

Detecting Host Based intrusion using Machine Learning

Edwin Lin

AGENDA

- 01 Intro to IDE
- 02 Problem Definition
- 03 Host-Based dataset
- 04 Machine Learning Models
- 05 Application and results
- 06 Conclusion

intrusion detection system (IDS)

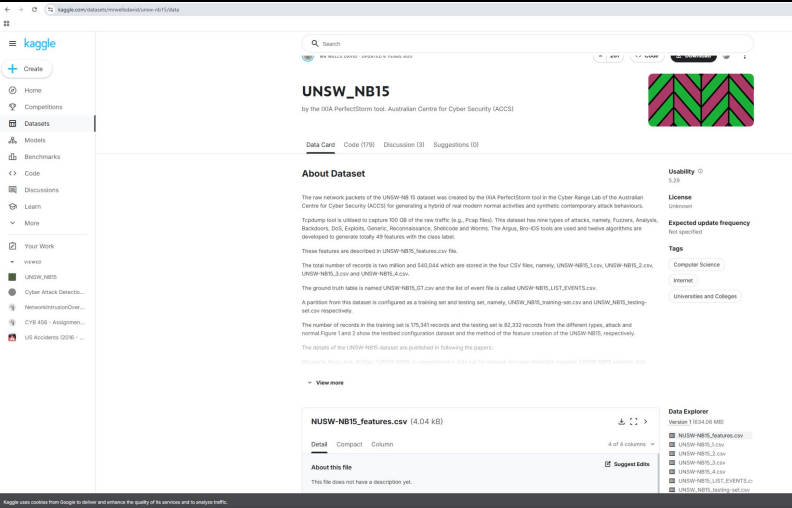
Misuse-based methods: detect known attacks by using signatures of those attacks

Anomaly-based techniques: model the normal network and system behavior and identify anomalies as deviations from normal behavior.

Hybrid techniques: combine misuse and anomaly detection

Problem Statement

Using a hybrid detection method on a UNSW_NB15 dataset for host based intrusion detection using machine learning



What is Machine Learning?

Machine Learning is like if you combine statistics, coding, and telling it to predict

Albert Einstein: Insanity Is Doing the Same Thing Over and Over Again and Expecting Different Results

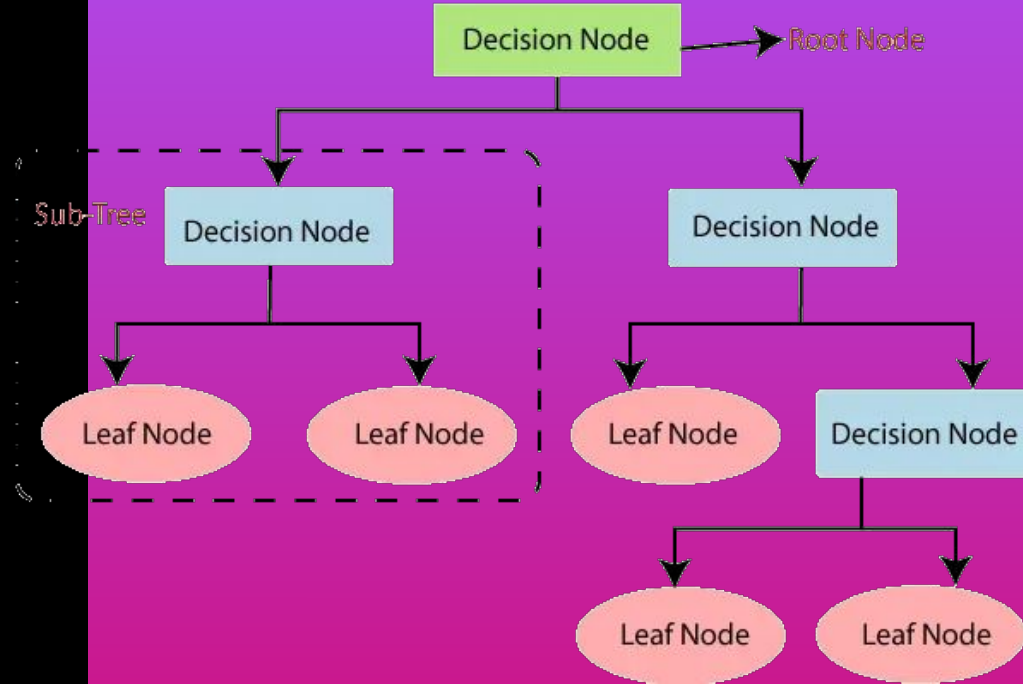
Machine learning:



