

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: • **Computing methodologies** → **Supervised learning by classification**; **Unsupervised learning**; *Natural language processing*; *Machine learning*; *Machine learning approaches*; • **Security and privacy** → Intrusion detection systems.

Additional Key Words and Phrases: Machine learning, supervised learning, unsupervised learning, language models, text classification, clustering, intent classification, SSH shell attacks, security log analysis

CONTENTS

Abstract	1
Contents	1
1 INTRODUCTION	1
2 BACKGROUND	2
3 DATA EXPLORATION AND PRE-PROCESSING	3
4 SUPERVISED LEARNING - CLUSTERING	6
5 UNSUPERVISED LEARNING - CLUSTERING	9
6 LANGUAGE MODEL EXPLORATION	13
7 CONCLUSION	16

1 INTRODUCTION

1.1 Motivation

Security logs play a crucial role in understanding and mitigating cyber attacks, particularly in the domain of network and system security. With the increasing sophistication of cyber threats, analyzing and interpreting

*The authors collaborated closely in developing this project.

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.

security logs has become paramount for detecting, preventing, and responding to potential security breaches. Unix shell attacks, especially those executed through SSH, represent a significant vector for potential system compromises.

The complexity of security log analysis stems from several key challenges:

- Logs are often unstructured and contain ambiguous or malformed text.
- Manual parsing and interpretation of logs is time-consuming and error-prone.
- The sheer volume of log data makes comprehensive manual review impractical.

These challenges underscore the need for automated, intelligent approaches to log analysis that can efficiently extract meaningful insights and identify potential security threats.

1.2 Objective

The primary objective of this research is to develop and evaluate machine learning techniques for automatic analysis and classification of SSH shell attack logs. Specifically, we aim to:

- Explore and preprocess a large dataset of Unix shell attack sessions.
- Apply supervised learning techniques to classify attack tactics based on session characteristics.
- Utilize unsupervised learning methods to discover patterns and clusters in attack sessions.
- Investigate the potential of advanced language models in understanding and categorizing attack intents.

By leveraging machine learning approaches, we seek to:

- Automate the process of log analysis and intent classification.
- Provide security professionals with insights into attack strategies.
- Develop a framework for understanding and categorizing SSH shell attacks using the MITRE ATT&CK tactics as a reference.

The significance of this research lies in its potential to enhance cybersecurity threat detection and response capabilities by transforming complex, unstructured log data into actionable intelligence.

2 BACKGROUND

2.1 Security Logs and Attack Analysis

Security logs represent a critical source of information for understanding system vulnerabilities and potential cyber attacks. These logs capture detailed records of system events, network interactions, and user activities, providing valuable insights into potential security breaches.

In the context of SSH shell attacks, logs document the sequence of commands executed during a malicious session, enabling security researchers to analyze attacker behaviors, techniques, and potential system impacts. However, the manual analysis of these logs is challenging due to their volume, complexity, and often non-standard formatting.

2.2 MITRE ATT&CK Framework

The MITRE ATT&CK (Adversarial Tactics, Techniques, and Common Knowledge) framework provides a comprehensive knowledge base of adversary tactics and techniques observed in real-world cyber attacks. This framework serves as a standardized methodology for understanding and categorizing attack strategies.

For our research, we focus on seven key intents derived from the MITRE ATT&CK framework:

- **Persistence:** Techniques used by adversaries to maintain system access across restarts or credential changes.
- **Discovery:** Methods for gathering information about the target system and network environment.
- **Defense Evasion:** Strategies to avoid detection by security mechanisms.

- **Execution:** Techniques for running malicious code on target systems.
- **Impact:** Actions aimed at manipulating, interrupting, or destroying systems and data.
- **Other:** Less common tactics including Reconnaissance, Resource Development, Initial Access, etc.
- **Harmless:** Non-malicious code or actions.

