

PROJET COMPLET : SYSTÈME DE DÉTECTION D'INTRUSIONS AVEC ML

Guide Étape par Étape - De Zéro à GitHub

Auteur : Ton Nom

Date : Janvier 2025

Temps total : 10-15 heures (sur 3-5 jours)

Prérequis : VSCode, Python, pandas (que tu as déjà ✓)

TABLE DES MATIÈRES

- [Étape 0 : Setup et Vérification](#)
 - [Étape 1 : Télécharger le Dataset](#)
 - [Étape 2 : Exploration des Données](#)
 - [Étape 3 : Prétraitement](#)
 - [Étape 4 : Sélection des Features](#)
 - [Étape 5 : Entraînement du Modèle](#)
 - [Étape 6 : Évaluation](#)
 - [Étape 7 : Visualisations](#)
 - [Étape 8 : GitHub et Documentation](#)
-

ÉTAPE 0 : SETUP ET VÉRIFICATION

Temps : 10 minutes

0.1 Vérifier Python et pandas

Ouvre VSCode et le terminal intégré (Ctrl + `)

```
bash

# Vérifier Python
python --version
# Ou
python3 --version

# Vérifier pandas
python -c "import pandas as pd; print(pd.__version__)"
```

Résultat attendu :

```
Python 3.x.x  
2.x.x (version de pandas)
```

0.2 Installer les bibliothèques manquantes

```
bash  
  
# Installer tout d'un coup  
pip install numpy matplotlib seaborn scikit-learn jupyter  
  
# Vérifier  
pip list | grep -E "numpy|matplotlib|seaborn|scikit"
```

0.3 Crée la structure du projet

Windows (PowerShell ou cmd) :

```
cmd  
  
cd Desktop  
mkdir IDS-ML-Project  
cd IDS-ML-Project  
mkdir data src results notebooks docs
```

Linux/Mac :

```
bash  
  
cd ~/Desktop  
mkdir IDS-ML-Project  
cd IDS-ML-Project  
mkdir data src results notebooks docs
```

Structure finale :

```
IDS-ML-Project/  
├── data/      # Datasets  
├── src/       # Code Python  
├── results/   # Graphiques, métriques  
├── notebooks/ # Jupyter notebooks  
├── docs/      # Documentation  
└── README.md  # (on va créer)
```

0.4 Ouvrir dans VSCode

```
bash  
# Dans le dossier IDS-ML-Project  
code .
```

 **Checkpoint :** VSCode ouvert avec la structure du projet visible

ÉTAPE 1 : TÉLÉCHARGER LE DATASET

Temps : 30-60 minutes (selon connexion)

1.1 Aller sur le site officiel

 **Lien :** <https://www.unb.ca/cic/datasets/ids-2017.html>

1.2 Télécharger les fichiers

Pour commencer, télécharge ces 3 fichiers :

1. **Monday-WorkingHours.pcap_ISCX.csv** (~450 MB)
 - Trafic NORMAL uniquement
2. **Friday-WorkingHours-Afternoon-DDoS.pcap_ISCX.csv** (~600 MB)
 - Attaques DDoS
3. **Friday-WorkingHours-Afternoon-PortScan.pcap_ISCX.csv** (~300 MB)
 - Scans de ports

1.3 Déplacer dans le dossier data/

Une fois téléchargés, déplace-les :

Windows :

```
De : C:\Users\TonNom\Downloads\  
Vers : C:\Users\TonNom\Desktop\IDS-ML-Project\data\
```

Linux/Mac :

```
bash  
mv ~/Downloads/*.csv ~/Desktop/IDS-ML-Project/data/
```

1.4 Vérifier

```
bash  
ls data/  
# Ou Windows  
dir data
```

Tu devrais voir :

```
Monday-WorkingHours.pcap_ISCX.csv  
Friday-WorkingHours-Afternoon-DDoS.pcap_ISCX.csv  
Friday-WorkingHours-Afternoon-PortScan.pcap_ISCX.csv
```

Checkpoint : 3 fichiers CSV dans le dossier data/

ÉTAPE 2 : EXPLORATION DES DONNÉES

Temps : 30 minutes

2.1 Créer le script d'exploration

Dans VSCode, crée `src/01_explore_data.py` :

```
python
```

....

