

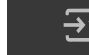
dash ce projet j'ai fait plusieurs étapes clés pour évaluer et améliorer la robustesse des modèles de machine learning face aux attaques adversariales que je les ai trouvé sur les slides :

1. Modélisation et Prédiction Initiales
2. Attaques Adversariales
3. Test du Modèle
4. Visualisation
5. Entrainement Adversarial

j'ai travailler avec deux methode d'attack, le fgsm et le pgd sur les donnes FashionMnist

▼ FGSM

```
1 !pip install torch
2
```

 Afficher la sortie masquée


Double-cliquez (ou appuyez sur Entrée) pour modifier

```
1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 import torchvision
5 import torchvision.transforms as transforms
6 import matplotlib.pyplot as plt
7
8 # Configuration du dispositif
9 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
```

```
1 # Chargement des datasets
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
```

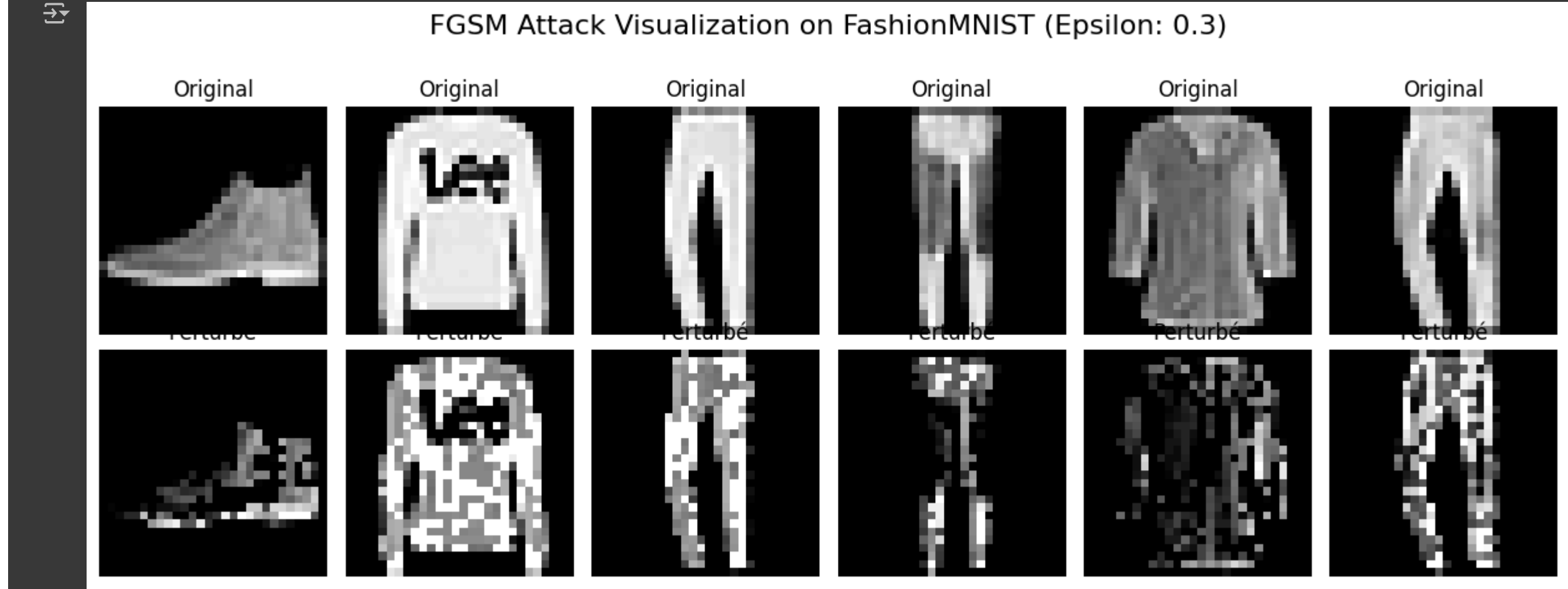
```
1
2
3 # Modele simple de réseau de neurones
4 class SimpleCNN28(nn.Module):
5     """Modele pour les images de taille 28x28 (MNIST, FashionMNIST)"""
6     def __init__(self, input_channels=1):
7         super(SimpleCNN28, self).__init__()
8         self.conv1 = nn.Conv2d(input_channels, 16, kernel_size=5, stride=1, padding=2)
9         self.conv2 = nn.Conv2d(16, 32, kernel_size=5, stride=1, padding=2)
10        self.fc1 = nn.Linear(32 * 7 * 7, 128)
11        self.fc2 = nn.Linear(128, 10)
12
13    def forward(self, x):
14        x = torch.relu(self.conv1(x))
15        x = torch.max_pool2d(x, 2)
16        x = torch.relu(self.conv2(x))
17        x = torch.max_pool2d(x, 2)
18        x = x.view(x.size(0), -1)
19        x = torch.relu(self.fc1(x))
20        x = self.fc2(x)
21        return x
22
23 # Initialisation de la fonction de perte et de l'optimizer
24 criterion = nn.CrossEntropyLoss()
25
26 # Fonction d'attaque FGSM
27 def fgsm_attack(image, epsilon, data_grad):
28     sign_data_grad = data_grad.sign()
