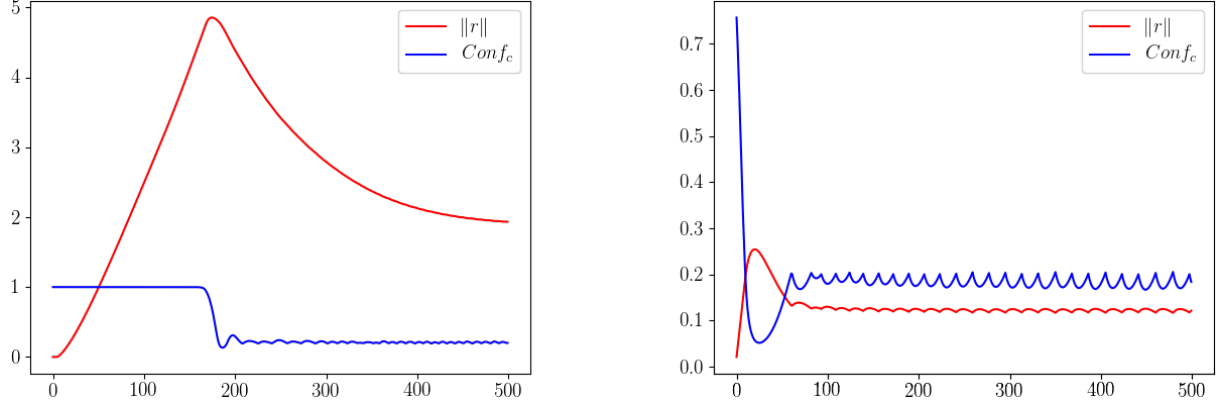


FIGURE 1 – évolution de  $\text{Conf}_c$  (en bleu) et de  $\|r\|$  (en rouge) au cours d'attaques adversaires

L'image de gauche peut être qualifiée de "difficile à attaquer" : un grand nombre d'étapes ont été nécessaires pour casser la prédiction du réseau ; il a été fallu pour cela d'atteindre une norme de  $r$  très élevée ; et la norme finale de la perturbation est importante.

L'image de droite peut au contraire être qualifiée de "facile à attaquer" : il a fallu peu d'étapes pour casser la prédiction du réseau ; le pic de  $\|r\|$  de très faible amplitude ; et la norme finale de  $r$  basse.

Les images de l'**Annexe C** se classifient aisément de la sorte entre "faciles" et "difficiles" à attaquer. Résumons les principales différences lors du déroulement de l'attaque entre les deux types d'images :

	Images "difficiles"	Images "faciles"
Étapes nécessaires	plus de 200	moins de 50
Pic de $\ r\ $	haut	absent ou faible
Norme de $r$ finale	élevée	faible

*Remarque* : Il n'y aucune différence visuelle entre les images faciles et difficiles à attaquer.

### 2.2 Quantification de la résistance à une attaque

Pour quantifier plus précisément cette difficulté à attaquer une image, introduisons le concept de *résistance*. Motivés par les observations de la partie précédente, on utilisera les trois expressions suivantes :  $\text{Res}_N$  la norme de la perturbation obtenue après  $N$  étapes ;  $\text{Res}_{\max}$  la plus haute norme de la perturbation au cours de l'attaque (qui correspond à la hauteur du pic) ; et  $\text{Res}_{\min}$  le nombre d'étapes qu'il a fallu pour abaisser  $\text{Conf}_c$  à 0.2. On pose donc :

- $\text{Res}_N(\text{img}) = \|\text{Pert}_N(\text{img})\|$
- $\text{Res}_{\max}(\text{img}) = \max \{ \|\text{Pert}_N(\text{img})\| \text{ tel que } N \in \mathbb{N} \}$
- $\text{Res}_{\min}(\text{img}) = \min \{ N \in \mathbb{N} \text{ tel que } \text{Conf}_c(\text{Pert}_N(\text{img})) < 0.2 \}$

Dans les cas où l'attaque échoue, on prendra systématiquement  $Res = +\infty$  dans les trois expressions. Dans tout ce qui suit, on prendra  $N = 500$  pour le calcul de  $Res_N$ .

*Remarque* : Une définition optimale de la résistance serait la norme de la perturbation minimale mettant en échec le réseau. Cette expression n'est cependant que d'un faible intérêt pratique, car incalculable.

## 2.3 Une corrélation avec la justesse de la prédiction

Les images attaquées dans l'**Annexe C** n'ont pas été choisies au hasard : les premières correspondent à des classifications correctes, et les suivantes à des erreurs de classification. Étudions la généralisation de ces résultats en observant la répartition des valeurs de la résistance sur des images correctement classifiées (notées **V**), et incorrectement classifiées (notées **F**) : pour **MNIST** avec 500 images dans **V** et les 62 erreurs dans **F** ; et pour **FashionMNIST** avec 500 images dans **V** et dans **F**.

<b>MNIST</b>	$Res_N$	$Res_{\max}$	$Res_{\min}$	<b>FashionMNIST</b>	$Res_N$	$Res_{\max}$	$Res_{\min}$
90% de <b>V</b>	> 0.97	> 2.8	> 109	80% de <b>V</b>	> 0.28	> 0.544	> 26
90% de <b>F</b>	< 0.57	< 1.3	< 58	80% de <b>F</b>	< 0.29	< 0.543	< 25

TABLE 1 – Répartition des résistances des images de **V** et **F**

Selon que les images sont correctement classifiées ou non, la répartition des résistances est alors très inégale : on trouve des valeurs des résistances qui discriminent de part et d'autre respectivement 90% (voire 95%) des images **V** et **F** dans le cas de **MNIST**, et tout juste 80% pour **FashionMNIST**. Une nette corrélation se dessine donc entre la résistance et la justesse de la prédiction du réseau : une résistance élevée est souvent associée à une prédiction juste, et une résistance faible à une erreur de classification.

On pourrait alors y voir une méthode pour améliorer la précision d'un réseau de neurones. Cependant, le nombre de faux positifs (images correctement classifiées mais identifiées comme des erreurs) et de faux négatifs (erreurs identifiées comme correctement classifiées) est bien supérieur à l'erreur totale commise par le réseau dans les deux cas. Une telle méthode semble ainsi peu pertinente : elle réduirait la précision du réseau. J'ai essayé d'affiner cette séparation en m'intéressant d'un seul coup à toutes les valeurs prises par  $\text{Conf}_c$  et  $\|r\|$  au cours d'une attaque (c'est-à-dire les deux courbes complètes), puis en construisant un réseau de neurones "discriminateur", entraîné à faire la différence entre images correctement classifiées et incorrectement classifiées à partir de ces données, mais ici encore, le taux d'erreur est resté supérieur à celui du réseau initial.

## 2.4 Une méthode de détection des exemples adversaires

On observe des résultats similaires dans le cas des attaques adversaires : les exemples adversaires sont en général plus faciles à attaquer que les "vraies" images. Étudions la validité de ce résultat.

La partie 1 présente une méthode efficace de génération d'exemple adversaire. On souhaite cependant se prémunir contre les meilleures attaques existantes, et c'est pourquoi j'ai confié la génération d'exemples adversaires à la bibliothèque **Foolbox** [8] (dont j'ai dû légèrement modifier le code source pour l'adapter aux réseaux de neurones utilisés ici). On utilisera les attaques suivantes :

- 1 **GradientAttack** qui procède par descente de gradient sur l'image,
- 2 **ContrastReductionAttack** qui modifie le contraste de l'image,
- 3 **PointwiseAttack** qui modifie fortement un faible nombre de pixels.

On considèrera qu'un exemple adversaire est satisfaisant quand l'assurance du réseau sur la mauvaise catégorie est supérieure à 95%.

Des exemples adversaires correspondant aux trois attaques utilisées sont présentés en **Annexe D** pour MNIST et **Annexe E** pour FashionMNIST. Dans le cas de MNIST, les attaques **GradientAttack** et **Contrast ReductionAttack** sont très difficiles à réaliser : en moyenne l'attaque ne réussit que sur deux images sur 100, rendant une collecte de données difficiles ; et **PointwiseAttack**, même si elle réussit presque à chaque fois, propose des images qui s'éloignent fortement de l'image originale, au point d'en perturber la lecture par l'œil : l'image en annexe est suffisamment déformée pour que l'on puisse la confondre avec un 8. Le cas de FashionMNIST ne présente pas ces difficultés, et on se concentrera donc sur cette base de données pour étudier la résistance des exemples adversaires.

Étudions donc la répartition des valeurs de la résistance sur 500 images non altérées de FashionMNIST (notées **V**), et sur trois lots de 200 exemples adversaires correspondant aux trois méthodes ci-dessus, notées **A1**, **A2** et **A3**.

