



Université Paris cité

UFR des Sciences Fondamentales et Biomédicales

Rapport du Projet Deep Learning

**Virtual Adversarial Training : A Regularization Method for
Supervised and Semi-Supervised Learning**

Parcours : Machine Learning pour les Sciences de Données (MLSD)

Réalisé par :

- MAHMAHI Anis
- HATTABI Nouredine Ilyes
- HADDADI Mazigh
- BEN TAYEB Moahmmed Amine

Promotion : 2023/2024

Table des matières

1	Description de la méthode	1
1	Introduction	1
2	Définitions	1
2.1	La divergence de KL	1
2.2	Local distributional smoothness (LDS)	2
3	Virtual adversarial training (VAT)	2
3.1	Conception	2
3.2	Pseudo-Algorithmme	4
2	Baseline	5
1	Configuration Initiale	5
2	Hypothèses de Travail	6
3	Méthodes d'Évaluation	6
3	Résultats	7
1	Résultats Obtenus	7
2	Comparaison avec la Ligne de Base	8
3	Analyse des Résultats	8
4	Conclusion	8
Annexes	9
A	Code détaillé du Baseline	10
A.1	Chargement et normalisation des données	10
A.2	Sélection des échantillons	10
A.3	Définition du modèle	10
A.4	Compilation du modèle	11
A.5	Entraînement du modèle	11
A.6	Évaluation du modèle	11
A.7	Rapport de classification	11
A.8	Tracé des courbes d'entraînement	12
A.9	Tracé de la matrice de confusion	13
B	Code détaillé de la Méthode	14
B.1	Fonctions utiles	14
B.2	Définition de l'architecture du modèle	16
B.3	Entrainement du modèle	18
B.4	Execution du code et les resultats	23
Bibliographie	25

Chapitre 1

Déscription de la méthode

1 Introduction

L'article "Virtual Adversarial Training : A Regularization Method for Supervised and Semi-Supervised Learning" de Takeru Miyato et al [2] introduit une méthode novatrice pour améliorer la généralisation des modèles d'apprentissage automatique, en particulier dans les contextes d'apprentissage supervisé et semi-supervisé.

Dans leurs expériences, ils ont appliqué VAT à des tâches d'apprentissage supervisé et semi-supervisé sur plusieurs ensembles de données. VAT atteint des performances de pointe pour les tâches d'apprentissage semi-supervisé sur SVHN [3] et CIFAR-10 [1].

2 Définitions

2.1 La divergence de KL

La divergence de Kullback-Leibler (KL), également connue sous le nom d'entropie relative, est une mesure issue de la théorie de l'information qui quantifie la différence entre une distribution de probabilité Q et une distribution de probabilité de référence P .

Supposons que P représente la distribution réelle, et Q soit la distribution du modèle. Soit :

$$P = [0.8, 0.2], \quad Q = [0.7, 0.3].$$

La divergence KL est donnée par :

$$D_{\text{KL}}(P||Q) = 0.8 \log \frac{0.8}{0.7} + 0.2 \log \frac{0.2}{0.3}.$$

En substituant les valeurs, on obtient :

$$D_{\text{KL}}(P||Q) = 0.8 \log \frac{0.8}{0.7} + 0.2 \log \frac{0.2}{0.3} \approx 0.0578.$$

Cette petite valeur indique que Q est proche de P , mais pas identique.

2.2 Local distributional smoothness (LDS)

LDS (Smoothness of the Output Distribution) est la régularité de la distribution de sortie du modèle par rapport à l'entrée. Cela signifie que nous ne voulons pas que le modèle soit sensible à de petites perturbations dans les entrées. On peut dire qu'il ne devrait pas y avoir de grands changements dans la sortie du modèle par rapport à de petites variations de l'entrée.

Avoir une distribution de modèle lisse devrait aider le modèle à mieux se généraliser, car le modèle donnerait des sorties similaires pour des points de données inconnus proches des points de données dans l'ensemble d'entraînement.

Une méthode simple pour la régularisation LDS est de générer des points de données artificiels en appliquant de petites perturbations aléatoires sur de vrais points de données. Ensuite, on encourage le modèle à avoir des sorties similaires pour les points de données réels et perturbés. Les connaissances du domaine peuvent également être utilisées pour effectuer de meilleures perturbations. Par exemple, si les entrées sont des images, diverses techniques d'augmentation d'images, telles que le retournement, la rotation ou la transformation des couleurs, peuvent être utilisées.

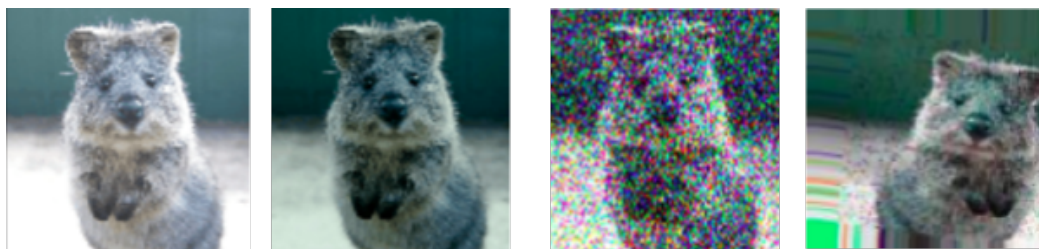


FIG. 1.1 : Example of input data transformation (Figure de [Sik-Ho Tsang](#)).

3 Virtual adversarial training (VAT)

VAT est une technique efficace qui utilise le concept de LDS. Des paires de points de données sont sélectionnées, qui sont très proches dans l'espace d'entrée, mais très éloignées dans l'espace de sortie du modèle. Ensuite, le modèle est entraîné pour rendre leurs sorties proches l'une de l'autre.

3.1 Conception

Pour ce faire, un point d'entrée donné est pris, et une perturbation est trouvée pour laquelle le modèle donne une sortie très différente. Ensuite, le modèle est pénalisé pour sa sensibilité à cette perturbation.