2.3 Dataset Overview

Our research utilizes a comprehensive dataset of SSH shell attacks collected from a honeypot deployment. Key characteristics of the dataset include:

- Approximately 230,000 unique Unix shell attack sessions
- Recorded after SSH login
- Stored in Parquet format for efficient data processing
- Columns include:
 - Session ID
 - Full session text (command sequences)
 - Timestamp
 - Intent labels based on MITRE ATT&CK tactics

It is important to note that the dataset's labels are automatically generated through a research project and may contain potential classification errors. This inherent uncertainty adds an additional layer of complexity to our analysis and highlights the importance of robust machine learning techniques.

2.4 Research Approach

Our research employs a multi-faceted approach to SSH shell attack log analysis:

- Comprehensive data exploration and preprocessing
- Supervised learning for attack intent classification
- Unsupervised learning for attack pattern discovery
- Advanced language model exploration

By combining these techniques, we aim to develop a comprehensive framework for understanding and categorizing SSH shell attacks, ultimately contributing to improved cybersecurity threat detection and response strategies.

3 DATA EXPLORATION AND PRE-PROCESSING

3.1 Introduction

This section details the data exploration and pre-processing steps undertaken to analyze SSH shell attack logs. The tasks include dataset preparation, temporal analysis, feature extraction, common words analysis, intent distribution, and text representation.

3.2 Dataset Preparation

The dataset used in this research is loaded from a Parquet file (`ssh_attacks.parquet`) into a Pandas DataFrame. The initial inspection involves checking the dataset's structure, identifying missing values, and detecting duplicate rows. The dataset contains columns such as `Session ID`, `Full session text`, `Timestamp`, and `Intent labels`.

Listing 1. Load and inspect the dataset

```
# Load the dataset
SSH_Attacks = pd.read_parquet("../data/processed/ssh_attacks_decoded.parquet")
```

```
# Inspect the dataset structure
print(SSH_Attacks.info())

# Check for missing values
print(SSH_Attacks.isnull().sum())

# Check for duplicate rows
print(SSH_Attacks.duplicated().sum())
```

The initial inspection revealed that the dataset is well-structured with columns that are essential for our analysis. However, it is important to handle any missing values and duplicates to ensure the integrity of the data. The following steps were taken to address these issues:

- **Missing Values**: We identified and handled missing values by either imputing them with appropriate values or removing the affected rows.
- **Duplicate Rows**: Duplicate rows were detected and removed to avoid redundancy in the analysis.

Placeholder for Dataset Structure Table

3.3 Temporal Analysis

The temporal analysis examines when the attacks were performed. The `first_timestamp` column is converted to a datetime format to analyze attack frequencies over time, including hourly, daily, and monthly trends.

Listing 2. Convert timestamps and analyze frequencies

```
# Convert first_timestamp to datetime format
SSH_Attacks['first_timestamp'] = pd.to_datetime(SSH_Attacks['first_timestamp'])

# Analyze attack frequencies over time
temporal_series = (
    SSH_Attacks.groupby(SSH_Attacks['first_timestamp'].dt.date)
    .size()
    .reset_index(name='attack_count')
)
```

The analysis includes plotting the number of attacks per hour, month, and year to identify patterns and trends. This helps in understanding the temporal distribution of attacks and identifying any periodic patterns or anomalies.

Placeholder for Temporal Analysis Plot

The temporal analysis revealed that the frequency of attacks varies significantly over time. By examining the hourly, daily, and monthly trends, we can gain insights into the behavior of attackers and the times when systems are most vulnerable.

3.4 Feature Extraction

Feature extraction involves identifying and extracting relevant features from the attack sessions. This includes analyzing the distribution of classes (intents) and visualizing the data using bar plots.

Listing 3. Extract and visualize class distribution

```
# Extract and count occurrences of each class
all_classes = SSH_Attacks['Set_Fingerprint'].explode().str.strip()
class_counts = all_classes.value_counts()

# Plot the distribution of classes
sns.barplot(x=class_counts.index, y=class_counts.values, palette='viridis')
```

The distribution of classes provides valuable information about the types of attacks and their prevalence. By visualizing this distribution, we can identify the most common attack intents and focus our analysis on these areas.

