

# Characterizing and Modeling Threat Feeds for Patch Management

Ashton Woiwood, University of Iowa

PI: Zubair Shafiq

## To patch or not to patch?

In consumer-grade systems if a new security patch is available most users have very little reason not to install a new patch. However, for mission critical systems, new security patches can put a significant strain of resources and safety due to potential downtime. Thus, it is important to know which patches need to be deployed immediately and which can wait.

To decide this, one must quantify the usefulness of various threat indicators from a multitude of threat feeds. Organizations spend millions of dollars to purchase threat information (Tounsi, Wiem, et al). These “threat feeds” may contain valuable information pertaining to future exploits, however processing the data to find this valuable information is costly and time consuming.

### Problem Statement :

- Threat intelligence is often confused with threat data.
  - Threat intelligence is the result of time consuming data mining, analysis, and how it is operationalized.
  - Threat feeds usually provide the data and the organization must figure out how to use it.
- Threat feeds differ in many aspects:
  - temporal delay, exclusivity, time scale, quality of textual description, domain specialization (Thomas, et al)
- How should an organization perform a comparative analysis of threat feeds to assess their quality in terms of exploit information?

### Related Research:

Researchers have compared the quality of data from different threat feeds (Li, Vector Guo, et al). Their key conclusions are:

- Threat Intelligence (TI) feeds collect data in ways that vary from each TI vendor. Most of the data collected did not have a clear type of category and was labeled with ambiguous terms such as “Malicious”.
- There is a low overlap between threat feeds, and the methods that TI vendors use may cause a sampling bias depending what the TI vendor is looking for.
- Many exploit events are experienced differently amongst different sectors. So an organization must consider what threats may be more relevant to their security policies.

### Methods:

We use a supervised machine learning approach to identify whether threat events contain relevant exploit.

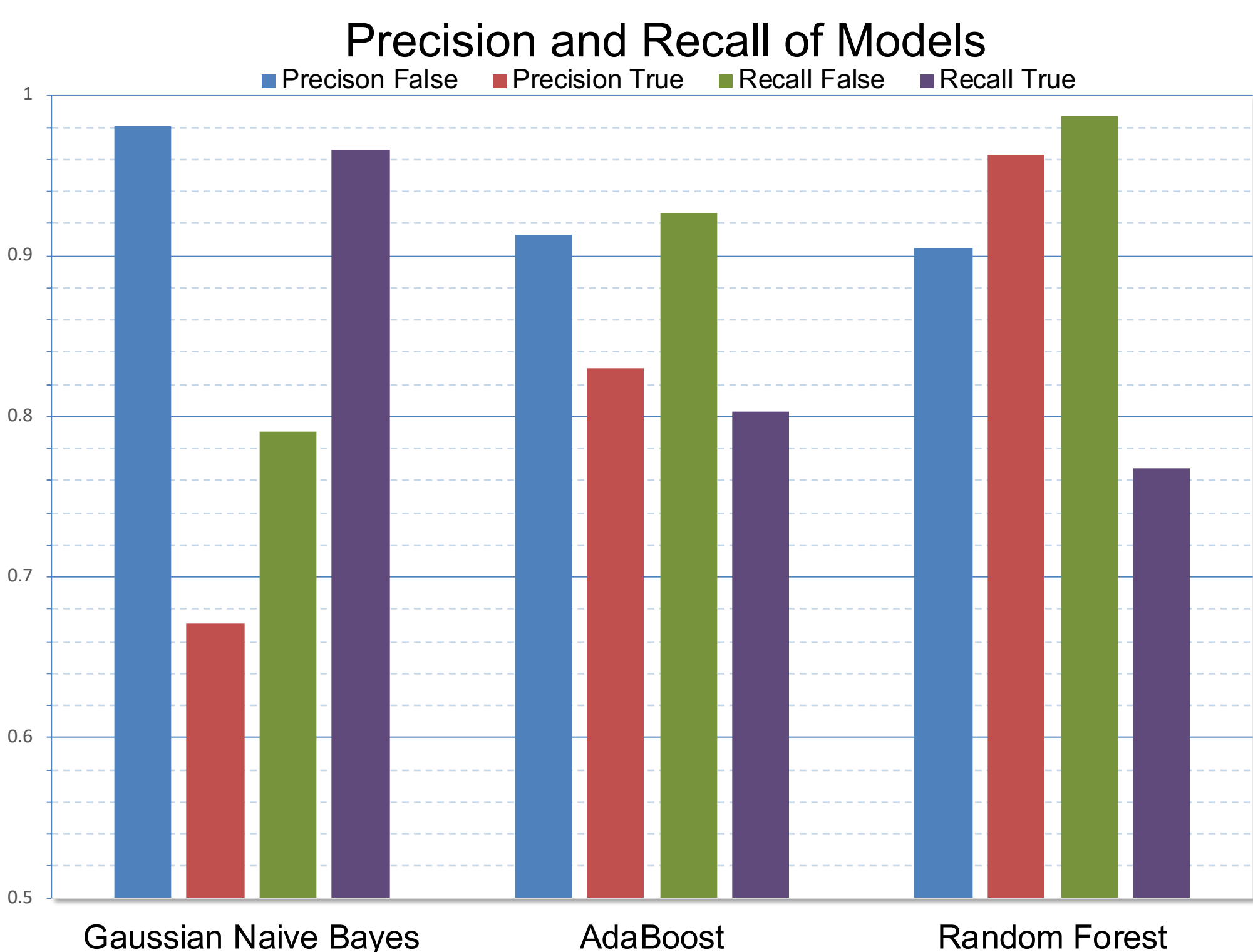
**Ground Truth:** We manually established ground truth as exploit or not by analyzing the textual description of a total of 4000 threat events.