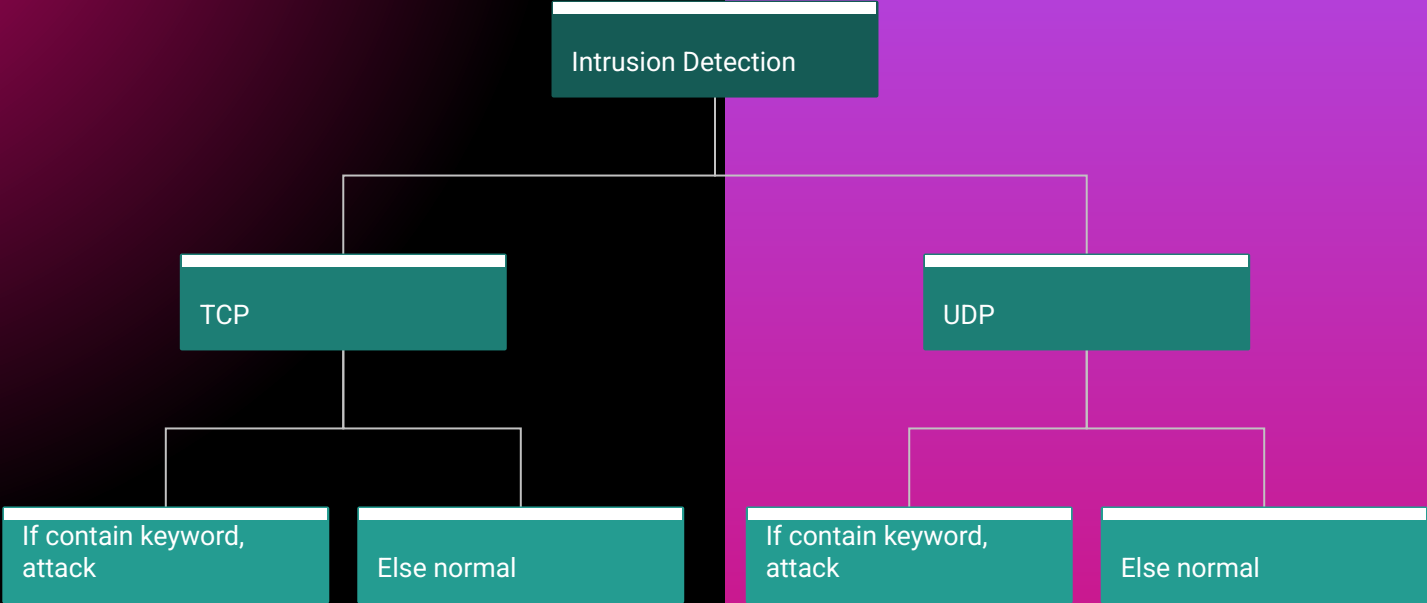
Decision Tree

- Our Misused based method
- Supervised Learning
- 2 Labeling options:
Normal = 0
Attack = 1

Each node level represents a feature, based on that feature, makes a decision



Example:



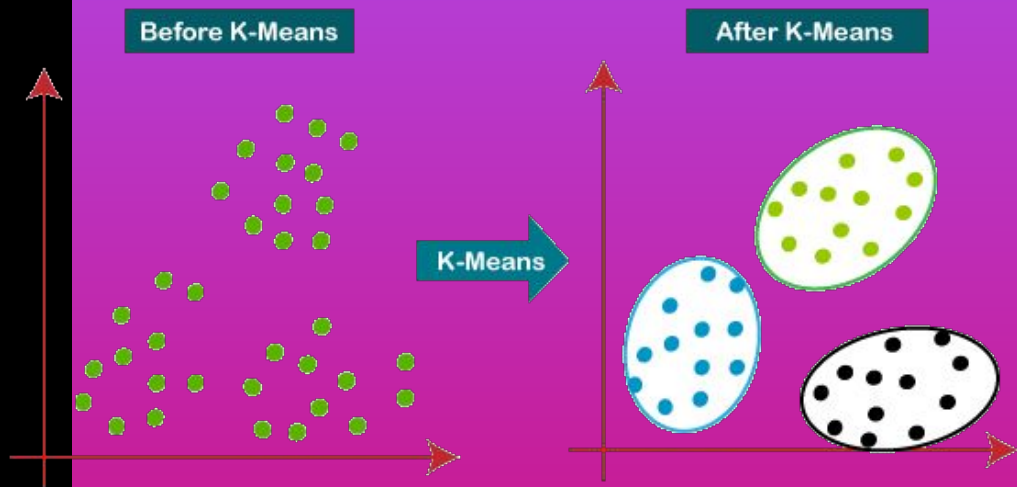
K-Means clustering

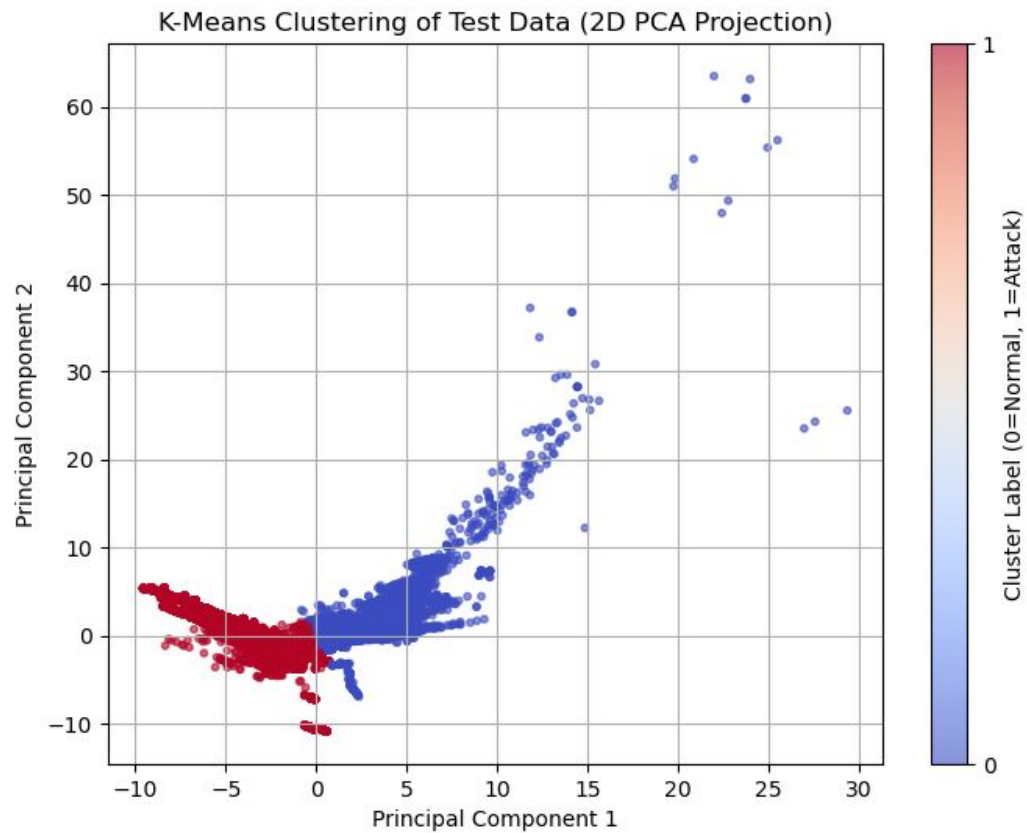
Anomaly-based approach

Unsupervised Learning

Find interesting patterns for
Anomalies

Centroid clusters



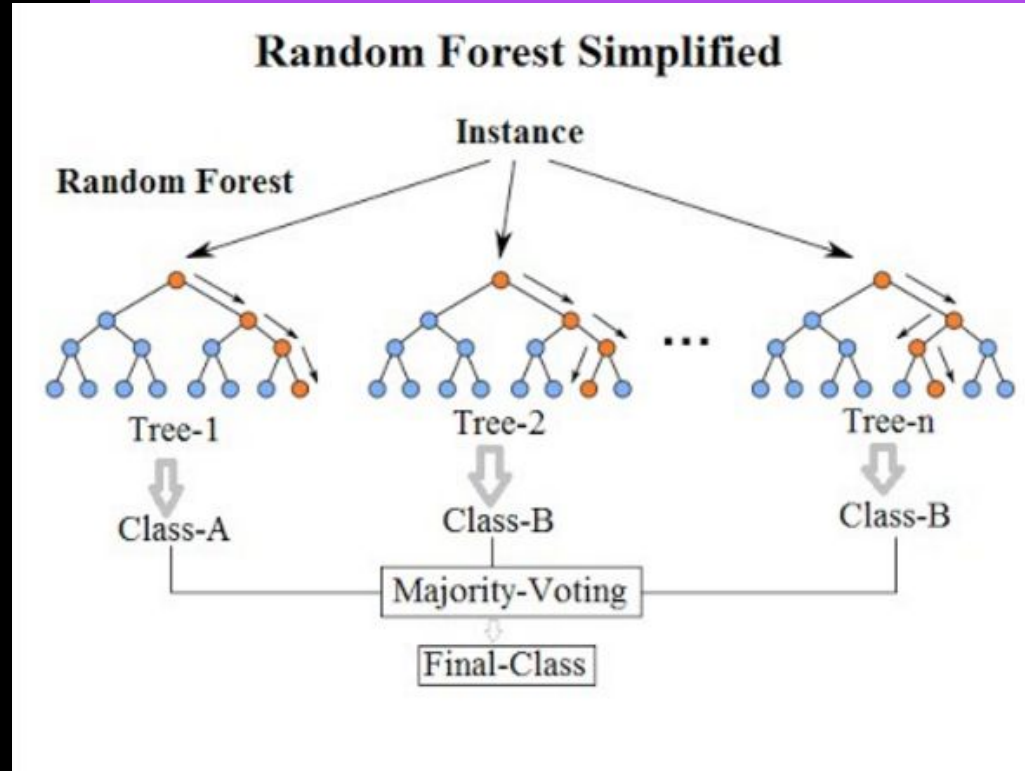


Random Forest

Hybrid approach

Multiple decision trees

Ensembled ML model



Performance Analysis

K means performance the worst as expected