29     perturbed_image = image + epsilon * sign_data_grad
30     return torch.clamp(perturbed_image, 0, 1)
31
32 # Fonction de test avec FGSM
33 def test_fgsm(model, testloader, epsilon):
34     correct = 0
35     adv_examples = []
36     model.eval()
37
38     for data, target in testloader:
39         data, target = data.to(device), target.to(device)
40         data.requires_grad = True
41
42         output = model(data)
43         init_pred = output.max(1, keepdim=True)[1]
44         loss = criterion(output, target)
45         model.zero_grad()
46         loss.backward()
47
48         data_grad = data.grad.data
49         perturbed_data = fgsm_attack(data, epsilon, data_grad)
50
51         output = model(perturbed_data)
52         final_pred = output.max(1, keepdim=True)[1]
53         correct += (final_pred == target).sum().item()
54
55     final_acc = correct / float(len(testloader.dataset))
56     print(f'Epsilon: {epsilon}\tTest Accuracy = {final_acc * 100:.2f}%')
57
58 # Entrainement et test pour chaque dataset
59 def train_and_evaluate(dataset_name, trainloader, testloader, model, epsilon=0.3):
60     print(f"\nTraining and evaluating FGSM on {dataset_name} dataset")
61     model.to(device)
62     optimizer = optim.Adam(model.parameters(), lr=0.001)
63
64     # Entrainement du modèle
65     model.train()
66     for epoch in range(3): # Entraîner pour 3 époques pour cet exemple
67         for images, labels in trainloader:
68             images, labels = images.to(device), labels.to(device)
69             optimizer.zero_grad()
70             outputs = model(images)
71             loss = criterion(outputs, labels)
72             loss.backward()
73             optimizer.step()
74
75     # Test FGSM
76     test_fgsm(model, testloader, epsilon)
77
78
79
80
```

```
1 # Tester sur chaque dataset
2 model_fashion = SimpleCNN28(input_channels=1).to(device)
3 train_and_evaluate("FashionMNIST", trainloader_fashion, testloader_fashion, model=model_fashion, epsilon=0.3)
4
```