FashionMNIST	$Res_N$	$Res_{\max}$	$Res_{\min}$
94 % de <b>V</b>	> 0.16	> 0.30	> 12
94 % de <b>A1</b>	< 0.14	< 0.26	< 8
94 % de <b>A2</b>	< 0.14	< 0.25	< 7
94 % de <b>A3</b>	< 0.13	< 0.25	< 11

TABLE 2 – Répartition des résistances des images de **V**, **A1**, **A2** et **A3**

La corrélation est cette fois-ci très forte : une valeur bien choisie à laquelle on comparerait les résistances obtenues permet de discriminer avec 96% de précision les vraies images des exemples adversaires, et ce dans les trois attaques différentes envisagées.

Cette étude n'est évidemment pas exhaustive (plusieurs paramètres peuvent jouer, en particulier la consigne de devoir tromper le réseau à 95%) mais une corrélation aussi forte est un bon présage quant à la généralisation de ces résultats.

Enfin, il est envisageable que des exemples adversaires puissent être conçus spécifiquement pour tromper cette méthode de détection. Je n'ai pas eu le temps de créer de tels exemples, mais on peut conjecturer que cette condition supplémentaire les conduira à être plus visibles à l'œil humain, donc moins bons.

## Conclusion

Ce travail a ainsi permis de mettre en évidence la corrélation entre la résistance d'une image donnée et la justesse de classification du réseau, toutefois de manière trop imprécise pour y voir une méthode d'amélioration des performances d'un réseau classificateur.

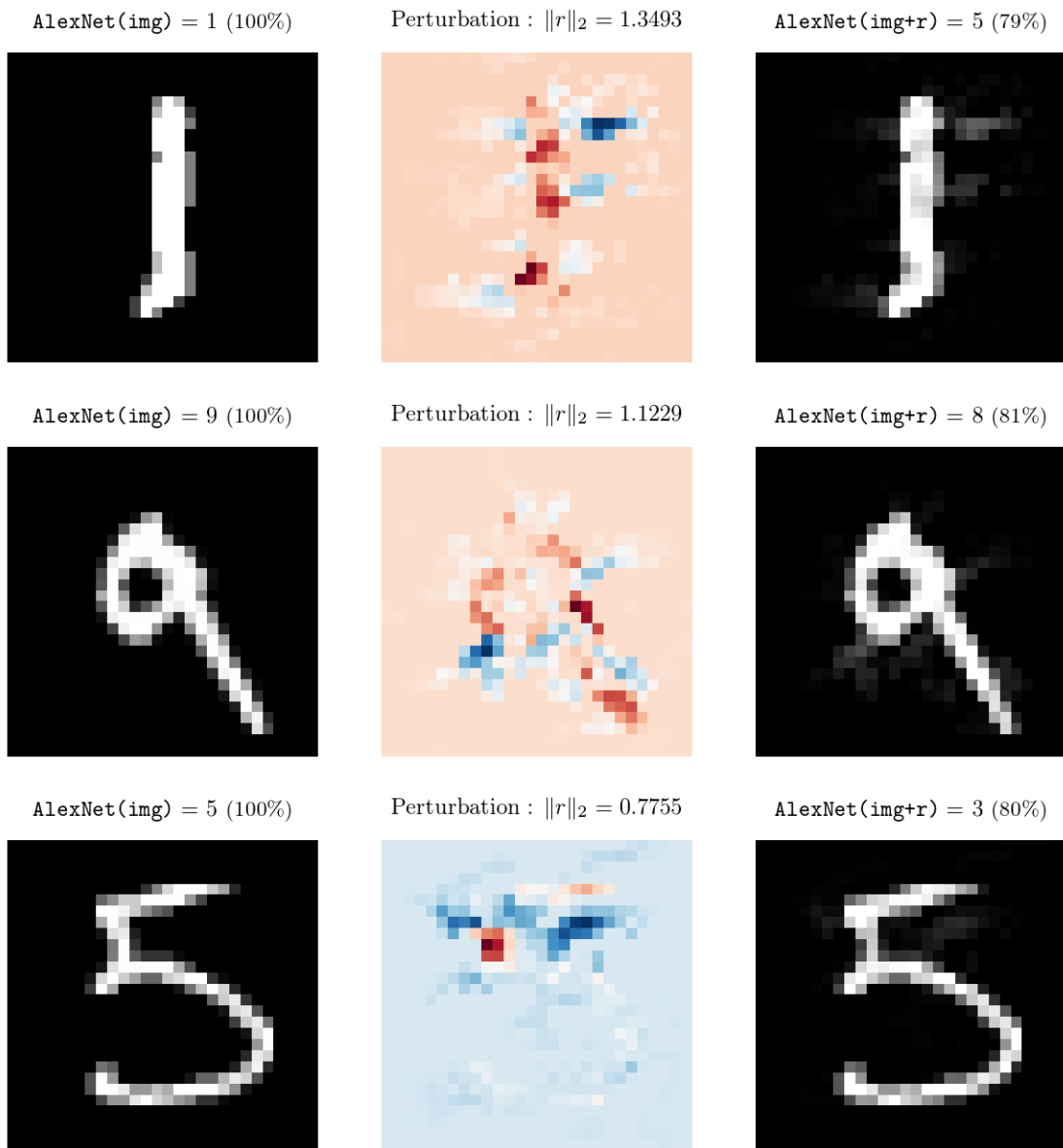
En revanche, la corrélation avec le caractère "vraie image"/"exemple adversaire" d'une image est plus intéressante, car plus précise. Dans un cadre où les exemples adversaires dépassent en fréquence le nombre de faux positifs ou faux négatifs (6% pour FashionMNIST), cette méthode peut être efficace pour se prémunir contre les attaques adversaires.

Ce dernier résultat est ainsi encourageant, mais reste à relativiser cependant : il serait utile de mener une étude plus variée (tant dans les méthodes et seuils d'attaques adversaires, dans les architectures de réseaux et dans les bases de données étudiées) de l'efficacité de cette méthode ; et d'étudier les manières de lutter contre cette méthode de détection, avant de pouvoir conclure quant à sa robustesse.

## Références

- [1] Jiawei Su, Danilo Vasconcellos Vargas, and Sakurai Kouichi. One pixel attack for fooling deep neural networks. *arXiv preprint arXiv :1710.08864*, 2017.
- [2] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv :1312.6199*, 2013.
- [3] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv :1412.6572*, 2014.
- [4] Diederik P Kingma and Jimmy Ba. Adam : A method for stochastic optimization. *arXiv preprint arXiv :1412.6980*, 2014.
- [5] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [6] Yann LeCun. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- [7] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST : a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv :1708.07747*, 2017.
- [8] Jonas Rauber, Wieland Brendel, and Matthias Bethge. Foolbox v0. 8.0 : A Python toolbox to benchmark the robustness of machine learning models. *arXiv preprint arXiv :1707.04131*, 2017.

## Annexe A. Résultats d'attaques adversaires



*(Le bleu est positif, le rouge négatif.)*

## Annexe B. Structure et entraînement des réseaux classificateurs

Les deux réseaux ont la même architecture (structure), inspirée de celle du réseau AlexNet :

- Entrée :  $28 \times 28$
- Convolution : 32 couches, noyau  $5 \times 5$
- ReLU
- MaxPool : noyau  $2 \times 2$
- Convolution : 64 couches, noyau  $3 \times 3$
- ReLU
- MaxPool : noyau  $2 \times 2$
- Dropout
- Couche complète :  $64 \times 5 \times 5 \rightarrow 120$
- ReLU
- Dropout
- Couche complète :  $120 \rightarrow 10$
- Softmax