Étape 2 : Exploration du Dataset CICIDS2017

....

```
import pandas as pd
import numpy as np

print("==" * 80)
print("EXPLORATION DU DATASET CICIDS2017")
print("==" * 80)

# Charger les datasets
print("\n📁 Chargement des fichiers...")

df_monday = pd.read_csv('data/Monday-WorkingHours.pcap_ISCX.csv')
print(f"✓ Monday (Normal) : {df_monday.shape}")

df_ddos = pd.read_csv('data/Friday-WorkingHours-Afternoon-DDoS.pcap_ISCX.csv')
print(f"✓ Friday DDoS : {df_ddos.shape}")

df_portscan = pd.read_csv('data/Friday-WorkingHours-Afternoon-PortScan.pcap_ISCX.csv')
print(f"✓ Friday PortScan : {df_portscan.shape}")

# Informations de base
print("\n" + "==" * 80)
print("📊 INFORMATIONS GÉNÉRALES")
print("==" * 80)

print(f"\nMonday (Trafic Normal) :")
print(f" Lignes : {len(df_monday)}")
print(f" Colonnes : {len(df_monday.columns)}")
print(f" Mémoire : {df_monday.memory_usage(deep=True).sum() / 1024**2:.1f} MB")

print(f"\nFriday DDoS :")
print(f" Lignes : {len(df_ddos)}")
print(f" Mémoire : {df_ddos.memory_usage(deep=True).sum() / 1024**2:.1f} MB")

print(f"\nFriday PortScan :")
print(f" Lignes : {len(df_portscan)}")
print(f" Mémoire : {df_portscan.memory_usage(deep=True).sum() / 1024**2:.1f} MB")

# Colonnes
print("\n" + "==" * 80)
print("📋 COLONNES DISPONIBLES")
print("==" * 80)
```

```

print(f"\nNombre total de colonnes : {len(df_monday.columns)}")
print("\nPremières 10 colonnes :")
for i, col in enumerate(df_monday.columns[:10], 1):
    print(f" {i:2d}. {col}")
print(f" ...")
print(f" {len(df_monday.columns)}. {df_monday.columns[-1]}")

# Distribution des labels
print("\n" + "=" * 80)
print("👉 DISTRIBUTION DES LABELS")
print("=" * 80)

print("\nMonday :")
print(df_monday['Label'].value_counts())

print("\nFriday DDoS :")
print(df_ddos['Label'].value_counts())

print("\nFriday PortScan :")
print(df_portscan['Label'].value_counts())

# Valeurs manquantes
print("\n" + "=" * 80)
print("❓ VALEURS MANQUANTES")
print("=" * 80)

print(f"\nMonday : {df_monday.isnull().sum().sum()} valeurs manquantes")
print(f"DDoS : {df_ddos.isnull().sum().sum()} valeurs manquantes")
print(f"PortScan : {df_portscan.isnull().sum().sum()} valeurs manquantes")

# Statistiques
print("\n" + "=" * 80)
print("📈 STATISTIQUES SUR QUELQUES COLONNES")
print("=" * 80)

colonnes_clés = ['Flow Duration', 'Total Fwd Packets', 'Total Backward Packets']
print("\nStatistiques pour DDoS :")
print(df_ddos[colonnes_clés].describe())

# Sauvegarder résumé
with open('docs/01_exploration_summary.txt', 'w', encoding='utf-8') as f:
    f.write("RÉSUMÉ EXPLORATION - CICIDS2017\n")
    f.write("=" * 80 + "\n\n")
    f.write(f"Monday : {len(df_monday):,} lignes\n")
    f.write(f"DDoS : {len(df_ddos):,} lignes\n")
    f.write(f"PortScan : {len(df_portscan):,} lignes\n")
    f.write(f"\nTotal : {len(df_monday) + len(df_ddos) + len(df_portscan):,} connexions\n")

```

```
print("\n✓ Résumé sauvegardé dans 'docs/01_exploration_summary.txt'")  
print("\n" + "=" * 80)  
print("✓ EXPLORATION TERMINÉE !")  
print("=" * 80)
```

2.2 Lancer le script

Dans le terminal VSCode :

```
bash  
python src/01_explore_data.py
```

2.3 Résultat attendu

Tu devrais voir :

EXPLORATION DU DATASET CICIDS2017

📁 Chargement des fichiers...
✓ Monday (Normal) : (529918, 79)
✓ Friday DDoS : (225745, 79)
✓ Friday PortScan : (286467, 79)

📊 INFORMATIONS GÉNÉRALES

Monday (Trafic Normal) :

Lignes : 529,918

Colonnes : 79

Mémoire : 318.6 MB

... (etc.)

✓ **Checkpoint** : Script exécuté sans erreur, fichier summary créé

ÉTAPE 3 : PRÉTRAITEMENT (NETTOYAGE)

Temps : 1 heure

3.1 Créer le script de prétraitement

Crée `src/02_preprocess_data.py` :

```
python
```

....

Étape 3 : Prétraitement et Nettoyage des Données

....

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import LabelEncoder
import pickle

print("==" * 80)
print("ÉTAPE 3 : PRÉTRAITEMENT DES DONNÉES")
print("==" * 80)

# =====
# 3.1 CHARGER ET COMBINER LES DATASETS
# =====
print("\n3.1 - Chargement et combinaison...")

df_monday = pd.read_csv('data/Monday-WorkingHours.pcap_ISCX.csv')
df_ddos = pd.read_csv('data/Friday-WorkingHours-Afternoon-DDoS.pcap_ISCX.csv')
df_portscan = pd.read_csv('data/Friday-WorkingHours-Afternoon-PortScan.pcap_ISCX.csv')

print(f"Monday : {len(df_monday)} lignes")
print(f"DDoS : {len(df_ddos)} lignes")
print(f"PortScan : {len(df_portscan)} lignes")

# Combiner
df = pd.concat([df_monday, df_ddos, df_portscan], ignore_index=True)
print(f"\n✓ Dataset combiné : {len(df)} lignes")

# =====
# 3.2 NETTOYER LES NOMS DE COLONNES
# =====
print("\n" + "==" * 80)
print("3.2 - Nettoyage des noms de colonnes...")

