



Supplementary Data

Edit View Run Kernel Settings Help

+ ✂ 📄 📌 ▶ ■ ↺ ▶▶ Code ▾

```
4 # ## Explainable Intrusion Detection System (X-IDS)
5 # ### Integrating ML with SHAP/LIME for Transparent Cybersecurity
6
7 import pandas as pd
8 import numpy as np
9 import streamlit as st
10 import matplotlib.pyplot as plt
11 import seaborn as sns
12 from sklearn.model_selection import train_test_split, GridSearchCV
13 from sklearn.preprocessing import MinMaxScaler, OneHotEncoder
14 from sklearn.ensemble import RandomForestClassifier
15 from sklearn.svm import SVC
16 from sklearn.tree import DecisionTreeClassifier
17 from sklearn.linear_model import LogisticRegression
18 from sklearn.metrics import accuracy_score, f1_score, roc_auc_score
19 from imblearn.over_sampling import SMOTE
20 from sklearn.feature_selection import mutual_info_classif
21 import shap
22 import lime
23 import lime.lime_tabular
24 import tensorflow as tf
25 from tensorflow.keras.models import Sequential
26 from tensorflow.keras.layers import Dense, Dropout
27 import time
28 import json
29
30 # Configuration
31 pd.set_option('display.max_columns', None)
32 np.random.seed(42)
33 tf.random.set_seed(42)
```

```
Edit View Run Kernel Settings Help
- ✂ 📄 📌 ▶ ■ ⌂ ⏪ Code ▾

34 |
35 # --- Data Preparation ---
36 def load_and_preprocess_data(dataset_name='NSL-KDD'):
37     """Load and preprocess dataset with feature selection"""
38     if dataset_name == 'NSL-KDD':
39         # Load NSL-KDD dataset (replace with actual path)
40         df = pd.read_csv('KDDTrain+.txt', header=None)
41         cols = [
42             'duration', 'protocol_type', 'service', 'flag', 'src_bytes', 'dst_bytes',
43             'land', 'wrong_fragment', 'urgent', 'hot', 'num_failed_logins',
44             'logged_in', 'num_compromised', 'root_shell', 'su_attempted',
45             'num_root', 'num_file_creations', 'num_shells', 'num_access_files',
46             'num_outbound_cmds', 'is_host_login', 'is_guest_login', 'count',
47             'srv_count', 'serror_rate', 'srv_serror_rate', 'rerror_rate',
48             'srv_rerror_rate', 'same_srv_rate', 'diff_srv_rate',
49             'srv_diff_host_rate', 'dst_host_count', 'dst_host_srv_count',
50             'dst_host_same_srv_rate', 'dst_host_diff_srv_rate',
51             'dst_host_same_src_port_rate', 'dst_host_srv_diff_host_rate',
52             'dst_host_serror_rate', 'dst_host_srv_serror_rate',
53             'dst_host_rerror_rate', 'dst_host_srv_rerror_rate', 'label'
54         ]
55         df.columns = cols
56
57         # Convert labels to binary (normal = 0, attack = 1)
58         df['label'] = df['label'].apply(lambda x: 0 if x == 'normal' else 1)
59
60         # Select relevant features based on mutual information
61         X = df.drop('label', axis=1)
62         y = df['label']
63         cat_cols = ['protocol_type', 'service', 'flag']
64         num_cols = [col for col in X.columns if col not in cat_cols]
65
```

```
65
66     else: # CICIDS2017
67         # Load sample of CICIDS2017 (replace with actual path)
68         df = pd.read_csv('CICIDS2017_sample.csv')
69         df['label'] = df['Label'].apply(lambda x: 0 if 'BENIGN' in x else 1)
70         df = df.drop('Label', axis=1)
71
72         # Feature selection
73         X = df.drop('label', axis=1)
74         y = df['label']
75         cat_cols = []
76         num_cols = X.columns.tolist()
77
78     # Preprocessing pipeline
79     # One-hot encoding for categorical features
80     if cat_cols:
81         ohe = OneHotEncoder(handle_unknown='ignore', sparse=False)
82         X_cat = pd.DataFrame(ohe.fit_transform(X[cat_cols]))
83         X_cat.columns = ohe.get_feature_names_out(cat_cols)
84         X = pd.concat([X.drop(cat_cols, axis=1), X_cat], axis=1)
85
86     # Normalization
87     scaler = MinMaxScaler()
88     X = pd.DataFrame(scaler.fit_transform(X), columns=X.columns)
89
90     # Feature selection using mutual information
91     mi_scores = mutual_info_classif(X, y)
92     mi_df = pd.DataFrame({'Feature': X.columns, 'MI_Score': mi_scores})
93     top_features = mi_df.sort_values('MI_Score', ascending=False).head(20)['Feature'].tolist()
94     X = X[top_features]
95
96     return X, y, top_features
97
```

```
Edit View Run Kernel Settings Help
✂ 📄 🗑 ▶ ■ 🔄 ▶▶ Code ▼

97
98 # --- Model Training ---
99 def train_models(X_train, y_train):
100     """Train and tune multiple ML models"""
101     models = {}
102
103     # Random Forest
104     rf_params = {'n_estimators': [100, 300, 500], 'max_depth': [5, 10, 15]}
105     rf = GridSearchCV(RandomForestClassifier(), rf_params, cv=3, scoring='f1')
106     rf.fit(X_train, y_train)
107     models['Random Forest'] = rf.best_estimator_
108
109     # SVM
110     svm_params = {'C': [0.1, 1, 10], 'gamma': [0.01, 0.1, 1]}
111     svm = GridSearchCV(SVC(probability=True), svm_params, cv=3, scoring='f1')
112     svm.fit(X_train, y_train)
113     models['SVM'] = svm.best_estimator_
114
115     # Decision Tree
116     dt_params = {'max_depth': [5, 10, 15], 'criterion': ['gini', 'entropy']}
117     dt = GridSearchCV(DecisionTreeClassifier(), dt_params, cv=3, scoring='f1')
118     dt.fit(X_train, y_train)
119     models['Decision Tree'] = dt.best_estimator_
120
121     # Logistic Regression
122     lr = LogisticRegression(max_iter=1000, C=0.1)
123     lr.fit(X_train, y_train)
124     models['Logistic Regression'] = lr
125
126     # Deep Neural Network
127     model = Sequential([
128         Dense(128, activation='relu', input_shape=(X_train.shape[1],)),
129         Dropout(0.3),
130         Dense(64, activation='relu'),
131         Dense(32, activation='relu'),
132         Dense(1, activation='sigmoid')
133     ])
134     model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
135     model.fit(X_train, y_train, epochs=20, batch_size=64, validation_split=0.1, verbose=0)
136     models['DNN'] = model
137
138     return models
139
```

