# Using Machines to Exploit Machines

## Harnessing AI to Accelerate Exploitation

Guy Barnhart-Magen
Ezra Caltum

@barnhartguy    @aCaltum

# Legal Notice and Disclaimers

~~@barnhartguy~~   @aCaltum

# $ ID

**Guy Barnhart-Magen**

@barnhartguy

BSidesTLV Chairman and CTF Lead



**Ezra Caltum**

@acaltum

BSidesTLV Co-Founder

DC9723 Lead

BSIDES LAS VEGAS

CRYPTO + PRIVACY VILLAGE

44CON

t2.fi INFOSEC

LKSAT SKYTA TALKS ATSKY

# OUR PROBLEM

**Fuzz Testing**
Literally thousands
of crashes to analyze

(good problem to
have?)

1

# OUR PROBLEM

**Automation**

Might miss something
important, but helps
reduce from thousands
to hundreds of results

2

1

**Fuzz Testing**
Literally thousands
of crashes to analyze

(good problem to
have?)

@barnhartguy    @aCaltum

# OUR PROBLEM

**Manual Analysis**
Can only do a limited
amount with limited
researchers time

**3**

**Automation**

Might miss something
important, but helps
reduce from thousands
to hundreds of results

**2**

**Fuzz Testing**
Literally thousands
of crashes to analyze

**1**

(good problem to
have?)

@barnhartguy    @aCaltum

# EFFORT BALANCE

Build the
Model

@barnhartguy     @aCaltum

# EFFORT BALANCE



Gather Data

Build the
Model

@barnhartguy    @aCaltum

# EFFORT BALANCE

Keep <u>Good</u> Data

Build the Model

@barnhartguy    @aCaltum

# PROBLEM STATEMENT

# PROBLEM STATEMENT

## What is Australia?

# PROBLEM STATEMENT

Can we create an ML model that can triage crashes and help us focus on the exploitable ones?

(we got a lot of crashes from AFL)

@barnhartguy    @aCaltum

# REVISED PROBLEM STATEMENT

Can we create an ML model that can outperform <u>exploitable</u>, based on the same data?

it should perform at least as well as

<u>exploitable</u>

# FULL DISCLOSURE

Limited dataset - but we tried anyway (no DL today)

We want to focus on the methodology

We can't trust this results, but they are worth sharing

@barnhartguy    @aCaltum

# MACHINE LEARNING

See our previous talks on hacking machine learning systems
:-)

# MACHINE LEARNING

What it isn't:

- **Magic**
- A solution to every problem
- Difficult or Complex
- One of the holy VC buzzwords:
  - Blockchain
  - Cyber
  - Zero Trust

@barnhartguy     @aCaltum

# THE DIFFERENCE BETWEEN ML AND AI

If it is written in Python, it's probably Machine Learning

If it is written in PowerPoint, it's probably AI

@barnhartguy    @aCaltum

# EXAMPLE

EXAMPLE

Using Machines to Exploit Machines

Harnessing **AI** to Accelerate Exploitation

# Everyone Confuses

# "AI" with "ML"

So do We
Sorry

# WHAT IS IT GOOD FOR?

Finding patterns in a lot of data, patterns you did not expect (counter intuitive)

# WHAT IS IT GOOD FOR?

Finding patterns in a lot of data, patterns you did not expect (counter intuitive)

Correlating different inputs you suspect are related somehow

# WHAT IS IT GOOD FOR?

Finding patterns in a lot of data, patterns you did not expect (counter intuitive)

Correlating different inputs you suspect are related somehow

Abstracting a problem and throwing it at an algorithm, hoping for the best (e.g. being lazy)

@barnhartguy    @aCaltum

# PREDICTIONS

ML makes predictions based on previously seen data

Your data quality is important! (data is not information)

# WHAT DO YOU GET?

How is this new sample I am testing now similar to all the other samples I've seen in the past?

Testing - extracting and then comparing features against your model

@barnhartguy    @aCaltum

# Windows

An error has occured.We dont even know what is it.So cant fix it ,
and you have restart your computer.By the way if you restart your ,
computer you will ▮▮▮▮▮▮▮▮▮▮▮▮ in all open applications.
In the other hand ▮▮▮▮▮▮▮▮▮▮ on :)

Press Enter to return to Windows(It wont work), or

Press CTRL+ALT+DEL to restart your computer.

