# Finding the Needle in the Haystack

## Anomaly Detection in the Cybersecurity Industry

## AT A GLANCE

Capstone Company: Rapid7 Location: Boston, USA

Problem Statement: Develop and implement an algorithm that reduces the amount of time necessary for analysts to detect intrusions of bad actors in client networks and facilitates the detection of new intrusions not discovered before.

Timeline & Milestor	nes					
February	March	April	May	June	July	August
Feb 6: Initial kick off meeting structure	하나 있다. 이 경우 전에 되면 보다는 것이 나는 없는 것은 전에 되면 보다는 것이다.	pr 3: Exploratory resentation due	May 8: Phase 1 prese report, summer work			5: Final report, Showcase tion, Project poster due
Feb 20: Project charter and workplan due Mar	<b>4</b> : Operational ata usage	<b>Apr 25</b> : Unders Rapid7 presen	standing analytics @ tation due		<b>Jul 14</b> : Mid-summer project status due	Aug 23: Final Capstone Presentation

#### The Team:

MBAn



Jocelyn Johnny Beauchesne

Roy

Hodgman

Rapid7

Vasudha Shivamoggi

**Faculty Advisors** 



Rahul Mazumder

Hussein Hazimeh

### IMPACT

1. Direct Labor Savings 2. Avoidance Savings

\$1M+

\$35M+

3. Novel ML Tools

## **Impact to Security Analysts:**

Created machine learning models to classify 90% of the data as "normal", significantly reducing manual review of client network data at a projected \$1M+

#### Impact to Clients:

By laying the foundation for modern machine learning in cybersecurity and driving initial findings, Rapid7 and their clients can avoid costs of at least \$35M+ annually

## Impact to Rapid7:

Packaged flexible, scalable, and auto-tuned machine learning pipeline to be used on any data sets for future anomaly detection use cases

patentable research

## DATASET

Rapid7 deploys "hunts" on client computer networks, which are deep downloads of computer behavior in a two week period; we used this data to conduct our anomaly detection analysis:

#### **Hunt Data Statistics**

- 4TB+ of data from 300+ clients
- Unstructured raw data
- 100% unlabeled without prior examples of intrusions

## TWO-PART METHODOLOGY

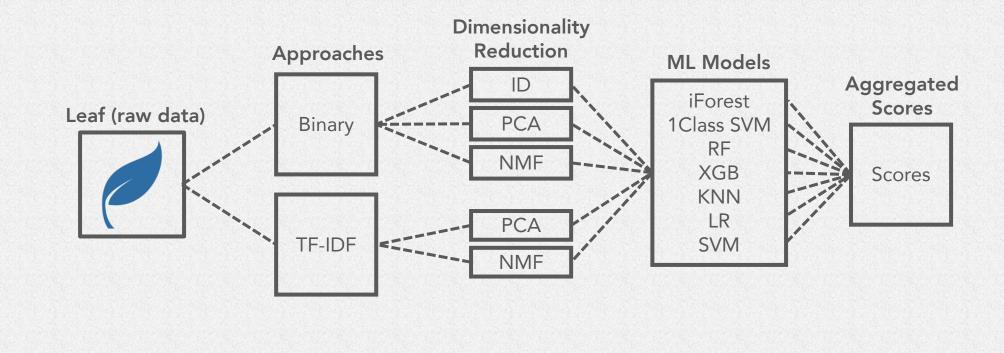
Tree Approach to Contextualize Data

To contextualize data, we used paths to create a tree grouping similar processes together into leaves

Path	Command Line	Signing Details	MD5	Modules	Network	
chrome.exe	option = private navigation	valid	2jljljiwldk24830	<ol> <li>Network Driver</li> <li>GIF Loader</li> <li>Ad Blocker</li> </ol>	IP[country = Russia]]= 12.35.534.34	
1 Tree Structure Using Paths			2 Confidence Interval With Outlier			
				x x x x	X	

## **Machine Learning Pipeline**

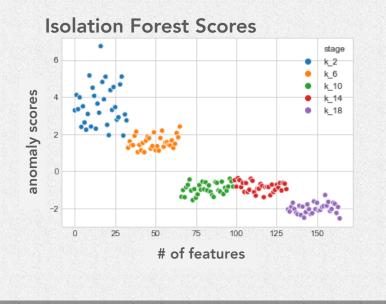
Within each leaf, we processed the data to reduce dimensionality and executed machine learning models automatically tuned with Bayesian optimization to identify anomalies and designate anomaly scores

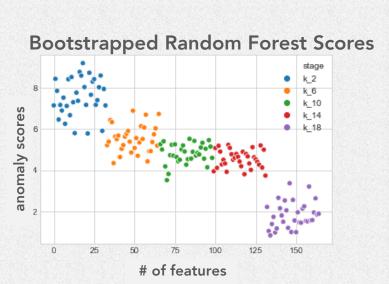


## KEY RESULTS

#### Synthetic Data

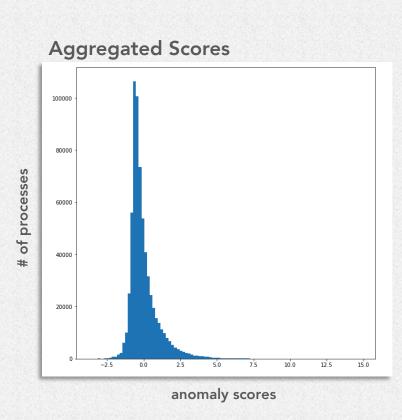
To test model efficacy, we ran multiple analyses on synthetic data; 85% accuracy on synthetic anomalies





#### **Client Data**

The anomaly distribution score is narrow: 95% of client data lies within  $\pm 2.25\sigma$  of the mean



#### **Examples of Anomalous Activities**



Unregistered remote control software installed on user files



#### INTERPRETABILITY

In order to add a layer of interpretability to our models, we devised two simple approaches to tie feature importance to our anomaly scores:

- Interpretable regression model against anomaly scores on the original features
- Individual feature anomaly scores for each process

#### FUTURE DIRECTION

- Confirm results with additional analyses using synthetic data
- Develop cross feature interpretability
- Apply machine learning pipeline on other datasets for future research
- Create database of labeled intrusions and hacks for improved machine learning
- Continue developing key relationships with security analysts for feedback