 jupyter XAI Last Checkpoint: 1 hour ago

File Edit View Run Kernel Settings Help

 Code

```
139
140 # --- XAI Functions ---
141 def generate_shap_explanations(model, X_sample, model_type='tree'):
142     """Generate SHAP explanations for a model"""
143     if model_type == 'tree':
144         explainer = shap.TreeExplainer(model)
145     elif model_type == 'dnn':
146         explainer = shap.DeepExplainer(model, X_sample.values)
147     else:
148         explainer = shap.KernelExplainer(model.predict_proba, X_sample)
149
150     shap_values = explainer.shap_values(X_sample)
151     return explainer, shap_values
152
153 def generate_lime_explanation(model, X_train, instance, feature_names):
154     """Generate LIME explanation for a single instance"""
155     explainer = lime.lime_tabular.LimeTabularExplainer(
156         training_data=X_train.values,
157         feature_names=feature_names,
158         class_names=['Normal', 'Attack'],
159         mode='classification'
160     )
161     exp = explainer.explain_instance(
162         data_row=instance.values,
163         predict_fn=model.predict_proba,
164         num_features=10
165     )
166     return exp
```

```
jupyter XAI Last Checkpoint: 1 hour ago
e Edit View Run Kernel Settings Help
+ 🔍 📄 ▶ ■ ↻ ▶▶ Code ▾ JupyterLab

167
168 # --- Dashboard Interface ---
169 def main():
170     st.set_page_config(layout="wide")
171     st.title("🔍 Explainable Intrusion Detection System (X-IDS)")
172     st.markdown("""
173     Machine Learning-powered IDS with SHAP & LIME explanations for transparent threat detection
174     """)
175
176     # Sidebar configuration
177     st.sidebar.header("Configuration")
178     dataset = st.sidebar.selectbox("Dataset", ["NSL-KDD", "CICIDS2017"])
179     model_type = st.sidebar.selectbox("Model", ["Random Forest", "SVM", "Decision Tree", "Logistic Regression", "DNN"])
180
181     # Load data
182     with st.spinner("Loading and preprocessing data..."):
183         X, y, feature_names = load_and_preprocess_data(dataset)
184         X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
185
186         # Handle imbalance with SMOTE
187         sm = SMOTE(random_state=42)
188         X_train, y_train = sm.fit_resample(X_train, y_train)
189
190     # Train models
191     with st.spinner(f"Training {model_type} model..."):
192         if 'models' not in st.session_state:
193             st.session_state.models = train_models(X_train, y_train)
194         model = st.session_state.models[model_type]
195
196     # Make predictions
197     if model_type == 'DNN':
198         y_pred = (model.predict(X_test) > 0.5).astype(int)
199     else:
200         y_pred = model.predict(X_test)
201
202     # Calculate metrics
203     accuracy = accuracy_score(y_test, y_pred)
204     f1 = f1_score(y_test, y_pred)
205     roc_auc = roc_auc_score(y_test, y_pred)
206
207     # Display performance metrics
208     st.subheader("Model Performance")
209     col1, col2, col3 = st.columns(3)
210     col1.metric("Accuracy", f"{accuracy:.4f}")
211     col2.metric("F1 Score", f"{f1:.4f}")
212     col3.metric("ROC AUC", f"{roc_auc:.4f}")
213
```

```
213 |
214 # Explanation section
215 st.subheader("Explainability Dashboard")
216 explanation_type = st.radio("Explanation Type", ["Global (SHAP)", "Local (SHAP)", "Local (LIME)"])
217
218 if explanation_type == "Global (SHAP)":
219     st.info("Global Feature Importance using SHAP")
220     with st.spinner("Generating global explanations..."):
221         # Sample data for faster computation
222         X_sample = X_train.sample(100, random_state=42)
223
224         # Generate SHAP values
225         model_for_shap = model
226         shap_model_type = 'tree' if model_type != 'DNN' else 'dnn'
227         explainer, shap_values = generate_shap_explanations(
228             model_for_shap,
229             X_sample,
230             model_type=shap_model_type
231         )
232
233         # Plot summary
234         fig, ax = plt.subplots()
235         shap.summary_plot(shap_values, X_sample, plot_type="bar", show=False)
236         st.pyplot(fig)
237
238         # Feature importance dataframe
239         shap_df = pd.DataFrame({
240             'Feature': feature_names,
241             'Importance': np.abs(shap_values).mean(0)
242         }).sort_values('Importance', ascending=False).head(10)
```