Error : 0E : 016F : BFF9B3D4

Crash Triage

# A COMMON MORNING IN MY LIFE

- I start a fuzzing process overnight and go home

# A COMMON MORNING IN MY LIFE

- I start a fuzzing process overnight and go home

- At first light in the morning (11:00) I drink a cup of coffee

@barnhartguy     @aCaltum

# A COMMON MORNING IN MY LIFE

- I start a fuzzing process overnight and go home

- At first light in the morning (11:00) I drink a cup of coffee

- I analyze the data from the crash dump with the help of a debugger

# A COMMON MORNING IN MY LIFE

- I start a fuzzing process overnight and go home

- At first light in the morning (11:00) I drink a cup of coffee

- I analyze the data from the crash dump with the help of a debugger

- Based on my experience, and the output of some plugins, I classify the crashes as either exploitable or not

@barnhartguy     @aCaltum

# A COMMON MORNING IN MY LIFE

- I start a fuzzing process overnight and go home

- At first light in the morning (11:00) I drink a cup of coffee

- I analyze the data from the crash dump with the help of a debugger

- Based on my experience, and the output of some plugins, I classify the crashes as either exploitable or not

- I start developing a POC for the exploitable crashes.

@barnhartguy    @aCaltum

# A COMMON MORNING IN MY LIFE

- I start a fuzzing process overnight and go home
- At first light in the morning (11:00) I drink a cup of coffee
- I analyze the data from the crash dump with the help of a debugger
- Based on my experience, and the output of some plugins, I classify the crashes as either exploitable or not
- I start developing a POC for the exploitable crashes.

➔ No need for sleep for our AI overlords

➔ No need for coffee for our AI overlords

➔ Preprocessing phase prepares the data for the ML analysis

➔ ML analyzes the data, based on its experience (training data), emits predictions (human intuition or heuristics)

➔ Human minions will develop a PoC for the overlords

@barnhartguy    @aCaltum

Our Data Set

# DARPA CYBER GRAND CHALLENGE

We have 632 test cases that we know are exploitable

We ran `exploitable` against them and got:

- 607 were definitely exploitable
- 12 were probably exploitable
- 13 were unknown - the tool couldn't reach a decision

# SO, WHAT DOES A CRASH GIVE US?

**EAX, EBX, ECX, EDX** - general purpose (values, addresses)

**ESP, EBP** - Stack pointers

**ESI, EDI** - Source and Destination Index (for string operations)

**EIP** - Instruction pointer

**eflags** - metadata (wasn't actually useful at all, empty values)

**CS, SS, DS, ES, FS, GS** - Segment registers

Also a whole lot of other things which we didn't look at

@barnhartguy    @aCaltum

# OUR PROCESS

| Creating Crashes | Crash Analysis | Feature Extracting |
|---|---|---|

**Running tests** against a ~600 programs with known crashes, collecting the crash dumps

**Analyzing the crash dumps** using `exploitable`, collecting the **stack** and **register values**

Converting the data collected from the `exploitable` output to a canonical representation, **extracting the features** we cared about

@barnhartguy     @aCaltum

# PROBLEM

Register values are discrete and unrelated to each other

What can we learn from specific register values?



@barnhartguy     @aCaltum

# CLASSIFYING DATA

We tried breaking the values of the registers into three groups:

- High address range (kernel)
- Low address range (userland)
- Values

Bad results - data distribution not uniform :-(

# BINNING

Dividing the values to evenly spaced bins

10 bins total, evenly distributed between [min_val, max_val]

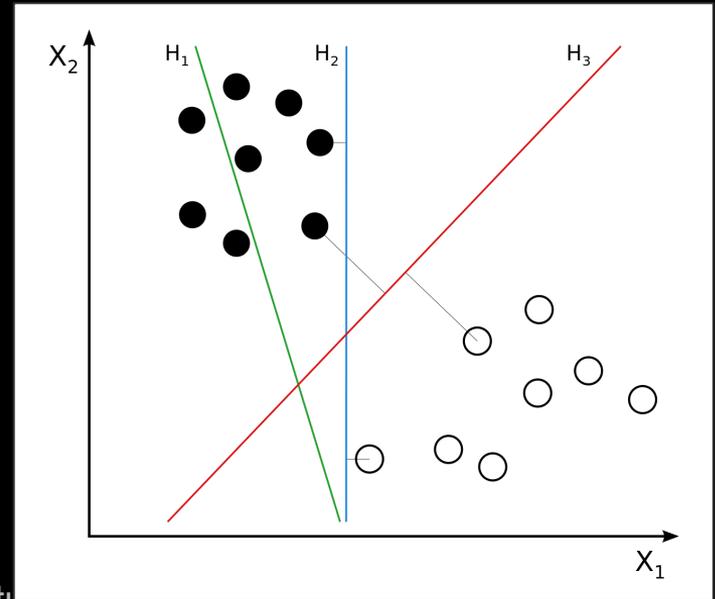This helps the model ignore specific values, and look at them as ranges

Good results :-)

# OneClassSVM

Train your major class (609 records, EXPLOITABLE)

