

Abstract

Computer network system administrators need to inspect and analyze network traffic and detect malicious communications, monitor system performance, and provide operational services. However, identifying threats contained within encrypted network traffic, which has become increasingly prevalent, poses a unique set of challenges. It is imperative to monitor this traffic or threats and malware but do so in a way that maintains privacy. This project aims to develop a machine learning-based system that can accurately detect malware communication in this setting.

Introduction

We were fortunate to enough to have been matched with Dr. Liang Zhou to work on a web development project to inspect and analyze network traffic and detect malicious communications, monitor system performance, and provide operational services. It is imperative to monitor this traffic or threats while maintaining privacy. This project's goal was to develop a machine learning-based system that accurately detect malware communications.

Editing Content w/ Python

```

# Importing datasets
malicious_dataset = pd.read_csv('malicious_flow.csv')
benign_dataset = pd.read_csv('benign_flow.csv')

# Removing duplicated rows from benign_dataset (GDSP row removed)
benign_dataset = benign_dataset[benign_dataset.duplicated(keep=False) == False]

# Combining both datasets together
all_data = pd.concat([malicious_dataset, benign_dataset])

# Reducing the size of the dataset to reduce the amount of time taken in training model
reduced_dataset = all_data.sample(25000)

# Removing non values
df = reduced_dataset.drop(reduced_dataset.columns[reduced_dataset.isnull().any()], axis=1)
reduced_dataset = df

# Dropping information
reduced_dataset.info()

# Class 'random_forest_classifier'
import pandas as pd
from sklearn import svm, linear_model, tree, ensemble

# Splitting datasets into training and test data
x_train, x_test, y_train, y_test = train_test_split(reduced_x, reduced_y, test_size=0.2)

# Training random forest classifier
rf = RandomForestClassifier(n_estimators=100)
rf.fit(x_train, y_train)

# Predicting on the test set
y_pred = rf.predict(x_test)

# Calculating accuracy score
print('Random Forest Classifier Accuracy score: ', accuracy_score(y_test, y_pred))

# Function to plot most important features of random forest model
def plot_feature_importance(importance, names, model_type):
    # Create arrays for feature importance and feature names
    feature_importance = pd.Series(importance)
    feature_names = pd.Series(names)

    # Create a DataFrame using a Dictionary
    df = pd.DataFrame({'Importance': feature_importance, 'Feature': feature_names})

    # Sort the DataFrame in order decreasing feature importance
    df = df.sort_values('Importance', ascending=False, inplace=True)

    # Create size of the plot
    plt.figure(figsize=(10,8))
    plt.bar(df['Feature'], df['Importance'], yerr=df['Importance'])
    plt.xticks(rotation=90)
    plt.xlabel('FEATURE IMPORTANCE')
    plt.ylabel('FEATURE NAMES')

    # Plot Feature Importance
    plt.figure(figsize=(10,8))
    plt.bar(df['Feature'], df['Importance'])
    plt.xticks(rotation=90)
    plt.xlabel('FEATURE IMPORTANCE')
    plt.ylabel('FEATURE NAMES')
    
```

