# Classifying Computer Processes in the DARPA OpTC dataset

Andrew Veal aveal@acm.org

XCS229ii – 003 Project

Introduce yourself:

Over 30 years in development/research/leadership in UK Government

Head of Data Mining Research, Head of Innovation, Future Technology Officer

Programme Committee member of ACM KDD Gov and Industry track 2010-2012

Evangelist for Andrew Ng's courses since the very first run of Stanford ML MOOC

# AI for Cybersecurity
## Hype or Reality?

The New York Times | https://nyti.ms/352Bp5W

### As Understanding of Russian Hacking Grows, So Does Alarm

Those behind the widespread intrusion into government and corporate networks exploited seams in U.S. defenses and gave away nothing to American monitoring of their systems.

By David E. Sanger, Nicole Perlroth and Julian E. Barnes

Published Jan. 2, 2021 Updated Jan. 5, 2021

On Election Day, General Paul M. Nakasone, the nation's top cyberwarrior, reported that the battle against Russian interference in the presidential campaign had posted major successes and exposed the other side's online weapons, tools and tradecraft.

"We've broadened our operations and feel very good where we're at right now," he told journalists.

Eight weeks later, General Nakasone and other American officials responsible for cybersecurity are now consumed by what they missed for at least nine months: a hacking, now believed to have affected upward of 250 federal agencies and businesses, that Russia aimed not at the election system but at the rest of the United States government and many large American corporations.

Three weeks after the intrusion came to light, American officials are still trying to understand whether what the Russians pulled off was simply an espionage operation inside the systems of the American bureaucracy or something more sinister, inserting "backdoor" access into government agencies, major corporations, the electric grid and laboratories developing and transporting new generations of nuclear weapons.

At a minimum it has set off alarms about the vulnerability of government and private sector networks in the United States to attack and raised questions about how and why the nation's cyberdefenses failed so spectacularly.

Those questions have taken on particular urgency given that the breach was not detected by any of the government agencies that share responsibility for cyberdefense — the military's Cyber Command and the National Security Agency, both of which are run by General Nakasone, and the Department of Homeland Security — but by a private cybersecurity company, FireEye.

**Andrew Veal**
16 Jan

AI better get a whole lot smarter ... If this prediction is to come true ...

ZDNET.COM
AI set to replace humans in cybersecurity by 2030, says Tr...

Like     Comment     Share

---

* "As Understanding of Russian Hacking Grows, So Does Alarm – Those behind the widespread intrusion into government and corporate networks exploited seams in U.S. defenses and gave away nothing to Americans monitoring of their systems."
The New York Times - January 2, 2021  - David Sanger, Nicole Perlroth, Julian Barnes
https://www.nytimes.com/2021/01/02/us/politics/russian-hacking-government.html

* NCSC Cyber UK ONLINE: "Cyber Threat: Oh that was clever! When even jaded incident responders are impressed"
NCSC's Tech Director for Incident Management for a tour of some of the interesting technical aspects that have exercised (and perhaps grudgingly impressed) our incident management team over the last year. Unsurprisingly, there will be a lot of discussion of UNC2452 this year.
https://youtu.be/ppXOt8f5H8Q

• The NCSC, CISA, FBI and NSA publish advice on detection and mitigation of SVR activity following the attribution of the SolarWinds compromise.
https://www.ncsc.gov.uk/news/joint-advisory-further-ttps-associated-with-svr-cyber-actors

**DARPA OpTC dataset** – is this the new "MNIST" benchmark dataset for cybersecurity researchers?

- DARPA OpTC dataset is a new open-source dataset that appears to have the requisite scale, richness and class imbalance to drive AI research.
- Can we take a high-level summary view of process activity, aggregating frequency counts of (object, action) events to form feature vectors for Machine Learning?
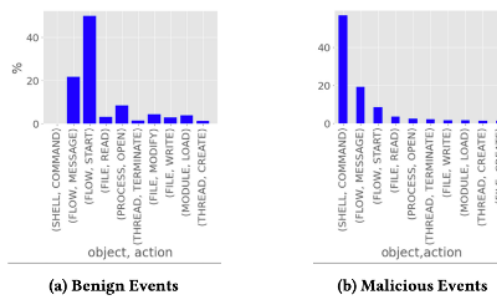
(a) Benign Events   (b) Malicious Events

Figure 2: Distribution of Benign and Malicious events

---

 *"The lack of diverse and useful data sets for cyber security research continues to play a profound and limiting role within the relevant research communities and their resulting published research"*

Melissa J. M. Turcotte, Alexander D. Kent and Curtis Hash. 2018. Unified Host and Network Data Set. In *Data Science for Cyber-Security*, Chapter 1 (November 2018), 1-22. World Scientific DOI: https://doi.org/10.1142/9781786345646_001