Test your data against similarity to the model {-1,1}

+1 = very similar to the model

-1 = very not similar to the model



@barnhartguy    @aCaltum

# RESULTS - OneClassSVM

Anomaly detection using OneClassSVM: 23 records (from 25) are successfully recognized as belong to "exploit" class

- 23 records recognized as a major class:
  - 13 records previously labeled as "unknown"
  - 10 records previously labeled as "probably exploitable"
- 2 "probably exploitable" records identified as outliers

| Class | 1ClsSVM |
|---|---|
| Exploitable | +23 |
| Probably Exploitable | 2 |
| Unknown | |

@barnhartguy     @aCaltum

# COSINE SIMILARITY

- Cluster our data (609 records, EXPLOITABLE)
- Measure similarity between each data point (24 records) to the cluster
- We also used binning and not the actual register values



@barnhartguy    @aCaltum

# RESULTS - Cosine Similarity

We tried comparing using linear or centroid methods

Started with 9 register values, then adding the rest (15 register values, using binning)

~65% using values of 9 registers

~87% using values of 15 discretized registers

| Class | CosSim Linear | CosSim Centroid |
|---|---|---|
| Exploitable | +16 | +22 |
| Probably Exploitable | | |
| Unknown | | |

@barnhartguy    @aCaltum

# XGBoost

"Tree" that is built using the most contributing features

Very easy to explain how decisions are made, good for insights

Select 80% of the data (evenly sample from each group) for training, 20% for testing

# RESULTS - XGBoost
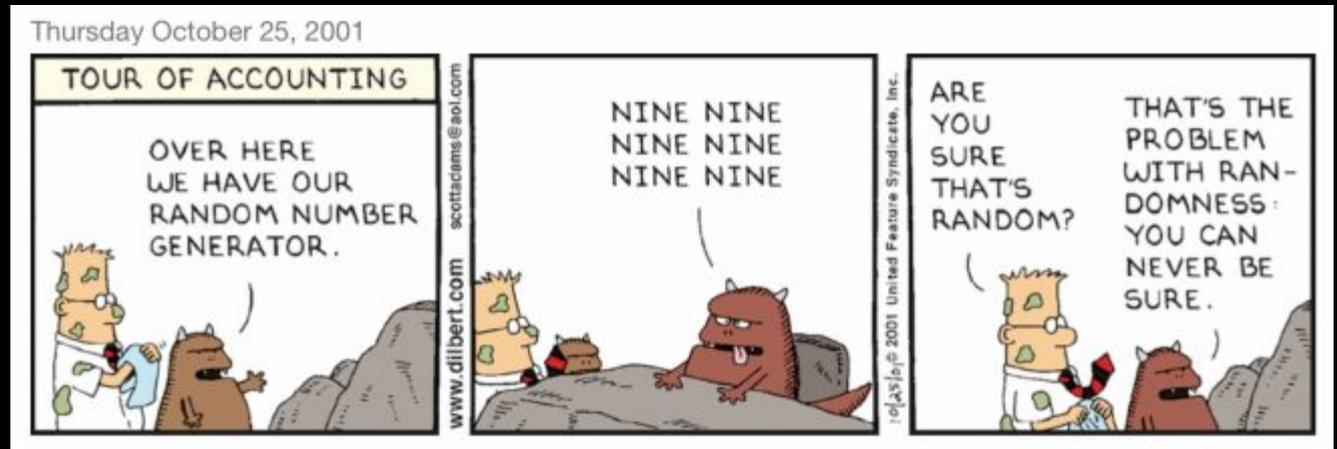
95-99% accuracy

# RESULTS - XGBoost

95-99% accuracy

This is not very good, you can get very high success rate guessing EXPLOITABLE all the time - be correct 96% of your guesses