 Training and evaluating FGSM on FashionMNIST dataset
Epsilon: 0.3 Test Accuracy = 643.39%

```
1 import matplotlib.pyplot as plt
2
3 def visualize_attack(model, testloader, epsilon, input_channels=1, dataset_name="Dataset"):
4     """
5     Visualise l'effet de l'attaque FGSM sur un modèle en affichant des images originales et perturbées.
6
7     Arguments :
8     - model : Le modèle entraîné sur lequel effectuer l'attaque.
9     - testloader : DataLoader contenant le jeu de test.
10    - epsilon : Magnitude de la perturbation pour FGSM.
11    - input_channels : Nombre de canaux d'entrée (1 pour les images en niveaux de gris, 3 pour les images en couleur).
12    - dataset_name : Nom du dataset, utilisé pour distinguer les images dans les annotations.
13    """
14
15    # Extraire un batch de données
16    dataiter = iter(testloader)
17    images, labels = next(dataiter)
18    images, labels = images.to(device), labels.to(device)
19
20    # Générer une attaque FGSM
21    images.requires_grad = True
22    output = model(images)
23    loss = criterion(output, labels)
24    model.zero_grad()
25    loss.backward()
26    data_grad = images.grad.data
27    perturbed_data = fgsm_attack(images, epsilon, data_grad)
28
29    # Créer la figure
30    fig, axes = plt.subplots(2, 6, figsize=(12, 5))
31    fig.suptitle(f"FGSM Attack Visualization on {dataset_name} (Epsilon: {epsilon})", fontsize=16)
32
33    for i in range(6):
34        # Image originale
35        ax = axes[0, i]
36        original_img = images[i].detach().cpu().numpy()
37        if input_channels == 1:
38            ax.imshow(original_img.squeeze(), cmap="gray")
39        else:
40            ax.imshow(original_img.transpose(1, 2, 0))
41        ax.axis('off')
42        ax.set_title("Original")
43
44        # Image perturbée
45        ax = axes[1, i]
46        perturbed_img = perturbed_data[i].detach().cpu().numpy()
47        if input_channels == 1:
48            ax.imshow(perturbed_img.squeeze(), cmap="gray")
49        else:
50            ax.imshow(perturbed_img.transpose(1, 2, 0))
51        ax.axis('off')
52        ax.set_title("Perturbed")
53
54    plt.tight_layout(rect=[0, 0.03, 1, 0.95]) # Ajuster l'espace pour le titre principal
55    plt.show()
56
```

```
1 # Visualiser l'attaque FGSM sur MNIST
2 visualize_attack(model_fashion, testloader_fashion, epsilon=0.3, input_channels=1, dataset_name="FashionMNIST")
3
```



```
1 def certify_robustness(model, testloader, epsilon, input_channels=1):
2     """
3     Certifie que les prédictions du modele sont robustes sous une perturbation FGSM de magnitude epsilon.
4
5     Arguments :
6     - model : le modèle à tester.
7     - testloader : le DataLoader pour le jeu de test.
8     - epsilon : la magnitude de la perturbation.
9     - input_channels : Le nombre de canaux d'entrée (1 pour les images en niveaux de gris, 3 pour les images en couleur).
10    """
11    model.eval()
12    correct_certified = 0
13    total = 0
14
15    for data, target in testloader:
16        data, target = data.to(device), target.to(device)
17        total += target.size(0)
18
19        # Prédiction initiale
20        output = model(data)
21        init_pred = output.max(1, keepdim=True)[1]
22
23        # Bornes inférieure et supérieure pour les perturbations
24        lower_bound = torch.clamp(data - epsilon, 0, 1)
25        upper_bound = torch.clamp(data + epsilon, 0, 1)
26
27        # Prédictions pour les bornes
28        lower_output = model(lower_bound)
29        upper_output = model(upper_bound)
30        lower_pred = lower_output.max(1, keepdim=True)[1]
31        upper_pred = upper_output.max(1, keepdim=True)[1]
32
33        # Certifier si les prédictions restent stables
34        is_certified = (lower_pred == init_pred).all(dim=1) & (upper_pred == init_pred).all(dim=1)
35        correct_certified += is_certified.sum().item()
36
37    # Calcul du taux de robustesse certifié
38    final_certified_acc = correct_certified / total
39    print(f'Certified robustness for epsilon {epsilon}: {final_certified_acc * 100:.2f}% over {total} samples.')
40
```

```
1 # Exécuter la certification de robustesse pour FashionMNIST
2 certify_robustness(model_fashion, testloader_fashion, epsilon=0.3, input_channels=1)
3
```