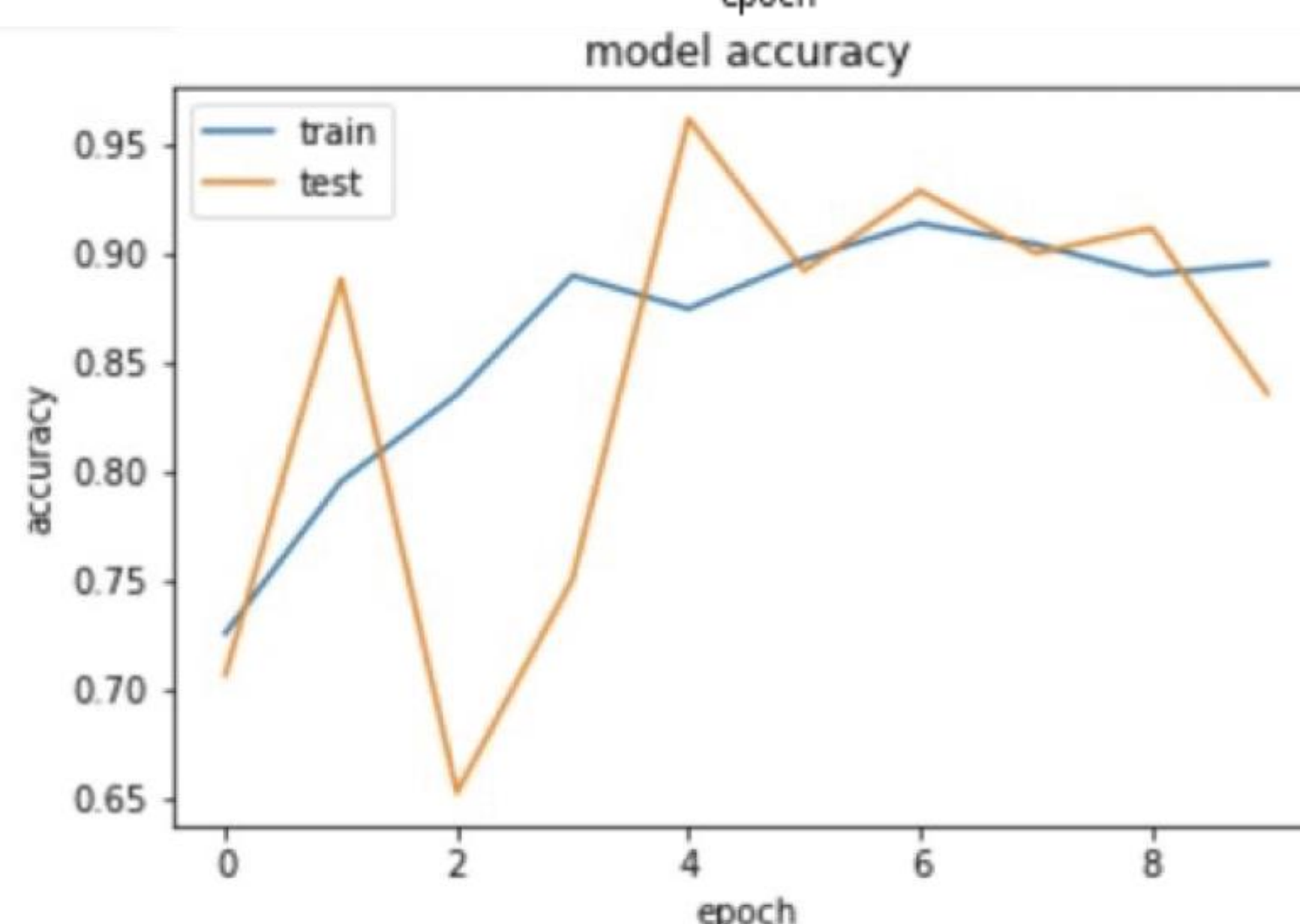
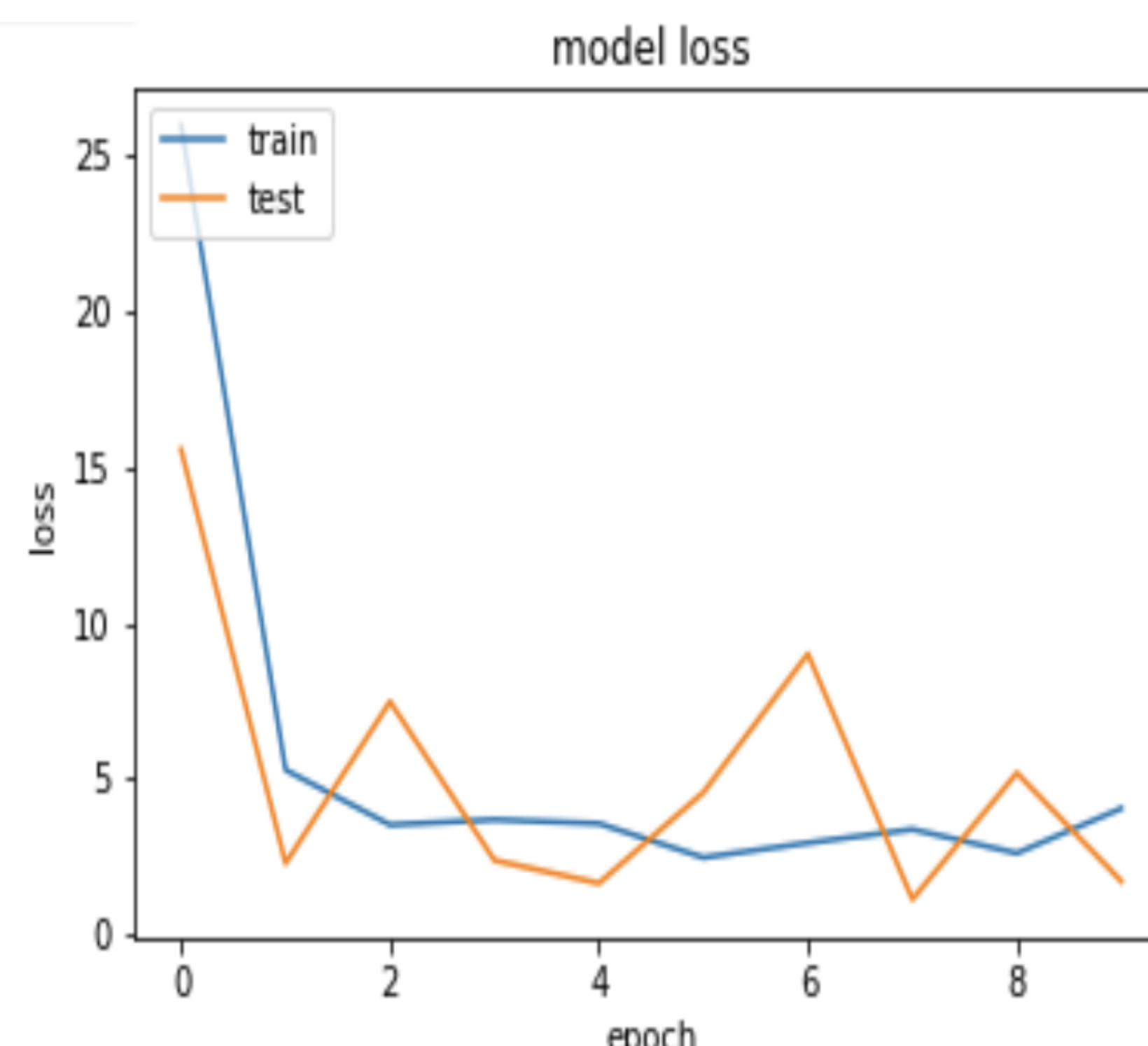
Technologies