# Supprimer espaces
df.columns = df.columns.str.strip()
print("✓ Espaces supprimés des noms de colonnes")

# =====
# 3.3 GÉRER LES VALEURS INFINIES
# =====
print("\n" + "==" * 80)
print("3.3 - Gestion des valeurs infinies...")
```

```
print(f"Valeurs infinies avant : {np.isinf(df.select_dtypes(include=[np.number])).sum().sum()}"
```

```
# Remplacer inf par NaN
```

```
df.replace([np.inf, -np.inf], np.nan, inplace=True)
```

```
print(f"Valeurs infinies après : {np.isinf(df.select_dtypes(include=[np.number])).sum().sum()}"
```

```
print("✓ Valeurs infinies remplacées par NaN")
```

```
# =====
```

```
# 3.4 GÉRER LES VALEURS MANQUANTES
```

```
# =====
```

```
print("\n" + "=" * 80)
```

```
print("3.4 - Gestion des valeurs manquantes...")
```

```
print(f"Valeurs manquantes avant : {df.isnull().sum().sum()}"
```

```
print(f"Nombre de lignes avant suppression : {len(df)}")
```

```
# Supprimer lignes avec NaN
```

```
df.dropna(inplace=True)
```

```
print(f"\nValeurs manquantes après : {df.isnull().sum().sum()}"
```

```
print(f"Nombre de lignes après suppression : {len(df)}")
```

```
print(f"Nombre de lignes supprimées : {529918 + 225745 + 286467 - len(df)}")
```

```
# =====
```

```
# 3.5 CRÉER LABEL BINAIRE (NORMAL vs ATTAQUE)
```

```
# =====
```

```
print("\n" + "=" * 80)
```

```
print("3.5 - Création du label binaire...")
```

```
print("\nDistribution des labels originaux :")
```

```
print(df['Label'].value_counts())
```

```
# 0 = BENIGN (Normal), 1 = Attaque
```

```
df['Attack'] = (df['Label'] != 'BENIGN').astype(int)
```

```
print("\nDistribution du label binaire :")
```

```
print(df['Attack'].value_counts())
```

```
print(f" 0 (Normal) : {(df['Attack'] == 0).sum();}"
```

```
print(f" 1 (Attaque) : {(df['Attack'] == 1).sum();}"
```

```
# =====
```

```
# 3.6 SUPPRIMER COLONNES INUTILES
```

```
# =====
```

```
print("\n" + "=" * 80)
```

```
print("3.6 - Suppression des colonnes inutiles...")
```

```

colonnes_a_supprimer = [
    'Flow ID', # Identifiant unique (pas utile pour ML)
    'Source IP', # Adresses IP (pas généralisables)
    'Destination IP',
    'Timestamp', # Horodatage (pas utile)
    'Label' # On a 'Attack' maintenant
]

print(f"Colonnes avant : {len(df.columns)}")

for col in colonnes_a_supprimer:
    if col in df.columns:
        df.drop(col, axis=1, inplace=True)
    print(f"✓ Supprimé : {col}")

print(f"Colonnes après : {len(df.columns)}")

# =====
# 3.7 SAUVEGARDER LE DATASET NETTOYÉ
# =====

print("\n" + "=" * 80)
print("3.7 - Sauvegarde du dataset nettoyé...")

df.to_csv('data/dataset_clean.csv', index=False)
print(f"✓ Dataset nettoyé sauvégarde : {len(df):,} lignes, {len(df.columns)} colonnes")

# Sauvegarder un échantillon (pour tests rapides)
df_sample = df.sample(n=50000, random_state=42)
df_sample.to_csv('data/dataset_sample.csv', index=False)
print(f"✓ Échantillon sauvégarde : 50,000 lignes")

# =====
# 3.8 RÉSUMÉ DU PRÉTRAITEMENT
# =====

print("\n" + "=" * 80)
print("📊 RÉSUMÉ DU PRÉTRAITEMENT")
print("=" * 80)

summary = f"""

DONNÉES INITIALES:
- Monday (Normal) : 529,918 lignes
- DDoS : 225,745 lignes
- PortScan : 286,467 lignes
- Total : 1,042,130 lignes

NETTOYAGE:
- Valeurs infinies remplacées
"""

```

- Lignes avec NaN supprimées
- Colonnes inutiles supprimées

RÉSULTAT FINAL:

- Lignes finales : `{len(df):}`
- Colonnes : `{len(df.columns)}`
- Normal : `{(df['Attack'] == 0).sum():}`
- Attaques : `{(df['Attack'] == 1).sum():}`

FICHIERS CRÉÉS:

- data/dataset_clean.csv (complet)
- data/dataset_sample.csv (50K lignes pour tests)
-

