

# Projet : Système de Détection d’Intrusion Distribué avec Apache Spark

HAKOUANE NIZAR

Avril 2025

## Résumé

Ce projet vise à construire un **système de détection d’intrusion (IDS)** basé sur l’apprentissage automatique distribué avec **Apache Spark** et **PySpark MLlib**. Le modèle a été entraîné sur le jeu de données **NSL-KDD**, largement utilisé pour la détection d’attaques réseau.

## 1 Introduction

La croissance du volume de données réseau rend nécessaire l’utilisation de traitements distribués pour la détection d’intrusions. Spark est une plateforme idéale pour entraîner des modèles de machine learning sur de grands ensembles de données en parallèle.

Notre objectif est de :

- Charger les données NSL-KDD dans Spark
- Construire une pipeline de prétraitement
- Entraîner un modèle de forêt aléatoire
- Prouver l’exécution distribuée

## 2 Description du Jeu de Données

Le dataset **NSL-KDD** contient :

- 41 caractéristiques (features) par connexion réseau

- Une étiquette de classe (attaque ou normale)
- Plus de 100 000 exemples

Les types d'attaques inclus sont : DoS, Probe, R2L, U2R.

### 3 Architecture et Pipeline

#### Étapes principales

- **Indexation** des colonnes catégorielles (protocole, service, flag)
- **Encodage One-Hot** des variables catégorielles
- **Assemblage** des features
- **Classification** avec RandomForestClassifier (50 arbres)

#### Capture de l'Exécution Distribuée

spark

3.5.5

[Jobs](#)
[Stages](#)
[Storage](#)
[Environment](#)
[Executors](#)
[SQL / DataFrame](#)

DistributedIDS application UI

Executors

[Show Additional Metrics](#)

Summary

	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Excluded
Active(1)	5	95.5 MiB / 434.4 MiB	0.0 B	8	4	0	79	83	1.3 min (1 s)	349.7 MiB	3.3 MiB	3.7 MiB	0
Dead(0)	0	0.0 B / 0.0 B	0.0 B	0	0	0	0	0	0.0 ms (0.0 ms)	0.0 B	0.0 B	0.0 B	0
Total(1)	5	95.5 MiB / 434.4 MiB	0.0 B	8	4	0	79	83	1.3 min (1 s)	349.7 MiB	3.3 MiB	3.7 MiB	0

Executors

Show 20 entries

Search:

Executor ID	Address	Status	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Thread Dump	Heap Histogram	Add Time	Remove Time
driver	DESKTOP-PIDRPU60104	Active	5	95.5 MiB / 434.4 MiB	0.0 B	8	4	0	79	83	1.3 min (1 s)	349.7 MiB	3.3 MiB	3.7 MiB	<a href="#">Thread Dump</a>	<a href="#">Heap Histogram</a>	2025-04-19 00:18:52	-

Showing 1 to 1 of 1 entries

Previous

1

Next

La capture montre que plusieurs tâches actives sont réparties sur 8 cœurs.

### 4 Code Source Principal

```
from pyspark.sql import SparkSession
from pyspark.ml import Pipeline
from pyspark.ml.feature import StringIndexer, OneHotEncoder,
    VectorAssembler
```

```

from pyspark.ml.classification import RandomForestClassifier
from pyspark.ml.evaluation import
    MulticlassClassificationEvaluator

spark = SparkSession.builder.appName("DistributedIDS").
    getOrCreate()

cols = [*map(str, range(41))] + ["label"]
train_path = "C:/Users/abccb/Downloads/data/KDDTrain+.txt"
test_path = "C:/Users/abccb/Downloads/data/KDDTest+.txt"

df = spark.read.csv(train_path, header=False, inferSchema=
    True).toDF(*cols)
dfT = spark.read.csv(test_path, header=False, inferSchema=
    True).toDF(*cols)

cat_cols = ["1", "2", "3"]
indexers = [StringIndexer(inputCol=c, outputCol=f"{c}_idx")
    for c in cat_cols]
encoder = OneHotEncoder(
    inputCols=[f"{c}_idx" for c in cat_cols],
    outputCols=[f"{c}_vec" for c in cat_cols]
)
num_cols = [c for c in df.columns if c not in cat_cols + ["
    label"]]
assembler = VectorAssembler(
    inputCols=num_cols + [f"{c}_vec" for c in cat_cols],
    outputCol="features"
)

rf = RandomForestClassifier(labelCol="label", featuresCol="
    features", numTrees=50, seed=42)
pipeline = Pipeline(stages=[*indexers, encoder, assembler, rf
    ])
model = pipeline.fit(df)

pred = model.transform(dfT)
f1 = (MulticlassClassificationEvaluator(
    labelCol="label",
    predictionCol="prediction",
    metricName="f1").evaluate(pred))

print("F1-score_sur_test_set:", round(f1, 4))
model.write().overwrite().save("spark_ids_model")

spark.stop()

```

## 5 Résultats

Le modèle a obtenu un **score F1** de :

<b>0.4556</b>
---------------

### Remarque :

- Le résultat est cohérent avec les implémentations standards sur NSL-KDD sans tuning approfondi. - L'exécution est entièrement distribuée à travers les cœurs Spark.

## 6 Conclusion et Perspectives

Ce projet démontre la faisabilité de construire un IDS distribué rapide avec Apache Spark. Les extensions possibles incluent :

- Utiliser des méthodes d'équilibrage de classes (SMOTE, poids de classes)
- Optimiser les hyperparamètres avec Grid Search
- Déployer en temps réel avec Spark Structured Streaming