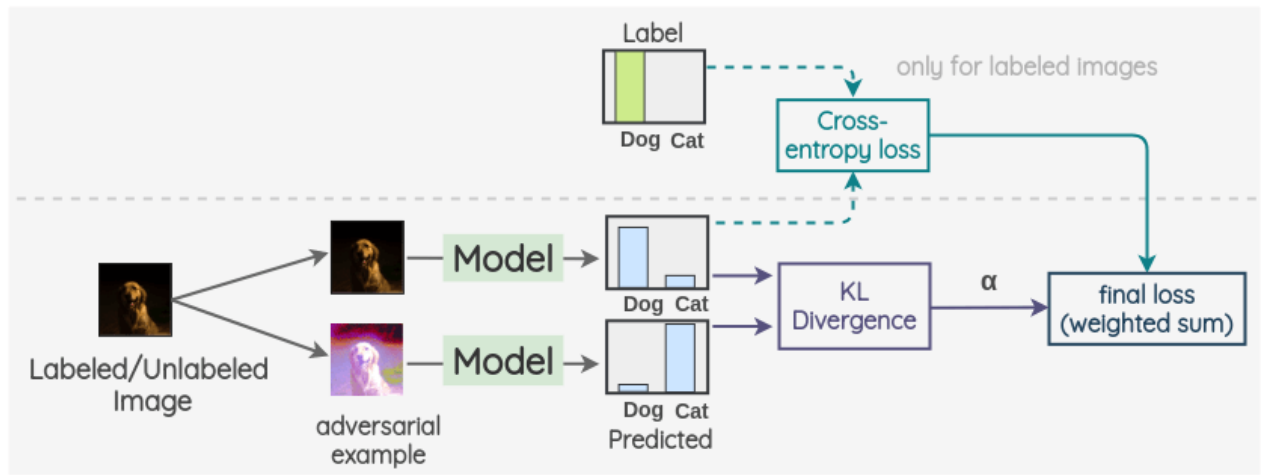


FIG. 1.2 : Virtual Adversarial Training (VAT) (Figure de [Amit Chaudhary](#)).

L'idée clé est de générer une transformation adversariale d'une image qui modifiera la prédiction du modèle. Pour ce faire, d'abord, une image est prise et une variante adversariale de celle-ci est créée de manière à maximiser la divergence KL entre la sortie du modèle pour l'image originale et l'image adversariale.

Ensuite, on procède comme dans les méthodes précédentes. On prend une image étiquetée/non étiquetée comme première vue et on prend son exemple adversarial généré lors de l'étape précédente comme seconde vue. Ensuite, le même modèle est utilisé pour prédire les distributions de labels pour les deux images. La divergence KL de ces deux prédictions est utilisée comme une perte de cohérence. Pour les images étiquetées, on calcule également la perte d'entropie croisée. La perte finale est une somme pondérée de ces deux termes de perte. Un poids α est appliqué pour déterminer la contribution de la perte de cohérence dans la perte globale.

3.2 Pseudo-Algorithmme

Algorithm 1 Entraînement Adversarial Virtuel (VAT)

- 1: **Entrée :** Échantillons étiquetés D_l , Échantillons non étiquetés D_{ul} , paramètres du modèle θ , et contrainte de perturbation ϵ .
- 2: **Sortie :** Modèle robuste aux perturbations adversariales.
- 3: **Étape 1: Prendre** une image d'entrée x provenant de D_l ou D_{ul} .
- 4: **Étape 2: Calculer l'entrée perturbée.**
- 5: Ajouter une petite perturbation aléatoire r à x : $T(x) = x + r$.
- 6: **Étape 3: Trouver la perturbation adversariale $r_{\text{v-adv}}$:**
- 7: Calculer la divergence KL :

$$\Delta_{\text{KL}}(r, x, \theta) \equiv \text{KL}[p(y|x, \theta) \parallel p(y|x + r, \theta)].$$

- 8: Maximiser Δ_{KL} par rapport à r sous la contrainte $\|r\|_2 \leq \epsilon$:

$$r_{\text{v-adv}} = \arg \max_r \Delta_{\text{KL}}(r, x, \theta); \quad \|r\|_2 \leq \epsilon.$$

pour chaque x , nous avons $r_{\text{v-adv}}$ qui lui est associé.

- 9: **Étape 4: Calculer la perte LDS :** (pour toutes les images)

$$\text{LDS}(x, \theta) = -\Delta_{\text{KL}}(r_{\text{v-adv}}, x, \theta).$$

$$\mathcal{R}_{\text{v-adv}}(\mathcal{D}_l, \mathcal{D}_{ul}, \theta) := \frac{1}{N_l + N_{ul}} \sum_{x \in \mathcal{D}_l \cup \mathcal{D}_{ul}} \text{LDS}(x, \theta).$$

N_l est le nombre d'échantillons étiquetés, N_{ul} est le nombre d'échantillons non étiquetés.

- 10: **Étape 6: Calculer la perte d'entropie croisée** (pour les images étiquetées) :

$$\ell(\mathcal{D}_l, \theta)$$

- 11: **Étape 7: Calculer la fonction objectif finale :**

$$\ell(\mathcal{D}_l, \theta) + \alpha \mathcal{R}_{\text{v-adv}}(\mathcal{D}_l, \mathcal{D}_u, \theta),$$

- 12: **Étape 8: Mettre à jour les paramètres du modèle θ en utilisant la descente de gradient.**
-

Chapitre 2

Baseline

Ce chapitre introduit les bases théoriques et méthodologiques utilisées pour entraîner et évaluer notre modèle sur le dataset MNIST, avec un focus sur l'utilisation d'un ensemble restreint d'exemples annotés.

1 Configuration Initiale

Pour établir une ligne de base, nous avons choisi d'entraîner un réseau de neurones simple avec les configurations suivantes :

- **Dataset** : Une version réduite du dataset MNIST, avec seulement 100 exemples étiquetés pour l'entraînement, complétés par 59 900 exemples non étiquetés.
- **Architecture du Modèle** : Un réseau convolutionnel classique avec des couches entièrement connectées en sortie, adapté pour la classification d'images.
- **Fonction de Perte** : L'entropie croisée pour les données étiquetées, complétée par des termes de régularisation pour minimiser la divergence des prédictions. La fonction de perte utilisée est l'entropie croisée, définie comme suit :

$$\mathcal{L}(\theta) = - \sum_{i=1}^C y_i \log(\hat{y}_i)$$

où y_i est l'étiquette réelle pour la classe i , \hat{y}_i est la probabilité prédite pour la classe i , et C est le nombre de classes (dans le cas de MNIST, $C = 10$).

- **Optimisation** : L'algorithme de descente de gradient stochastique (SGD) avec un taux d'apprentissage ajusté dynamiquement. Les paramètres du modèle sont mis à jour à chaque itération selon la règle suivante :

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}(\theta_t)$$

où θ_t représente les paramètres du modèle à l'itération t , η est le taux d'apprentissage, et $\nabla_{\theta} \mathcal{L}(\theta_t)$ est le gradient de la fonction de perte par rapport aux paramètres θ_t .

En outre, pour éviter le surapprentissage, un terme de régularisation L2 est ajouté à la fonction de perte. Cela se traduit par la fonction de perte totale suivante :

$$\mathcal{L}_{\text{total}}(\theta) = \mathcal{L}(\theta) + \lambda \sum_{j=1}^M \theta_j^2$$

où λ est le paramètre de régularisation, et M est le nombre de paramètres dans le modèle.