Certified robustness for epsilon 0.3: 47.70% over 10000 samples.

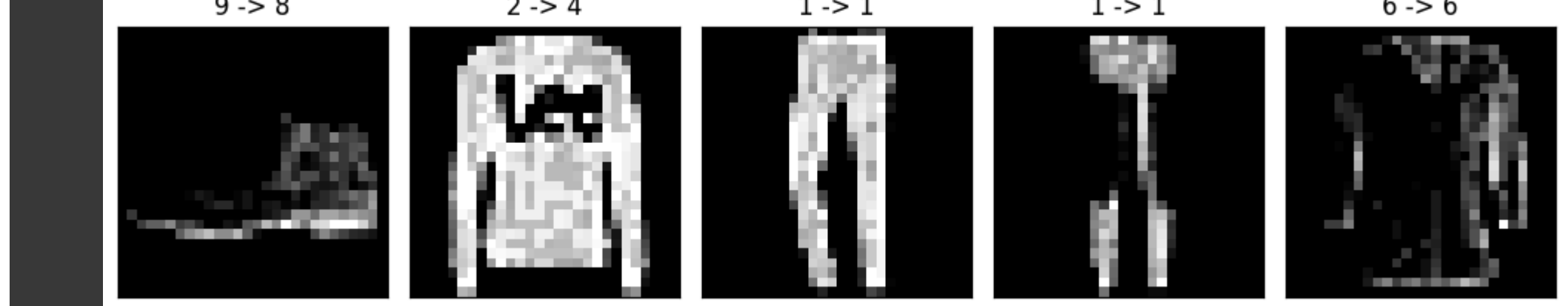
```
1 def adversarial_training(model, train_loader, epsilon, num_epochs=10):
2     # Placer le modèle en mode entraînement
3     model.train()
4
5     for epoch in range(num_epochs):
6         for images, labels in train_loader:
7             images, labels = images.to(device), labels.to(device)
8
9             # Zéro gradients
10            optimizer.zero_grad()
11
12            # Calculer les perturbations adversariales sur les images d'entraînement
13            images.requires_grad = True
14            outputs = model(images)
15            loss = criterion(outputs, labels)
16            model.zero_grad()
17            loss.backward()
18            data_grad = images.grad.data
19
20            # Créer des images adversariales en appliquant FGSM
21            perturbed_images = fgsm_attack(images, epsilon, data_grad)
22
23            # Ré-entraînement sur les images perturbées
24            outputs_adv = model(perturbed_images)
25            loss_adv = criterion(outputs_adv, labels)
26            loss_adv.backward()
27            optimizer.step()
28
29            # Affichage de la perte moyenne après chaque époque
30            print(f'Epoch {epoch+1}, Loss: {loss.item()}')
31
```

```
1 def test_model(model, test_loader, epsilon):
2     correct = 0
3     total = 0
4     adv_examples = []
5
6     model.eval()
7     for images, labels in test_loader:
8         images, labels = images.to(device), labels.to(device)
9         images.requires_grad = True
10
11        # Forward pass the data through the model
12        outputs = model(images)
13        _, init_pred = outputs.max(1) # Get the index of the max log-probability
14
15        # Loop over the batch
16        for idx in range(images.size(0)):
17            original_image = images[idx]
18            init_prediction = init_pred[idx]
19            label = labels[idx]
20
21            if init_prediction == label:
22                # Calculate loss only for correct predictions to generate adversarial examples
23                loss = criterion(outputs[idx].squeeze(0), labels[idx].unsqueeze(0))
24                model.zero_grad()
25                loss.backward(retain_graph=True)
26                data_grad = images.grad.data[idx]
27
28                # Generate adversarial example
29                perturbed_data = fgsm_attack(original_image.unsqueeze(0), epsilon, data_grad.unsqueeze(0))
30                output = model(perturbed_data)
31                final_pred = output.max(1)[1]
32
33                if final_pred.item() == label.item():
34                    correct += 1
35                if len(adv_examples) < 5: # Save some examples to visualize
36                    adv_ex = perturbed_data.squeeze().detach().cpu().numpy()
37                    adv_examples.append((init_prediction.item(), final_pred.item(), adv_ex))
38
39            total += images.size(0)
40
41        # Calculate final accuracy
42        final_acc = correct / float(total)
43        print(f'Epsilon: {epsilon}, Test Accuracy: {final_acc * 100:.2f}%')
44        return adv_examples
45
```

```
1 def visualize_results(adv_examples):
2     plt.figure(figsize=(10, 5))
3     for i, (init_pred, final_pred, ex) in enumerate(adv_examples):
4         plt.subplot(1, len(adv_examples), i+1)
5         plt.xticks([], [])
6         plt.yticks([], [])
7         plt.title(f'{init_pred} -> {final_pred}')
8         plt.imshow(ex, cmap="gray")
9     plt.tight_layout()
10    plt.show()
11
```

```
1 # Assume epsilon = 0.1 has been set
2 adv_examples = test_model(model_fashion, testloader_fashion, epsilon=0.1)
3 visualize_results(adv_examples)
4
```

Epsilon: 0.1, Test Accuracy: 35.50%



PGD

```
1 import torch
2
3 def pgd_attack(model, images, labels, epsilon, alpha, num_iter):
4     """
5     Effectue une attaque PGD sur les images.
6
7     Arguments :
8     - model : le modèle à attaquer.
9     - images : les images d'entrée.
10    - labels : les étiquettes correctes des images.
11    - epsilon : le maximum de perturbation.
12    - alpha : la taille de chaque étape.
13    - num_iter : le nombre d'itérations de PGD.
14
15    Retourne :
16    - perturbed_images : les images perturbées.
17    """
18    perturbed_images = images.clone().detach().to(images.device)
19    perturbed_images.requires_grad = True
20
21    for i in range(num_iter):
22        outputs = model(perturbed_images)
23        loss = criterion(outputs, labels)
24        model.zero_grad()
25        loss.backward()
26
27        # Calculer la perturbation et mettre à jour les images
28        grad = perturbed_images.grad.data
29        perturbed_images = perturbed_images + alpha * grad.sign()
30
31        # Projeter les images perturbées pour rester dans l'intervalle epsilon autour des images d'origine
32        perturbed_images = torch.max(torch.min(perturbed_images, images + epsilon), images - epsilon)
33        perturbed_images = torch.clamp(perturbed_images, 0, 1)
34        perturbed_images = perturbed_images.detach()
35        perturbed_images.requires_grad = True
36
37    return perturbed_images
38
```