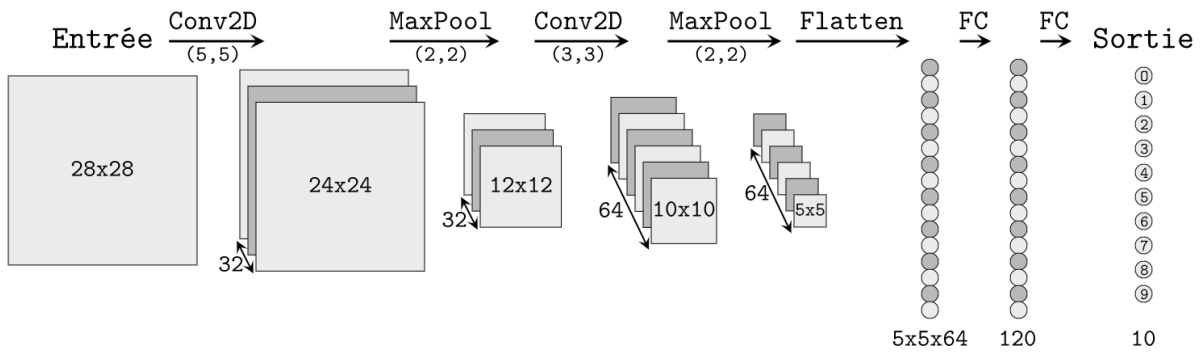


FIGURE 2 – Schéma de l'architecture utilisée



## Annexe C. Évolutions de $\|r\|$ et $Conf_c$ au cours d'attaques

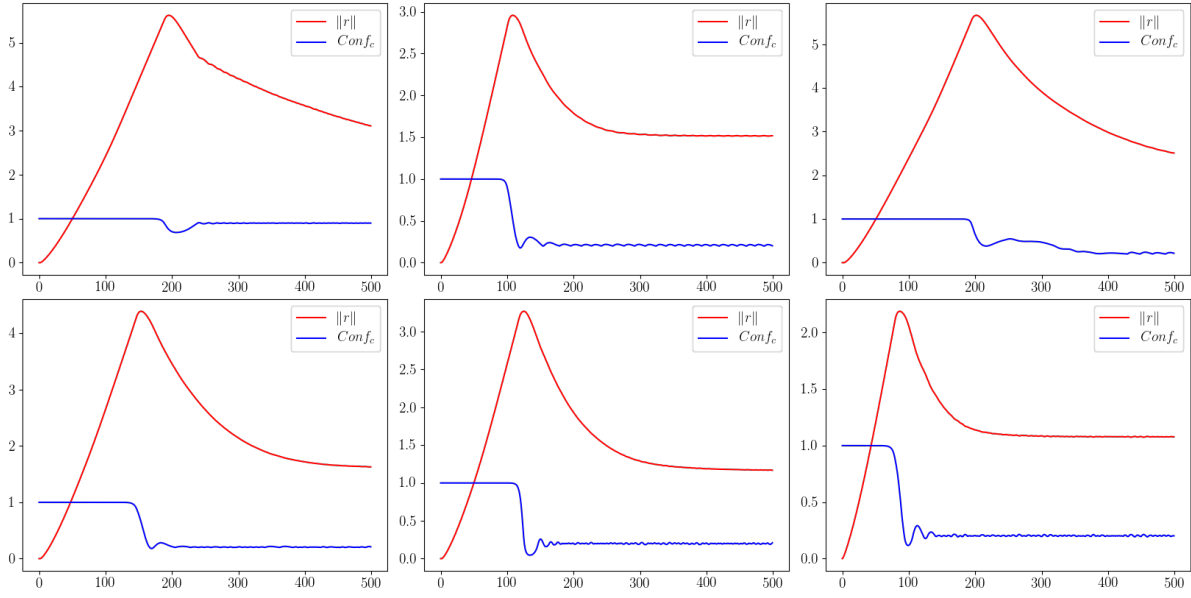


FIGURE 3 – Valeurs prises par  $\|r\|$  et  $Conf_c$  au cours de l'attaque de 6 images "difficiles" à attaquer

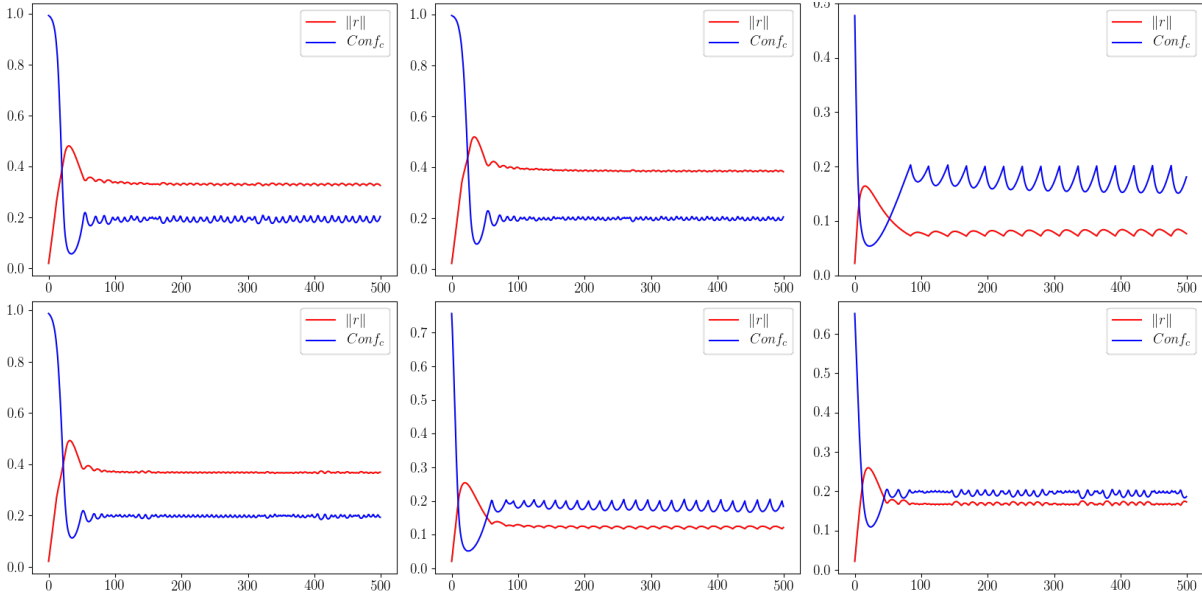


FIGURE 4 – valeurs prises par  $\|r\|$  et  $Conf_c$  au cours de l'attaque de 6 images "faciles" à attaquer

