

# **Rapport sur la Défense contre les Exemples Adversariaux**

Kouamé Gérard Kra

Mercredi, 15 Janvier 2025

# 1. Introduction

Les exemples adversariaux représentent une des principales menaces pour les systèmes d'apprentissage automatique. Ces exemples sont des entrées artificiellement modifiées de manière imperceptible pour l'œil humain, mais suffisamment perturbées pour induire en erreur les modèles de classification. Par exemple, une image de chiffre manuscrit classée correctement par un modèle peut, après une petite perturbation, être catégorisée de manière erronée. Ces attaques adversariales exploitent les vulnérabilités des modèles de deep learning, soulevant des questions cruciales en matière de sécurité et de fiabilité des systèmes. Par exemple, les exemples adversariaux, tels que décrits par *Goodfellow et al. (2015)* dans leur travail sur le **Fast Gradient Sign Method (FGSM)**, représentent une menace significative pour les systèmes de deep learning. Ces exemples exploitent les gradients du modèle pour créer des perturbations imperceptibles qui peuvent tromper les modèles.

Avec l'adoption croissante de l'apprentissage automatique dans des domaines critiques comme la vision par ordinateur, les véhicules autonomes et la cybersécurité, il devient impératif de garantir la robustesse des modèles face à ces attaques. Ce projet s'inscrit dans cette problématique et vise à concevoir un mécanisme de défense capable de contrer des attaques adversariales variées tout en maintenant de bonnes performances sur des exemples propres. En combinant des stratégies comme l'entraînement adversarial et l'utilisation de multiples types d'attaques pendant l'apprentissage, ce travail explore comment renforcer les modèles et les rendre plus résilients face à ces menaces.

## 2. Méthodologie

Le jeu de données MNIST, constitué d'images de chiffres manuscrits, a été utilisé.

### 2.1. Architecture du Modèle

Le modèle utilise l'architecture suivante :

- Deux couches convolutionnelles (32 et 64 filtres respectivement).
- Deux couches de max pooling.
- Deux couches linéaires (64 et 512 neurones).
- Une couche de sortie avec 10 neurones pour la classification. Le modèle a d'abord été entraîné sur des données propres, puis affiné grâce à un entraînement adversarial.

### 2.2. Attaques Adversariales

Trois types d'attaques ont été implémentés :

1. **FGSM (Fast Gradient Sign Method)** : Une attaque rapide basée sur le gradient. Développée par *Goodfellow et al. (2015)*, cette méthode exploite le gradient de la fonction de perte pour générer une perturbation en un seul pas.
2. **BIM (Basic Iterative Method)** : Une version itérative de FGSM.
3. **PGD (Projected Gradient Descent)** : Une attaque puissante qui constitue une généralisation de BIM. Selon *Madry et al. (2017)*, le Projected Gradient Descent est une extension plus puissante et itérative de FGSM, largement reconnue pour tester la robustesse des modèles.

### 2.3. Entraînement Adversarial

*Carlini et al. (2019)* ont proposé des critères stricts pour tester l'efficacité des défenses contre les attaques adversariales, ce qui guide notre méthodologie.

Ainsi, pour améliorer la robustesse, le modèle a été entraîné sur un mélange d'exemples propres et adversariaux. Les attaques FGSM, BIM, et PGD ont été utilisées avec des paramètres variés pour garantir une diversité et éviter le surapprentissage à une attaque spécifique.

## 3. Résultats et Analyse

### 3.1. Performance du Modèle Baseline

Le modèle baseline a été évalué sur des exemples propres et adversariaux. Les résultats sont présentés dans le tableau ci-dessous :

Type d'Exemple	Précision (%)
Exemples Propres	98.75
FGSM ( $\epsilon = 0.3$ )	9.12
BIM ( $\epsilon = 0.3$ , 40 itérations)	0.00
PGD ( $\epsilon = 0.3$ , 40 itérations)	0.00

TABLE 1 – Performance du modèle baseline.

### 3.2. Performance Après Entraînement Adversarial

Après entraînement adversarial, les précisions sur des exemples propres et adversariaux sont :

Type d'Exemple	Précision (%)
Exemples Propres	99.03
FGSM ( $\epsilon = 0.3$ )	84.44
BIM ( $\epsilon = 0.3$ , 40 itérations)	56.85
PGD ( $\epsilon = 0.3$ , 40 itérations)	58.80

TABLE 2 – Performance après entraînement adversarial.