Placeholder for Class Distribution Bar Plot

The feature extraction process also involves creating new features that can enhance the analysis. For example, we can extract the length of each session, the number of commands executed, and other relevant metrics.

3.5 Common Words Analysis

The common words analysis identifies the most frequent words used in the attack sessions. This is achieved using word clouds and other text analysis techniques.

Listing 4. Generate a word cloud from session text

```
# Generate a word cloud for the session text
wordcloud = WordCloud(width=800, height=400, background_color='white').generate(
    ' '.join(SSH_Attacks['Full_session_text']))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis('off')
plt.show()
```

The word cloud visualization highlights the most common words used in the attack sessions, providing insights into the attackers' behavior and strategies. This can help in identifying common commands and patterns used in the attacks.

Placeholder for Word Cloud Visualization

Additionally, we can create bar plots to show the frequency of the top 10 most common words, which can further aid in understanding the textual characteristics of the attack sessions.

Placeholder for Common Words Bar Plot

3.6 Intent Distribution

The intent distribution analysis examines the distribution of intents over time. This involves grouping the data by date and intent to count occurrences and visualize the trends.

Listing 5. Group attacks by fingerprint and date

```
# Group by Set_Fingerprint and date to count occurrences
grouped_SSH_Attacks = (
    SSH_Attacks.explode('Set_Fingerprint')
    .groupby([SSH_Attacks['first_timestamp'].dt.date, 'Set_Fingerprint'])
    .size()
    .reset_index(name='attack_count')
)
```

By analyzing the distribution of intents over time, we can identify trends and patterns in the attackers' behavior. This can help in understanding how different types of attacks evolve and vary over time.

Placeholder for Intent Distribution Plot

The intent distribution analysis also helps in identifying any seasonal or periodic patterns in the attacks, which can be crucial for developing effective defense strategies.

3.7 Text Representation

Text representation converts the session text into numerical representations using techniques such as Bag of Words (BoW) and Term Frequency-Inverse Document Frequency (TF-IDF). These representations are used for further analysis and machine learning tasks.

Listing 6. Convert text into numerical representations

```
# Convert text into numerical representations using Bag of Words (BoW)
from sklearn.feature_extraction.text import CountVectorizer
bow_vectorizer = CountVectorizer()
X_bow = bow_vectorizer.fit_transform(SSH_Attacks['Full_session_text'])

# Convert text into numerical representations using TF-IDF
from sklearn.feature_extraction.text import TfidfVectorizer
tfidf_vectorizer = TfidfVectorizer()
X_tfidf = tfidf_vectorizer.fit_transform(SSH_Attacks['Full_session_text'])
```

The resulting numerical representations from both BoW and TF-IDF are normalized and used for subsequent analysis and modeling. These representations are essential for applying machine learning algorithms to classify and predict attack intents.

Placeholder for Text Representation Table

The text representation techniques help in transforming the unstructured session text into structured numerical data, which can be used for various analytical and predictive tasks. By comparing the performance of different representation techniques, we can select the most effective method for our analysis.

4 SUPERVISED LEARNING - CLUSTERING

4.1 Introduction

This section provides an overview of the supervised learning task and its objectives. The goal is to classify attack session tactics based on the provided dataset. We will implement and evaluate various machine learning models to determine the most effective approach for this classification task.

4.2 Data Splitting

The first step in the supervised learning process is to split the dataset into training and test sets. This ensures that we can evaluate the performance of our models on unseen data.

Data Loading: The dataset is loaded from a Parquet file into a Pandas DataFrame.

Listing 7. Load the dataset

```
# Load the dataset
SSH_Attacks = pd.read_parquet("../data/processed/ssh_attacks_decoded.parquet")
```

Data Splitting: We split the dataset into training and test sets, ensuring a 70/30 split while maintaining reproducibility.

Listing 8. Split the dataset into training and test sets

```
# Split the dataset into training and test sets
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, random_state=42
)
```

Placeholder for Data Splitting Summary Table

The data splitting process ensures that the training set is used to train the models, while the test set is used to evaluate their performance.