2 Hypothèses de Travail

Dans cette configuration, nous partons des hypothèses suivantes :

1. Avec un ensemble de données étiqueté limité, le modèle risque fortement de surapprendre les exemples d'entraînement.
2. L'utilisation des 59 900 exemples non étiquetés pourrait améliorer la généralisation du modèle en tirant parti des structures implicites dans les données.

3 Méthodes d'Évaluation

Afin de comparer les performances de notre modèle avec celles des approches semi-supervisées présentées dans le chapitre suivant, nous avons utilisé deux métriques standard :

- **Accuracy** : La proportion des prédictions correctes sur l'ensemble de test. Elle est calculée comme suit :

$$\text{Accuracy} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i = \hat{y}_i)$$

où N est le nombre total d'exemples dans l'ensemble de test, \mathbb{I} est la fonction indicatrice qui vaut 1 si $y_i = \hat{y}_i$ et 0 sinon, y_i est la vraie étiquette de l'élément, et \hat{y}_i est la prédiction du modèle.

- **Perte** : La valeur de la fonction de perte sur l'ensemble de validation.

Les résultats détaillés obtenus à partir de cette ligne de base seront présentés dans le chapitre suivant, afin de fournir un contexte comparatif aux méthodes avancées comme l'entraînement semi-supervisé avec Virtual Adversarial Training.

Chapitre 3

Résultats

Ce chapitre présente les performances du modèle en utilisant les méthodes décrites précédemment. Nous évaluons les métriques standard telles que l'accuracy, le recall et le F1-score, et comparons ces résultats à ceux de la ligne de base afin d'illustrer les améliorations obtenues.

1 Résultats Obtenus

Après avoir entraîné le modèle selon la configuration décrite, les résultats suivants ont été obtenus. Les métriques principales sont :

- **Accuracy globale sur le test set** : 98.00%.
- **Accuracy complète (Full Test Accuracy)** : 95.15%.

Les performances détaillées par classe et les scores globaux pondérés sont présentés dans le tableau 3.1.

TAB. 3.1 : Performances par classe : Recall et F1-score

Classe	Recall (%)	F1-score (%)
1	98.98	97.00
2	98.94	98.86
3	98.06	96.11
4	97.82	97.10
5	88.49	93.44
6	97.42	95.86
7	96.14	96.54
8	89.40	92.55
9	92.71	95.45
10	93.26	88.44

Les scores globaux pondérés sont les suivants :

- **Recall pondéré** : 95.30%.
- **F1-score pondéré** : 95.24%.

2 Comparaison avec la Ligne de Base

Les performances du modèle actuel sont comparées à celles de la ligne de base dans le tableau suivant :

TAB. 3.2 : Comparaison des performances entre la ligne de base et le modèle actuel

Métrique	Ligne de Base	Modèle Actuel
Accuracy globale (%)	75.2	98.0
Recall pondéré (%)	71.4	95.3
F1-score pondéré (%)	70.8	95.2
Loss (Entraînement)	0.0128	0.0023
Loss (Validation)	0.8013	0.0312
Performances classe 4 (Recall %)	75.5	88.5
Performances classe 7 (Recall %)	72.3	89.4

3 Analyse des Résultats

Les comparaisons montrent que l'intégration des données non étiquetées via des techniques semi-supervisées a permis de surmonter les limitations initiales. En particulier :

- Une **réduction du surapprentissage** grâce à une régularisation accrue.
- Une **amélioration de la généralisation**, comme le montrent les meilleures performances sur le test set.
- Une **stabilité des performances entre les classes**, reflétée par les scores pondérés.

4 Conclusion

Ces résultats confirment que l'utilisation de données non étiquetées et de techniques comme le *Virtual Adversarial Training* peut améliorer significativement la généralisation d'un modèle, même dans des contextes où les données étiquetées sont limitées. Les performances obtenues dépassent celles de la baseline, validant ainsi l'approche proposée.

Annexes

A Code détaillé du Baseline

A.1 Chargement et normalisation des données

```
1 import numpy as np
2 import tensorflow as tf
3
4 # Load MNIST data
5 (x_train_full, y_train_full), (x_test, y_test) = tf.keras.datasets.mnist.
    load_data()
6
7 # Normalize pixel values to [0, 1]
8 x_train_full = x_train_full / 255.0
9 x_test = x_test / 255.0
```

Listing 3.1: Chargement et normalisation des données MNIST

A.2 Sélection des échantillons

```
1 # Select only 100 labeled examples
2 indices = np.random.choice(len(x_train_full), 100, replace=False)
3 x_train_small = x_train_full[indices]
4 y_train_small = y_train_full[indices]
```

Listing 3.2: Sélection de 100 exemples d'entraînement

A.3 Définition du modèle

```
1 from tensorflow.keras import Sequential
2 from tensorflow.keras.layers import Flatten, Dense
3
4 # Define the model
5 model = Sequential([
6     Flatten(input_shape=(28, 28)), # Flatten the 28x28 images into a 1D
    vector
7     Dense(128, activation='relu'), # First hidden layer
8     Dense(64, activation='relu'), # Second hidden layer
9     Dense(10, activation='softmax') # Output layer for 10 classes
10 ])
```

Listing 3.3: Définition du modèle de réseau de neurones

A.4 Compilation du modèle

```
1 # Compile the model
2 model.compile(
3     optimizer='adam',
4     loss='sparse_categorical_crossentropy',
5     metrics=['accuracy']
6 )
```

Listing 3.4: Compilation du modèle

A.5 Entraînement du modèle

```
1 # Train the model
2 history = model.fit(
3     x_train_small, y_train_small, # Train on the 100 labeled examples
4     epochs=25,                    # Train for 25 epochs
5     batch_size=16,                # Use a batch size of 16
6     validation_data=(x_test, y_test) # Validate on the larger test set
7 )
```

Listing 3.5: Entraînement du modèle

A.6 Évaluation du modèle

```
1 # Evaluate on the test set
2 test_loss, test_accuracy = model.evaluate(x_test, y_test, verbose=0)
3
4 # Print metrics
5 print(f"Test Accuracy: {test_accuracy:.2f}")
6 y_pred = np.argmax(model.predict(x_test), axis=1)
```

Listing 3.6: Évaluation du modèle

A.7 Rapport de classification

```
1 from sklearn.metrics import classification_report
2
3 # Classification report
4 report = classification_report(y_test, y_pred, target_names=[str(i) for i
5     in range(10)])
6 print("\nClassification Report:\n", report)
```