## Annexe D. Attaques utilisées (MNIST)

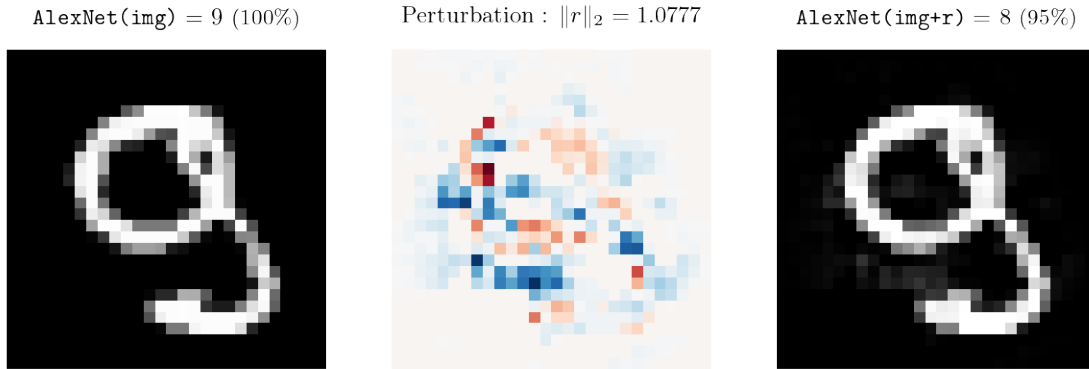


FIGURE 5 – `foolbox.attacks.GradientAttack`

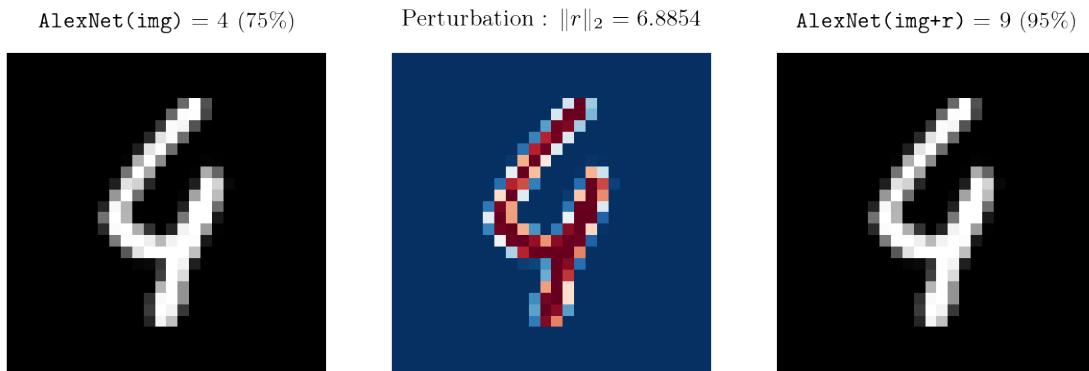


FIGURE 6 – `foolbox.attacks.ContrastReductionAttack`

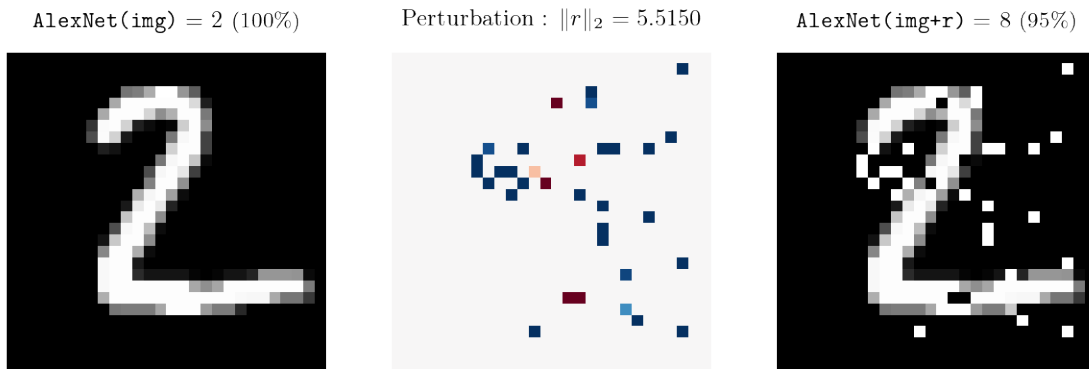


FIGURE 7 – `foolbox.attacks.PointwiseAttack`

## Annexe E. Attaques utilisées (FashionMNIST)

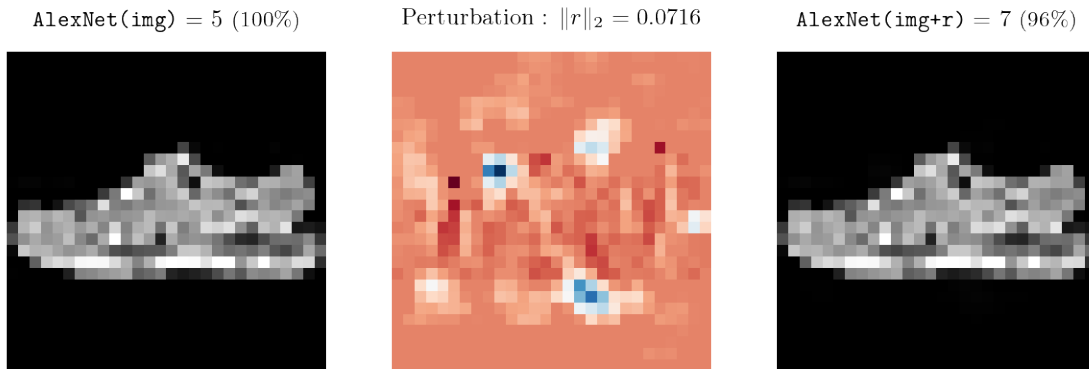


FIGURE 8 – `foolbox.attacks.GradientAttack`

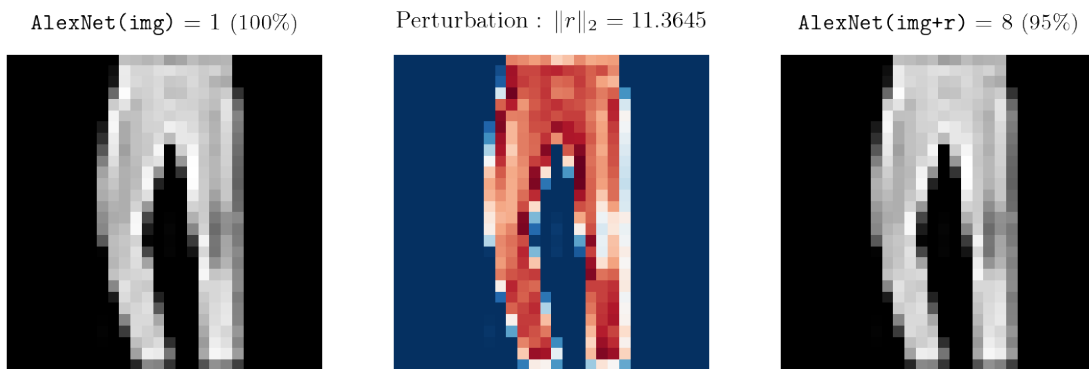


FIGURE 9 – `foolbox.attacks.ContrastReductionAttack`

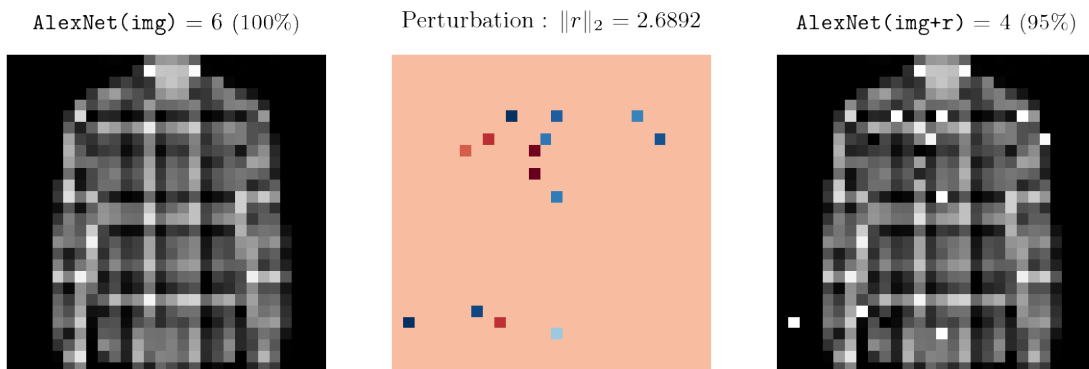


FIGURE 10 – `foolbox.attacks.PointwiseAttack`

## Annexe F. Code source

Tout le code est présent à l'adresse <https://github.com/maxdinech/resistance>.

On utilise la librairie **PyTorch** (qui permet des calculs efficaces de gradients sur des tenseurs) dans sa version 0.3 (la plus récente).

Le code est divisé en six fichiers :

- `basics.py`
- `data_loader.py`
- `architectures.py`
- `train.py`
- `accuracy.py`
- `attack.py`
- `plot.py`

On utilisera les modules et bibliothèques Python suivantes :

- `torch` (PyTorch)
- `torchvision`
- `foolbox`
- `matplotlib`
- `tkinter`
- `tqdm`
- `numpy`
- `cv2`

Ainsi que `texlive` pour la trace de courbes avec des titres en  $\text{\LaTeX}$ .

Ce code est fonctionnel sous Linux, et devrait normalement être compatible avec Windows (ce que je n'ai pas néanmoins pas vérifié en pratique).

### a) `basics.py`

Ce fichier contient quelques fonctions utilisées dans les autres fichiers : une fonction de chargement d'une architecture (c'est-à-dire une classe qui définit une architecture d'un réseau de neurones), et une fonction de chargement de *modèle*, c'est-à-dire un réseau de neurones déjà entraîné (donc son architecture + ses poids calculés lors de l'entraînement).

```

1  """
2  Basic PyTorch functions used in most of the other files.
3
4  """
5
6
7  import os
8  import warnings
9
10 import torch
11
```

```

12 import architectures
13
14
15 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
16
17
18 def load_architecture(model_name):
19     model = getattr(architectures, model_name)()
20     model = model.to(device)
21     return model
22
23
24 def load_model(dataset, model_name):
25     try:
26         path = os.path.join("../", "models", dataset, model_name + ".pt")
27         with warnings.catch_warnings():
28             # Ignores the compatibility warning between pytorch updates
29             from torch.serialization import SourceChangeWarning
30             warnings.simplefilter('ignore', SourceChangeWarning)
31             model = torch.load(path, map_location=lambda storage, loc: storage)
32             model = model.to(device)
33         return model
34     except FileNotFoundError:
35         raise ValueError('No trained model found.')

