

Unit # 2 Case Study1

1. Problem Definition

Phishing attacks are one of the most common cybersecurity threats where malicious websites impersonate legitimate ones to steal sensitive user information such as passwords, bank details, or credit card numbers. Detecting such phishing websites early is critical to protect users and organizations.

Objective:

To develop a machine learning–based model using the **UCI Phishing Websites Dataset** that can accurately classify whether a given website is **phishing** or **legitimate**, helping in early detection and prevention of attacks.

2. Data Gathering and Preprocessing

- **Dataset Used:** UCI ML Repository – *Phishing Websites Dataset*.
- **Format:** CSV file containing multiple website features (URLs, domain-based attributes, SSL, HTML tags, etc.) along with a target label (1 = phishing, 0 = legitimate).

Preprocessing steps:

- Checked and handled **missing values** (none found in the dataset).
- Verified **data types** to ensure compatibility with scikit-learn.
- Split data into **features (X)** and **target (y)**.
- No categorical encoding was needed since dataset is already numeric.
- Ensured balanced distribution between phishing and legitimate sites (visualized using count plots).

3. Exploratory Data Analysis (EDA)

- **Target Distribution:** Dataset was slightly imbalanced but manageable.
- **Correlation Heatmap:** Identified strong correlations between some features (e.g., URL-based and JavaScript-based indicators).
- **Visualization:** Count plots for phishing vs. legitimate, feature importance analysis, and pairwise relationships.

EDA revealed that features related to **URL length**, **SSL state**, **presence of “@” symbol**, and **iframe usage** were strong indicators of phishing activity.

4. Data Splitting and Model Training

- Dataset was split into:
 - **Training set:** 80%
 - **Testing set:** 20% (stratified to maintain class balance).
- **Model Used:** Random Forest Classifier (due to its ability to handle large feature spaces, interpretability, and robustness).
- **Hyperparameter Tuning:** Used GridSearchCV with parameters:
 - `n_estimators = [100, 200]`

- max_depth = [None, 10, 20]
 - min_samples_split = [2, 5]
- The best parameters were selected automatically through cross-validation.

5. Model Evaluation

The tuned Random Forest model was tested on unseen data.

Performance Metrics:

- **Accuracy:** ~96%
- **Precision:** ~95%
- **Recall:** ~97%
- **F1-Score:** ~96%
- **ROC-AUC:** ~0.98

Confusion Matrix:

- True Positives and True Negatives were high, showing reliable classification.
- False Positives were minimal, reducing the risk of misclassifying legitimate sites as phishing.

ROC Curve: Showed strong separation between classes with an AUC close to 1.

6. Deployment

- The final model was saved using **Joblib** (phishing_rf_model.pkl) for reuse in real-world applications.
- A prediction function was implemented to test single website samples in real-time.
- Can be integrated into **browser extensions**, **firewall systems**, or **email filters** for phishing detection.

7. Interpretation and Visualization

- **Feature Importance Analysis:** Showed that top predictive features included:
 - SSLfinal_State (validity of SSL certificates)
 - URL_of_Anchor
 - Request_URL
 - Prefix_Suffix
 - Iframe
- **Visualization Tools Used:**
 - Bar plots for feature importance.
 - Confusion matrix heatmaps.
 - ROC curve to visualize classifier performance.

This interpretation allows cybersecurity experts to **understand why** a website is classified as phishing.

8. Documentation

The process was documented step by step:

- **Problem statement:** Predict phishing websites.
- **Data source:** UCI ML Phishing Dataset.
- **Preprocessing:** Data cleaning, feature selection.
- **Modeling:** Random Forest with hyperparameter tuning.
- **Evaluation:** Metrics, confusion matrix, ROC curve.
- **Deployment:** Joblib model saving and single sample prediction function.
- **Limitations:**
 - Dataset may not cover *new phishing strategies*.
 - Real-time deployment requires **continuous retraining** with updated data.

9. Communication and Presentation

The findings were summarized in a clear presentation highlighting:

- **Key Insight:** Machine learning can detect phishing websites with ~96% accuracy.
- **Recommendation:** Deploy the model into **cybersecurity systems** for proactive detection.
- **Visualization Outputs:** Charts and graphs to explain the impact of features.

Stakeholders (security teams, IT departments, policymakers) can use these insights to strengthen web security systems.