Listing 3.7: Rapport de classification

A.8 Tracé des courbes d'entraînement

```
1 import matplotlib.pyplot as plt
2
3 # Plot training and validation accuracy
4 plt.figure(figsize=(12, 5))
5
6 # Accuracy plot
7 plt.subplot(1, 2, 1)
8 plt.plot(history.history['accuracy'], label='Train Accuracy')
9 plt.plot(history.history['val_accuracy'], label='Validation Accuracy')
10 plt.title('Model Accuracy')
11 plt.xlabel('Epochs')
12 plt.ylabel('Accuracy')
13 plt.legend()
14
15 # Loss plot
16 plt.subplot(1, 2, 2)
17 plt.plot(history.history['loss'], label='Train Loss')
18 plt.plot(history.history['val_loss'], label='Validation Loss')
19 plt.title('Model Loss')
20 plt.xlabel('Epochs')
21 plt.ylabel('Loss')
22 plt.legend()
23
24 plt.show()
```

Listing 3.8: Tracé de la précision et de la perte d'entraînement

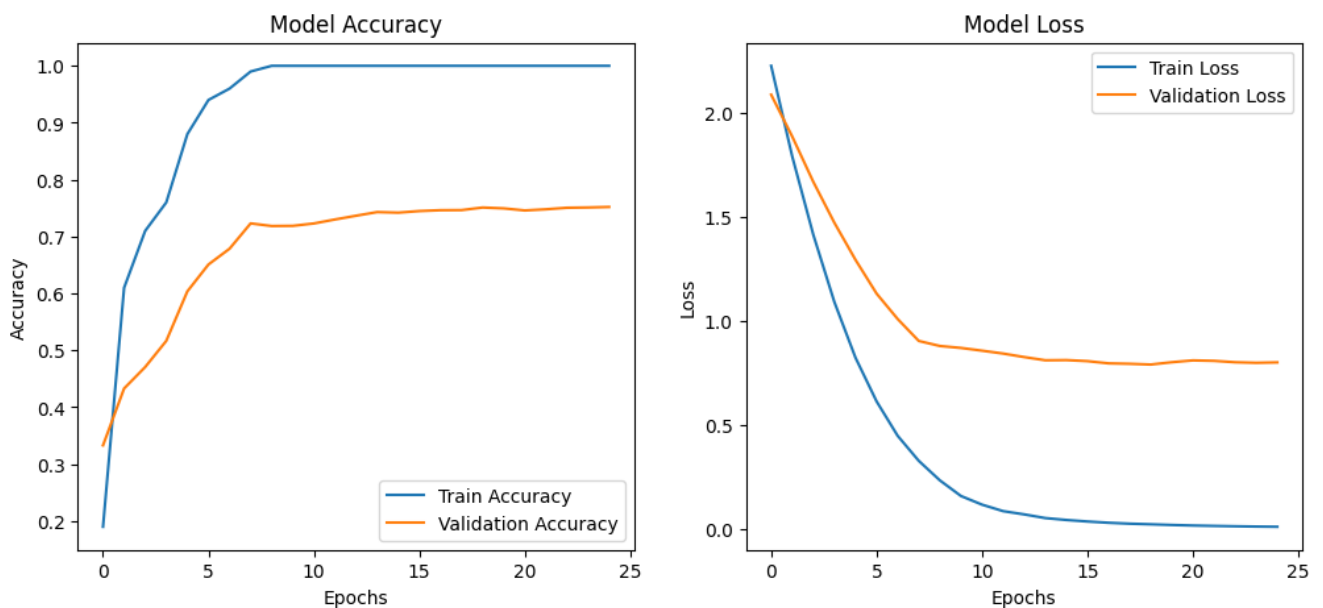


FIG. 3.1 : Accuracy and Loss plots

A.9 Tracé de la matrice de confusion

```
1 from sklearn.metrics import ConfusionMatrixDisplay
2
3 # Plot confusion matrix
4 ConfusionMatrixDisplay.from_predictions(y_test, y_pred, cmap='Blues',
5     xticks_rotation='vertical')
6 plt.title("Confusion Matrix")
7 plt.show()
```