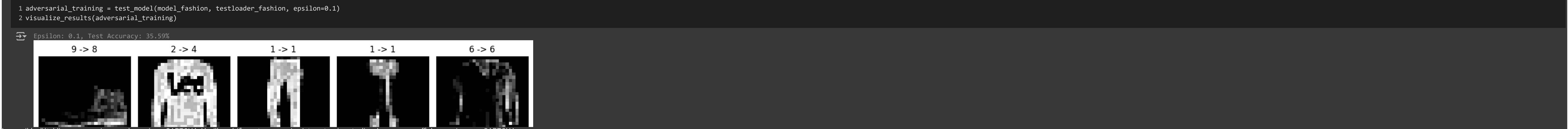
```
1 def visualize_attack_pgdl(model, testloader, epsilon, alpha, num_iter, input_channels=1, dataset_name="Dataset"):  
2     """  
3     Visualise l'effet de l'attaque PGD sur un modèle en affichant des images originales et perturbées.  
4     """  
5     # Extraire un batch de données  
6     dataiter = iter(testloader)  
7     images, labels = next(dataiter)  
8     images, labels = images.to(device), labels.to(device)  
9  
10    # Générer une attaque PGD  
11    perturbed_data = pgd_attack(model, images, labels, epsilon, alpha, num_iter)  
12  
13    # Créer la figure  
14    fig, axes = plt.subplots(2, 6, figsize=(12, 5))  
15    fig.suptitle(f"PGD Attack Visualization on {dataset_name} (Epsilon: {epsilon}, Alpha: {alpha}, Iterations: {num_iter})", fontsize=16)  
16  
17    for i in range(6):  
18        # Image originale  
19        ax = axes[0, i]  
20        original_img = images[i].detach().cpu().numpy()  
21        if input_channels == 1:  
22            ax.imshow(original_img.squeeze(), cmap="gray")  
23        else:  
24            ax.imshow(original_img.transpose(1, 2, 0))  
25        ax.axis("off")  
26        ax.set_title("Original")  
27  
28        # Image perturbée  
29        ax = axes[1, i]  
30        perturbed_img = perturbed_data[i].detach().cpu().numpy()  
31        if input_channels == 1:  
32            ax.imshow(perturbed_img.squeeze(), cmap="gray")  
33        else:  
34            ax.imshow(perturbed_img.transpose(1, 2, 0))  
35        ax.axis("off")  
36        ax.set_title("Perturbée")  
37  
38    plt.tight_layout(rect=[0, 0.03, 1, 0.95]) # Ajuster l'espace pour le titre principal  
39    plt.show()  
40
```



```
1 def certify_robustness_pgdl(model, testloader, epsilon, alpha, num_iter, input_channels=1):  
2     """  
3     Certifie que les prédictions du modèle sont robustes sous une perturbation PGD de magnitude epsilon.  
4     """  
5     model.eval()  
6     correct_certified = 0  
7     total = 0  
8  
9     for data, target in testloader:  
10        data, target = data.to(device), target.to(device)  
11        total += target.size(0)  
12  
13        # Générer des images perturbées avec PGD  
14        perturbed_data = pgd_attack(model, data, target, epsilon, alpha, num_iter)  
15  
16        # Prédiction pour les images perturbées  
17        output = model(perturbed_data)  
18        final_pred = output.max(1, keepdim=True)[1]  
19        correct_certified += (final_pred == target.view_as(final_pred)).sum().item()  
20  
21    # Calcul du taux de robustesse certifié  
22    final_certified_acc = correct_certified / total  
23    print(f"Certified robustness for PGD attack (epsilon={epsilon}, alpha={alpha}, iterations={num_iter}): {final_certified_acc * 100:.2f}%")  
24
```



```
1 def adversarial_training_pgdl(model, train_loader, epsilon, alpha, num_iter, num_epochs=10):  
2     model.train()  
3     for epoch in range(num_epochs):  
4         for images, labels in train_loader:  
5             images, labels = images.to(device), labels.to(device)  
6  
7             optimizer.zero_grad()  
8             perturbed_images = pgd_attack(model, images, labels, epsilon, alpha, num_iter)  
9             outputs_adv = model(perturbed_images)  
10            loss_adv = criterion(outputs_adv, labels)  
11            loss_adv.backward()  
12            optimizer.step()  
13  
14            print(f"Epoch {epoch+1}, Loss: {loss_adv.item()}")  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100
```



Impossible d'établir une connexion avec le service reCAPTCHA. Veuillez vérifier votre connexion Internet, puis actualiser la page pour afficher une image reCAPTCHA.