4.3 Baseline Model Implementation

In this subsection, we implement and evaluate baseline models to establish a performance benchmark. We will use Logistic Regression, Random Forest, and Support Vector Machine (SVM) as our baseline models.

Logistic Regression:

Listing 9. Train Logistic Regression model

```
# Initialize and train Logistic Regression model
model = LogisticRegression(max_iter=1000, random_state=42)
model.fit(X_train_tfidf, y_train_binary)
```

Random Forest:

Listing 10. Train Random Forest model

```
# Initialize and train Random Forest model
model = RandomForestClassifier(n_estimators=100, random_state=42)
model.fit(X_train_tfidf, y_train_binary)
```

Support Vector Machine (SVM):

Listing 11. Train SVM model

```
# Initialize and train SVM model
model = SVC(kernel='linear', random_state=42)
model.fit(X_train_tfidf, y_train_binary)
```

Placeholder for Baseline Model Performance Table

The baseline model implementation provides a reference point for evaluating the performance of more advanced models.

4.4 Hyperparameter Tuning

Hyperparameter tuning involves optimizing the parameters of the models to improve their performance. We use GridSearchCV to perform an exhaustive search over specified parameter values.

Logistic Regression Hyperparameter Tuning:

Listing 12. Parameter grid for Logistic Regression

```
# Define parameter grid for Logistic Regression
param_grid = {'C': [0.1, 1, 10, 100]}
grid_search = GridSearchCV(LogisticRegression(max_iter=1000, random_state=42),
    param_grid, cv=5)
grid_search.fit(X_train_tfidf, y_train_binary)
```

Random Forest Hyperparameter Tuning:

Listing 13. Parameter grid for Random Forest

```
# Define parameter grid for Random Forest
param_grid = {'n_estimators': [50, 100, 200]}
grid_search = GridSearchCV(RandomForestClassifier(random_state=42), param_grid,
    cv=5)
grid_search.fit(X_train_tfidf, y_train_binary)
```

Placeholder for Hyperparameter Tuning Results Table

Hyperparameter tuning helps in finding the best parameters for each model, thereby improving their performance.

4.5 Result Analysis

In this subsection, we analyze the results of the models for each intent. We use metrics such as accuracy, precision, recall, and F1-score to evaluate the performance.

Classification Report:

Listing 14. Generate classification report

```
# Generate classification report
report = classification_report(y_test_binary, y_pred, zero_division=0)
print(report)
```

Confusion Matrix:

Listing 15. Generate confusion matrix

```
# Generate confusion matrix
cm = confusion_matrix(y_test_binary, y_pred)
sns.heatmap(cm, annot=True, fmt='d', cmap='coolwarm')
plt.show()
```

Placeholder for Classification Report and Confusion Matrix Plots

The result analysis provides insights into the performance of the models and helps in identifying areas for improvement.

4.6 Feature Experimentation

Feature experimentation involves exploring different feature combinations and their impact on model performance. We experiment with various text representation techniques such as Bag of Words (BoW) and Term Frequency-Inverse Document Frequency (TF-IDF).

Bag of Words (BoW):

Listing 16. Convert text using Bag of Words (BoW)

```
# Convert text into numerical representations using Bag of Words (BoW)
bow_vectorizer = CountVectorizer()
X_train_bow = bow_vectorizer.fit_transform(X_train)
X_test_bow = bow_vectorizer.transform(X_test)
```

TF-IDF:

Listing 17. Convert text using TF-IDF

```
# Convert text into numerical representations using TF-IDF
tfidf_vectorizer = TfidfVectorizer()
X_train_tfidf = tfidf_vectorizer.fit_transform(X_train)
X_test_tfidf = tfidf_vectorizer.transform(X_test)
```

Placeholder for Feature Experimentation Results Table

By experimenting with different features, we can identify the most effective representation techniques for our classification task.

