

Fiche Technique - Attaques Adversariales sur Modeles ML

Equipe de Cybersecurite

18 novembre 2025

Introduction

Cette fiche explique comment tester la robustesse de nos modeles ML (Random Forest, KNN, Decision Tree) face aux attaques adversariales dans le contexte de detection d'intrusion.

1 Concepts Fondamentaux

1.1 Attaque Adversariale : Qu'est-ce que c'est ?

Analogie : Un intrus qui modifie legerement son apparence pour tromper un systeme de securite.

Exemple concret : Notre modele detecte une attaque DDoS normalement. On modifie 2-3 features (duree, nombre de packets) de facon subtile. Le modele pense maintenant que c'est du trafic normal.

1.2 Pourquoi nos modeles sont vulnerables ?

- **Random Forest** : Sensible aux perturbations sur les features importantes
- **KNN** : Peut etre trompe en deplacant les points pres des frontieres
- **Decision Tree** : Une petite modification peut changer la feuille de decision

1.3 Choix des methodes d'attaque - par Copilot

J'ai choisi ces deux exemples (HopSkipJump et surrogate+FGSM) car ils sont les plus adaptes aux modeles scikit-learn (arbres, KNN, random forest) pour les raisons suivantes :

HopSkipJump (decision-based, black-box)

- Fonctionne sans acces au gradient, donc compatible avec les modeles non-differentiables (arbres, KNN)
- C'est une attaque generique, applicable a tout classificateur qui donne une prediction (score ou classe)
- Represente le scenario "attaquant externe" qui ne connait pas l'architecture du modele

Surrogate + FGSM (transfert)

- Les attaques par gradient (FGSM, PGD) ne marchent que sur des modeles differentiables (reseaux de neurones)
- Pour attaquer un modele non-differentiable, on entraine un "surrogate" (reseau Keras) sur les memes donnees ou sur les predictions du modele cible, puis on genere des adversariaux sur le surrogate et on teste leur effet sur le modele cible
- Cette approche est tres utilisee en recherche pour tester la robustesse des modeles non-differentiables

Pourquoi pas d'autres attaques ?

- Les attaques directes par gradient (FGSM, PGD, Carlini-Wagner) ne fonctionnent pas sur scikit-learn trees/KNN
- Les attaques par heuristique (Square, Boundary, NES) sont aussi possibles, mais HopSkipJump est la plus connue et la plus efficace pour le "decision-based"
- Les attaques par perturbation aleatoire ou optimisation evolutive sont moins efficaces ou plus lentes

2 Strategies d'Attaque Implementees

2.1 1. Attaque HopSkipJump (Black-Box)

Idee : Methode generique qui fonctionne avec n'importe quel modele sans connaitre ses internes.

Avantage : Simple a mettre en œuvre, universelle

Comment ca marche : Teste plein de petites modifications jusqu'a trouver celles qui trompent le modele.

Listing 1 – Exemple minimal HopSkipJump

```
1 # Exemple minimal - HopSkipJump (black-box / decision-based)
2 import numpy as np
3 from art.estimators.classification import SklearnClassifier
4 from art.attacks.evasion import HopSkipJump
5 # from sklearn.ensemble import RandomForestClassifier # votre modele deja
6 # entraine
7
8 # Supposons : model = RandomForestClassifier(...); model.fit(X_train,
9 # y_train)
10 # et X_test, y_test disponibles
11
12 # Definir clip_values selon vos features (ex: (0., 1.) ou (X_min, X_max))
13 clip_values = (X_train.min(axis=0).astype(float), X_train.max(axis=0).astype
14 # (float))
15 # ART attend un tuple (min, max). Si vous preferez uniformiser : (0., 1.)
16
17 art_clf = SklearnClassifier(model=model, clip_values=(0.0, 1.0)) # ajuster
18 clip_values
```

```

15
16 attack = HopSkipJump(classifier=art_clf, max_iter=50, max_eval=1000,
17   init_eval=10)
# Pour un petit test, utiliser un sous-ensemble
18 X_sub = X_test[:100].astype(np.float32)
19
20 # Generer adversariaux (peut prendre du temps selon max_iter)
21 X_adv = attack.generate(x=X_sub)
22
23 # Evaluer transfert/effet sur le modele
24 preds_clean = model.predict(X_sub)
25 preds_adv = model.predict(X_adv)
26
27 from sklearn.metrics import accuracy_score
28 print("Accuracy_clean:", accuracy_score(y_test[:100], preds_clean))
29 print("Accuracy_adv:", accuracy_score(y_test[:100], preds_adv))

```

2.2 2. Attaque par Transfert (Surrogate)

Idee :

1. Entrainer un reseau neuronal qui imite notre modele cible
2. Generer des attaques sur ce reseau (plus facile)
3. Tester si ces attaques fonctionnent sur le modele original

