# Network Data and Cybersecurity

**Problem Statement:**

**Goal:** Develop a machine learning model that can accurately detect and classify network attacks using the CIC-UNSW-NB15 dataset.

**Problem:** Given a set of network traffic logs with features describing connection metadata and behavior, predict whether each connection instance represents a normal (benign) event or a malicious (attack) event — and ideally identify the type of attack.

**Context:**

With the exponential growth of internet-connected systems, cyberattacks have become more sophisticated and frequent. Network administrators need automated and intelligent systems to detect and prevent malicious activity in real time. Traditional intrusion detection systems (IDS) rely heavily on static rules and signatures, which cannot keep up with new or evolving attack patterns.

Machine learning offers a dynamic alternative: by training models on labeled network traffic data, we can detect patterns that distinguish normal traffic from various types of attacks.

The CIC-UNSW-NB15 dataset, created by the University of New South Wales and the Canadian Institute for Cybersecurity, provides a realistic and modern dataset for evaluating intrusion detection systems.

**Criteria for Success:**

A successful project will:

- Be efficient at predicting attacks, due to the nature of attacks happening fast. The model has to perform faster or at least just as fast.
- Have an accuracy of 90% or better. Whether predicting an attack or benign behavior, to avoid false alarms.
- Have a scale on attacks to determine the right course of action to take to stop or fix the problem. As well as damage migrations steps.

**Scope:**

- Data preprocessing (handling missing data, scaling, encoding, feature selection)
- Binary classification: Attack vs. Normal behavior
- Model comparison: Logistic Regression, Random Forest, XGBoost, Neural Networks
- Evaluation using confusion matrix, ROC-AUC, F1-score, coefficient of determination

**Out of Scope for now:**

- Real time detection system

- Continuous learning (online learning models)
- Deep packet inspection or raw network traffic capture

**Constraints:**

- The dataset is large (~2 million records), which may require efficient sampling or computing resources.
- Class imbalance, some attacks are much rarer than others which could create a bias in the model.
- Privacy and ethical considerations: The model must be trained on public data only.
- Time limitations for project
- Limited access to real-world network environments. (Deployment/testing)

| Stakeholders: | Role/ Interest: |
|---|---|
| Network Security Analysts | Use the model to flag suspicious network activities . |
| IT Infrastructure Teams | Integrate the model into monitoring systems for real-time alerts. |
| Management / Executives | Interested in reducing risk and improving security posture. |
| Data Scientist | Design, train, and validate the model; communicate findings. |

**Data Sources:**

Dataset: CIC-UNSW-NB15

Source Organizations:

- Canadian Institute for Cybersecurity (CIC)
- University of New South Wales (UNSW)

Citation:

H. Mohammadian, A. H. Lashkari, A. Ghorbani. "Poisoning and Evasion: Deep Learning-Based NIDS under Adversarial Attacks," 21st Annual International Conference on Privacy, Security and Trust (PST), 2024.

**Approach:**

1. Data Acquisition: download from the official CIC website.
2. Data Cleaning and Preprocessing, handle missing values, categorical encoding, and normalization
3. Exploratory Data Analysis (EDA): Analyze feature distributions. Visualize correlations and class imbalance.
4. Feature Selection: Use correlation analysis or feature importance from a tree-based model
5. Model Development: Train multiple models (Logistic Regression, Random Forest, XGBoost, Neural Network). Then use cross-validation for robust results.
6. Model Evaluation: confusion matrix, ROC-AUC, F1-score, coefficient of determination, accuracy, and precision. Analyze false positives/negatives
7. Interpretability: Use SHAP or feature impotence plots to explain the model.
8. Deliverables: Github Repository, Final project report, Slide deck