5 UNSUPERVISED LEARNING - CLUSTERING

5.1 Introduction

This section provides an overview of the clustering task and its objectives. The goal is to group similar attack sessions into clusters based on their characteristics. By identifying clusters, we can gain insights into common patterns and behaviors in the attack data. We will use various clustering techniques and evaluate their effectiveness.

5.2 Determine the Number of Clusters

Determining the optimal number of clusters is a crucial step in the clustering process. We use methods like the elbow method and silhouette analysis to identify the appropriate number of clusters.

Elbow Method: The elbow method involves plotting the sum of squared distances (inertia) for a range of cluster numbers and identifying the point where the inertia starts to decrease more slowly (the "elbow").

Listing 18. Elbow Method for k-Means Clustering

```
# Elbow Method
n_cluster_list = []
inertia_list = []
for n_clusters in range(3, 17):
    kmeans = KMeans(n_clusters=n_clusters, n_init=10, random_state=42)
    kmeans.fit(X)
    inertia_list.append(kmeans.inertia_)
    n_cluster_list.append(n_clusters)
```

```
# Plot Elbow Method
plt.figure(figsize=(5, 3.5))
plt.plot(n_cluster_list, inertia_list, marker='o', markersize=5, color='blue')
plt.xlabel('Number_of_clusters')
plt.ylabel('k-Means_clustering_error')
plt.title('Elbow_Method')
plt.show()
```

Placeholder for Elbow Method Plot

Silhouette Analysis: Silhouette analysis involves calculating the silhouette score for each potential cluster number. The silhouette score measures how similar an object is to its own cluster compared to other clusters.

Listing 19. Silhouette Analysis for k-Means Clustering

```
# Silhouette Analysis
silhouette_list = []
for n_clusters in range(3, 17):
    kmeans = KMeans(n_clusters=n_clusters, n_init=10, random_state=42)
    labels = kmeans.fit_predict(X)
    silhouette_score_value = silhouette_score(X, labels)
    silhouette_list.append(silhouette_score_value)

# Plot Silhouette Analysis
plt.figure(figsize=(5, 3.5))
plt.plot(n_cluster_list, silhouette_list, marker='o', markersize=5, color='blue')
plt.xlabel('Number_of_clusters')
plt.ylabel('Silhouette_Score')
plt.title('Silhouette_Analysis')
plt.show()
```

Placeholder for Silhouette Analysis Plot

Based on the elbow method and silhouette analysis, we select the optimal number of clusters for further analysis.

5.3 Hyperparameter Tuning

Hyperparameter tuning involves optimizing the parameters of the clustering algorithms to improve their performance. We use GridSearchCV to perform an exhaustive search over specified parameter values.

K-Means Hyperparameter Tuning:

Listing 20. Grid Search for k-Means Clustering

```
# Define parameter grid for K-Means
param_grid_kmeans = {
    'init': ['k-means++', 'random'],
    'n_init': list(range(10, 21, 2)),
    'max_iter': list(range(50, 200, 50)),
}
```

```

# Create KMeans object
kmeans = KMeans(n_clusters=10, random_state=42)

# Create GridSearchCV object
grid_search_kmeans = GridSearchCV(kmeans, param_grid=param_grid_kmeans, cv=5)

# Fit the grid search to the data
grid_search_kmeans.fit(X)

# Get the best parameters
best_params_kmeans = grid_search_kmeans.best_params_
print("Best parameters:", best_params_kmeans)

```

Placeholder for K-Means Hyperparameter Tuning Results

Gaussian Mixture Model (GMM) Hyperparameter Tuning:

Listing 21. Grid Search for Gaussian Mixture Model (GMM)

```

# Define parameter grid for GMM
param_grid_gmm = {
    'init_params': ['kmeans'],
    'covariance_type': ['full', 'spherical'],
    'tol': [1e-3, 1e-4, 1e-5],
    'max_iter': list(range(50, 300, 50)),
}

# Create GaussianMixture object
gmm = GaussianMixture(n_components=10, random_state=42)

# Create GridSearchCV object
grid_search_gmm = GridSearchCV(gmm, param_grid=param_grid_gmm, cv=5, scoring=
    silhouette_scorer)

# Fit the grid search to the data
grid_search_gmm.fit(X)

# Get the best parameters
best_params_gmm = grid_search_gmm.best_params_
print("Best parameters:", best_params_gmm)