High False Positives as expected for normal

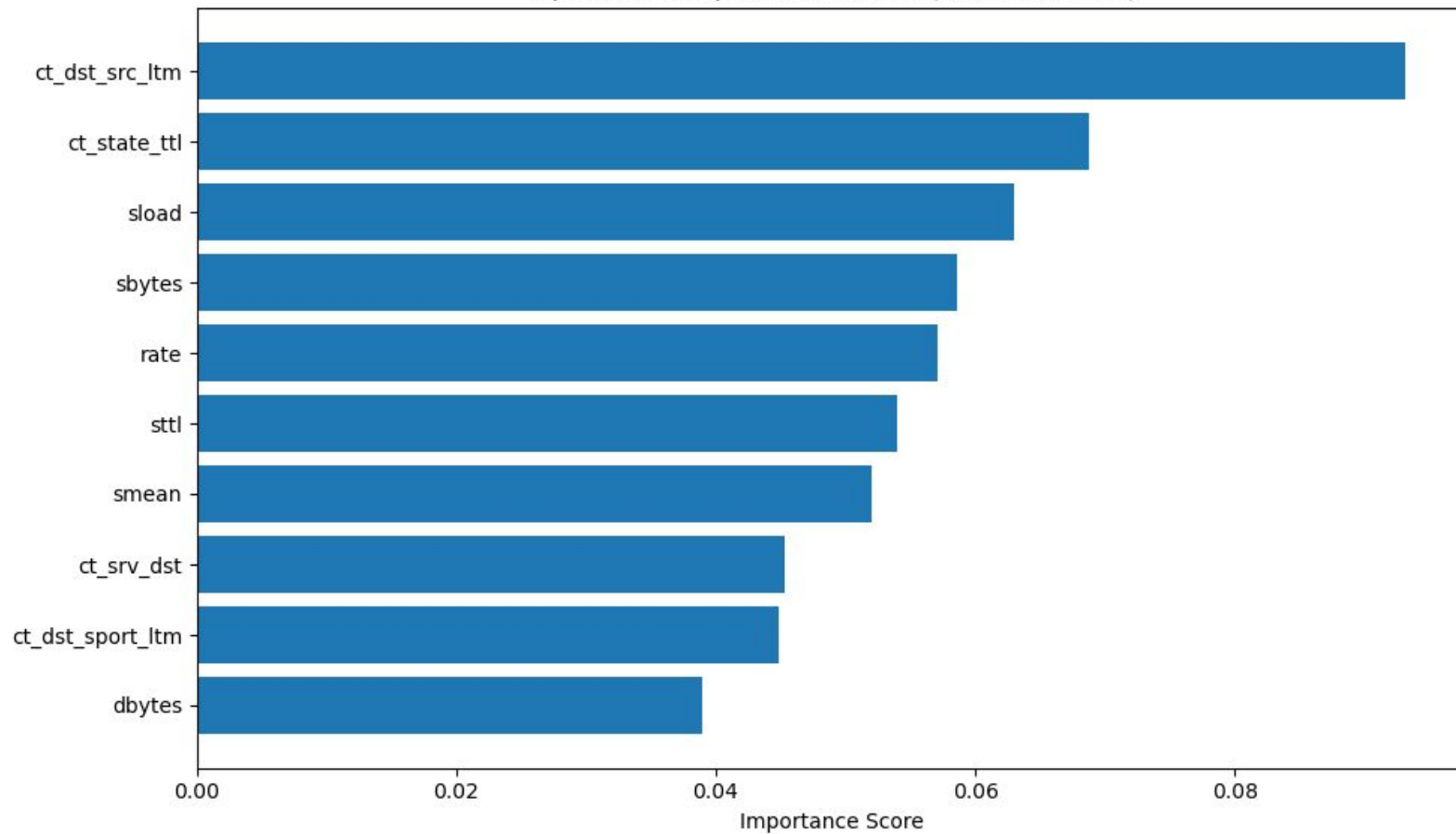
99% accuracy in finding attacks
Great performance.

Decision Tree Report:					
	precision	recall	f1-score	support	
0	0.77	0.97	0.86	56000	
1	0.98	0.86	0.92	119341	
accuracy			0.90	175341	
macro avg	0.88	0.91	0.89	175341	
weighted avg	0.91	0.90	0.90	175341	

K-Means Report (Anomaly Detection):				
	precision	recall	f1-score	support
0	0.54	0.88	0.67	56000
1	0.92	0.65	0.76	119341
accuracy			0.72	175341
macro avg	0.73	0.77	0.72	175341
weighted avg	0.80	0.72	0.73	175341

Random Forest Report (Hybrid Detection):				
	precision	recall	f1-score	support
0	0.77	0.98	0.86	56000
1	0.99	0.86	0.92	119341
accuracy			0.90	175341
macro avg	0.88	0.92	0.89	175341
weighted avg	0.92	0.90	0.90	175341

Top 10 Most Important Features (Random Forest)



THANK YOU

Any questions?