# Using Splunk ML for Threat Hunting

**Joe Partlow, ReliaQuest**
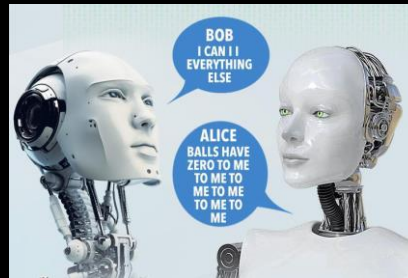
# >whoami

## Joe Partlow – CTO, ReliaQuest

Joe has been in the IT and informaxtion Security industry for 20+ years and most recently been working with simulated attack & defense networks, security analytics, building big data platforms and machine learning.  Reliaquest partners with the worlds largest enterprise splunk customers performing analyst, engineering/architecture and content development functions.

# Inspired by gibson

# What ml/ai is not

▸ Algos will not replace good analysts, just another resource

- But when can I have wintermute in my soc???

▸ Many times, successful hunt campaigns achievable with just effective searches

▸ Many good products/models fail because of poor and incomplete data

▸ Haven't even scratched the surface of basic ML, let alone "ai" …. yet

▸ Gone too soon:



splunk> .conf18

# Ensuring success

▸ Clean data will make or break your training models

- Include filtering, black/white lists to remove large datasets that could skew results

▸ Data normalization across the various sources

▸ Know your data! (supervised and unsupervised learning)

▸ Numeric data shouldn't always be treated as such (ie. port numbers)

▸ Incorporate red teamers with your data science team for better domain knowledge

▸ Continuously retrain your models as the environment changes

splunk> .conf18

# (most) Relevant algorithms

▸ YMMV but below are some algorithms well-suited for common hunt campaigns:

- Detect Numeric Outliers – Useful for determining weird status codes or abnormal event IDs

- Detect Categorical Outliers – Useful for finding weird DNS queries or abnormal user-agents/page requests

- Cluster Numeric Events – Helpful for finding outlying host login counts or application usage counts

- Prediction and Forecasting algorithms might be better suited until the environment is stable and baselined

# Web attack use case

# Web attack use case

# Web attack use case

# Web attack use case

# Breaking the models

▸ Already proven for image classification deep learning*

▸ Similar to SIEM issues, overload the data ingestion with enough noise that the "abnormal becomes normal"

▸ Attackers are already good at blending in (living off the land, pivoting, etc.)

*\* https://blog.openai.com/adversarial-example-research/*



splunk> .conf18

# Future enhancements

▸ Field is progressing amazingly fast. Just because something isn't possible now, give it 6 months!

▸ Move towards stacked/ensemble learning to avoid "jack of all trades, master of none" algos

▸ Build up and better utilize belief networks

• Attempts to increase accuracy by adding conditional dependencies

• Excellent blackhat talk by raffael marty - https://www.slideshare.net/zrlram

# THANK YOU

**Questions?** jpartlow@reliaquest.com

## Don't forget to rate this session in the .conf18 mobile app

.conf18

splunk>