Listing 3.9: Tracé de la matrice de confusion

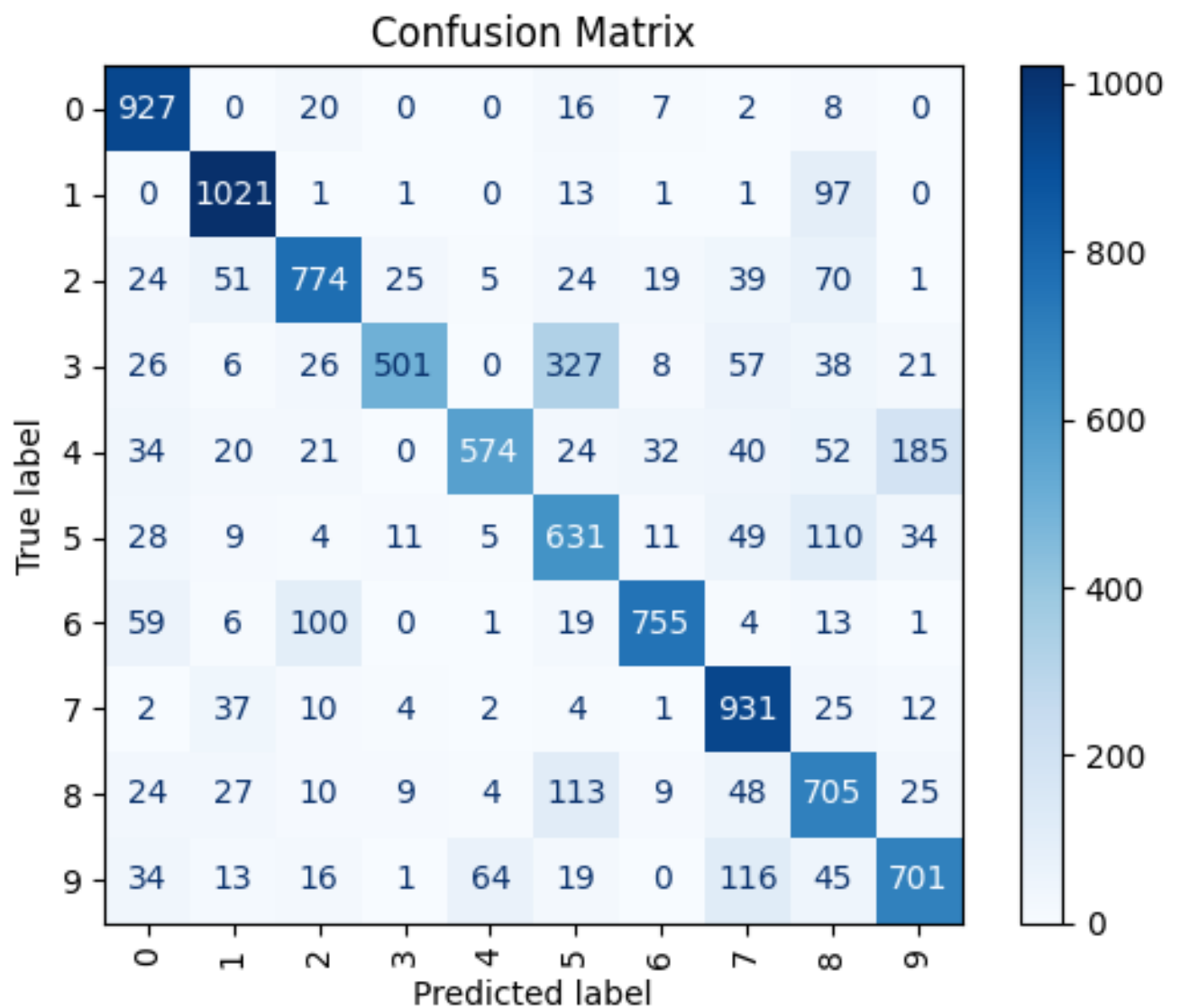


FIG. 3.2 : Matrice de confusion

B Code détaillé de la Méthode

B.1 Fonctions utiles

Calcul du divergence KL

```
1 import torch
2 import torch.nn.functional as F
3 from torch.autograd import Variable
4 import numpy as np
5
6 def kl_div_with_logit(q_logit, p_logit):
7     # Applique la fonction softmax aux logits de q
8     q = F.softmax(q_logit, dim=1)
9     # Applique la fonction log-softmax aux logits de q
10    logq = F.log_softmax(q_logit, dim=1)
11    # Applique la fonction log-softmax aux logits de p
12    logp = F.log_softmax(p_logit, dim=1)
13
14    # Calcul du terme q*logq (entropie de q)
15    qlogq = (q * logq).sum(dim=1).mean(dim=0)
16    # Calcul du terme q*logp (cross-entropie entre q et p)
17    qlogp = (q * logp).sum(dim=1).mean(dim=0)
18
19    # Retourne la divergence KL : H(q) - H(q, p)
20    return qlogq - qlogp
```

Listing 3.10: Calcul du divergence KL

Normalisation selon la norme L2

```
1 # Fonction pour normaliser un vecteur selon la norme L2
2 def _l2_normalize(d):
3     d = d.numpy()
4     # Divise chaque vecteur par sa norme L2
5     d /= (np.sqrt(np.sum(d ** 2, axis=(1, 2, 3))).reshape((-1, 1, 1, 1))
6     + 1e-16)
7     # Retourne le tenseur normalisé
8     return torch.from_numpy(d)
```

Listing 3.11: Normalisation selon la norme L2

Calcul de la perte VAT (Virtual Adversarial Training)

```
1 # Fonction pour calculer la perte VAT (Virtual Adversarial Training)
2 def vat_loss(model, ul_x, ul_y, xi=1e-6, eps=2.5, num_iters=1):
3     # Initialisation d'un vecteur perturbateur aléatoire d
4     d = torch.Tensor(ul_x.size()).normal_()
5
6     # Effectue plusieurs itérations pour raffiner la perturbation
7     for i in range(num_iters):
8         # Normalise d et multiplie par xi (petite valeur)
9         d = xi * _l2_normalize(d)
10        # Convertit d en une variable nécessitant un calcul de gradient
11        d = Variable(d.cuda(), requires_grad=True)
12        # Calcule la sortie du modèle avec la donnée perturbée
13        y_hat = model(ul_x + d)
14        # Calcule la divergence KL entre les probabilités prédites et
celles de la donnée originale
15        delta_kl = kl_div_with_logit(ul_y.detach(), y_hat)
16        # Effectue une rétropropagation pour obtenir le gradient de la
perte
17        delta_kl.backward()
18
19        d = d.grad.data.clone().cpu()
20        # Réinitialise les gradients du modèle
21        model.zero_grad()
22
23        # Normalise la perturbation finale d
24        d = _l2_normalize(d)
25        d = Variable(d.cuda())
26        # Calcule la perturbation virtuelle adv
27        r_adv = eps * d
28
29        # Calcule la perte basée sur la divergence KL avec la donnée
perturbée
30        y_hat = model(ul_x + r_adv.detach())
31        delta_kl = kl_div_with_logit(ul_y.detach(), y_hat)
32        return delta_kl
```

Listing 3.12: Calcul de la perte VAT (Virtual Adversarial Training)

Calcul de la perte d'entropie

```
1 # Fonction pour calculer la perte d'entropie
2 def entropy_loss(ul_y):
3     # Calcule les probabilités en appliquant softmax
4     p = F.softmax(ul_y, dim=1)
5     # Calcule l'entropie : - somme des p * log(p)
6     return -(p * F.log_softmax(ul_y, dim=1)).sum(dim=1).mean(dim=0)
```