```
print(summary)
```

```
with open('docs/02_preprocessing_summary.txt', 'w', encoding='utf-8') as f:  
    f.write(summary)
```

```
print("\n✓ Résumé sauvegardé dans 'docs/02_preprocessing_summary.txt'")  
print("=" * 80)  
print("✓ PRÉTRAITEMENT TERMINÉ !")  
print("=" * 80)
```

3.2 Lancer le script

```
bash
```

```
python src/02_preprocess_data.py
```

Temps d'exécution : 5-10 minutes (selon ton PC)

Checkpoint : Fichier `dataset_clean.csv` créé dans `data/`

ÉTAPE 4 : SÉLECTION DES FEATURES

Temps : 30 minutes

4.1 Créer le script de sélection

Crée `src/03_select_features.py` :

```
python
```

....

Étape 4 : Sélection des Features Importantes

....

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.ensemble import RandomForestClassifier
```

```
print("==" * 80)
print("ÉTAPE 4 : SÉLECTION DES FEATURES")
print("==" * 80)
```

```
# Charger dataset nettoyé (échantillon pour aller plus vite)
print("\n📁 Chargement du dataset...")
df = pd.read_csv('data/dataset_sample.csv')
print(f"✓ Dataset chargé : {df.shape}")
```

```
# Séparer X et y
X = df.drop('Attack', axis=1)
y = df['Attack']
```

```
print(f"\nFeatures (X) : {X.shape}")
print(f"Target (y) : {y.shape}")
```

```
# -----
# 4.1 ENTRAÎNER UN RANDOM FOREST POUR IMPORTANCE
# -----
print("\n" + "==" * 80)
print("4.1 - Calcul de l'importance des features...")
```

```
print("Entraînement d'un Random Forest...")
rf = RandomForestClassifier(n_estimators=100, random_state=42, n_jobs=-1)
rf.fit(X, y)
print("✓ Random Forest entraîné")
```

```
# Importance des features
importances = rf.feature_importances_
feature_importance_df = pd.DataFrame({
    'Feature': X.columns,
    'Importance': importances
}).sort_values('Importance', ascending=False)
```

```
print("\n📊 Top 20 Features les plus importantes :")
print(feature_importance_df.head(20))
```

```

# =====
# 4.2 SÉLECTIONNER LES MEILLEURES FEATURES
# =====
print("\n" + "=" * 80)
print("4.2 - Sélection des meilleures features...")

# Garder features avec importance > 0.01
selected_features = feature_importance_df[feature_importance_df['Importance'] > 0.01]['Feature'].tolist()

print(f"\nFeatures sélectionnées : {len(selected_features)}")
print("\nListe des features sélectionnées :")
for i, feat in enumerate(selected_features, 1):
    imp = feature_importance_df[feature_importance_df['Feature'] == feat]['Importance'].values[0]
    print(f" {i:2d}. {feat}: {imp:.4f} (Importance: {imp:.4f})")

# =====
# 4.3 VISUALISER L'IMPORTANCE
# =====
print("\n" + "=" * 80)
print("4.3 - Visualisation de l'importance...")

plt.figure(figsize=(12, 8))
top_n = 20
top_features = feature_importance_df.head(top_n)

plt.barh(range(top_n), top_features['Importance'], color='steelblue', edgecolor='black')
plt.yticks(range(top_n), top_features['Feature'])
plt.xlabel('Importance', fontsize=12, fontweight='bold')
plt.title(f'Top {top_n} Features Importantes pour la Détection d'Intrusions',
          fontsize=14, fontweight='bold', pad=20)
plt.gca().invert_yaxis()
plt.grid(axis='x', alpha=0.3)
plt.tight_layout()
plt.savefig('results/feature_importance.png', dpi=300, bbox_inches='tight')
print("✓ Graphique sauvégarde : results/feature_importance.png")
plt.show()

# =====
# 4.4 SAUVEGARDER LA LISTE DES FEATURES
# =====
print("\n" + "=" * 80)
print("4.4 - Sauvegarde de la liste des features...")

import pickle

with open('data/selected_features.pkl', 'wb') as f:
    pickle.dump(selected_features, f)

```

```
print("✓ Liste sauvegardée : {len(selected_features)} features")  
  
with open('docs/03_selected_features.txt', 'w', encoding='utf-8') as f:  
    f.write(f'FEATURES SÉLECTIONNÉES ({len(selected_features)})\n')  
    f.write("=" * 80 + "\n\n")  
    for i, feat in enumerate(selected_features, 1):  
        imp = feature_importance_df[feature_importance_df['Feature'] == feat]['Importance'].values[0]  
        f.write(f'{i:2d}. {feat}: {imp:.4f}\n')  
  
print("✓ Liste sauvegardée : docs/03_selected_features.txt")  
  
print("\n" + "=" * 80)  
print("✓ SÉLECTION DES FEATURES TERMINÉE !")  
print("=" * 80)
```

4.2 Lancer le script

bash

```
python src/03_select_features.py
```

Checkpoint : Graphique d'importance créé, liste de features sauvegardée

ÉTAPE 5 : ENTRAÎNEMENT DU MODÈLE

Temps : 1 heure

5.1 Crée le script d'entraînement

Crée `src/04_train_model.py` :

python

....

