

# SSH Shell Attacks

ANDREA BOTTICELLA\*, Politecnico di Torino, Italy  
ELIA INNOCENTI\*, Politecnico di Torino, Italy  
RENATO MIGNONE\*, Politecnico di Torino, Italy  
SIMONE ROMANO\*, Politecnico di Torino, Italy

This paper introduces a comprehensive machine learning framework to analyze SSH shell attack sessions, leveraging both supervised and unsupervised learning techniques. Using a dataset of 230,000 unique Unix shell attack sessions, the framework aims to classify attacker tactics based on the MITRE ATT&CK framework and uncover latent patterns through clustering. The key contributions of this work are:

- Development of a robust pre-processing pipeline to analyze temporal trends, extract numerical features, and evaluate intent distributions from large-scale SSH attack session data.
- Implementation of supervised classification models to accurately predict multiple attacker tactics, supported by hyperparameter tuning and feature engineering for enhanced performance.
- Application of unsupervised clustering techniques to uncover hidden patterns in attack behaviors, leveraging visualization tools and cluster analysis for fine-grained categorization.
- Exploration of advanced language models, such as BERT and Doc2Vec, for representation learning and fine-tuning to improve intent classification and session interpretation.

CCS Concepts: • **Security and privacy** → **Intrusion detection systems; Malware and its mitigation;** • **Computing methodologies** → *Supervised learning by classification; Unsupervised learning; Natural language processing.*

Additional Key Words and Phrases: SSH shell attacks, machine learning, supervised learning, unsupervised learning, language models, security logs, intrusion detection.

## CONTENTS

Abstract	1
Contents	1
1 INTRODUCTION	1
2 BACKGROUND	2
3 DATA EXPLORATION AND PRE-PROCESSING	2
4 SUPERVISED LEARNING - CLUSTERING	3
5 UNSUPERVISED LEARNING - CLUSTERING	3
6 LANGUAGE MODEL EXPLORATION	4
7 CONCLUSION	4

## 1 INTRODUCTION

This section introduces the topic of the project, provides background information, and outlines the objectives.

### 1.1 Motivation

Provide an explanation of why this topic is important and relevant.

\*All authors contributed equally to this research.

Authors' Contact Information: Andrea Botticella, andrea.botticella@studenti.polito.it, Politecnico di Torino, Turin, Italy; Elia Innocenti, elia.innocenti@studenti.polito.it, Politecnico di Torino, Turin, Italy; Renato Mignone, renato.mignone@studenti.polito.it, Politecnico di Torino, Turin, Italy; Simone Romano, simone.romano@studenti.polito.it, Politecnico di Torino, Turin, Italy.

## 1.2 Objective

Clearly state the objectives and what the project aims to accomplish.

## 2 BACKGROUND

This section ...

### 2.1 Subsection 1

...

### 2.2 Subsection 2

...

## 3 DATA EXPLORATION AND PRE-PROCESSING

This section ...

### 3.1 Introduction

Brief introduction to the data exploration and pre-processing tasks.

...

### 3.2 Dataset Preparation

Loading the dataset and initial inspection.

...

### 3.3 Temporal Analysis

Analysis of when the attacks were performed.

...

### 3.4 Feature Extraction

Extracting features from the attack sessions.

...

### 3.5 Common Words Analysis

Identifying the most common words in the sessions.

...

### 3.6 Intent Distribution

Analyzing the distribution of intents over time.

...

### 3.7 Text Representation

Converting text into numerical representations (BoW, TF-IDF).

...

## 4 SUPERVISED LEARNING - CLUSTERING

This section ...

### 4.1 Introduction

Overview of the supervised learning task and its objectives.

...

### 4.2 Data Splitting

Splitting the dataset into training and test sets.

...

### 4.3 Baseline Model Implementation

Implementing and evaluating baseline models.

...

### 4.4 Hyperparameter Tuning

Tuning hyperparameters and evaluating performance.

...

### 4.5 Result Analysis

Analyzing the results for each intent.

...

### 4.6 Feature Experimentation

Exploring different feature combinations and their impact on performance.

...

## 5 UNSUPERVISED LEARNING - CLUSTERING

This section ...

### 5.1 Introduction

Overview of the clustering task and its objectives.

...

### 5.2 Determine the Number of Clusters

Using methods like the elbow method or silhouette analysis.

...

### 5.3 Hyperparameter Tuning

Tuning other hyperparameters, if any.

...

### 5.4 Cluster Visualization

Visualizing the clusters through t-SNE.

...

## 5.5 Cluster Analysis

Analyzing the characteristics of each cluster.

...

## 5.6 Intent Homogeneity

Assessing if clusters reflect intent division.

...

## 5.7 Specific Attack Categories

Associating clusters with specific attack categories.

...

# 6 LANGUAGE MODEL EXPLORATION

This section ...

## 6.1 Introduction

Overview of the language models task and its objectives.

...

## 6.2 Pretraining

Pretraining Doc2Vec or using a pretrained Bert model.

...

## 6.3 Model Fine-tuning

Fine-tuning the last layer of the network.

...

## 6.4 Learning Curves

Plotting learning curves and determining the optimal number of epochs.

...

# 7 CONCLUSION

This section ...

## 7.1 Subsection 1

...

## 7.2 Subsection 2

...