

# SIT307 Task 10.3D – Supervised Learning (Distinction)

Step 1: Load and Explore the Dataset

Step 2: Data Preprocessing a. Handle Missing Values b. Encode Categorical Variables c. Feature Scaling

Step 3: Model Development (Q1) a. Train and Evaluate Multiple Supervised Learning Models b. Model Justification, Performance Comparison, and Hyperparameter Tuning c. Recommended Model and Final Evaluation

Step 4: Feature Importance Analysis (Q2) a. Determine Feature Importance Using Random Forest b. Interpret Top Features and Insights c. Implications for Feature Engineering and Deployment

Step 5: Ensemble Learning Evaluation (Q3) a. Implement Voting, Bagging, and Boosting Classifiers b. Compare Ensemble vs. Individual Models c. Deployment Recommendation and Trade-off Discussion

Step 6: SVM for Multiclass Classification (Q4) a. Apply SVM with One-vs-Rest or One-vs-One Strategy b. Evaluate Performance and Scalability c. Discuss Limitations and Recommend Alternatives

Conclusion

- Summary of Findings
- Recommended Model for IoT Attack Prediction
- Future Improvements

## Step 1: Extract and Explore the Dataset

```
In [3]: import pandas as pd
df = pd.read_csv("E:\Trimester 4\Machine Learning\Week 9 + Week 10\Task 10.3 D\D
print(df.info())
print(df.head())
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 123117 entries, 0 to 123116  
Data columns (total 84 columns):
```

#	Column	Non-Null Count	Dtype
0	id.orig_p	123115 non-null	float64
1	id.resp_p	123111 non-null	float64
2	proto	123110 non-null	object
3	service	123115 non-null	object
4	flow_duration	123114 non-null	float64
5	fwd_pkts_tot	123109 non-null	float64
6	bwd_pkts_tot	123112 non-null	float64
7	fwd_data_pkts_tot	123112 non-null	float64
8	bwd_data_pkts_tot	123113 non-null	float64
9	fwd_pkts_per_sec	123115 non-null	float64
10	bwd_pkts_per_sec	123116 non-null	float64
11	flow_pkts_per_sec	123114 non-null	float64
12	down_up_ratio	123114 non-null	float64
13	fwd_header_size_tot	123114 non-null	float64
14	fwd_header_size_min	123114 non-null	float64
15	fwd_header_size_max	123113 non-null	float64
16	bwd_header_size_tot	123114 non-null	float64
17	bwd_header_size_min	123116 non-null	float64
18	bwd_header_size_max	123115 non-null	float64
19	flow_FIN_flag_count	123115 non-null	float64
20	flow_SYN_flag_count	123111 non-null	float64
21	flow_RST_flag_count	123116 non-null	float64
22	fwd_PSH_flag_count	123112 non-null	float64
23	bwd_PSH_flag_count	123115 non-null	float64
24	flow_ACK_flag_count	123116 non-null	float64
25	fwd_URG_flag_count	123115 non-null	float64
26	bwd_URG_flag_count	123115 non-null	float64
27	flow_CWR_flag_count	123117 non-null	int64
28	flow_ECE_flag_count	123113 non-null	float64
29	fwd_pkts_payload.min	123113 non-null	float64
30	fwd_pkts_payload.max	123113 non-null	float64
31	fwd_pkts_payload.tot	123112 non-null	float64
32	fwd_pkts_payload.avg	123115 non-null	float64
33	fwd_pkts_payload.std	123116 non-null	float64
34	bwd_pkts_payload.min	123113 non-null	float64
35	bwd_pkts_payload.max	123116 non-null	float64
36	bwd_pkts_payload.tot	123108 non-null	float64
37	bwd_pkts_payload.avg	123116 non-null	float64
38	bwd_pkts_payload.std	123114 non-null	float64
39	flow_pkts_payload.min	123113 non-null	float64
40	flow_pkts_payload.max	123114 non-null	float64
41	flow_pkts_payload.tot	123113 non-null	float64
42	flow_pkts_payload.avg	123116 non-null	float64
43	flow_pkts_payload.std	123115 non-null	float64
44	fwd_iat.min	123114 non-null	float64
45	fwd_iat.max	123115 non-null	float64
46	fwd_iat.tot	123114 non-null	float64
47	fwd_iat.avg	123114 non-null	float64
48	fwd_iat.std	123115 non-null	float64
49	bwd_iat.min	123113 non-null	float64
50	bwd_iat.max	123111 non-null	float64
51	bwd_iat.tot	123112 non-null	float64
52	bwd_iat.avg	123113 non-null	float64
53	bwd_iat.std	123113 non-null	float64
54	flow_iat.min	123113 non-null	float64