**Feature Selection:** We extracted 20 million unique word tokens as features from the textual descriptions. We then performed unsupervised feature selection using variance threshold.

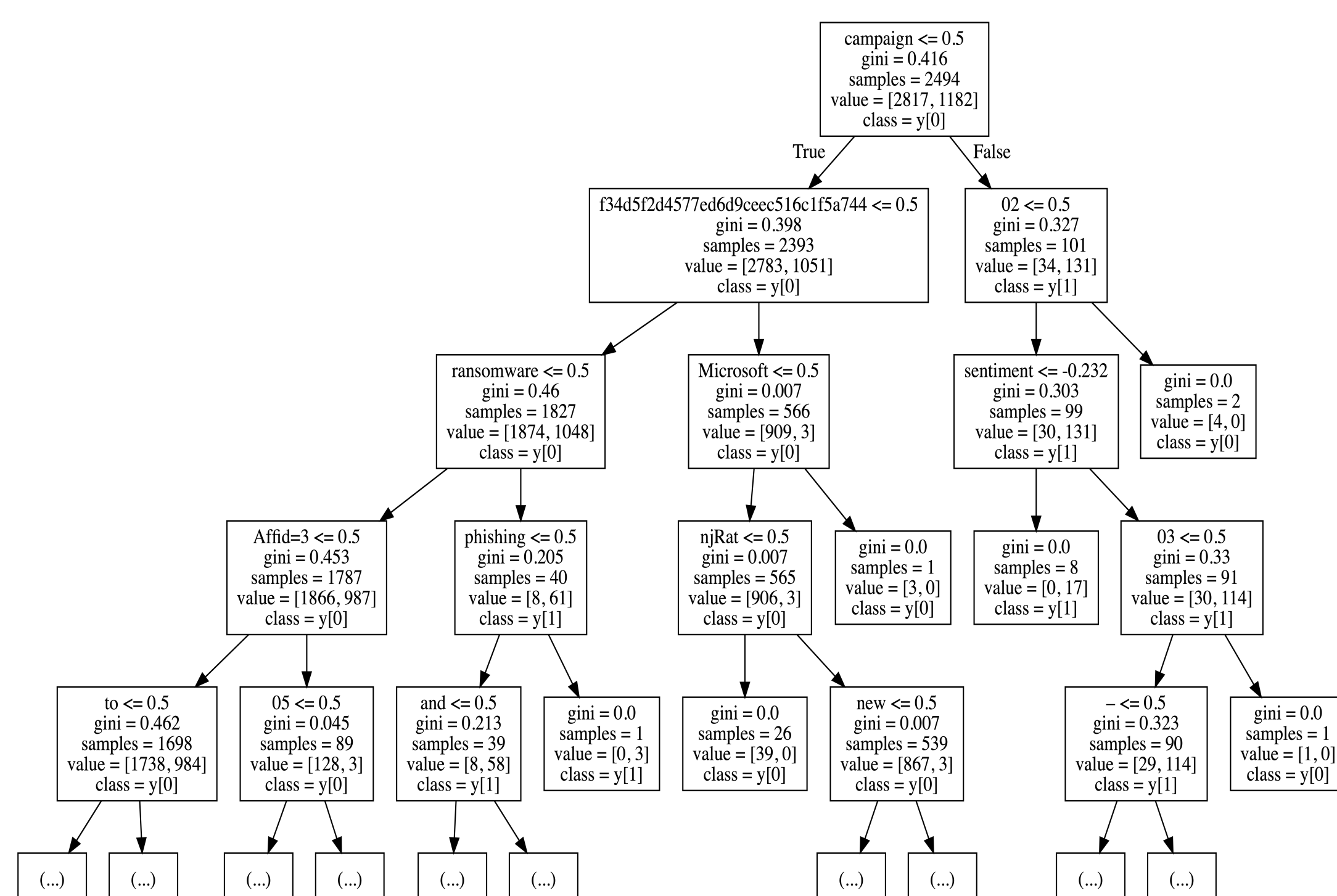
**Model Training:** We then trained 3 supervised machine learning models, Gaussian Naïve Bayes (GNB), Random Forest, and AdaBoost.

**Results:** 10-fold cross validation accuracy

- Gaussian Naïve Bayes 75.47%
- Random Forest: 78.69%
- AdaBoost: 82.29%



Example Decision Tree Model



### Future work and Impact:

- Improving prediction accuracy using better features (e.g. TF-IDF or word embeddings) and deep neural network classifiers.
- Comparison of different threat feeds in terms of which contain more/less information for specific exploit events.
- Use trained models to forecast risk for security exploits to better inform patch management
- Our research will help organizations make more informed data-driven security decisions

Works cited: THOMAS, K., AMIRA, R., BEN-YOASH, A., FOLGER, O., HARDON, A., BERGER, A., BURSZEIN, E., ANDBAILEY, M. The abuse sharing economy: Understanding the limits of threat exchanges. International Symposium on Research in Attacks, Intrusions, and Defenses(2016), Springer.

Tounsi, Wiem, Threat intelligence market analysis by solution, services, deployment, application and segment forecast, 2018 - 2025. <https://www.grandviewresearch.com/industry-analysis/threat-intelligence-market>.