### 3.3. Pourquoi le modèle est-il robuste ?

La robustesse est démontrée en suivant les recommandations de *Madry et al. (2017)* pour diversifier les exemples adversariaux pendant l'entraînement, et les critères stricts définis par *Carlini et al. (2019)* pour évaluer les performances face à des attaques itératives comme PGD. Plus précisément, le modèle devient robuste grâce aux facteurs suivants :

- Diversité des exemples adversariaux** : L'entraînement intègre des attaques FGSM, BIM, et PGD avec des paramètres variés, couvrant un large spectre de perturbations adversariales.
- Mélange équilibré d'exemples** : En mélangeant des exemples propres et adversariaux dans chaque lot d'entraînement, le modèle apprend à bien généraliser tout en renforçant sa résistance aux attaques.
- Entraînement progressif** : L'augmentation graduelle des paramètres adversariaux (comme epsilon et alpha) pendant l'entraînement permet au modèle de s'adapter sans surapprentissage.
- Approche itérative** : L'utilisation d'attaques itératives comme BIM et PGD lors de l'entraînement prépare le modèle à des attaques similaires, qui sont souvent plus puissantes que FGSM.

## 4. Conclusion

Les résultats de ce projet montrent une avancée significative dans la conception d'un modèle robuste face aux exemples adversariaux. Grâce à l'utilisation d'un entraînement adversarial diversifié, le modèle a réussi à maintenir une précision élevée sur des données propres (99,03%) tout en améliorant considérablement sa résistance aux attaques FGSM, BIM et PGD. La précision atteinte sur FGSM (84,44%) et les améliorations notables sur BIM (56,87%) et PGD (58,80%) démontrent l'efficacité des stratégies mises en œuvre.

Ces performances témoignent de l'importance d'une approche équilibrée et progressive, mêlant des exemples adversariaux variés et des données propres, pour garantir une généralisation robuste. De plus, l'intégration de techniques d'attaques itératives a permis de mieux préparer le modèle aux perturbations plus complexes, tout en minimisant le compromis sur les performances globales.

En conclusion, ce projet met en évidence l'efficacité de l'entraînement adversarial comme stratégie défensive clé contre les attaques adversariales. Les résultats obtenus ouvrent également la voie à des recherches futures, telles que l'exploration d'autres types d'attaques ou la combinaison avec des régularisations supplémentaires pour encore renforcer la robustesse.

## 5. Références

1. Madry et al, (2017). *Towards deep learning models resistant to adversarial attacks* .
2. Goodfellow et al (2015). *Explaining and harnessing adversarial examples*.
3. N. Carlini et al., *On evaluating adversarial robustness*, 2019.

## Annexe

Implémentation en python avec Tensorflow :

```
import tensorflow as tf
import numpy as np
from tensorflow.keras.datasets import mnist
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Conv2D, MaxPooling2D, Flatten, Dense, Input
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.losses import SparseCategoricalCrossentropy

# Load and preprocess the MNIST dataset
(x_train, y_train), (x_test, y_test) = mnist.load_data()

x_train = x_train.astype('float32') / 255.0
x_test = x_test.astype('float32') / 255.0

x_train = np.expand_dims(x_train, axis=-1)
x_test = np.expand_dims(x_test, axis=-1)

# Model architecture
def build_model():
    model = Sequential([
        Input(shape=(28, 28, 1)),
        Conv2D(32, (5, 5), activation='relu', use_bias=False),
        MaxPooling2D(pool_size=(2, 2), strides=(2, 2)),
```

```

        Conv2D(64, (3, 3), activation='relu', use_bias=False),
        MaxPooling2D(pool_size=(2, 2), strides=(2, 2)),
        Flatten(),
        Dense(64, activation='relu', use_bias=False),
        Dense(512, activation='relu', use_bias=False),
        Dense(10, activation='softmax')
    ])
return model

# Compile the model
model = build_model()
model.compile(optimizer=Adam(learning_rate=0.001),
              loss=SparseCategoricalCrossentropy(),
              metrics=['accuracy'])

# Train the baseline model
print("Training the baseline model...")
history = model.fit(x_train, y_train, epochs=5, batch_size=64, validation_split=0.1)

# Evaluate baseline model accuracy
baseline_accuracy = model.evaluate(x_test, y_test, verbose=0)
print(f"Baseline model accuracy on clean examples: {baseline_accuracy[1] * 100:.2f}%")

# Define adversarial attack functions
def fgsm_attack(model, x, y, epsilon):
    x = tf.convert_to_tensor(x)
    y = tf.convert_to_tensor(y)
    with tf.GradientTape() as tape:
        tape.watch(x)
        prediction = model(x)
        loss = tf.keras.losses.sparse_categorical_crossentropy(y, prediction)
    gradient = tape.gradient(loss, x)
    perturbation = epsilon * tf.sign(gradient)
    x_adv = tf.clip_by_value(x + perturbation, 0, 1)
    return x_adv

def bim_attack(model, x, y, epsilon, alpha, iterations):
    x = tf.convert_to_tensor(x)
    y = tf.convert_to_tensor(y)
    x_adv = tf.identity(x)
    for i in range(iterations):
        x_adv = fgsm_attack(model, x_adv, y, alpha)
        perturbation = tf.clip_by_value(x_adv - x, -epsilon, epsilon)
        x_adv = tf.clip_by_value(x + perturbation, 0, 1)
    return x_adv

def pgd_attack(model, x, y, epsilon, alpha, iterations):
    x = tf.convert_to_tensor(x)
    y = tf.convert_to_tensor(y)
    x_adv = x + tf.random.uniform(tf.shape(x), minval=-epsilon, maxval=epsilon)