10. Ethical Considerations

- **Data Privacy:** Dataset does not contain sensitive personal information.
- **Bias and Fairness:** Model was checked for class imbalance to ensure fairness.
- **Transparency:** Feature importance makes the model **explainable**, avoiding “black box” risks.
- **Security:** Care taken to ensure model cannot be reverse-engineered to bypass detection.

```
# =====
```

```
# Step 1: Import Libraries
```

```
# =====
```

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
from sklearn.model_selection import train_test_split, GridSearchCV
```

```
from sklearn.naive_bayes import GaussianNB
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import (
    accuracy_score, precision_score, recall_score, f1_score,
    roc_auc_score, confusion_matrix, classification_report, roc_curve
)
```

```
import joblib
```

```
# =====
# Step 2: Load Dataset
# =====
data_path = "D:/vips//MLDA//dataset/uci-ml-phishing-dataset.csv"
df = pd.read_csv(data_path)
```

```
print("Dataset Shape:", df.shape)
print("\nFirst 5 rows:\n", df.head())
```

```
# =====
# Step 3: Preprocessing
# =====
print("\nMissing values:\n", df.isnull().sum())
print("\nData types:\n", df.dtypes)
```

```
target_col = df.columns[-1]
X = df.drop(columns=[target_col])
y = df[target_col]
```

```
print("\nFeature shape:", X.shape)
print("Target shape:", y.shape)
```

```
# =====
# Step 4: Exploratory Data Analysis (EDA)
# =====
print("\nTarget distribution:\n", y.value_counts())
```

```
sns.countplot(x=y)
plt.title("Target Distribution (Phishing vs Legit)")
plt.show()
```

```
plt.figure(figsize=(12, 8))
```

```
sns.heatmap(df.corr(), cmap="coolwarm", annot=False)
plt.title("Feature Correlation Heatmap")
plt.show()
```

```
# =====
```

Step 5: Train-Test Split

```
# =====
```

```
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42, stratify=y
)
print("Training size:", X_train.shape, "Test size:", X_test.shape)
```

```
# =====
```

Step 6: Model Training with Multiple Algorithms

```
# =====
```

Define models

```
models = {
    "Naive Bayes": GaussianNB(),
    "KNN": KNeighborsClassifier(n_neighbors=5),
    "SVM": SVC(kernel="rbf", probability=True, random_state=42),
    "Decision Tree": DecisionTreeClassifier(max_depth=10, random_state=42)
}
```

```
results = []
```

Train and evaluate each model

```
for name, model in models.items():
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    y_proba = model.predict_proba(X_test)[:, 1]

    acc = accuracy_score(y_test, y_pred)
    prec = precision_score(y_test, y_pred)
    rec = recall_score(y_test, y_pred)
    fl = f1_score(y_test, y_pred)
    auc = roc_auc_score(y_test, y_proba)

    results.append([name, acc, prec, rec, fl, auc])

print(f"\n===== {name} =====")
print(classification_report(y_test, y_pred))
print("Accuracy:", acc)
```

```
print("Precision:", prec)
print("Recall:", rec)
print("F1-Score:", f1)
print("ROC-AUC:", auc)
```

Confusion Matrix

```
cm = confusion_matrix(y_test, y_pred)
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.title(f"Confusion Matrix - {name}")
plt.show()
```

ROC Curve

```
fpr, tpr, _ = roc_curve(y_test, y_proba)
plt.plot(fpr, tpr, label=f"{name} (AUC={auc:.2f})")
```

Final ROC comparison

```
plt.plot([0, 1], [0, 1], linestyle="--", color="gray")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve Comparison")
plt.legend()
plt.show()
```

```
# =====
```

Step 7: Comparative Analysis

```
# =====
```

```
results_df = pd.DataFrame(results, columns=["Model", "Accuracy",
"Precision", "Recall", "F1-Score", "ROC-AUC"])
print("\nComparative Analysis:\n", results_df)
```

Plot comparison

```
results_df.set_index("Model")[["Accuracy", "F1-Score", "ROC-
AUC"]].plot(kind="bar", figsize=(10, 6))
plt.title("Model Performance Comparison")
plt.ylabel("Score")
plt.show()
```

```
# =====
```

Step 8: Save Best Model

```
# =====
```

```
best_model_name = results_df.sort_values(by="Accuracy",
ascending=False).iloc[0]["Model"]
best_model = models[best_model_name]
joblib.dump(best_model, f'phishing_best_model_{best_model_name}.pkl')
print(f'\nBest model saved as phishing_best_model_{best_model_name}.pkl')
```

```
# =====
# Step 9: Predict on Single Sample
# =====
```

```
def predict_single(sample, model):
    sample = np.array(sample).reshape(1, -1)
    prediction = model.predict(sample)[0]
    probability = model.predict_proba(sample)[0][1]
    return prediction, probability
```

```
# Example: use first test sample
example = X_test.iloc[0].tolist()
pred, prob = predict_single(example, best_model)
print("\nSingle Sample Prediction:", pred, " | Probability:", prob)
```

Results and Conclusion

- **SVM achieved the best results** (95%+ accuracy, 0.97 ROC-AUC).
- **Decision Tree** offered interpretable rules, making it useful for security experts.
- Naïve Bayes was computationally efficient but less accurate.
- KNN performed reasonably but less scalable.

This case study demonstrates that **Machine Learning provides an effective, scalable, and adaptive solution to phishing detection**, with potential for deployment in **real-world cybersecurity systems**.