Avantage : Plus efficace, permet d'attaquer des modeles non-differentiables

Listing 2 – Exemple minimal Surrogate + FGSM

```

1 # Exemple minimal - surrogate (Keras) + FGSM
2 import numpy as np
3 from tensorflow.keras.models import Sequential
4 from tensorflow.keras.layers import Dense, InputLayer
5 from art.estimators.classification import KerasClassifier
6 from art.attacks.evasion import FastGradientMethod
7 from sklearn.metrics import accuracy_score
8
9 # 1) creer les labels "surrogate" comme predictions du target
10 y_surrogate = model.predict(X_train) # si classification multi-
11   classe, shape (n,)
12 # si model.predict renvoie classes, transformer en one-hot si besoin pour
13   Keras
14 from tensorflow.keras.utils import to_categorical
15 num_classes = len(np.unique(y_surrogate))
16 y_surrogate_oh = to_categorical(y_surrogate, num_classes)
17
18 # 2) petit reseau Keras comme surrogate
19 surrogate = Sequential([
20     InputLayer(input_shape=(X_train.shape[1],)),
21     Dense(128, activation='relu'),
22     Dense(64, activation='relu'),
23     Dense(num_classes, activation='softmax')
24 ])
25 surrogate.compile(optimizer='adam', loss='categorical_crossentropy', metrics
26  =['accuracy'])
27
28 # 3) Normaliser / clip_values correctement (ici on suppose features deja
29   dans [0,1])

```

```

26 surrogate.fit(X_train.astype(np.float32), y_surrogate_oh, epochs=10,
27   batch_size=64, verbose=1)
28
29 # 4) wrapper ART
30 keras_clf = KerasClassifier(model=surrogate, clip_values=(0.0, 1.0))
31
32 # 5) creer adversariaux via FGSM
33 attack = FastGradientMethod(estimator=keras_clf, eps=0.05) # eps =
34   magnitude
35 X_sub = X_test[:500].astype(np.float32)
36 X_adv = attack.generate(x=X_sub)
37
38 # 6) tester l'effet sur vos modeles scikit-learn
39 preds_clean = model.predict(X_sub)
40 preds_adv = model.predict(X_adv)
41 print("Acc_clean:", accuracy_score(y_test[:len(X_sub)], preds_clean))
42 print("Acc_adv:", accuracy_score(y_test[:len(X_sub)], preds_adv))

```

3 Guide d'Implementation

3.1 Code Complet - Demo d'Attaques

Listing 3 – Script complet demo_attack.py

```

1 """ Demo script: entrainer vos modeles et effectuer deux attaques
2   adversariales
3
4 - Attack 1: HopSkipJump (decision-based / black-box) via ART
5 - Attack 2: Surrogate (Keras) + FGSM puis test de transferabilite
6
7 Usage:
8   python scripts/attacks/demo_attack.py
9
10 Pre-requis:
11   pip install adversarial-robustness-toolbox tensorflow scikit-learn numpy
12 """
13 import warnings
14 warnings.filterwarnings('ignore')
15
16 import numpy as np
17 from sklearn.datasets import make_classification
18 from sklearn.model_selection import train_test_split
19 from sklearn.preprocessing import MinMaxScaler
20
21 from scripts.models.random_forest import train_random_forest
22 from scripts.models.knn import train_knn
23 from scripts.models.decision_tree import train_decision_tree
24
25 try:
26     from art.estimators.classification import SklearnClassifier,
27       KerasClassifier
28     from art.attacks.evasion import HopSkipJump, FastGradientMethod
29 except Exception as e:
30     raise ImportError(
31         "This demo requires 'adversarial-robustness-toolbox'.\n"
32         "Install with: pip install adversarial-robustness-toolbox tensorflow"
33     )
34     ) from e