```

---

## b) data\_loader.py

Ce fichier permet de travailler simplement avec les bases de données MNIST et FashionMNIST en les téléchargeant automatiquement et en permettant un accès facilité.

```
1  """
2  Automatically creates and loads the MNIST and FashionMNIST datasets.
3  """
4
5
6  import os
7  import shutil
8
9  import torch
10 from torchvision import datasets, transforms
11
12
13 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
14
15
16 # Creates `train.pt` and `test.pt` from the specified dataset.
17 def create(dataset):
18     root = os.path.join '..', 'data'
19     if dataset == 'MNIST':
20         datasets.MNIST(root=root, train=True,
21                        transform=transforms.ToTensor(),
22                        download=True)
23     elif dataset == 'FashionMNIST':
24         datasets.FashionMNIST(root=root, train=True,
25                               transform=transforms.ToTensor(),
26                               download=True)
27     os.mkdir(os.path.join(root, dataset))
28     shutil.move(os.path.join(root, 'processed', 'training.pt'),
29               os.path.join(root, dataset, 'train.pt'))
30     shutil.move(os.path.join(root, 'processed', 'test.pt'),
31               os.path.join(root, dataset, 'test.pt'))
32     shutil.rmtree(os.path.join(root, 'raw'))
33     shutil.rmtree(os.path.join(root, 'processed'))
34
35
36 # Loads a subset from a dataset.
37 def load(dataset, subset, num_elements=None):
38     root = os.path.join '..', 'data'
39     if dataset in ['MNISTnorms', 'FashionMNISTnorms',
40                  'MNISTconfs', 'FashionMNISTconfs']:
41         file_name = subset + '_' + dataset[-5:] + '.pt'
42         path = os.path.join(root, dataset[:-5], file_name)
43         values, labels = torch.load(path)
44         return values, labels.long()
45     elif dataset in ['MNIST', 'FashionMNIST']:
46         path = os.path.join(root, dataset, subset + '.pt')
```

```

47     if not os.path.exists(path):
48         create(dataset)
49     images, labels = torch.load(path)
50     if num_elements:
51         images = images[:num_elements].clone()
52         labels = labels[:num_elements].clone()
53     images = images.float() / 255
54     labels = labels.long()
55     images = images.to(device)
56     labels = labels.to(device)
57     images = images.view(len(images), 1, 28, 28) # Channels first
58     return images, labels
59 else:
60     raise ValueError('Unknown dataset')
61
62
63 # Loads ((1 - val_split) * num_images) from `train` starting at position 0
64 def train(dataset, val_split, num_images=None):
65     images, labels = load(dataset, 'train', num_images)
66     num_images = num_images if num_images else len(images)
67     num_train = round((1-val_split) * num_images)
68     train_images = images[:num_train].clone()
69     train_labels = labels[:num_train].clone()
70     return train_images, train_labels
71
72
73 # Loads (val_split * num_images) from `train` starting at position `num_train`
74 def val(dataset, val_split, num_images=None):
75     images, labels = load(dataset, 'train', num_images)
76     num_images = num_images if num_images else len(images)
77     num_train = round((1-val_split) * num_images)
78     train_images = images[num_train:num_images].clone()
79     train_labels = labels[num_train:num_images].clone()
80     return train_images, train_labels
81
82
83 # Loads the test images
84 # The unused second argument gives the same type to the three functions.
85 def test(dataset, _, num_test=None):
86     return load(dataset, 'test', num_test)

```

---

### c) architectures.py

Ce fichier contient les descriptions d'architectures formelles de réseaux de neurones, auxquelles on a ajouté des paramètres d'entraînement (nombre d'époques, taux d'apprentissage, taille de mini-batch) efficaces, afin de pouvoir passer plus simplement d'une architecture à l'autre pendant la phase d'expérimentation.

```
1  """ architectures.py
2  Network architectures.
3
4  Default hyperparameters (learning rate, number of epochs, batch size,
5  optimizer) are included in the definition of the architectures for more
6  flexibility.
7  """
8
9
10 import torch
11 from torch import nn
12
13
14 class AlexNet(nn.Module):
15     def __init__(self):
16         super(AlexNet, self).__init__()
17         # Training hyperparameters
18         self.lr = 5e-4
19         self.epochs = 100
20         self.batch_size = 64
21         self.features = nn.Sequential(
22             nn.Conv2d(1, 32, kernel_size=5),
23             nn.ReLU(inplace=True),
24             nn.MaxPool2d(kernel_size=2),
25             nn.Conv2d(32, 64, kernel_size=3),
26             nn.ReLU(inplace=True),
27             nn.MaxPool2d(kernel_size=2)
28         )
29         self.classifier = nn.Sequential(
30             nn.Dropout(),
31             nn.Linear(64 * 5 * 5, 120),
32             nn.ReLU(inplace=True),
33             nn.Dropout(),
34             nn.Linear(120, 10),
35             nn.Softmax(dim=1)
36         )
37         # Optimizer and loss function
38         self.optimizer = torch.optim.Adam(self.parameters(), lr=self.lr)
39         self.loss_fn = nn.CrossEntropyLoss()
40
41     def forward(self, x):
42         x = self.features(x)
43         x = x.view(x.size(0), -1) # Flatten
44         x = self.classifier(x)
45         return x
```



```

46
47
48 class AlexNet_bn(nn.Module):
49     def __init__(self):
50         super(AlexNet_bn, self).__init__()
51         # Training hyperparameters
52         self.lr = 5e-4
53         self.epochs = 100
54         self.batch_size = 64
55         self.features = nn.Sequential(
56             nn.Conv2d(1, 32, kernel_size=5),
57             nn.ReLU(inplace=True),
58             nn.MaxPool2d(kernel_size=2),
59             nn.BatchNorm2d(32),
60             nn.Conv2d(32, 64, kernel_size=3),
61             nn.ReLU(inplace=True),
62             nn.MaxPool2d(kernel_size=2),
63             nn.BatchNorm2d(64)
64         )
65         self.classifier = nn.Sequential(
66             nn.Linear(64 * 5 * 5, 120),
67             nn.ReLU(inplace=True),
68             nn.BatchNorm1d(120),
69             nn.Linear(120, 10),
70             nn.Softmax(dim=1)
71         )
72         # Optimizer and loss function
73         self.optimizer = torch.optim.Adam(self.parameters(), lr=self.lr)
74         self.loss_fn = nn.CrossEntropyLoss()
75
76     def forward(self, x):
77         x = self.features(x)
78         x = x.view(x.size(0), -1) # Flatten
79         x = self.classifier(x)
80         return x
81
82
83 class VGG(nn.Module):
84     def __init__(self):
85         super(VGG, self).__init__()
86         # Training hyperparameters
87         self.lr = 1e-4
88         self.epochs = 40
89         self.batch_size = 32
90         self.features = nn.Sequential(
91             nn.Conv2d(1, 64, kernel_size=3),
92             nn.ReLU(inplace=True),
93             nn.Conv2d(64, 64, kernel_size=3),
94             nn.ReLU(inplace=True),
95             nn.MaxPool2d(kernel_size=2, stride=2),
96             nn.Conv2d(64, 128, kernel_size=3),
97             nn.ReLU(inplace=True),
98             nn.Conv2d(128, 128, kernel_size=3),

```

```

99         nn.ReLU(inplace=True),
100         nn.MaxPool2d(kernel_size=2, stride=2)
101     )
102     self.classifier = nn.Sequential(
103         nn.Linear(128 * 4 * 4, 4096),
104         nn.ReLU(True),
105         nn.Dropout(),
106         nn.Linear(4096, 4096),
107         nn.ReLU(True),
108         nn.Dropout(),
109         nn.Linear(4096, 10),
110         nn.Softmax(dim=1)
111     )
112     # Optimizer and loss function
113     self.optimizer = torch.optim.Adam(self.parameters(), lr=self.lr)
114     self.loss_fn = nn.CrossEntropyLoss()
115
116     def forward(self, x):
117         x = self.features(x)
118         x = x.view(x.size(0), -1) # Flatten
119         x = self.classifier(x)
120         return x
121
122
123 class VGG_bn(nn.Module):
124     def __init__(self):
125         super(VGG_bn, self).__init__()
126         # Training hyperparameters
127         self.lr = 1e-4
128         self.epochs = 40
129         self.batch_size = 32
130         self.features = nn.Sequential(
131             nn.Conv2d(1, 64, kernel_size=3),
132             nn.BatchNorm2d(64),
133             nn.ReLU(inplace=True),
134             nn.Conv2d(64, 64, kernel_size=3),
135             nn.BatchNorm2d(64),
136             nn.ReLU(inplace=True),
137             nn.MaxPool2d(kernel_size=2, stride=2),
138             nn.Conv2d(64, 128, kernel_size=3),
139             nn.BatchNorm2d(128),
140             nn.ReLU(inplace=True),
141             nn.Conv2d(128, 128, kernel_size=3),
142             nn.BatchNorm2d(128),
143             nn.ReLU(inplace=True),
144             nn.MaxPool2d(kernel_size=2, stride=2)
145         )
146         self.classifier = nn.Sequential(
147             nn.Linear(128 * 4 * 4, 4096),
148             nn.ReLU(True),
149             nn.Dropout(),
150             nn.Linear(4096, 4096),
151             nn.ReLU(True),