Étape 5 : Entraînement du Modèle de Détection d'Intrusions

....

```
import pandas as pd
import numpy as np
import pickle
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
import time

print("=" * 80)
print("ÉTAPE 5 : ENTRAÎNEMENT DU MODÈLE")
print("=" * 80)

# =====
# 5.1 CHARGER LES DONNÉES ET FEATURES SÉLECTIONNÉES
# =====
print("\n5.1 - Chargement des données...")

df = pd.read_csv('data/dataset_clean.csv')
print(f"✓ Dataset : {df.shape}")

with open('data/selected_features.pkl', 'rb') as f:
    selected_features = pickle.load(f)

print(f"✓ Features sélectionnées : {len(selected_features)}")

# =====
# 5.2 PRÉPARER X ET Y
# =====
print("\n" + "=" * 80)
print("5.2 - Préparation des données...")

X = df[selected_features]
y = df['Attack']

print(f"X shape : {X.shape}")
print(f"y shape : {y.shape}")
print(f"\nDistribution des classes :")
print(f" Normal (0) : {(y == 0).sum()}/{(y == 0).sum() / len(y) * 100:.1f}%")
print(f" Attaque (1) : {(y == 1).sum()}/{(y == 1).sum() / len(y) * 100:.1f}%")
```

```

# =====
# 5.3 SPLIT TRAIN/TEST
# =====
print("\n" + "=" * 80)
print("5.3 - Split train/test (80/20)...")

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42, stratify=y
)

print(f"Train set : {X_train.shape[0]} échantillons ({X_train.shape[0] / len(X) * 100:.1f}%)")
print(f"Test set : {X_test.shape[0]} échantillons ({X_test.shape[0] / len(X) * 100:.1f}%)"
# =====
# 5.4 NORMALISATION (StandardScaler)
# =====
print("\n" + "=" * 80)
print("5.4 - Normalisation des features...")

scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

print("✓ Features normalisées")
print(f"Moyenne des features après normalisation : {X_train_scaled.mean():.6f} (≈ 0)")
print(f"Écart-type après normalisation : {X_train_scaled.std():.6f} (≈ 1)")

# Sauvegarder le scaler
with open('data/scaler.pkl', 'wb') as f:
    pickle.dump(scaler, f)
print("✓ Scaler sauvégarde")

# =====
# 5.5 ENTRAÎNER DECISION TREE
# =====
print("\n" + "=" * 80)
print("5.5 - Entraînement Decision Tree...")

start_time = time.time()

dt = DecisionTreeClassifier(max_depth=10, random_state=42)
dt.fit(X_train_scaled, y_train)

dt_time = time.time() - start_time

print(f"✓ Decision Tree entraîné en {dt_time:.2f} secondes")
print(f" Profondeur : {dt.get_depth()}")

```

```

print(f" Feuilles : {dt.get_n_leaves()}")


# Prédictions
y_pred_dt = dt.predict(X_test_scaled)

# Métriques
dt_accuracy = accuracy_score(y_test, y_pred_dt)
dt_precision = precision_score(y_test, y_pred_dt)
dt_recall = recall_score(y_test, y_pred_dt)
dt_f1 = f1_score(y_test, y_pred_dt)

print(f"\n📊 Performances Decision Tree :")
print(f" Accuracy : {dt_accuracy:.4f} ({dt_accuracy * 100:.2f}%)")
print(f" Precision : {dt_precision:.4f}")
print(f" Recall : {dt_recall:.4f}")
print(f" F1-Score : {dt_f1:.4f}")

# =====
# 5.6 ENTRAÎNER RANDOM FOREST
# =====

print("\n" + "=" * 80)
print("5.6 - Entraînement Random Forest...")

start_time = time.time()

rf = RandomForestClassifier(n_estimators=100, max_depth=10, random_state=42, n_jobs=-1)
rf.fit(X_train_scaled, y_train)

rf_time = time.time() - start_time

print(f"✓ Random Forest entraîné en {rf_time:.2f} secondes")
print(f" Arbres : {rf.n_estimators}")

# Prédictions
y_pred_rf = rf.predict(X_test_scaled)

# Métriques
rf_accuracy = accuracy_score(y_test, y_pred_rf)
rf_precision = precision_score(y_test, y_pred_rf)
rf_recall = recall_score(y_test, y_pred_rf)
rf_f1 = f1_score(y_test, y_pred_rf)

print(f"\n📊 Performances Random Forest :")
print(f" Accuracy : {rf_accuracy:.4f} ({rf_accuracy * 100:.2f}%)")
print(f" Precision : {rf_precision:.4f}")
print(f" Recall : {rf_recall:.4f}")
print(f" F1-Score : {rf_f1:.4f}")