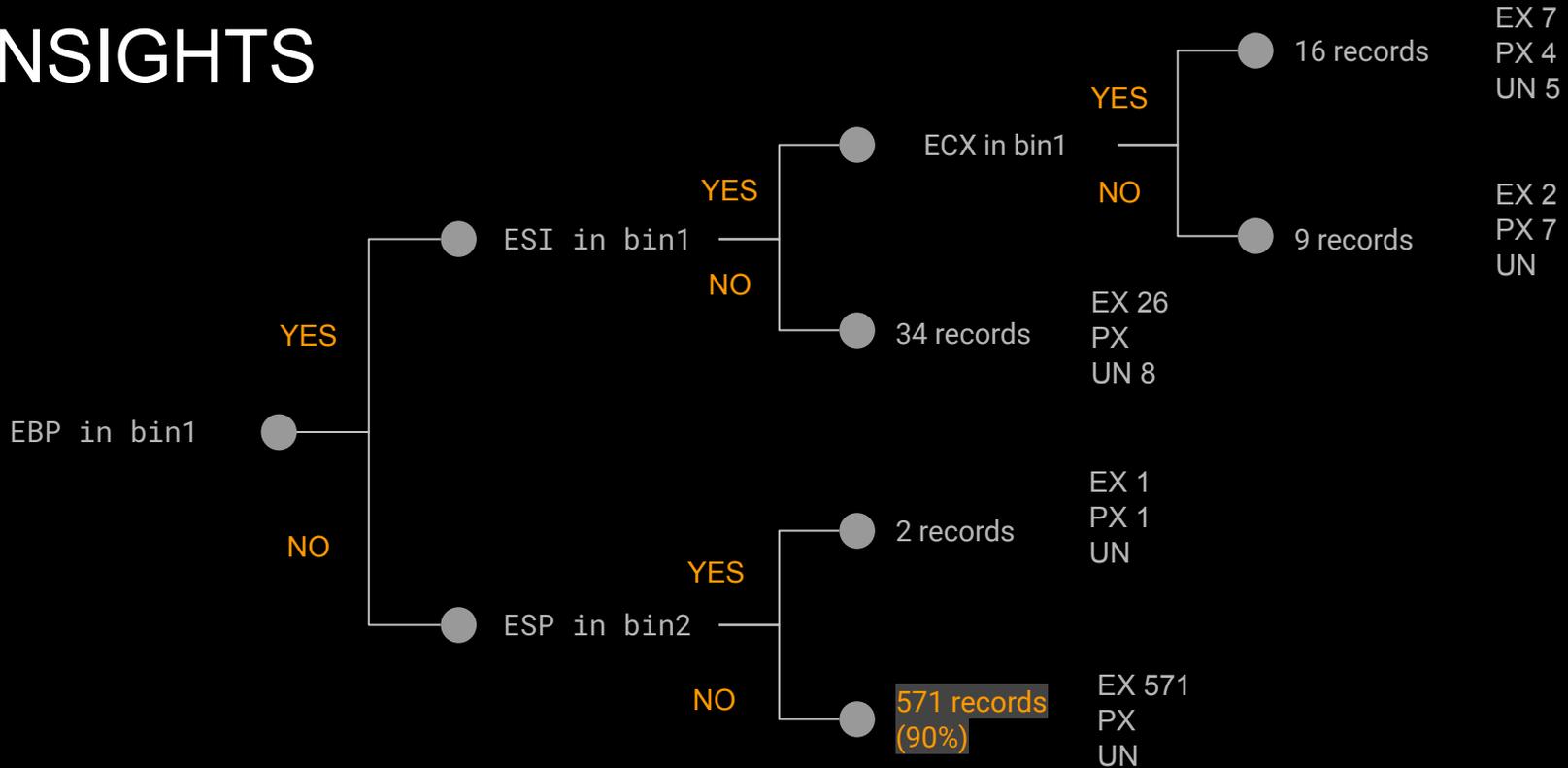# RESULTS - XGBoost

95-99% accuracy

This is not very good, you can get very high success rate guessing EXPLOITABLE all the time - be correct 96% of your guesses



@barnhartguy    @aCaltum

# INSIGHTS



EBP in bin1

YES

ESI in bin1

YES

ECX in bin1

YES

16 records

EX 7
PX 4
UN 5

NO

9 records

EX 2
PX 7
UN

NO

34 records

EX 26
PX
UN 8

NO

ESP in bin2

YES

2 records

EX 1
PX 1
UN

NO

571 records (90%)

EX 571
PX
UN

@barnhartguy    @aCaltum

# RULE OF THUMB?

For 571 (90%) of our records, it is enough to test:
!(EBP in bin1) & !(ESP in bin2) to classify it as EXPLOITABLE

Does this make any sense?

Will this remain true with more data?

# COMPARISON AGAINST exploitable

Built and tested against a set of heuristics - works very well

Out method shows that we can perform as well or better against the same data set

However, we need more data to give any certainty to these claims

# HOW TO BUILD THIS YOURSELF

We released a <span style="color:orange">whitepaper</span> to explain our methodology and results

*https://www.productsecurity.info/files/Whitepaper_SAS19.pdf*

More research, and especially more data is needed!

@barnhartguy     @aCaltum

# CONCLUSIONS

ML is only as good as your dataset, you're answering "how similar"

This is still a work in progress.

We don't have enough non-exploitable crashes to test against

The insights we gathered are interesting, and merit a deeper look when more data is available

# WHERE CAN WE USE THIS?

Feedback for bug trackers (impact/importance)

Feedback for vuln hunters - focus areas

Feedback for fuzzers - where to focus

# MORE INSIGHTS

Data science is an <span style="color:orange">art</span>

We need to talk with people from <span style="color:orange">different disciplines</span> than us

# ACKNOWLEDGEMENTS

@barnhartguy    @aCaltum

# Thank You!

@barnhartguy
@aCaltum