```

Placeholder for GMM Hyperparameter Tuning Results

5.4 Cluster Visualization

Visualizing the clusters helps in understanding the distribution and characteristics of the clusters. We use t-SNE to reduce the dimensionality of the data and create clear visual representations of the clusters.

t-SNE Visualization:

Listing 22. t-SNE Visualization of Clusters

```

# Apply t-SNE to reduce the number of components
tsne = TSNE(n_components=2, random_state=42).fit_transform(X)
df_tsne = pd.DataFrame(tsne, columns=['x1', 'x2'])

# K-Means Clusters
df_tsne['cluster_kmeans'] = kmeans_tuned.labels_
sns.scatterplot(data=df_tsne, x='x1', y='x2', hue='cluster_kmeans', palette='
    viridis')
plt.title('t-SNE_Visualization_of_K-Means_Clusters')
plt.show()

# GMM Clusters
df_tsne['cluster_gmm'] = gmm_tuned.predict(X)
sns.scatterplot(data=df_tsne, x='x1', y='x2', hue='cluster_gmm', palette='
    viridis')
plt.title('t-SNE_Visualization_of_GMM_Clusters')
plt.show()

```

Placeholder for t-SNE Visualization of K-Means Clusters

Placeholder for t-SNE Visualization of GMM Clusters

5.5 Cluster Analysis

Analyzing the characteristics of each cluster helps in understanding the common patterns and behaviors within the clusters. We examine the distribution of features and intents within each cluster.

Cluster Feature Analysis:

Listing 23. Feature Distribution Analysis by Cluster

```

# Analyze the distribution of features within each cluster
for cluster in range(10):
    cluster_data = df_tsne[df_tsne['cluster_kmeans'] == cluster]
    print(f"Cluster_{cluster}_Feature_Distribution:")
    print(cluster_data.describe())

```

Placeholder for Cluster Feature Distribution Tables

5.6 Intent Homogeneity

Assessing if clusters reflect intent division involves examining the homogeneity of intents within each cluster. We calculate the proportion of each intent within the clusters to determine if the clusters are homogeneous.

Intent Homogeneity Analysis:

Listing 24. Intent Proportions Analysis by Cluster

```

# Calculate the proportion of each intent within the clusters
for cluster in range(10):
    cluster_data = df_tsne[df_tsne['cluster_kmeans'] == cluster]
    intent_proportions = cluster_data['intent'].value_counts(normalize=True)
    print(f"Cluster_{cluster}_Intent_Proportions:")
    print(intent_proportions)

```

Placeholder for Intent Proportion Tables

5.7 Specific Attack Categories

Associating clusters with specific attack categories involves identifying the common attack patterns within each cluster. We analyze the most frequent attack categories within the clusters to understand their characteristics.

Attack Category Analysis:

Listing 25. Attack Categories Analysis by Cluster

```
# Analyze the most frequent attack categories within the clusters
for cluster in range(10):
    cluster_data = df_tsne[df_tsne['cluster_kmeans'] == cluster]
    attack_categories = cluster_data['attack_category'].value_counts()
    print(f"Cluster_{cluster}_Attack_Categories:")
    print(attack_categories)
```

Placeholder for Attack Category Distribution Tables

6 LANGUAGE MODEL EXPLORATION

6.1 Introduction

This section provides an overview of the language models task and its objectives. The goal is to leverage advanced language models to classify attack session tactics based on the provided dataset. We will explore the use of pretrained models such as BERT and Doc2Vec, fine-tune them for our specific task, and analyze their performance.

Language models have revolutionized natural language processing (NLP) by enabling transfer learning, where models pretrained on large datasets can be fine-tuned on specific tasks. This approach allows us to benefit from the knowledge encoded in these models and achieve better performance with less training data.