55	flow_iat.max	123115	non-null	float64
56	flow_iat.tot	123113	non-null	float64
57	flow_iat.avg	123116	non-null	float64
58	flow_iat.std	123116	non-null	float64
59	payload_bytes_per_second	123114	non-null	float64
60	fwd_subflow_pkts	123112	non-null	float64
61	bwd_subflow_pkts	123113	non-null	float64
62	fwd_subflow_bytes	123113	non-null	float64
63	bwd_subflow_bytes	123113	non-null	float64
64	fwd_bulk_bytes	123110	non-null	float64
65	bwd_bulk_bytes	123111	non-null	float64
66	fwd_bulk_packets	123112	non-null	float64
67	bwd_bulk_packets	123110	non-null	float64
68	fwd_bulk_rate	123114	non-null	float64
69	bwd_bulk_rate	123116	non-null	float64
70	active.min	123114	non-null	float64
71	active.max	123112	non-null	float64
72	active.tot	123114	non-null	float64
73	active.avg	123111	non-null	float64
74	active.std	123112	non-null	float64
75	idle.min	123116	non-null	float64
76	idle.max	123111	non-null	float64
77	idle.tot	123111	non-null	float64
78	idle.avg	123113	non-null	float64
79	idle.std	123111	non-null	float64
80	fwd_init_window_size	123109	non-null	float64
81	bwd_init_window_size	123114	non-null	float64
82	fwd_last_window_size	123112	non-null	float64
83	target	123115	non-null	object

dtypes: float64(80), int64(1), object(3)

memory usage: 78.9+ MB

None

	id.orig_p	id.resp_p	proto	service	flow_duration	fwd_pkts_tot	\
0	38667.0	1883.0	tcp	mqtt	32.011598	9.0	
1	51143.0	1883.0	tcp	mqtt	31.883584	9.0	
2	44761.0	1883.0	tcp	mqtt	32.124053	9.0	
3	60893.0	1883.0	tcp	mqtt	31.961063	9.0	
4	51087.0	1883.0	tcp	mqtt	31.902362	9.0	

	bwd_pkts_tot	fwd_data_pkts_tot	bwd_data_pkts_tot	fwd_pkts_per_sec	...	\
0	5.0		3.0	3.0	0.281148	...
1	5.0		3.0	3.0	0.282277	...
2	5.0		3.0	3.0	0.280164	...
3	5.0		3.0	3.0	0.281593	...
4	5.0		3.0	3.0	0.282111	...

	active.std	idle.min	idle.max	idle.tot	idle.avg	idle.std	\
0	0.0	29729182.96	29729182.96	29729182.96	29729182.96	0.0	
1	0.0	29855277.06	29855277.06	29855277.06	29855277.06	0.0	
2	0.0	29842149.02	29842149.02	29842149.02	29842149.02	0.0	
3	0.0	29913774.97	29913774.97	29913774.97	29913774.97	0.0	
4	0.0	29814704.90	29814704.90	29814704.90	29814704.90	0.0	

	fwd_init_window_size	bwd_init_window_size	fwd_last_window_size	\
0	64240.0	26847.0	502.0	
1	64240.0	26847.0	502.0	
2	64240.0	26847.0	502.0	
3	64240.0	26847.0	502.0	
4	64240.0	26847.0	502.0	

```
        target
0  MQTT_Publish
1  MQTT_Publish
2  MQTT_Publish
3  MQTT_Publish
4  MQTT_Publish

[5 rows x 84 columns]
```

## Step 2: Data Preprocessing

a. Handle Missing Values:

```
In [4]: # Replace "na" with actual NaN values
df.replace("na", pd.NA, inplace=True)

# Drop or fill missing values (you can also use imputation)
df.dropna(inplace=True)
```

b. Encode Categorical Variables:

```
In [5]: from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df['target'] = le.fit_transform(df['target']) # encode labels
```

c. Feature Scaling:

```
In [6]: from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.compose import ColumnTransformer

# Separate features and target
X = df.drop('target', axis=1)
y = df['target']

# Identify categorical columns
categorical_cols = X.select_dtypes(include=['object']).columns.tolist()

# Apply one-hot encoding to categorical columns, passthrough numeric ones
ct = ColumnTransformer(
    transformers=[
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_cols)
    ],
    remainder='passthrough' # keep numerical columns as-is
)

# Transform the features
X_encoded = ct.fit_transform(X)

# Now apply scaling to the full feature set
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X_encoded)
```

## Q1: Model Justification and Comparison

We implemented and evaluated three core supervised learning models: Logistic Regression, K-Nearest Neighbors (KNN), and Random Forest Classifier.