```

```

# =====
# 5.7 COMPARAISON DES MODÈLES
# =====
print("\n" + "=" * 80)
print("5.7 - Comparaison des modèles...")

comparison = pd.DataFrame({
    'Modèle': ['Decision Tree', 'Random Forest'],
    'Accuracy': [dt_accuracy, rf_accuracy],
    'Precision': [dt_precision, rf_precision],
    'Recall': [dt_recall, rf_recall],
    'F1-Score': [dt_f1, rf_f1],
    'Temps (s)': [dt_time, rf_time]
})

print("\n" + comparison.to_string(index=False))

# Meilleur modèle
best_model_name = comparison.loc[comparison['Accuracy'].idxmax(), 'Modèle']
best_accuracy = comparison['Accuracy'].max()

print(f"\n🏆 Meilleur modèle : {best_model_name} (Accuracy: {best_accuracy:.4f})"

# =====
# 5.8 SAUVEGARDER LES MODÈLES
# =====
print("\n" + "=" * 80)
print("5.8 - Sauvegarde des modèles... ")

with open('data/model_decision_tree.pkl', 'wb') as f:
    pickle.dump(dt, f)
print("✓ Decision Tree sauvégarde")

with open('data/model_random_forest.pkl', 'wb') as f:
    pickle.dump(rf, f)
print("✓ Random Forest sauvégarde")

# Sauvegarder métriques
metrics = {
    'decision_tree': {
        'accuracy': dt_accuracy,
        'precision': dt_precision,
        'recall': dt_recall,
        'f1': dt_f1
    },
    'random_forest': {

```

```
'accuracy': rf_accuracy,
'precision': rf_precision,
'recall': rf_recall,
'f1': rf_f1
}
}

with open('data/metrics.pkl', 'wb') as f:
    pickle.dump(metrics, f)
print("✓ Métriques sauvegardées")

print("\n" + "=" * 80)
print("✓ ENTRAÎNEMENT TERMINÉ !")
print("=" * 80)
```

5.2 Lancer le script

```
bash
python src/04_train_model.py
```

Temps d'exécution : 10-30 minutes (selon ton PC et taille du dataset)

Checkpoint : 2 modèles entraînés et sauvegardés

ÉTAPE 6 : ÉVALUATION DES PERFORMANCES

Temps : 30 minutes

6.1 Créer le script d'évaluation

Crée `src/05_evaluate_model.py` :

```
python
```

....

Étape 6 : Évaluation Détailée des Modèles

....

```
import pandas as pd
import numpy as np
import pickle
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix, classification_report

print("==" * 80)
print("ÉTAPE 6 : ÉVALUATION DES MODÈLES")
print("==" * 80)

# =====
# 6.1 CHARGER DONNÉES ET MODÈLES
# =====
print("\n6.1 - Chargement...")

df = pd.read_csv('data/dataset_clean.csv')

with open('data/selected_features.pkl', 'rb') as f:
    selected_features = pickle.load(f)

with open('data/scaler.pkl', 'rb') as f:
    scaler = pickle.load(f)

with open('data/model_decision_tree.pkl', 'rb') as f:
    dt = pickle.load(f)

with open('data/model_random_forest.pkl', 'rb') as f:
    rf = pickle.load(f)

print("✓ Tout chargé")

# Préparer données
X = df[selected_features]
y = df['Attack']

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42, stratify=y
)

X_test_scaled = scaler.transform(X_test)
```

```

# =====
# 6.2 CONFUSION MATRIX - DECISION TREE
# =====
print("\n" + "=" * 80)
print("6.2 - Confusion Matrix - Decision Tree...")

y_pred_dt = dt.predict(X_test_scaled)
cm_dt = confusion_matrix(y_test, y_pred_dt)

plt.figure(figsize=(8, 6))
sns.heatmap(cm_dt, annot=True, fmt='d', cmap='Blues',
            xticklabels=['Normal', 'Attaque'],
            yticklabels=['Normal', 'Attaque'],
            cbar_kws={'label': 'Count'})
plt.title('Confusion Matrix - Decision Tree', fontsize=14, fontweight='bold', pad=20)
plt.ylabel('Vraie Classe', fontsize=12)
plt.xlabel('Classe Prédite', fontsize=12)

# Ajouter pourcentages
for i in range(2):
    for j in range(2):
        pct = cm_dt[i, j] / cm_dt[i].sum() * 100
        plt.text(j + 0.5, i + 0.7, f'{pct:.1f}%', ha='center', va='center', fontsize=10, color='red')

plt.tight_layout()
plt.savefig('results/confusion_matrix_dt.png', dpi=300, bbox_inches='tight')
print("✓ Sauvegardé : results/confusion_matrix_dt.png")
plt.show()

# =====
# 6.3 CONFUSION MATRIX - RANDOM FOREST
# =====
print("\n" + "=" * 80)
print("6.3 - Confusion Matrix - Random Forest...")

y_pred_rf = rf.predict(X_test_scaled)
cm_rf = confusion_matrix(y_test, y_pred_rf)

plt.figure(figsize=(8, 6))
sns.heatmap(cm_rf, annot=True, fmt='d', cmap='Greens',
            xticklabels=['Normal', 'Attaque'],
            yticklabels=['Normal', 'Attaque'],
            cbar_kws={'label': 'Count'})
plt.title('Confusion Matrix - Random Forest', fontsize=14, fontweight='bold', pad=20)
plt.ylabel('Vraie Classe', fontsize=12)