6.2 Pretraining

Pretraining involves using a pretrained language model or training a model from scratch on a large corpus of text. For this task, we will use a pretrained BERT model from HuggingFace's Transformers library.

Installing Dependencies:

Listing 26. Install required packages

```
!pip install transformers torch
```

Loading the Dataset:

Listing 27. Load dataset and print its size

```
import pandas as pd

# Load the dataset
df = pd.read_parquet("../data/processed/ssh_attacks_sampled_decoded.parquet")
print(f"Dataset_size:_{df.shape[0]}_rows")
```

Preprocessing:

Listing 28. Preprocess 'Set_Fingerprint' column

```

from sklearn.preprocessing import MultiLabelBinarizer

# Preprocess Set_Fingerprint column
df['Set_Fingerprint'] = df['Set_Fingerprint'].apply(lambda x: [intent.strip()
    for intent in x.split(',')])
mlb = MultiLabelBinarizer()
y = mlb.fit_transform(df['Set_Fingerprint'])
print(f"Classes_identified:_{mlb.classes_}")

```

Placeholder for Dataset Summary Table

6.3 Model Fine-tuning

Fine-tuning involves training the last layer of the pretrained model on our specific dataset. We will use BERT for sequence classification and fine-tune it on the SSH attack sessions.

Tokenization:

Listing 29. Tokenize text data using BERT tokenizer

```

from transformers import BertTokenizer

# Tokenize the text data
tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
train_encodings = tokenizer(list(train_texts.fillna("")).astype(str), truncation=
    =True, padding=True, max_length=128)
val_encodings = tokenizer(list(val_texts.fillna("")).astype(str), truncation=
    True, padding=True, max_length=128)

```

Model Initialization:

Listing 30. Initialize BERT model for sequence classification

```

from transformers import BertForSequenceClassification, AdamW

# Initialize the BERT model for sequence classification
model = BertForSequenceClassification.from_pretrained('bert-base-uncased',
    num_labels=y.shape[1])
model.to(device)

# Optimizer and Loss
optimizer = AdamW(model.parameters(), lr=5e-5)
criterion = torch.nn.BCEWithLogitsLoss()

```

Training Loop:

Listing 31. Fine-tune BERT model

```

train_loss_list, val_loss_list = [], []

```

```

for epoch in range(5): # Fine-tune for 5 epochs
    model.train()
    total_loss = 0

    for batch in train_loader:
        optimizer.zero_grad()
        input_ids, attention_mask, labels = (
            batch['input_ids'].to(device),
            batch['attention_mask'].to(device),
            batch['labels'].to(device),
        )
        outputs = model(input_ids=input_ids, attention_mask=attention_mask)
        loss = criterion(outputs.logits, labels)
        loss.backward()
        optimizer.step()
        total_loss += loss.item()

    train_loss_list.append(total_loss / len(train_loader))

# Validation
model.eval()
val_loss = 0
with torch.no_grad():
    for batch in val_loader:
        input_ids, attention_mask, labels = (
            batch['input_ids'].to(device),
            batch['attention_mask'].to(device),
            batch['labels'].to(device),
        )
        outputs = model(input_ids=input_ids, attention_mask=attention_mask)
        loss = criterion(outputs.logits, labels)
        val_loss += loss.item()
    val_loss_list.append(val_loss / len(val_loader))

```

Placeholder for Training and Validation Loss Table

6.4 Learning Curves

Plotting learning curves helps in understanding the model's performance over epochs and determining the optimal number of epochs for training.

Plotting Learning Curves:

Listing 32. Plot learning curves

```

import matplotlib.pyplot as plt

# Plot learning curves
plt.plot(range(1, 6), train_loss_list, label="Training_Loss")
plt.plot(range(1, 6), val_loss_list, label="Validation_Loss")
plt.xlabel("Epochs")

```

```
plt.ylabel("Loss")  
plt.legend()  
plt.show()
```

Placeholder for Learning Curves Plot

By analyzing the learning curves, we can determine the optimal number of epochs to stop training and avoid overfitting. The point where the validation loss stops decreasing or starts increasing indicates the optimal stopping point.