```

```
152         nn.Dropout(),
153         nn.Linear(4096, 10),
154         nn.Softmax(dim=1)
155     )
156     # Optimizer and loss function
157     self.optimizer = torch.optim.Adam(self.parameters(), lr=self.lr)
158     self.loss_fn = nn.CrossEntropyLoss()
159
160     def forward(self, x):
161         x = self.features(x)
162         x = x.view(x.size(0), -1) # Flatten
163         x = self.classifier(x)
164         return x
```

---

#### d) train.py

Ce fichier permet l'entraînement d'un réseau de neurones en gardant la plus grande souplesse possible : on passe en argument l'architecture utilisée, la base de donnée sur laquelle on travaille, le taux de validation (c'est-à-dire le pourcentage d'images utilisées pour l'entraînement, les autres servant à valider la généralisation des résultats), la taille du mini-batch, le nombre d'epochs...

```
----- train.py -----
1  """
2  Trains a specified architecture to classify a specified dataset.
3  The networks are defined in architectures.py
4
5  ---
6
7  usage: python3 train.py [-num NUM] [-split SPLIT] [-lr LR] [-e E] [-bs BS]
8  [-t T] [-S] model dataset
9
10 positional arguments:
11     model           Network architecture (defined in architectures.py)
12     dataset        Dataset used for training
13
14 optional arguments:
15     -num NUM       Number of images used (default: all)
16     -split SPLIT   Images proportion in val (default: 1/6)
17     -lr LR         Learning rate (default: value in the model class)
18     -e E           Num. of epochs (default: value in the model class)
19     -bs BS         batch size (default: value in the model class)
20     -k K           Top-k error metric (default: 1)
21     -S, --save     Saves the trained model (default: not saved)
22
23 Note: the -split s argument is used to divide the training images into `train`
24 and `val`, with proportions respectively `s` and `(1-s)`.
25
26  """
27
28
29  import os
30  import sys
31  import argparse
32
33  import torch
34  from torch.utils.data import DataLoader, TensorDataset
35  import matplotlib.pyplot as plt
36  from tqdm import tqdm
37
38  from basics import load_architecture
39  import data_loader
40  import plot
41
42
43  # Parameters parsing
44  parser = argparse.ArgumentParser()
```

```

45 parser.add_argument("model", type=str,
46                     help="Network architecture (defined in architectures.py)")
47 parser.add_argument("dataset", type=str,
48                     help="Dataset used for training")
49 parser.add_argument("-num", type=int,
50                     help="Number of images used (default: all)")
51 parser.add_argument("-split", type=str, default="1/6",
52                     help="Images proportion in val (default: 1/6)")
53 parser.add_argument("-lr", type=float,
54                     help="Learning rate (default: value in the model class)")
55 parser.add_argument("-e", type=int,
56                     help="Num. of epochs (default: value in the model class)")
57 parser.add_argument("-bs", type=int,
58                     help="batch size (default: value in the model class)")
59 parser.add_argument("-k", type=int, default=1,
60                     help="Top-k error metric (default: 1)")
61 parser.add_argument("-S", "--save", action="store_true",
62                     help="Saves the trained model (default: not saved)")
63 args = parser.parse_args()
64
65 model_name = args.model
66 dset_name = args.dataset
67 num_img = args.num
68 val_split = eval(args.split) # Allows to pass fractions in parameters
69 k = args.k
70 save_model = args.save
71
72
73 # Model instantiation
74 model = load_architecture(model_name)
75
76
77 # Loads model hyperparameters (if not specified in args)
78 batch_size = args.bs if args.bs else model.batch_size
79 lr = args.lr if args.lr else model.lr
80 epochs = args.e if args.e else model.epochs
81
82
83 # Loads model functions
84 loss_fn = model.loss_fn
85 optimizer = model.optimizer
86
87
88 # Loads the training database, and splits it in into `train` and `val`.
89 train_images, train_labels = data_loader.train(dset_name, val_split, num_img)
90 val_images, val_labels = data_loader.val(dset_name, val_split, num_img)
91 num_train = len(train_images)
92 num_val = len(val_images)
93
94
95 # DataLoader of the `train` images
96 train_loader = DataLoader(TensorDataset(train_images, train_labels),
97                           batch_size=batch_size,

```

```

98                 shuffle=True)
99
100 num_batches = len(train_loader)
101
102
103 # Computes the Top-k accuracy of the model.
104 # (computing the accuracy mini-batch after mini-batch avoids memory overload)
105 def accuracy(images, labels, k=1):
106     data = TensorDataset(images, labels)
107     loader = DataLoader(data, batch_size=10, shuffle=False)
108     count = 0
109     for (x, y) in loader:
110         y_pred = model.eval()(x)
111         y_pred_k = y_pred.topk(k, 1, True, True)[1]
112         count += sum(sum((y_pred_k.t() == y).double()))
113         # .double(): ByteTensor sums are limited at 255.
114     return 100 * count / len(images)
115
116
117 # Computes the loss of the model.
118 # (computing the loss mini-batch after mini-batch avoids memory overload)
119 def big_loss(images, labels):
120     data = TensorDataset(images, labels)
121     loader = DataLoader(data, batch_size=100, shuffle=False)
122     count = 0
123     for (x, y) in loader:
124         y_pred = model.eval()(x)
125         count += len(x) * loss_fn(y_pred, y).item()
126     return count / len(images)
127
128
129 # NETWORK TRAINING
130 # -----
131
132 # Prints the hyperparameters before the training.
133 print(f"Train on {num_train} samples, val on {num_val} samples.")
134 print(f"Epochs: {epochs}, batch size: {batch_size}")
135 optimizer_name = type(optimizer).__name__
136 print(f"Optimizer: {optimizer_name}, learning rate: {lr}")
137 num_parameters = sum(param.numel() for param in model.parameters())
138 print(f"Parameters: {num_parameters}")
139 print(f"Save model : {save_model}\n")
140
141
142 # Custom progress bar.
143 def bar(data, e):
144     epoch = f"Epoch {e+1}/{epochs}"
145     left = "{desc}: {percentage:3.0f}%"
146     right = "{elapsed} - ETA:{remaining} - {rate_fmt}"
147     bar_format = left + " |{bar}| " + right
148     return tqdm(data, desc=epoch, ncols=74, unit='b', bar_format=bar_format)
149
150

```

```

151 train_accs, val_accs = [], []
152 train_losses, val_losses = [], []
153
154 try:
155     # Main loop over each epoch
156     for e in range(epochs):
157
158         # Secondary loop over each mini-batch
159         for (x, y) in bar(train_loader, e):
160
161             # Computes the network output
162             y_pred = model.train()(x)
163             loss = loss_fn(y_pred, y)
164
165             # Optimizer step
166             model.zero_grad()
167             loss.backward()
168             optimizer.step()
169
170             # Calculates accuracy and loss on the train database.
171             train_acc = accuracy(train_images, train_labels, k)
172             train_loss = big_loss(train_images, train_labels)
173             train_accs.append(train_acc)
174             train_losses.append(train_loss)
175
176             # Calculates accuracy and loss on the validation database.
177             val_acc = accuracy(val_images, val_labels, k)
178             val_loss = big_loss(val_images, val_labels)
179             val_accs.append(val_acc)
180             val_losses.append(val_loss)
181
182             # Prints the losses and accs at the end of each epoch.
183             print(f" |-> train_acc@{k}: {train_acc:5.2f}%",
184                   f" -- train_loss: {train_loss:6.4f}")
185             print(f" |-> val_acc@{k}: {val_acc:5.2f}%",
186                   f" -- val_loss: {val_loss:6.4f}\n")
187
188     # Allows to manually interrupt the training (early stopping).
189     except KeyboardInterrupt:
190         pass
191
192     # Saves the network if stated.
193     if save_model:
194         path = os.path.join("../", "models", dset_name, model_name + ".pt")
195         torch.save(model, path)
196         # Saves the accs history graph
197         plot.train_history(train_accs, val_accs)
198         plt.savefig(path + model_name + ".png", transparent=True)

```

---

## e) accuracy.py

Ce fichier, qui reprend la même syntaxe de paramètres que le précédent, sert à évaluer la performance d'un modèle (réseau entraîné).