```

```

x_adv = tf.clip_by_value(x_adv, 0, 1)
for i in range(iterations):
    x_adv = fgsm_attack(model, x_adv, y, alpha)
    perturbation = tf.clip_by_value(x_adv - x, -epsilon, epsilon)
    x_adv = tf.clip_by_value(x + perturbation, 0, 1)
return x_adv

# Test baseline model on adversarial examples
epsilon = 0.3
alpha = 0.01
iterations = 40

x_test_adv_fgsm = fgsm_attack(model, x_test, y_test, epsilon)
fgsm_accuracy = model.evaluate(x_test_adv_fgsm, y_test, verbose=0)
print(f"Baseline model accuracy on FGSM adversarial examples: {fgsm_accuracy[1] * 100:.2f}%")

x_test_adv_bim = bim_attack(model, x_test, y_test, epsilon, alpha, iterations)
bim_accuracy = model.evaluate(x_test_adv_bim, y_test, verbose=0)
print(f"Baseline model accuracy on BIM adversarial examples: {bim_accuracy[1] * 100:.2f}%")

x_test_adv_pgd = pgd_attack(model, x_test, y_test, epsilon, alpha, iterations)
pgd_accuracy = model.evaluate(x_test_adv_pgd, y_test, verbose=0)
print(f"Baseline model accuracy on PGD adversarial examples: {pgd_accuracy[1] * 100:.2f}%")

# Adversarial training
@tf.function
def train_step(model, optimizer, x, y, epsilon, alpha, iterations):
    # Generate a mix of adversarial and clean examples with varied attacks
    x_adv_pgd = pgd_attack(model, x, y, epsilon, alpha, iterations)
    x_adv_bim = bim_attack(model, x, y, epsilon, alpha / 2, iterations // 2)
    x_adv_fgsm = fgsm_attack(model, x, y, epsilon / 2)

    x_combined = tf.concat([x, x_adv_pgd, x_adv_bim, x_adv_fgsm], axis=0)
    y_combined = tf.concat([y, y, y, y], axis=0)

    with tf.GradientTape() as tape:
        predictions = model(x_combined)
        loss = tf.keras.losses.sparse_categorical_crossentropy(y_combined, predictions)
        gradients = tape.gradient(loss, model.trainable_variables)
        optimizer.apply_gradients(zip(gradients, model.trainable_variables))

# Create a new model for adversarial training
model_adv = build_model()
model_adv.compile(optimizer=Adam(learning_rate=0.001),
                  loss=SparseCategoricalCrossentropy(),
                  metrics=['accuracy'])

optimizer = Adam(learning_rate=0.001)

```

```

# Perform adversarial training
print("Adversarial training...")
for epoch in range(10):
    print(f"Epoch {epoch + 1}/10")
    for i in range(len(x_train) // 64):
        x_batch = x_train[i * 64:(i + 1) * 64]
        y_batch = y_train[i * 64:(i + 1) * 64]
        train_step(model_adv, optimizer, x_batch, y_batch, epsilon, alpha, iterations)

# Evaluate robust model accuracy
robust_accuracy = model_adv.evaluate(x_test, y_test, verbose=0)
print(f"Robust model accuracy on clean examples: {robust_accuracy[1] * 100:.2f}%")

x_test_adv_fgsm_robust = fgsm_attack(model_adv, x_test, y_test, epsilon)
robust_fgsm_accuracy = model_adv.evaluate(x_test_adv_fgsm_robust, y_test, verbose=0)
print(f"Robust model accuracy on FGSM adversarial examples:

{robust_fgsm_accuracy[1] * 100:.2f}%")

x_test_adv_bim_robust = bim_attack(model_adv, x_test, y_test, epsilon, alpha, iterations)
robust_bim_accuracy = model_adv.evaluate(x_test_adv_bim_robust, y_test, verbose=0)
print(f"Robust model accuracy on BIM adversarial examples:

{robust_bim_accuracy[1] * 100:.2f}%")

x_test_adv_pgd_robust = pgd_attack(model_adv, x_test, y_test, epsilon, alpha, iterations)
robust_pgd_accuracy = model_adv.evaluate(x_test_adv_pgd_robust, y_test, verbose=0)
print(f"Robust model accuracy on PGD adversarial examples:

{robust_pgd_accuracy[1] * 100:.2f}%")

```

---