Programming Assignment 5

# CYBERSECURITY ATTACK CLASSIFICATION USING RANDOM FOREST

You are a cybersecurity intern at SecureNet Corp. The company has been experiencing various network security breaches, and your team is tasked with developing a machine learning model to classify different types of cyberattacks.

The senior engineer hands you a dataset of simulated network traffic and says, *"We need a robust classifier to detect attacks like DDoS, Botnet, and others. Use your ML skills to build this. Oh, and don't forget— real-world data is messy; make sure your model can handle it."*

Your goal is to preprocess the data and train a Random Forest classifier to identify the type of attack.

## Dataset

The dataset cyber_attacks.csv contains the following features:
- protocol_type: Protocol (0: TCP, 1: UDP, 2: ICMP).
- duration: Connection duration (seconds).
- src_bytes: Bytes sent from source.
- dst_bytes: Bytes sent to destination.
- num_packets: Number of packets transmitted.
- num_connections: Number of connections.
- attack_type: Type of attack (Normal, Botnet, DDoS, PortScan, Phishing).

Download the dataset here: cyber_attacks.csv.

## General Guidelines
1. Load and perform exploratory the data analysis.
2. Perform data preprocessing.
3. Train a Random Forest Classifier using the preprocessed data to train the model.
4. Perform feature importance analysis.
5. Evaluate the model performance and report precision, recall, and F1-score.

## Guide Questions
Answer the following questions.
1. Preprocessing

    - How did you handle missing values in numerical features?
    - Why is there a need to encode protocol_type feature, and which encoding method can be used?

2. Model Training

    - How did you address class imbalance in the dataset?
    - What hyperparameters of the Random Forest did you tune, and why?
    - Why is a pipeline useful for this task?

3. Evaluation

    - Why is accuracy alone insufficient to evaluate this model?
    - Which attack type was hardest to classify, and why might this be?

- How would you improve the model's performance on minority classes?

4. Interpretation
   - Which features were most important for classification? Does this align with real-world attack patterns?
   - How would you explain the model's decision-making process to a non-technical team?

5. Application
   - What steps would you take to deploy this model in a real-time network monitoring system?
   - How would you handle new attack types not present in the training data?

## Requirements

- Ensure that your code is clean, well-commented, and organized.
- Use Python libraries such as `numpy` and `pandas` for data manipulation and `matplotlib` or `seaborn` for visualization.

## Submission

1. Submit your work as a Jupyter Notebook (.ipynb) file.
2. Upload your Jupyter Notebook to your GitHub repository. Ensure the notebook is well-documented with markdown cells explaining each step and the corresponding results.
3. Provide the link to your GitHub repository for grading.