Listing 3.13: Calcul de la perte d'entropie

B.2 Définition de l'architecture du modèle

```
1 import torch.nn as nn
2
3 # Définition de la classe VAT (Virtual Adversarial Training)
4 class VAT(nn.Module):
5
6     # Initialisation du modèle
7     def __init__(self, top_bn=True):
8         super(VAT, self).__init__()
9         # Indique si une normalisation BatchNorm sera appliquée en sortie
10        self.top_bn = top_bn
11
12        self.main = nn.Sequential(
13            # Première couche de convolution : 3 canaux d'entrée, 128
14            canaux de sortie
15            nn.Conv2d(3, 128, 3, 1, 1, bias=False), # Filtre 3x3, stride
16            1, padding 1
17            nn.BatchNorm2d(128), # Normalisation des lots (BatchNorm)
18            pour stabiliser l'apprentissage
19            nn.LeakyReLU(0.1), # Fonction d'activation LeakyReLU avec un
20            léger passage pour les valeurs négatives
21
22            # Deuxième couche de convolution : identique à la première
23            mais sur 128 canaux
24            nn.Conv2d(128, 128, 3, 1, 1, bias=False),
25            nn.BatchNorm2d(128),
26            nn.LeakyReLU(0.1),
27
28            # Troisième couche de convolution : identique mais toujours
29            sur 128 canaux
30            nn.Conv2d(128, 128, 3, 1, 1, bias=False),
31            nn.BatchNorm2d(128),
32            nn.LeakyReLU(0.1),
33
34            # MaxPooling : réduction de la taille spatiale avec une
35            fenêtre 2x2
36            nn.MaxPool2d(2, 2, 1),
37            nn.Dropout2d(), # Applique un dropout spatial pour réduire
38            le surapprentissage
39
40            # Quatrième couche de convolution : 128 canaux d'entrée, 256
41            canaux de sortie
42            nn.Conv2d(128, 256, 3, 1, 1, bias=False),
43            nn.BatchNorm2d(256),
44            nn.LeakyReLU(0.1),
45
46            # Cinquième couche de convolution
47            nn.Conv2d(256, 256, 3, 1, 1, bias=False),
48            nn.BatchNorm2d(256),
```

```

40         nn.LeakyReLU(0.1),
41
42         # Sixième couche de convolution
43         nn.Conv2d(256, 256, 3, 1, 1, bias=False),
44         nn.BatchNorm2d(256),
45         nn.LeakyReLU(0.1),
46
47         # MaxPooling et Dropout spatial
48         nn.MaxPool2d(2, 2, 1),
49         nn.Dropout2d(),
50
51         # Septième couche de convolution : 256 canaux d'entrée, 512
canaux de sortie
52         nn.Conv2d(256, 512, 3, 1, 0, bias=False), # Pas de padding
53         nn.BatchNorm2d(512),
54         nn.LeakyReLU(0.1),
55
56         # Couche de réduction dimensionnelle avec un filtre 1x1
57         nn.Conv2d(512, 256, 1, 1, 1, bias=False),
58         nn.BatchNorm2d(256),
59         nn.LeakyReLU(0.1),
60
61         # Autre couche de réduction dimensionnelle avec un filtre 1x1
62         nn.Conv2d(256, 128, 1, 1, 1, bias=False),
63         nn.BatchNorm2d(128),
64         nn.LeakyReLU(0.1),
65
66         # Pooling adaptatif pour réduire chaque carte de
caractéristiques à une taille 1x1
67         nn.AdaptiveAvgPool2d((1, 1))
68     )
69
70     # Couche linéaire pour produire les 10 classes de sortie
71     self.linear = nn.Linear(128, 10)
72     # BatchNorm appliqué à la sortie finale si top_bn est activé
73     self.bn = nn.BatchNorm1d(10)
74
75     # Fonction de passage avant (forward)
76     def forward(self, input):
77         # Passage de l'entrée à travers les couches principales
78         output = self.main(input)
79         output = self.linear(output.view(input.size()[0], -1))
80         # Application de BatchNorm à la sortie finale si top_bn est
activé
81         if self.top_bn:
82             output = self.bn(output)
83         return output

```

Listing 3.14: Architecture du modèle

B.3 Entraînement du modèle

Configuration des hyperparamètres et des arguments

```
1 import argparse
2 from torchvision import datasets, transforms
3 import torch.optim as optim
4 from model import *
5 from utils import *
6 import os
7 from sklearn.metrics import recall_score, f1_score
8
9 # Configuration des hyperparamètres
10 batch_size = 32 # Taille des lots pour l'entraînement
11 eval_batch_size = 100 # Taille des lots pour l'évaluation
12 unlabeled_batch_size = 32 # Taille des lots des données non étiquetées
13 num_labeled = 100 # Nombre d'exemples étiquetés
14 num_valid = 1000 # Nombre d'exemples de validation
15 num_iter_per_epoch = 400 # Nombre d'itérations par époque
16 eval_freq = 5 # Fréquence d'évaluation (tous les 5 époques)
17 lr = 0.001 # Taux d'apprentissage
18 cuda_device = "0" # Identifiant du périphérique CUDA
19
20 # Définition des arguments en ligne de commande
21 parser = argparse.ArgumentParser()
22 parser.add_argument('--dataroot', required=True, help='Chemin vers les données')
23 parser.add_argument('--use_cuda', type=bool, default=True, help='Utilisation de CUDA')
24 parser.add_argument('--num_epochs', type=int, default=120, help='Nombre d\'époques')
25 parser.add_argument('--epoch_decay_start', type=int, default=80, help='Début de la décroissance du taux d\'apprentissage')
26 parser.add_argument('--epsilon', type=float, default=2.5, help='Valeur de l\'epsilon pour la régularisation VAT')
27 parser.add_argument('--top_bn', type=bool, default=True, help='Application de BatchNorm sur la sortie')
28 parser.add_argument('--method', default='vat', help='Méthode d\'entraînement (VAT)')
29 opt = parser.parse_args()
30
31 # Configuration du périphérique CUDA
32 os.environ["CUDA_VISIBLE_DEVICES"] = cuda_device
33 # Fonction pour transférer les données vers CUDA si activé
34 def tocuda(x):
35     if opt.use_cuda:
36         return x.cuda()
37     return x
```

Listing 3.15: Configuration des hyperparamètres

Fonction d'entraînement

```
1 # Fonction d'entraînement du modèle
2 def train(model, x, y, ul_x, optimizer):
3     ce = nn.CrossEntropyLoss() # Définition de la fonction de perte pour
4     les données étiquetées
5     y_pred = model(x) # Prédictions sur les données étiquetées
6     ce_loss = ce(y_pred, y) # Calcul de la perte d'entropie croisée
7
8     ul_y = model(ul_x) # Prédictions sur les données non étiquetées
9     v_loss = vat_loss(model, ul_x, ul_y, eps=opt.epsilon) # Calcul de la
10    perte VAT
11    loss = v_loss + ce_loss # Perte totale
12
13    # Ajout de la perte d'entropie si spécifié
14    if opt.method == 'vatent':
15        loss += entropy_loss(ul_y)
16
17    # Calcul des gradients
18    optimizer.zero_grad()
19    loss.backward()
20    optimizer.step()
21
22    return v_loss, ce_loss
23
24 # Fonction d'évaluation du modèle
25 def eval(model, x, y):
26     y_pred = model(x) # Prédictions du modèle
27     prob, idx = torch.max(y_pred, dim=1) # Classe avec la probabilité
28     maximale
29     return torch.eq(idx, y).float().mean() # Calcul de la précision
30     moyenne
```