upyter XAI Last Checkpoint: 1 hour ago

Edit View Run Kernel Settings Help

- ✂ 📄 📄 ▶ ■ ↺ ▶▶ Code ▾

```
244
245     elif explanation_type.startswith("Local (SHAP)"):
246         st.info("Local Explanation for Individual Prediction using SHAP")
247         instance_idx = st.slider("Select instance", 0, len(X_test)-1, 50)
248         instance = X_test.iloc[[instance_idx]]
249         actual_label = "Attack" if y_test.iloc[instance_idx] == 1 else "Normal"
250         pred_label = "Attack" if y_pred[instance_idx] == 1 else "Normal"
251
252         st.write(f"**Actual:** {actual_label} | **Predicted:** {pred_label}")
253
254     with st.spinner("Generating local SHAP explanation..."):
255         # Generate SHAP values
256         model_for_shap = model
257         shap_model_type = 'tree' if model_type != 'DNN' else 'dnn'
258         explainer, shap_values = generate_shap_explanations(
259             model_for_shap,
260             instance,
261             model_type=shap_model_type
262         )
263
264         # Force plot
265         st.subheader("SHAP Force Plot")
266         fig, ax = plt.subplots()
267         shap.force_plot(
268             explainer.expected_value,
269             shap_values[0],
270             instance,
271             matplotlib=True,
272             show=False
273         )
274         st.pyplot(fig)
275
276         # Waterfall plot
277         st.subheader("Feature Contribution")
278         fig, ax = plt.subplots()
279         shap.plots.waterfall(shap_values[0], max_display=10, show=False)
280         st.pyplot(fig)
281
```

```
281
282     else: # LIME explanation
283         st.info("Local Explanation for Individual Prediction using LIME")
284         instance_idx = st.slider("Select instance", 0, len(X_test)-1, 100)
285         instance = X_test.iloc[[instance_idx]]
286         actual_label = "Attack" if y_test.iloc[instance_idx] == 1 else "Normal"
287         pred_label = "Attack" if y_pred[instance_idx] == 1 else "Normal"
288
289         st.write(f"**Actual:** {actual_label} | **Predicted:** {pred_label}")
290
291     with st.spinner("Generating LIME explanation..."):
292         # Generate LIME explanation
293         exp = generate_lime_explanation(
294             model,
295             X_train,
296             instance.iloc[0],
297             feature_names
298         )
299
300         # Show explanation
301         st.subheader("LIME Explanation")
302         st.write(exp.as_list())
303
304         # Visualize
305         fig = exp.as_pyplot_figure()
306         st.pyplot(fig)
307
308         # Show raw features
309         st.subheader("Instance Features")
310         st.dataframe(instance)
311
312 if __name__ == "__main__":
313     main()
```