7 CONCLUSION

7.1 Summary of Key Findings

In this project, we explored various techniques for analyzing and classifying SSH shell attack logs. The primary objectives were to preprocess the data, perform exploratory data analysis, implement supervised and unsupervised learning models, and leverage advanced language models for classification tasks. Here, we summarize the key findings from each section of the project.

Data Exploration and Pre-processing: We began by loading and inspecting the dataset, identifying missing values, and handling duplicates. Temporal analysis revealed significant variations in attack frequencies over time, with notable peaks during specific hours and months. Feature extraction and common words analysis provided insights into the most frequent commands and intents used in the attack sessions.

Supervised Learning - Classification: We implemented and evaluated several machine learning models, including Logistic Regression, Random Forest, and Support Vector Machine (SVM). Hyperparameter tuning improved the performance of these models, and the result analysis highlighted the strengths and weaknesses of each approach. Feature experimentation with different text representation techniques, such as Bag of Words (BoW) and TF-IDF, demonstrated the impact of feature selection on model performance.

Unsupervised Learning - Clustering: Clustering techniques, such as K-Means and Gaussian Mixture Models (GMM), were used to group similar attack sessions. The elbow method and silhouette analysis helped determine the optimal number of clusters. Cluster visualization using t-SNE provided a clear representation of the clusters, and cluster analysis revealed common patterns and behaviors within each group.

Language Model Exploration: We explored the use of advanced language models, such as BERT, for classifying attack session tactics. Fine-tuning the pretrained BERT model on our dataset improved classification performance. Learning curves indicated the optimal number of epochs for training, helping to avoid overfitting.

7.2 Challenges Faced

Throughout the project, we encountered several challenges that required careful consideration and problem-solving.

Data Quality and Preprocessing: Handling missing values, duplicates, and inconsistencies in the dataset was a critical step. Ensuring the data was clean and well-prepared for analysis required significant effort. Additionally, the unstructured nature of the session text posed challenges for text representation and feature extraction.

Model Selection and Tuning: Selecting appropriate machine learning models and tuning their hyperparameters was a complex task. Balancing model complexity with performance and avoiding overfitting required iterative experimentation and validation.

Computational Resources: Training advanced language models, such as BERT, required substantial computational resources. Efficiently managing these resources and optimizing the training process was essential to achieve timely results.

Interpretability of Results: Interpreting the results of clustering and classification models, especially in the context of cybersecurity, was challenging. Ensuring that the findings were meaningful and actionable required careful analysis and domain knowledge.

7.3 Future Work

Based on the findings and challenges encountered in this project, we propose several directions for future work.

Enhanced Feature Engineering: Further exploration of feature engineering techniques, such as incorporating domain-specific knowledge and using advanced text representation methods, could improve model performance. Experimenting with additional features, such as network metadata and contextual information, may provide deeper insights into attack patterns.

Advanced Model Architectures: Exploring more advanced model architectures, such as transformer-based models and deep neural networks, could enhance classification accuracy. Transfer learning with other pretrained models and ensemble methods may also yield better results.

Real-time Analysis and Detection: Implementing real-time analysis and detection systems for SSH shell attacks could provide immediate insights and responses to potential threats. Integrating the models developed in this project into a real-time monitoring framework would be a valuable extension.

Broader Dataset and Generalization: Expanding the dataset to include a wider range of attack types and sources would improve the generalizability of the models. Collaborating with other organizations to share data and insights could enhance the robustness and applicability of the findings.

7.4 Conclusion

This project demonstrated the potential of machine learning and advanced language models for analyzing and classifying SSH shell attack logs. By leveraging various techniques, we gained valuable insights into attack patterns and behaviors, which can inform cybersecurity strategies and defenses. Despite the challenges faced, the results highlight the importance of data-driven approaches in enhancing cybersecurity threat detection and response capabilities. Future work in this area holds promise for further advancements and practical applications in the field of cybersecurity.