Scale – 17 billion events
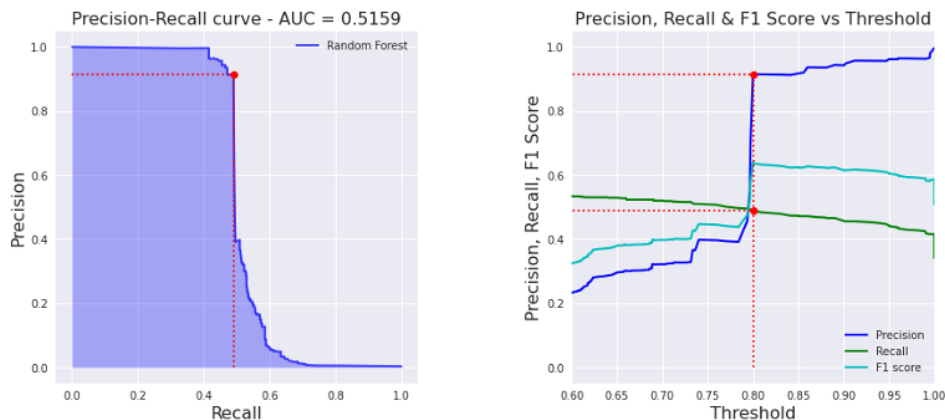Extreme class imbalance – 0.0016% of the events are malicious
Limited variety of attacks – the red team used related modus operandii (powershell)

Figure 2 reproduced from:

Md. Monowar Anjum, Shahrear Iqbal and Benoit Hamelin. 2021. Analysing the Usefulness of the DARPA OpTC Dataset in Cyber Threat Detection Research. arXiv:2103.03080v2. Retrieved from https://arxiv.org/abs/2103.03080

Accepted for *ACM Symposium on Access Control Models and Technologies (SACMAT)*, 16-18 June, 2021, Barcelona, Spain [virtual event]. ACM Inc., New York, NY. DOI: https://doi.org/10.1145/3450569.3463573

## Can we classify processes as benign or malicious using a Random Forest?

Our core hypothesis is that we can distinguish between malicious and benign processes using the frequency count of the (object, action) events associated with each process as a feature vector. If we examine the figures above we can see that we can recover 50% of the malicious processes (Recall) and of those 50%, 90% are true positives (Precision) without any tuning of hyperparameters – that suggests our hypothesis is reasonable.

In addition, by looking at the feature importance (not shown) we can identify the event types which are more important for the classification of malicious processes. This is important – decision trees provide explanatory rules and the results are explainable.

Why decision trees are likely to be good – they produce rules – if you do some Exploratory Data Analysis (EDA) you will see that the following code only selects malicious examples from the dataframe df – it selects 12 malicious processes:

```
> df[df['SHELL_COMMAND_'] > 2800][['label']].apply(lambda x: Counter(x))
{1: 12}
```

## Discussion – summary of progress so far

So far, we have:
- Created the ML dataset, done basic EDA and most of model selection
- Implemented a *simple* train/validation/test split strategy
- Got results for linear models (LR and LinearSVM) and Random Forest

We plan to:
- Try XGBoost, KNN and do hyperparameter tuning on selected model
- Examine misclassification errors (spoiler alert: low counts -> errors)
- Pay special attention to the choice of training, validation and test sets

Issues with our experimental protocol are:

- We are taking a summary view of process activity, aggregating frequency counts of high level (object, action) events only. We are not using the full richness of the dataset: each (object, action) event has detailed properties. For example, PROCESS.CREATE event table contains meta-data for timestamp, user, image path and command line.
- We need to pay special attention to the choice of training, validation and test sets, so that the distributions (as far as possible) reflect the data we expect to get in the future.

The real research challenge in the final week is to split the dataset in a way that avoids leakage of information from training into validation and test sets.

Andrew Ng. 2018. Machine Learning Yearning: Technical Strategy for AI Engineers, In the Era of Deep Learning. Draft Version. Retrieved May 6, 2021 from https://www.deeplearning.ai/programs/

"Colonial cyber attack is a warning of worse to come - A plague of ransomware will continue in every sector until the superpowers step in"

Hacking group tied to cyber attack on US pipeline said to have shut down

DarkSide's closure follows $5m ransom paid in bitcoin by Colonial Pipeline Company

A tank farm connected to the Colonial pipeline, which shut down for days after hackers invaded its information systems © Drew Angerer/Getty Images

Hannah Murphy in San Francisco, Myles McCormick in New York and Katrina Manson in Washington YESTERDAY

FTWeekend 15 May/16 May 2021

[1] Opinion – page 11 of print edition – "Colonial cyber attack is a warning of worse to come – A plague of ransomware will continue in every sector until the superpowers step in" by Misha Glenny

'The cyber attack on Colonial Pipeline, which transports 45 per cent of oil consumed on the east coast of the US, should be the event that finally wakes everyone up."

[2] Cyber attacks. Criminal gangs – page 4 of print edition – "Ransomware hackers stay one step ahead – Experts and governments debate best way to fight back after victims pay out billions" by Hannah Murphy

"Cybersecurity experts like to joke that the hackers who have turned ransomware into a multibillion-dollar industry are often more professional than even their biggest victims."