---

```
1  """
2  Computes the Top-k error of a trained model, over a given subset of a dataset.
3
4  ---
5
6  usage: python3 accuracy.py [-k K] model dataset subset
7
8  positional arguments:
9    model          Trained model to evaluate
10   dataset        Dataset used for training
11   subset         Subset to calculate the error on
12
13  optional arguments:
14    -split SPLIT  Images proportion in val (default: 1/6)
15    -k K          Top-k error metric (default: 1)
16
17  Note: to compute the accuracy on `train` or `val`, the -split argument is
18  needed in order to determine the right repartition of all training images into
19  `train` and `val`. To compute the accuracy over all training images at once,
20  use:
21      python3 accuracy.py [-k K] model dataset train -split 0
22
23  """
24
25
26  import sys
27  import argparse
28
29  import torch
30  from torch.utils.data import DataLoader, TensorDataset
31  from tqdm import tqdm
32
33  from basics import load_model
34  import data_loader
35
36
37  # Parameters parsing
38  parser = argparse.ArgumentParser()
39  parser.add_argument("model", type=str, help="Trained model to evaluate")
40  parser.add_argument("dataset", type=str, help="Dataset used for training")
41  parser.add_argument("subset", type=str, help="Subset to compute the error on")
42  parser.add_argument("-split", type=str, default="1/6",
43                      help="Images proportion in val (default: 1/6)")
44  parser.add_argument("-k", type=int, default=1,
45                      help="Top-k error metric (default: 1)")
46  args = parser.parse_args()
```



```

47
48 model_name = args.model
49 dset_name = args.dataset
50 subset = args.subset
51 val_split = eval(args.split) # Allows to pass fractions in parameters
52 k = args.k
53
54
55 # Loads the model
56 model = load_model(dset_name, model_name)
57
58
59 # Loads the specified subset from the dataset.
60 images, labels = getattr(data_loader, subset)(dset_name, val_split)
61
62
63 # Custom progress bar.
64 def bar(data):
65     bar_format = "{percentage:3.0f}% |{bar}| {elapsed} - ETA:{remaining}"
66     return tqdm(data, ncols=74, bar_format=bar_format)
67
68
69 # Computes the Top-k accuracy of the model.
70 # (computing the accuracy mini-batch after mini-batch avoids memory overload)
71 def accuracy(images, labels, k=1):
72     data = TensorDataset(images, labels)
73     loader = DataLoader(data, batch_size=100, shuffle=False)
74     count = 0
75     position = 0
76     for (x, y) in bar(loader):
77         position += len(x)
78         y_pred = model.eval()(x)
79         y_pred_k = y_pred.topk(k, 1, True, True)[1]
80         count += sum(sum((y_pred_k.t() == y).double()))
81         # .double(): ByteTensor sums are limited at 255.
82     return 100 * count / len(images)
83
84
85 # Prints the losses and accuracies at the end of each epoch.
86 print(f"Computing the Top-{k} error on the {len(images)} {subset} images...")
87 acc = accuracy(images, labels, k)
88 error = 100 - acc
89 print(f"Top-{k} error: {error:0.2f}%")

```

---

## f) attack.py

Ce fichier sert à exécuter les attaques adversaires par descente de gradient itérée, et calculer la résistance d'une image; ainsi que les attaques adversaires de la bibliothèque Foolbox.

```
1  """
2  Adversarial attacks
3
4  ---
5
6  usage: python3 -i attack.py model dataset
7
8  positional arguments:
9      model          Trained model to evaluate
10     dataset       Dataset used for training
11
12  """
13
14
15  import os
16  import sys
17  import shutil
18  import argparse
19
20  import torch
21  from torch import nn
22  import matplotlib.pyplot as plt
23  import foolbox
24  import numpy as np
25  import cv2
26
27  from basics import load_model
28  import data_loader
29  import plot
30
31
32  device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
33
34
35  # Parameters parsing
36  parser = argparse.ArgumentParser()
37  parser.add_argument("model", type=str, help="Trained model to evaluate")
38  parser.add_argument("dataset", type=str, help="Dataset used for training")
39  args = parser.parse_args()
40
41  model_name = args.model
42  dataset_name = args.dataset
43
44
45  # Loads the model
46  model = load_model(dataset_name, model_name).to(device)
```

```

47
48
49 # Loads the specified subset from the specified database
50 images, labels = data_loader.test(dataset_name, None)
51
52
53 # BASIC FUNCTIONS
54 # -----
55
56 # Loads an image from the webcam
57 def webcam():
58     cap = cv2.VideoCapture(0)
59     ret, frame = cap.read()
60     cap.release()
61     gray = cv2.cvtColor(frame, cv2.COLOR_BGR2GRAY)
62     x, y = len(gray[0]), len(gray)
63     d = int((x-y) / 2)
64     gray = np.array([line[d:y-d] for line in gray])
65     gray = cv2.resize(gray, (28, 28)) / 255
66     return torch.Tensor(gray).view(1, 1, 28, 28).to(device)
67
68
69 # Loads the #img_id image from the test database (torch.Tensor).
70 def load_image(img_id):
71     return images[img_id].view(1, 1, 28, 28).to(device)
72
73
74 # Loads the #img_id label from the test database (int).
75 def load_label(img_id):
76     return labels[img_id].item()
77
78
79 # Returns the label prediction of an image (float).
80 def prediction(image):
81     return model.eval()(image).max(1)[1].item()
82
83
84 # Returns the predictions of of the network on an image (ndarray).
85 def predictions(image):
86     return model.eval()(image).data[0].numpy()
87
88
89 # Returns the confidence of the network that the image is `digit`.
90 def confidence(image, category):
91     return model.eval()(image)[0, category].item()
92
93
94 # Yields the indices of the first n wrong predictions.
95 def errors(n=len(images)):
96     i = 0
97     l = len(images)
98     while i < l and n > 0:
99         image, label = load_image(i), load_label(i)

```

```

100         if prediction(image) != label:
101             yield i
102             n -= 1
103         i += 1
104
105
106     # Yields the indices of the first n correct predictions.
107     def not_errors(n=len(images)):
108         i = 0
109         l = len(images)
110         while i < l and n > 0:
111             image, label = load_image(i), load_label(i)
112             if prediction(image) == label:
113                 yield i
114                 n -= 1
115             i += 1
116
117
118     # ATTACK FUNCTIONS
119     # -----
120
121     class Attacker(nn.Module):
122         def __init__(self, p, lr):
123             super(Attacker, self).__init__()
124             self.p = p
125             self.r = nn.Parameter(torch.zeros(1, 1, 28, 28))
126             self.optimizer = torch.optim.Adam(self.parameters(), lr=lr)
127
128         def forward(self, x):
129             return (x + self.r).clamp(0, 1)
130
131         def loss_fn(self, image, digit):
132             adv = self.forward(image)
133             conf = model(adv)[0, digit]
134             norm = (adv - image).abs().pow(self.p).sum()
135             if conf < 0.2:
136                 return norm
137             elif conf < 0.9:
138                 return conf + norm
139             else:
140                 return conf - norm
141
142
143     def GDA(image, steps=500, p=2, lr=1e-3):
144         norms, confs = [], []
145         digit = prediction(image)
146         attacker = Attacker(p, lr)
147         attacker = attacker.to(device)
148         optim = attacker.optimizer
149         for i in range(steps):
150             # Training step
151             loss = attacker.loss_fn(image, digit)
152             attacker.zero_grad()

```

```

153     loss.backward()
154     optim.step()
155     # Prints results
156     adv = attacker.forward(image)
157     conf = confidence(adv, digit)
158     norm = (adv - image).norm(p).item()
159     print(f"Step {i:4} -- conf: {conf:0.4f}, L_{p}(r): {norm:0.20f}",
160           end='\r')
161     norms.append(norm)
162     confs.append(conf)
163     print()
164     success = prediction(adv) != digit
165     return(success, adv, norms, confs)
166
167
168 def GDA_break(image, max_steps=500, p=2, lr=1e-3):
169     norms, confs = [], []
170     digit = prediction(image)
171     attacker = Attacker(p, lr)
172     attacker = attacker.to(device)
173     optim = attacker.optimizer
174     adv = attacker.forward(image)
175     steps = 0
176     while confidence(adv, digit) >= 0.2 and steps < max_steps:
177         steps += 1
178         # Training step
179         loss = attacker.loss_fn(image, digit)
180         attacker.zero_grad()
181         loss.backward()
182         optim.step()
183         # Prints results
184         adv = attacker.forward(image)
185         conf = confidence(adv, digit)
186         norm = (adv - image).norm(p).item()
187         print(f"Step {steps:4} -- conf: {conf:0.4f}, L_{p}(r): {norm:0.10f}",
188               end='\r')
189         norms.append(norm)
190         confs.append(conf)
191     print()
192     return(steps, adv, norms, confs)
193
194
195 def GDA_graph(img, steps=500, p=2, lr=1e-3):
196     success, adv, norms, confs = GDA(img, steps, p, lr)
197     plot.attack_history(norms, confs)
198     plt.show()
199     path = os.path.join(".", "results", "last_attack_history.png")
200     plt.savefig(path, transparent=True)
201     confc = lambda i: confidence(i, prediction(img))
202     if success:
203         print("\nAttack succeeded")
204         img_pred = prediction(img)
205         img_conf = confidence(img, img_pred)