```

```

plt.xlabel('Classe Prédictive', fontsize=12)

for i in range(2):
    for j in range(2):
        pct = cm_rf[i, j] / cm_rf[i].sum() * 100
        plt.text(j + 0.5, i + 0.7, f'{pct:.1f}%', ha='center', va='center', fontsize=10, color='red')

plt.tight_layout()
plt.savefig('results/confusion_matrix_rf.png', dpi=300, bbox_inches='tight')
print("✓ Sauvegardé : results/confusion_matrix_rf.png")
plt.show()

# =====
# 6.4 CLASSIFICATION REPORT
# =====

print("\n" + "=" * 80)
print("6.4 - Classification Report...")

print("\nDECISION TREE :")
print(classification_report(y_test, y_pred_dt,
                            target_names=['Normal', 'Attaque']))

print("\nRANDOM FOREST :")
print(classification_report(y_test, y_pred_rf,
                            target_names=['Normal', 'Attaque']))

# =====
# 6.5 SAUVEGARDER RAPPORT
# =====

print("\n" + "=" * 80)
print("6.5 - Sauvegarde du rapport...")

with open('docs/05_evaluation_report.txt', 'w', encoding='utf-8') as f:
    f.write("RAPPORT D'ÉVALUATION - IDS ML\n")
    f.write("=" * 80 + "\n\n")

    f.write("DECISION TREE:\n")
    f.write("-" * 80 + "\n")
    f.write(classification_report(y_test, y_pred_dt, target_names=['Normal', 'Attaque']))
    f.write("\n\n")

    f.write("RANDOM FOREST:\n")
    f.write("-" * 80 + "\n")
    f.write(classification_report(y_test, y_pred_rf, target_names=['Normal', 'Attaque']))

print("✓ Rapport sauvegardé : docs/05_evaluation_report.txt")

```

```
print("\n" + "=" * 80)
print("✓ ÉVALUATION TERMINÉE !")
print("=" * 80)
```

6.2 Lancer le script

bash

```
python src/05_evaluate_model.py
```

✓ **Checkpoint :** Confusion matrices créées, rapport d'évaluation sauvegardé

ÉTAPE 7 : VISUALISATIONS FINALES

Temps : 20 minutes

7.1 Créer le script de visualisation

Crée `src/06_visualizations.py` :

python

....

Étape 7 : Visualisations Finales du Projet

....

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import pickle

print("==" * 80)
print("ÉTAPE 7 : VISUALISATIONS FINALES")
print("==" * 80)

# Charger métriques
with open('data/metrics.pkl', 'rb') as f:
    metrics = pickle.load(f)

# =====
# 7.1 COMPARAISON DES MODÈLES
# =====
print("\n7.1 - Graphique de comparaison...")

dt_metrics = metrics['decision_tree']
rf_metrics = metrics['random_forest']

metrics_df = pd.DataFrame({
    'Métrique': ['Accuracy', 'Precision', 'Recall', 'F1-Score'],
    'Decision Tree': [dt_metrics['accuracy'], dt_metrics['precision'],
                      dt_metrics['recall'], dt_metrics['f1']],
    'Random Forest': [rf_metrics['accuracy'], rf_metrics['precision'],
                      rf_metrics['recall'], rf_metrics['f1']]
})

fig, ax = plt.subplots(figsize=(10, 6))

x = range(len(metrics_df))
width = 0.35

bars1 = ax.bar([i - width/2 for i in x], metrics_df['Decision Tree'],
               width, label='Decision Tree', color='steelblue', edgecolor='black')
bars2 = ax.bar([i + width/2 for i in x], metrics_df['Random Forest'],
               width, label='Random Forest', color='forestgreen', edgecolor='black')

ax.set_xlabel('Métriques', fontsize=12, fontweight='bold')
ax.set_ylabel('Score', fontsize=12, fontweight='bold')
ax.set_title('Comparaison des Performances des Modèles', fontsize=14, fontweight='bold', pad=20)
```

```

ax.set_xticks(x)
ax.set_xticklabels(metrics_df['Métrique'])
ax.legend()
ax.set_xlim(0, 1.1)
ax.grid(axis='y', alpha=0.3)