Results

Models	Accuracy Score
Decision Tree	0.9996
Logistic Regression	0.939
Random Forest	0.9998
Neural Network	0.9998

		Predicted	
		0	1
Actual	0	2609	217
	1	98	2076



Conclusions

The final version of the product was delivered on November 17th, 2022, which includes a simple user interface that accurately displays the accuracy and loss from our train and test data. Although there were some complications and issues that led this project off track in the early stages, the client was overall satisfied by the product and all the requirements initially stated by the client were successfully met. Based upon our observations, Neural Network and Random Forest are most accurate.

Acknowledgments

1. Dr. Jack Zheng (Advisor)
2. Dr. Lian Zhou (Client)

Contact Information

<https://malicioustraffic.wixsite.com/malicious-traffic-de>

References

Dieng, W. (2020, August 14). *Walterdieng/TLS-malware-detection-with-machine-learning: Leveraging machine learning to detect TLS based malware in encrypted traffic without decryption*. GitHub. Retrieved November 14, 2022, from <https://github.com/WalterDieng/TLS-Malware-Detection-with-Machine-Learning>
 Johnson, Daniel. "Artificial Neural Network Tutorial with Tensorflow Ann Examples." *Guru99*, 17 Sept. 2022. <https://www.guru99.com/artificial-neural-network-tutorial.html>.