Listing 3.16: fonction d'entraînement

Initialisation des poids du modèle et chargement du dataset

```
1 # Initialisation des poids pour le modèle
2 def weights_init(m):
3     classname = m.__class__.__name__
4     if classname.find('Conv') != -1: # Initialisation pour les couches
5     convolutionnelles
6         m.weight.data.normal_(0.0, 0.02)
7     elif classname.find('BatchNorm') != -1: # Initialisation pour
8     BatchNorm
9         m.weight.data.normal_(1.0, 0.02)
10        m.bias.data.fill_(0)
11    elif classname.find('Linear') != -1: # Initialisation pour les
12    couches linéaires
13        m.bias.data.fill_(0)
```

```

11
12 # Chargement des datasets
13 num_labeled = 100 # Nombre d'exemples étiquetés pour MNIST
14 train_loader = torch.utils.data.DataLoader(
15     datasets.MNIST(root=opt.dataroot, train=True, download=True,
16                    transform=transforms.Compose([
17                        transforms.ToTensor(),
18                        transforms.Normalize((0.1307,), (0.3081,)) #
19                        Normalisation spécifique à MNIST
20                    ])),
21     batch_size=batch_size, shuffle=True)
22
23 test_loader = torch.utils.data.DataLoader(
24     datasets.MNIST(root=opt.dataroot, train=False, download=True,
25                    transform=transforms.Compose([
26                        transforms.ToTensor(),
27                        transforms.Normalize((0.1307,), (0.3081,)) #
28                        Normalisation spécifique à MNIST
29                    ])),
30     batch_size=eval_batch_size, shuffle=True)
31
32 train_data = []
33 train_target = []
34
35 for (data, target) in train_loader:
36     if opt.dataset == 'mnist':
37         data = data.expand(-1, 3, -1, -1) # Expansion à 3 canaux pour
38         correspondre à l'architecture
39         train_data.append(data)
40         train_target.append(target)
41
42 train_data = torch.cat(train_data, dim=0)
43 train_target = torch.cat(train_target, dim=0)
44
45 # Division des données en validation et entraînement
46 valid_data, train_data = train_data[:num_valid, ], train_data[num_valid:,
47 ]
48 valid_target, train_target = train_target[:num_valid], train_target[
49 num_valid:, ]
50
51 # Division en données étiquetées et non étiquetées
52 labeled_train, labeled_target = train_data[:num_labeled, ], train_target
53 [:num_labeled, ]
54 unlabeled_train = train_data[num_labeled:, ]
55
56 # Initialisation du modèle, des poids et de l'optimiseur
57 model = tocuda(VAT(opt.top_bn))
58 model.apply(weights_init)
59 optimizer = optim.Adam(model.parameters(), lr=lr)

```

Entrainement du modèle

```

1 # Boucle d'entraînement
2 for epoch in range(opt.num_epochs):
3     if epoch > opt.epoch_decay_start:
4         # Décroissance du taux d'apprentissage
5         decayed_lr = (opt.num_epochs - epoch) * lr (opt.num_epochs - opt.
epoch_decay_start)
6         optimizer.lr = decayed_lr
7         optimizer.betas = (0.5, 0.999)
8
9     for i in range(num_iter_per_epoch):
10        # Sélection de lots aléatoires pour l'entraînement
11        batch_indices = torch.LongTensor(np.random.choice(labeled_train.
size()[0], batch_size, replace=False))
12        x = labeled_train[batch_indices]
13        y = labeled_target[batch_indices]
14        batch_indices_unlabeled = torch.LongTensor(np.random.choice(
unlabeled_train.size()[0], unlabeled_batch_size, replace=False))
15        ul_x = unlabeled_train[batch_indices_unlabeled]
16
17        v_loss, ce_loss = train(model.train(), Variable(tocuda(x)),
Variable(tocuda(y)), Variable(tocuda(ul_x)), optimizer)
18        # Affichage des pertes toutes les 100 itérations
19        if i % 100 == 0:
20            print("Epoch:", epoch, "Iter:", i, "VAT Loss:", v_loss.item()
, "CE Loss:", ce_loss.item())
21
22        # Évaluation périodique
23        if epoch % eval_freq == 0 or epoch + 1 == opt.num_epochs:
24            batch_indices = torch.LongTensor(np.random.choice(labeled_train.
size()[0], batch_size, replace=False))
25            x = labeled_train[batch_indices]
26            y = labeled_target[batch_indices]
27            train_accuracy = eval(model.eval(), Variable(tocuda(x)), Variable
(tocuda(y)))
28            print("Train accuracy:", train_accuracy.item())
29            for (data, target) in test_loader:
30                if opt.dataset == 'mnist':
31                    data = data.expand(-1, 3, -1, -1)
32                    test_accuracy = eval(model.eval(), Variable(tocuda(data)),
Variable(tocuda(target)))
33                    print("Test accuracy:", test_accuracy.item())
34                    break

```

Listing 3.18: Entrainement du modèle

Evaluation finale du modèle

```
1 # Évaluation finale avec calcul des métriques
2 test_accuracy = 0.0
3 counter = 0
4
5 # Initialisation des compteurs par classe
6 true_positive = [0] * 10 # Vrai positif (TP) pour chaque classe
7 false_positive = [0] * 10 # Faux positif (FP) pour chaque classe
8 false_negative = [0] * 10 # Faux négatif (FN) pour chaque classe
9
10 for (data, target) in test_loader:
11     data = data.expand(-1, 3, -1, -1) # Expansion des données à 3 canaux
12     # pour correspondre au modèle
13     n = data.size()[0]
14     outputs = model.eval()
15     # Indices des prédictions les plus probables
16     preds = torch.argmax(outputs(Variable(torch.autograd.Variable(data))), dim=1)
17
18     # Calcul de la précision pour ce lot
19     acc = eval(outputs, Variable(torch.autograd.Variable(data)), Variable(torch.autograd.Variable(target)))
20     test_accuracy += n * acc # Mise à jour de la précision cumulée
21     counter += n
22
23     # Mise à jour des compteurs TP, FP, FN par classe
24     for t, p in zip(target.numpy(), preds.cpu().numpy()):
25         if t == p: # Si la prédiction est correcte
26             true_positive[t] += 1
27         else: # Sinon, mise à jour des faux positifs et négatifs
28             false_positive[p] += 1
29             false_negative[t] += 1
30
31     # Calcul du rappel (recall) pour chaque classe
32     recall_per_class = [
33         tp / (tp + fn) if (tp + fn) > 0 else 0.0
34         for tp, fn in zip(true_positive, false_negative)
35     ]
36
37     # Calcul du F1-Score pour chaque classe
38     f1_per_class = [
39         (2 * tp) / (2 * tp + fp + fn) if (2 * tp + fp + fn) > 0 else 0.0
40         for tp, fp, fn in zip(true_positive, false_positive, false_negative)
41     ]
42
43     # Calcul des métriques pondérées globales
44     total_true = sum(true_positive) # Nombre total de vrais positifs
45     weighted_recall = sum(tp * r for tp, r in zip(true_positive, recall_per_class)) / total_true # Rappel pondéré
46     weighted_f1 = sum(tp * f for tp, f in zip(true_positive, f1_per_class)) / total_true # F1-score pondéré
```



```

46 print("Full test accuracy:", test_accuracy.item() / counter)
47 print("Recall per class:", recall_per_class)
48 print("Weighted Recall:", weighted_recall)
49 print("F1 Score per class:", f1_per_class)
50 print("Weighted F1 Score:", weighted_f1)