# Ajouter valeurs sur les barres
for bars in [bars1, bars2]:
    for bar in bars:
        height = bar.get_height()
        ax.text(bar.get_x() + bar.get_width()/2., height,
                f'{height:.3f}',
                ha='center', va='bottom', fontsize=10, fontweight='bold')

plt.tight_layout()
plt.savefig('results/model_comparison.png', dpi=300, bbox_inches='tight')
print("✓ Sauvegardé : results/model_comparison.png")
plt.show()

print("\n" + "=" * 80)
print("✓ VISUALISATIONS TERMINÉES !")
print("=" * 80)

```

7.2 Lancer le script

```

bash
python src/06_visualizations.py

```

Checkpoint : Graphique de comparaison créé

ÉTAPE 8 : GITHUB ET DOCUMENTATION

Temps : 30 minutes

8.1 Créer README.md

Dans VSCode, crée `README.md` à la racine du projet :

```
markdown
```

🔒 Système de Détection d'Intrusions avec Machine Learning

📄 Description

Projet de détection d'intrusions réseau utilisant Machine Learning pour identifier les attaques DDoS et scans de ports.

🎯 Objectif

Développer un IDS (Intrusion Detection System) capable de classifier le trafic réseau en :

- Trafic normal (BENIGN)
- Attaques (DDoS, Port Scan)

📊 Dataset

CICIDS2017 - Canadian Institute for Cybersecurity

- Monday : Trafic normal (~530K connexions)
- Friday DDoS : Attaques DDoS (~226K connexions)
- Friday PortScan : Scans de ports (~286K connexions)

🔗 Source : <https://www.unb.ca/cic/datasets/ids-2017.html>

🛠 Technologies

- **Python 3.x**
- **pandas** - Manipulation de données
- **scikit-learn** - Machine Learning
- **matplotlib/seaborn** - Visualisations

📈 Résultats

| Modèle | Accuracy | Precision | Recall | F1-Score |

-----	-----	-----	-----	-----
Decision Tree	XX.XX%	X.XXX	X.XXX	X.XXX
Random Forest	XX.XX%	X.XXX	X.XXX	X.XXX

🚀 Utilisation

```bash

# 1. Cloner le repo

```
git clone https://github.com/tonpseudo/IDS-ML-Project.git
```

# 2. Installer dépendances

```
pip install -r requirements.txt
```

# 3. Lancer l'analyse

```
python src/01_explore_data.py
```

```
python src/02_preprocess_data.py
```

```
python src/03_select_features.py
python src/04_train_model.py
python src/05_evaluate_model.py
python src/06_visualizations.py
...
...
```

## ## 📁 Structure du Projet

IDS-ML-Project/

```
|── data/ # Datasets
|── src/ # Code source
|── results/ # Graphiques et résultats
|── docs/ # Documentation
└── README.md
```

## 💬 Auteur

\*\*Ton Nom\*\* - Étudiant en Réseaux et Sécurité

## 📄 Licence

Ce projet est à des fins éducatives.

## 8.2 Créer requirements.txt

```
txt

pandas==2.0.3
numpy==1.24.3
matplotlib==3.7.2
seaborn==0.12.2
scikit-learn==1.3.0
```

## 8.3 Initialiser Git

bash

```
cd ~/Desktop/IDS-ML-Project
```

```
Initialiser Git
```

```
git init
```

```
Créer .gitignore
```

```
echo "data/*.csv"
```

```
*.pkl
```

```
__pycache__ /
```

```
.vscode/" > .gitignore
```

```
Premier commit
```

```
git add .
```

```
git commit -m "Initial commit: IDS ML Project"
```

## 8.4 Mettre sur GitHub

1. Va sur <https://github.com>
2. Crée un nouveau repo : "IDS-ML-Project"
3. Dans le terminal :

```
bash
```

```
git remote add origin https://github.com/TonPseudo/IDS-ML-Project.git
```

```
git branch -M main
```

```
git push -u origin main
```

 **Checkpoint :** Projet sur GitHub !

---

## PROJET TERMINÉ !

### Récapitulatif

Tu as créé un projet complet avec :

- Dataset CICIDS2017 (1M+ connexions)
- Prétraitement et nettoyage
- Sélection de features
- 2 modèles ML entraînés
- Évaluation complète
- Visualisations professionnelles
- Documentation
- GitHub

## Prochaines Étapes

1. Ajoute plus de types d'attaques
2. Teste d'autres algorithmes (SVM, Neural Networks)
3. Optimise les hyperparamètres
4. Crée un dashboard interactif
5. Déploie le modèle en production

## Pour Ton CV

 Système de Détection d'Intrusions avec Machine Learning  
- Détection de DDoS et scans de ports avec 95%+ d'accuracy  
- Dataset CICIDS2017 (1M+ connexions réseau)  
- Random Forest et Decision Tree optimisés  
- GitHub: [github.com/TonPseudo/IDS-ML-Project](https://github.com/TonPseudo/IDS-ML-Project)

---

**BRAVO ! Tu as un projet ML complet en Cybersécurité ! 🎉**