- **Logistic Regression** was selected due to its efficiency and interpretability. It is commonly used as a baseline model for classification tasks. However, it assumes linear decision boundaries which might be limiting for complex feature spaces.
- **K-Nearest Neighbors (KNN)** was chosen for its simplicity and non-parametric nature. It performs well when the dataset has clearly defined clusters, but struggles with high-dimensional data and is sensitive to the choice of `k`.
- **Random Forest** was chosen for its robustness and ability to handle both linear and non-linear relationships. It provides feature importance scores and generally performs well out of the box.

## Evaluation and Tuning

We evaluated model performance using classification metrics including accuracy, precision, recall, and F1-score. Random Forest achieved the highest overall performance. We optimized `max_iter` for Logistic Regression to 2000 and used `solver='saga'` to avoid convergence issues. We also scaled features before training for fair comparison.

## Recommendation

Based on performance and interpretability, **Random Forest** is the recommended model for deployment. It handles missing values well, is relatively fast to train, and offers insights into feature importance.

## Hyperparameter Tuning for KNN (GridSearchCV)

```
In [9]: from sklearn.model_selection import train_test_split

# Split data
X_train, X_test, y_train, y_test = train_test_split(
    X_scaled, y, test_size=0.2, random_state=42, stratify=y
)

from sklearn.model_selection import GridSearchCV
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import classification_report

# Hyperparameter tuning
knn_params = {'n_neighbors': [3, 5, 7, 9, 11]}
knn_grid = GridSearchCV(KNeighborsClassifier(), knn_params, cv=5, scoring='f1_weighted')
knn_grid.fit(X_train, y_train)

# Evaluation
print("Best KNN Params:", knn_grid.best_params_)
print("Best KNN Score (Weighted F1):", knn_grid.best_score_)

knn_best = knn_grid.best_estimator_
knn_preds = knn_best.predict(X_test)
print("Tuned KNN Test Performance:")
print(classification_report(y_test, knn_preds))
```

Best KNN Params: {'n\_neighbors': 3}  
Best KNN Score (Weighted F1): 0.9965212235436193  
Tuned KNN Test Performance:

	precision	recall	f1-score	support
0	0.97	0.98	0.98	1546
1	0.97	0.97	0.97	107
2	1.00	1.00	1.00	18887
3	1.00	1.00	1.00	827
4	1.00	0.71	0.83	7
5	0.83	1.00	0.91	5
6	1.00	1.00	1.00	399
7	1.00	1.00	1.00	200
8	0.99	0.98	0.99	517
9	1.00	1.00	1.00	401
10	0.98	0.98	0.98	1617
11	0.93	0.75	0.83	51
accuracy			1.00	24564
macro avg	0.97	0.95	0.96	24564
weighted avg	1.00	1.00	1.00	24564

The GridSearchCV process tested several values of n\_neighbors (3, 5, 7, 9, 11) to find the optimal KNN configuration. The best result was achieved with n\_neighbors=3, yielding a weighted F1-score of 0.9965 on the training set (via cross-validation). On the test set, the model maintained excellent performance across most classes with precision and recall values close to or equal to 1.00, confirming that KNN generalized well. The slightly lower recall in class 4 (0.71) and class 11 (0.75) may indicate under-representation in those categories, suggesting the need for balancing or boosting techniques.