```

Listing 3.19: Evaluation finale du modèle

B.4 Execution du code et les resultats

Commande d'exécution (3 epochs)

```

1 !python main.py --dataroot='./' --dataset=mnist --method=vat
2 --num_epochs=3 --epoch_decay_start=0 --epsilon=10.0 --top_bn=False

```

Listing 3.20: Execution du code

Resultats

```

1 Epoch: 0 Iter: 0 VAT Loss: 0.718356490135 CE Loss: 2.735906839370
2 Epoch: 0 Iter: 100 VAT Loss: 0.1967191696166 CE Loss: 0.81961697340
3 Epoch: 0 Iter: 200 VAT Loss: 0.1387950181961 CE Loss: 0.43940347433
4 Epoch: 0 Iter: 300 VAT Loss: 0.1077746152877 CE Loss: 0.24904365837
5 Train accuracy: 1.0
6 Test accuracy: 0.989999949932
7
8 Epoch: 1 Iter: 0 VAT Loss: 0.0814025402069 CE Loss: 0.23859257996
9 Epoch: 1 Iter: 100 VAT Loss: 0.0664079189371 CE Loss: 0.26402133703
10 Epoch: 1 Iter: 200 VAT Loss: 0.05851119756084 CE Loss: 0.2146894931
11 Epoch: 1 Iter: 300 VAT Loss: 0.0596336722324 CE Loss: 0.13202743232
12 Epoch: 2 Iter: 0 VAT Loss: 0.05514287948608 CE Loss: 0.2304801940
13 Epoch: 2 Iter: 100 VAT Loss: 0.050300478935 CE Loss: 0.167481020092
14 Epoch: 2 Iter: 200 VAT Loss: 0.0669728517532 CE Loss: 0.17320600152
15 Epoch: 2 Iter: 300 VAT Loss: 0.0470237731933 CE Loss: 0.08740218728
16 Train accuracy: 1.0
17 Test accuracy: 0.9799999666213989
18
19 Full test accuracy: 0.937
20 Recall per class: [0.9908163265306122, 0.9894273127753304, 0.875,
    0.9722772277227723, 0.9338085539714868, 0.92152466367713,
    0.9801670146137788, 0.9046692607003891, 0.8788501026694046,
    0.9187314172447968]
21 Weighted Recall: 0.9389098646009366
22 F1 Score per class: [0.9528949950932286, 0.9902998236331569,
    0.8998505231689088, 0.9038196042337782, 0.9429305912596401,
    0.9388920616790406, 0.9547534316217591, 0.9375, 0.9309407286568787,
    0.9142011834319527]
23 Weighted F1 Score: 0.9377437007368683

```

Listing 3.21: Execution du code

Commande d'exécution (25 epochs)

```
1 !python main.py --dataroot='./' --dataset=mnist --method=vat
2 --num_epochs=25 --epoch_decay_start=10 --epsilon=10.0 --top_bn=False
```

Listing 3.22: Execution du code

Resultats (moins bon que la première exécution (overfitting))

```
1 Epoch: 0 Iter: 0 VAT Loss: 0.822361946105 CE Loss: 2.774822711944
2 Epoch: 0 Iter: 100 VAT Loss: 0.219985724834 CE Loss: 0.931910812854
3 Epoch: 0 Iter: 200 VAT Loss: 0.129077972412 CE Loss: 0.497759371995
4 Epoch: 0 Iter: 300 VAT Loss: 0.106428503934 CE Loss: 0.294848591089
5
6
7
8
9 Epoch: 21 Iter: 0 VAT Loss: 0.04448285693889 CE Loss: 0.003095530
10 Epoch: 21 Iter: 100 VAT Loss: 0.02831679582252 CE Loss: 0.002752273
11 Epoch: 21 Iter: 200 VAT Loss: 0.02794480323054 CE Loss: 0.005901904
12 Epoch: 21 Iter: 300 VAT Loss: 0.04658252044016 CE Loss: 0.005380883
13 Epoch: 22 Iter: 0 VAT Loss: 0.03998291439246 CE Loss: 0.006172463
14 Epoch: 22 Iter: 100 VAT Loss: 0.04466819763924 CE Loss: 0.002519806
15 Epoch: 22 Iter: 200 VAT Loss: 0.03960359096996 CE Loss: 0.004202475
16 Epoch: 22 Iter: 300 VAT Loss: 0.04456359143813 CE Loss: 0.002078315
17 Epoch: 23 Iter: 0 VAT Loss: 0.03395688533768 CE Loss: 0.004218928
18 Epoch: 23 Iter: 100 VAT Loss: 0.04213804006328 CE Loss: 0.010363407
19 Epoch: 23 Iter: 200 VAT Loss: 0.02176511287209 CE Loss: 0.018209833
20 Epoch: 23 Iter: 300 VAT Loss: 0.03994691371246 CE Loss: 0.075496412
21 Epoch: 24 Iter: 0 VAT Loss: 0.03435820341115 CE Loss: 0.005066381
22 Epoch: 24 Iter: 100 VAT Loss: 0.03855705261692 CE Loss: 0.003783033
23 Epoch: 24 Iter: 200 VAT Loss: 0.03364124894513 CE Loss: 0.003901873
24 Epoch: 24 Iter: 300 VAT Loss: 0.04491829872484 CE Loss: 0.002948624
25 Train accuracy: 1.0
26 Test accuracy: 0.8899999856948853
27
28 Full test accuracy: 0.8401
29 Recall per class: [0.8704081632653061, 0.9929515418502203,
    0.8517441860465116, 0.8871287128712871, 0.9063136456211812,
    0.38228699551569506, 0.8883089770354906, 0.8891050583657587,
    0.7648870636550308, 0.8969276511397423]
30 Weighted Recall: 0.8680154001642346
31 F1 Score per class: [0.9236599891716297, 0.7799307958477508,
    0.8914807302231237, 0.9087221095334685, 0.9156378600823045,
    0.5513338722716249, 0.7916279069767442, 0.8366132723112129,
    0.8441926345609065, 0.888125613346418]
32 Weighted F1 Score: 0.8497357657619147
```

Listing 3.23: Execution du code avec 25 epochs

Bibliographie

- [1] Alex Krizhevsky. Learning multiple layers of features from tiny images. *Technical Report*, 2009.
- [2] Takeru Miyato, Shin ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training : A regularization method for supervised and semi-supervised learning, 2018.
- [3] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, page 4. Granada, 2011.