```

```

33 from tensorflow.keras.models import Sequential
34 from tensorflow.keras.layers import Dense, InputLayer
35 from tensorflow.keras.utils import to_categorical
36
37
38 def make_scaled_dataset(n_samples=2000, n_features=20, n_classes=2,
39   random_state=0):
40   X, y = make_classification(
41     n_samples=n_samples,
42     n_features=n_features,
43     n_informative=int(n_features / 2),
44     n_redundant=0,
45     n_classes=n_classes,
46     random_state=random_state,
47   )
48   X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
49     =0.25, random_state=random_state)
50   scaler = MinMaxScaler(feature_range=(0.0, 1.0))
51   X_train = scaler.fit_transform(X_train)
52   X_test = scaler.transform(X_test)
53   return X_train.astype(np.float32), X_test.astype(np.float32), y_train.
54     astype(np.int64), y_test.astype(np.int64)
55
56
57 def train_models(X_train, y_train):
58   # Utilise vos fonctions d'entraînement
59   rf_model, rf_cv = train_random_forest(X_train, y_train, n_estimators=50,
60     max_depth=8)
61   knn_model, knn_cv = train_knn(X_train, y_train, n_neighbors=5)
62   dt_model, dt_cv = train_decision_tree(X_train, y_train, max_depth=8)
63   return rf_model, knn_model, dt_model
64
65
66 def hopskip_attack_and_eval(model, X_sub, y_sub, clip=(0.0, 1.0), max_iter
67 =20):
68   art_clf = SklearnClassifier(model=model, clip_values=clip)
69   attack = HopSkipJump(classifier=art_clf, max_iter=max_iter)
70   print("Generating HopSkipJump adversarials (this may take a while)...")
71   X_adv = attack.generate(x=X_sub)
72   preds_clean = model.predict(X_sub)
73   preds_adv = model.predict(X_adv)
74   acc_clean = (preds_clean == y_sub).mean()
75   acc_adv = (preds_adv == y_sub).mean()
76   return acc_clean, acc_adv
77
78
79 def surrogate_fgsm_transfer(rf_model, knn_model, dt_model, X_train, y_train,
80   X_sub, y_sub, clip=(0.0, 1.0), eps=0.05):
81   # Construire labels à partir du modèle cible (ici on s'appuie sur
82   # rf_model comme source des étiquettes)
83   surrogate_labels = rf_model.predict(X_train)
84   num_classes = len(np.unique(surrogate_labels))
85   surrogate_labels_oh = to_categorical(surrogate_labels, num_classes)
86
87   # Petit réseau Keras
88   surrogate = Sequential([
89     InputLayer(input_shape=(X_train.shape[1],)),
90     Dense(128, activation='relu'),
91     Dense(64, activation='relu'),
92     Dense(num_classes, activation='softmax')
93   ])

```

```

87     surrogate.compile(optimizer='adam', loss='categorical_crossentropy',
88                         metrics=['accuracy'])
89     surrogate.fit(X_train, surrogate_labels_oh, epochs=5, batch_size=64,
90                     verbose=0)
91
92     keras_clf = KerasClassifier(model=surrogate, clip_values=clip)
93     fgsm = FastGradientMethod(estimator=keras_clf, eps=eps)
94     print("Generating FGSM adversarials on surrogate...")
95     X_adv = fgsm.generate(x=X_sub)
96
97     # Evaluer la transferabilite sur les 3 modeles
98     results = {}
99     for name, model in (('RandomForest', rf_model), ('KNN', knn_model), ('DecisionTree', dt_model)):
100         preds_clean = model.predict(X_sub)
101         preds_adv = model.predict(X_adv)
102         acc_clean = (preds_clean == y_sub).mean()
103         acc_adv = (preds_adv == y_sub).mean()
104         results[name] = (acc_clean, acc_adv)
105
106     return results
107
108 def main():
109     print("Preparing dataset...")
110     X_train, X_test, y_train, y_test = make_scaled_dataset(n_samples=3000,
111                 n_features=20, n_classes=2)
112
113     print("Training models (this is fast with small data)...")
114     rf_model, knn_model, dt_model = train_models(X_train, y_train)
115
116     # Sous-ensemble pour attaques (plus rapide)
117     X_sub = X_test[:100]
118     y_sub = y_test[:100]
119
120     print('\n==== HopSkipJump attack on RandomForest ===')
121     acc_clean, acc_adv = hopskip_attack_and_eval(rf_model, X_sub, y_sub,
122             clip=(0.0, 1.0), max_iter=20)
123     print(f"RandomForest - clean acc: {acc_clean:.3f}, adv acc: {acc_adv:.3f}")
124
125     print('\n==== Surrogate (FGSM) transfer attack ===')
126     transfer_results = surrogate_fgsm_transfer(rf_model, knn_model, dt_model,
127             X_train, y_train, X_sub, y_sub, clip=(0.0, 1.0), eps=0.05)
128     for model_name, (acc_clean, acc_adv) in transfer_results.items():
129         print(f"{model_name} - clean acc: {acc_clean:.3f}, adv acc (transferred): {acc_adv:.3f}")

if __name__ == '__main__':
    main()

```

3.2 Pre-requis Techniques

```
1 pip install adversarial-robustness-toolbox tensorflow scikit-learn numpy
```

3.3 Structure du Projet

scripts/

```
|-- attacks/
|   |-- demo_attack.py          # Script principal
|-- models/
|   |-- random_forest.py        # Modeles existants
|   |-- knn.py
|   |-- decision_tree.py
```

3.4 Execution

```
1 python scripts/attacks/demo_attack.py
```

4 Metriques d’Evaluation

4.1 Calcul de robustesse

Success Rate = (Nb attaques réussies / Nb total) * 100
Accuracy Drop = Accuracy_avant - Accuracy_apres

4.2 Template de rapport

Test [Modele] - [Methode]

- Modele : Random Forest
- Methode : HopSkipJump
- Accuracy avant : 95% -> apres : 72%
- Taux de succes : 23%
- Features sensibles : [Src_bytes, Duration]
- Temps execution : 15 min

Conclusion

Ces deux methodes (HopSkipJump et Surrogate+FGSM) sont les standards pour tester la robustesse de vos modeles scikit-learn. Commencer par HopSkipJump pour une evaluation rapide, puis utiliser la methode surrogate pour des analyses plus poussées.