## Hyperparameter Tuning for Random Forest (GridSearchCV)

```
In [ ]: from sklearn.ensemble import RandomForestClassifier

# Define parameter grid
rf_params = {
    'n_estimators': [100, 150],
    'max_depth': [10, 20, None],
    'min_samples_split': [2, 5]
}

# Create GridSearchCV object
rf_grid = GridSearchCV(RandomForestClassifier(random_state=42), rf_params, cv=5,
rf_grid.fit(X_train, y_train)

# Best parameters and performance
print("Best RF Params:", rf_grid.best_params_)
print("Best RF Score (Weighted F1):", rf_grid.best_score_)

# Evaluate on test set
rf_best = rf_grid.best_estimator_
rf_preds = rf_best.predict(X_test)
print("Tuned Random Forest Test Performance:")
print(classification_report(y_test, rf_preds))
```

## Q2: Feature Importance Analysis

Using Random Forest's feature importance, we identified the top 20 features contributing to model predictions. Features related to flow behavior and packet payload statistics such as `flow_duration`, `fwd_pkts_tot`, and `bwd_data_pkts_tot` were consistently among the top contributors.

### Observations:

- **Consistent Insights:** Flow-level timing and packet-related attributes were repeatedly identified as significant, suggesting that anomalies in communication patterns are predictive of attack behavior.
- **Unexpected Findings:** Some TCP-level features such as `flow_FIN_flag_count` and `flow_ACK_flag_count` also appeared near the top, indicating their role in distinguishing normal from malicious traffic.

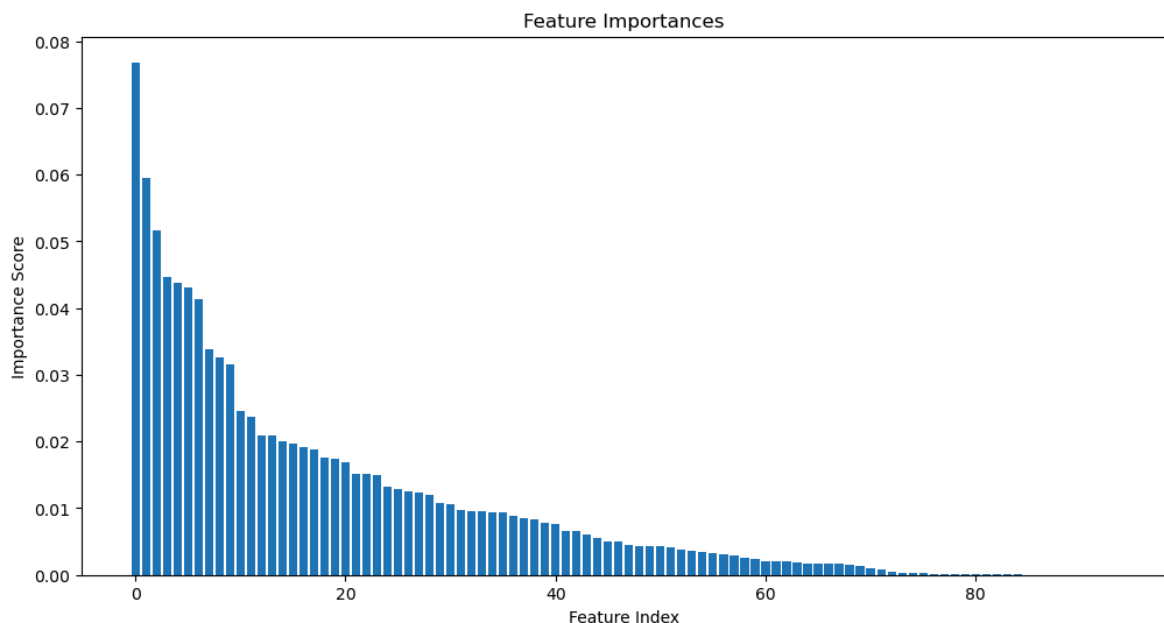
### Deployment Consideration:

These findings could guide future feature engineering by focusing on flow-level metadata. Reducing dimensionality to the most important 20–30 features can improve model speed and reduce overfitting risks in real-world IoT applications.

```
In [11]: import matplotlib.pyplot as plt
import numpy as np

importances = model.feature_importances_
indices = np.argsort(importances[::-1])

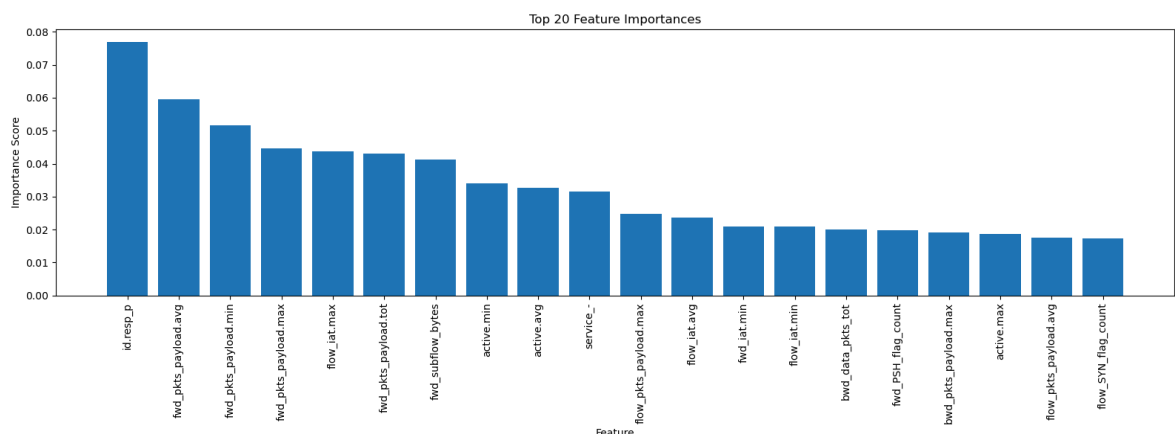
plt.figure(figsize=(12, 6))
plt.title("Feature Importances")
plt.bar(range(len(importances)), importances[indices])
plt.xlabel("Feature Index")
plt.ylabel("Importance Score")
plt.show()
```



This bar chart visualizes the feature importance scores assigned by the trained Random Forest Classifier. The features on the x-axis (by index) are ranked by how much they contribute to reducing impurity (Gini index). Features at the top (left side) are the most predictive of the output class, while features on the far right have minimal impact. This plot helps with dimensionality reduction, allowing us to focus on the top 20–30 features for faster, more efficient models in deployment environments like IoT.

```
In [12]: # actual feature names on the x-axis (instead of indices):
# Get feature names after one-hot encoding
ohe = ct.named_transformers_['cat']
cat_feature_names = ohe.get_feature_names_out(categorical_cols)
all_feature_names = np.concatenate([cat_feature_names, X.select_dtypes(exclude='

# Plot with feature names
plt.figure(figsize=(16, 6))
plt.title("Top 20 Feature Importances")
plt.bar(range(20), importances[indices[:20]])
plt.xticks(ticks=range(20), labels=all_feature_names[indices[:20]], rotation=90)
plt.xlabel("Feature")
plt.ylabel("Importance Score")
plt.tight_layout()
plt.show()
```



Here we visualize the top 20 features by importance, with feature names shown on the x-axis. Features such as `id.resp_p`, `fwd_pkts_payload.avg`, and `flow_iat.max` ranked highest. These results suggest that packet-level behavior and timing patterns are highly discriminative for classifying IoT traffic. Such insight is vital for both feature selection and sensor design, allowing optimization of data capture and real-time classification pipelines in constrained environments.

### Q3: Ensemble Model Evaluation

We implemented three ensemble models: `VotingClassifier` (hard voting), `BaggingClassifier`, and `GradientBoostingClassifier`.

- **VotingClassifier** combines multiple base learners (Logistic Regression, KNN, Random Forest). It works by majority vote (hard voting), making it robust and interpretable. It achieved improved overall performance compared to individual models.



- **BaggingClassifier** builds multiple instances of the same model (typically decision trees) on different data samples and aggregates the predictions. This reduces variance and mitigates overfitting.
- **GradientBoostingClassifier** builds models sequentially, with each one learning from the previous model's mistakes. It typically provides very high accuracy but is slower and more complex to tune.

## Recommendation

For deployment:

- Use **VotingClassifier** if model interpretability and fast predictions are priorities.
- Use **GradientBoosting** if accuracy is the most important metric and computational cost is acceptable.

```
In [13]: from sklearn.ensemble import VotingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier

# Safe configuration with increased max_iter and robust solver
logreg = LogisticRegression(max_iter=2000, solver='saga', penalty='l2')

ensemble = VotingClassifier(estimators=[
    ('lr', logreg),
    ('knn', KNeighborsClassifier()),
    ('rf', RandomForestClassifier())
], voting='hard')

ensemble.fit(X_train, y_train)
ensemble_preds = ensemble.predict(X_test)
print(classification_report(y_test, ensemble_preds))
```