```

```

206     adv_pred = prediction(adv)
207     adv_conf = confidence(adv, adv_pred)
208     plot.attack_result(model_name, p,
209                        img, img_pred, img_conf,
210                        adv, adv_pred, adv_conf)
211     path = os.path.join("../", "results", "last_attack_result.png")
212     plt.savefig(path, transparent=True)
213     plt.show()
214 else:
215     print("\nAttack failed")
216
217
218 def GDA_break_graph(img, max_steps=500, p=2, lr=1e-3):
219     success, adv, norms, confs = GDA_break(img, max_steps, p, lr)
220     plot.attack_history(norms, confs)
221     path = os.path.join("../", "results", "last_attack_history.png")
222     plt.savefig(path, transparent=True)
223     plt.show()
224     img_pred = prediction(img)
225     img_conf = confidence(img, img_pred)
226     adv_pred = prediction(adv)
227     adv_conf = confidence(adv, adv_pred)
228     plot.attack_result(model_name, p,
229                        img, img_pred, img_conf,
230                        adv, adv_pred, adv_conf)
231     path = os.path.join("../", "results", "last_attack_result.png")
232     plt.savefig(path, transparent=True)
233     plt.show()
234
235
236 # RESISTANCE FUNCTIONS
237 # -----
238
239 # A resistance value greater than 1000 is not possible.
240 # Which is why 10000 will represent infinity when the attack fails.
241
242 def resistance_N(image, steps=500):
243     success, _, norms, _ = GDA(image, steps)
244     if success:
245         return norms[-1]
246     return 10000
247
248
249 def resistance_max(image, steps=500):
250     success, _, norms, _ = GDA(image, steps)
251     if success:
252         return max(norms)
253     return 10000
254
255
256 def resistance_min(image, max_steps=500):
257     steps = attack_break(image, max_steps)[0]
258     if steps < max_steps:

```

```

259         return steps
260     return 10000
261
262
263 # Computes the N-resistance, max_resistance and min_resistance in a single pass
264 def resistances(image, steps=500):
265     success, _, norms, confs = GDA(image, steps)
266     if success:
267         res_N = norms[-1]
268         res_max = max(norms)
269         res_min = 1 + next((i for i, c in enumerate(confs) if c <= 0.2), steps)
270         return (res_N, res_max, res_min)
271     else:
272         return (10000, 10000, 10000)
273
274
275 # Computes the N-resistance, max_resistance and min_resistance of an image list
276 def resistances_list(images_list, steps=500):
277     L_res_N, L_res_max, L_res_min = [], [], []
278     i, l = 1, len(images_list)
279     for image in images_list:
280         print(f"{i}/{l} : ")
281         res_N, res_max, res_min = resistances(image, steps)
282         L_res_N += [res_N]
283         L_res_max += [res_max]
284         L_res_min += [res_min]
285         i += 1
286     return (L_res_N, L_res_max, L_res_min)
287
288
289 # FOOLBOX ATTACKS
290 # -----
291
292
293 fmodel = foolbox.models.PyTorchModel(model, (0, 1), num_classes=10,
294                                     channel_axis=1,
295                                     cuda=torch.cuda.is_available())
296
297
298 # Runs the specified attack on an image, ensuring that the confidence of the
299 # network's prediction on the adversarial example is above `p`.
300 def foolbox_attack(img, attack_name, p=0.95):
301     path = f"../results/{dataset_name}/{attack_name}/"
302     if not os.path.exists(path):
303         os.mkdir(path)
304     img_pred = prediction(img)
305     img_conf = confidence(img, img_pred)
306     try:
307         # Finds an adversarial example
308         attack = getattr(foolbox.attacks, attack_name)(fmodel)
309         np_adv = attack(np.array(img).reshape(1, 28, 28), img_pred)
310         adv = torch.Tensor(np_adv).view(1, 1, 28, 28).to(device)
311         adv_pred = prediction(adv)

```

```

312         # Increases its classification probability above p
313         criterion = foolbox.criteria.TargetClassProbability(adv_pred, p)
314         attack = getattr(foolbox.attacks, attack_name)(fmodel, criterion)
315         np_adv = attack(np.array(img).reshape(1, 28, 28), img_pred)
316         adv = torch.Tensor(np_adv).view(1, 1, 28, 28).to(device)
317         adv_pred = prediction(adv)
318         adv_conf = confidence(adv, adv_pred)
319         plot.attack_result(model_name, 2,
320                           img, img_pred, img_conf,
321                           adv, adv_pred, adv_conf)
322         plt.show()
323         shutil.move("../results/latest/attack_result.png",
324                     path + f"{img_id:04d}.png")
325         torch.save(adv, path + f"{img_id:04d}.pt")
326         return adv
327     except KeyboardInterrupt:
328         break
329     except:
330         pass
331
332
333 def foolbox_attacks_list(size, attack_name, p=0.95):
334     adv_list = []
335     for img_id in not_errors():
336         if len(adv_list) >= size or img_id >= 10000:
337             break
338         print(len(adv_list), "/", img_id, end='\r')
339         img = load_image(img_id)
340         adv = foolbox_attack(img, attack_name, p)
341         if adv is not None:
342             adv_list.append(adv)
343     return adv_list

```

---



g) plot.py

Ce fichier permet de tracer les différentes courbes de ce document.

```
1  """
2  Plotting functions used in attack.py
3  """
4
5
6  import matplotlib.pyplot as plt
7  from matplotlib import rcParams
8  from tqdm import tqdm
9
10
11  rcParams['text.usetex'] = True
12  rcParams['text.latex.unicode'] = True
13  rcParams['font.family'] = "serif"
14  rcParams['font.serif'] = "cm"
15
16
17  # Plots an image (Variable)
18  def plot_image(image):
19      plt.imshow(image.data.view(28, 28).numpy(), cmap='gray')
20
21
22  # Plots and saves the comparison graph of an adversarial image
23  def attack_result(model_name, p,
24                    img, img_pred, img_conf,
25                    adv, adv_pred, adv_conf):
26      model_name = model_name.replace('_', '\\_') # Escapes '_' characters
27      r = (adv - img)
28      norm = r.norm(p)
29      # Matplotlib settings
30      rcParams['axes.titlepad'] = 10
31      rcParams['font.size'] = 8
32      fig = plt.figure(figsize=(7, 2.5), dpi=180)
33      # Image
34      ax1 = fig.add_subplot(1, 3, 1)
35      ax1.imshow(img.data.view(28, 28).cpu().numpy(), cmap='gray')
36      plt.title(f"\\texttt{{{model_name}}}{img}} = {img_pred} \\small{{{100*img_conf:0.0f}\\%}}")
37      plt.axis('off')
38      # Perturbation
39      ax2 = fig.add_subplot(1, 3, 2)
40      ax2.imshow(r.data.view(28, 28).cpu().numpy(), cmap='RdBu')
41      plt.title(f"Perturbation : $\\Vert r \\Vert_{{{p}}} = {norm:0.4f}$")
42      plt.axis('off')
43      # Adversarial image
44      ax3 = fig.add_subplot(1, 3, 3)
45      ax3.imshow(adv.data.view(28, 28).cpu().numpy(), cmap='gray')
46      plt.title(f"\\texttt{{{model_name}}}{img+r}} = {adv_pred} \\small{{{100*adv_conf:0.0f}\\%}}")
47      plt.axis('off')
```

```

48     # Save and plot
49     fig.tight_layout(pad=1)
50     plt.subplots_adjust(left=0.05, right=0.95, top=0.80, bottom=0.05)
51     plt.savefig("../results/latest/attack_result.png", transparent=True)
52
53
54 # Plots the history of a model training
55 def train_history(train_accs, val_accs):
56     rcParams['font.size'] = 12
57     t = list(range(len(train_accs)))
58     plt.plot(t, train_accs, 'r')
59     plt.plot(t, val_accs, 'b')
60     plt.title("Network training history")
61     plt.legend(["train accuracy", "val accuracy"])
62
63
64 # Plots the history of an attack
65 def attack_history(norms, confs):
66     rcParams['font.size'] = 14
67     t = list(range(len(norms)))
68     plt.plot(t, norms, 'r')
69     plt.plot(t, confs, 'b')
70     plt.legend([" $\text{Vert } r$ ", " $\text{Conf}_c$ "])
71     plt.savefig("../results/latest/attack_history.png", transparent=True)

```

---