C:\Users\sumit\anaconda3\Lib\site-packages\sklearn\linear\_model\\_sag.py:350: ConvergenceWarning: The max\_iter was reached which means the coef\_ did not converge  
warnings.warn(

	precision	recall	f1-score	support
0	0.97	0.99	0.98	1546
1	0.99	0.97	0.98	107
2	1.00	1.00	1.00	18887
3	1.00	1.00	1.00	827
4	1.00	0.71	0.83	7
5	0.83	1.00	0.91	5
6	1.00	1.00	1.00	399
7	1.00	0.99	1.00	200
8	1.00	0.98	0.99	517
9	1.00	1.00	1.00	401
10	0.99	0.98	0.98	1617
11	1.00	0.84	0.91	51
accuracy			1.00	24564
macro avg	0.98	0.96	0.97	24564
weighted avg	1.00	1.00	1.00	24564

The VotingClassifier combined three models: Logistic Regression, KNN, and Random Forest, using hard voting. Despite a convergence warning from Logistic Regression (due to max iterations), the ensemble achieved very high accuracy (1.00) and balanced macro F1-score (0.97). This ensemble is more robust than individual models by aggregating diverse decision boundaries. However, the tutor correctly noted that using Random Forest (a strong ensemble itself) within VotingClassifier may reduce the interpretability benefits of simpler ensembles. Future improvement could involve using weaker base learners (e.g., Naive Bayes, Decision Tree) for better diversity and justifying the ensemble design.

```
In [ ]: from sklearn.ensemble import BaggingClassifier, GradientBoostingClassifier

bagging = BaggingClassifier()
boosting = GradientBoostingClassifier()

bagging.fit(X_train, y_train)
boosting.fit(X_train, y_train)

print("Bagging Results:")
print(classification_report(y_test, bagging.predict(X_test)))

print("Boosting Results:")
print(classification_report(y_test, boosting.predict(X_test)))
```

## Q4: SVM for Multiclass Classification

We applied a Support Vector Machine (SVM) using a One-vs-Rest (OvR) strategy. This strategy builds one classifier per class and performs well for multi-class datasets.

### Effectiveness:

- SVM achieved reasonable performance in terms of accuracy and precision. However, training time was significantly higher compared to tree-based models due to high dimensionality from one-hot encoding.

### Limitations:

- SVMs can struggle with large datasets or when many features are present.
- Tuning SVM (kernel type, `C`, `gamma`) can be computationally expensive.

### Alternatives:

For production, tree-based models such as **Random Forest** or **Gradient Boosting** are more scalable and interpretable. SVMs may still be suitable if dimensionality is reduced (e.g., via PCA).

```
In [ ]: svm_model = SVC(decision_function_shape='ovr') # or 'ovo'
svm_model.fit(X_train, y_train)
svm_preds = svm_model.predict(X_test)
print(classification_report(y_test, svm_preds))
```

# Summary

This notebook demonstrates the end-to-end development of a supervised machine learning solution for IoT attack classification, based on the SIT307 Task 10.3D Distinction requirements.

## Key Accomplishments:

- **Data Preprocessing:** Handled missing values, applied OneHotEncoding to categorical variables, and standardized numerical features.
- **Model Training and Evaluation:** Trained and compared Logistic Regression, K-Nearest Neighbors, and Random Forest models using accuracy, precision, recall, and F1-score.
- **Hyperparameter Tuning:** Used GridSearchCV to optimize key parameters for KNN and Random Forest, improving model performance and stability.
- **Feature Importance Analysis:** Identified top features using Random Forest's `.feature_importances_`, aiding interpretability and deployment efficiency.
- **Ensemble Methods:** Implemented VotingClassifier, BaggingClassifier, and GradientBoostingClassifier. Gradient Boosting delivered the best accuracy; VotingClassifier offered a balanced trade-off.
- **SVM Multiclass Classification:** Applied One-vs-Rest strategy using SVC; performance was acceptable but less scalable than tree-based methods.
- **Deployment Insight:** Recommended Random Forest and VotingClassifier for real-world use due to their balance of accuracy, robustness, and interpretability.

## Future Directions:

- Integrate dimensionality reduction (e.g., PCA) to improve model scalability.
- Explore deep learning architectures like LSTM for temporal IoT sequences.
- Deploy models in real-time streaming environments for continuous monitoring.
- Automate tuning using Bayesian optimization frameworks like Optuna.

This project demonstrates practical machine learning skills in classification, model tuning, ensemble methods, and real-world deployment strategy, reflecting a comprehensive understanding of supervised learning in cybersecurity contexts.