## 4. Detection

Refined and Predictive Alerts via Supervised Machine Learning

Copyright © SAS Institute Inc. All rights reserved.

### **Cybersecurity Data Science (CSDS)**





## Learning Objectives





#### Cybersecurity Data Science (CSDS) Lifecycle



Objectives for Cybersecurity Detection and Prediction Predictive Detection Using Supervised Machine Learning

Detection and prediction: targeted alerts using supervised machine learning

- Overview: the power of labeled data and identified incidents
- Apply supervised analytics machine learning algorithms to predict and to detect targeted anomalies
- Demo / Exercise: Enterprise Miner for supervised analytics



#### CSDS Process Unified Orchestration



#### Role of Algorithms





## Predictive (Supervised) Machine Learning



## **Machine Learning**

#### Descriptive (Unsupervised)

- Cluster analysis
- Factor analysis
- Self-Organizing Maps (SOMs)

#### Predictive (Supervised)

- K-Means
- Decision Trees (DT) (random forests, boosted trees)
- Naïve Bayes classifier
- Neural networks
- Support Vector Machine (SVM)
- Ensembles / Ensemble Learning







Support Vector Machines

#### Machine learning tasks

Regression

**Bayesian Statistics** 

Decision trees

Gradient boosting

Random forests

SVM

Gaussian processes

#### Supervised Learning

- Trained on labeled examples. We have a target we are predicting.
- Map inputs to desired output.
- Suitable for classification, regression, prediction.
  Considerations
- Getting labeled data for rare events can be a challenge
- Suspicion is not fraud
- Data is skewed 99-1







### **Econometric Forecasting**



#### **Econometric Diagnostics**

#### An individual data "points" are anomalous based on the data





**Bayesian Algorithms** 

Many **Predictive** Machine Learning Algorithms...



Artificial Neural Network Algorithms



Ensemble Algorithms

http://machinelearningmastery.com/a-tour-of-machine-learning-algorithms/

### **Predictive Machine Learning (Supervised)**





# **Predictive ML** Why and How for Cyber?



### **Predictive Cybersecurity Use Cases**

### WHAT?

1. Malware detection

2. Device behavior

3. Known incidents

## HOW?

- Recording of behavior (e.g. Wireshark communications)
- 2. CMDB listing of devices + recording of behavior (e.g. authentication + NetFlow + DNS)
- 3. E.g. password attack in Linux authentication logs; DDOS attack

### **Predictive ML for Malware**

- 1. <u>Collect</u> examples of malware and benignware. Use these examples to train ML algorithm to recognize malware.
- 2. <u>Extract</u> features from each training example to represent example as an array of numbers. This step also includes research to design good features that will help your machine learning system make accurate inferences.
- 3. <u>**Train**</u> the machine learning system to recognize malware using the features we have extracted.
- 4. <u>Test</u> the approach on some data not included in our training examples to see how well our detection system works.



Saxe, Joshua Malware Data Science <u>Simple Example</u>: Training Based on a Decision Boundary (2 feature Logistic Regression)



#### ...and 3 dimensional Logistic Regression feature space



**Reference**: Saxe, Joshua. Malware Data Science

#### Irregular and Disjointed Decision Boundaries: e.g. K-Nearest Neighbor

K-Nearest Neighbors



## Where do I get malware for analysis?

- Malware Zoos (e.g. theZoo) <u>https://github.com/ytisf/theZoo</u>
- Build your own zoo <a href="https://www.sans.org/reading-room/whitepapers/malicious/paper/33543">https://www.sans.org/reading-room/whitepapers/malicious/paper/33543</a>
- Metasploit
  - Metasploit Framework is a complex tool that can surveil targets with scans, select exploits to match vulnerabilities (~600), create payloads (~200), receive connections from backdoors, launch exploits

#### Where can I test and analyze?

- 'Sandbox' of test computers
- Honeypots (e.g. RasPis)
- VM lab



Building Virtual Machine Labs

A Hands-On Guide

Tony Robinson

Maar let op! Openbaar Ministerie beschouwt Metasploit en Nmap als malware

https://www.security.nl/posting/567852/Openbaar+Ministerie+beschouwt+Metasploit+en+Nmap+als+malware



Error



# **Model Diagnostics**

# **Confusion Matrix**



Metric	Formula	Interpretation	
Accuracy	$\frac{\mathrm{TP} + \mathrm{TN}}{\mathrm{TP} + \mathrm{TN} + \mathrm{FP} + \mathrm{FN}}$	Overall performance of model	
Precision	$\frac{\mathrm{TP}}{\mathrm{TP}+\mathrm{FP}}$	How accurate the positive predictions are	
Recall Sensitivity	$\frac{\mathrm{TP}}{\mathrm{TP}+\mathrm{FN}}$	Coverage of actual positive sample	
Specificity	$\frac{\mathrm{TN}}{\mathrm{TN}+\mathrm{FP}}$	Coverage of actual negative sample	
F1 score	$\frac{2 \mathrm{TP}}{2 \mathrm{TP} + \mathrm{FP} + \mathrm{FN}}$	Hybrid metric useful for unbalanced classes	

#### Model Diagnostics: Lift and Misclassification Rate



#### **Receiver Operating Characteristic (ROC curves)**



False-positive (FP) rate is also known as the fall-out or probability of false alarm and can be **calculated** as (1 - specificity, where specificity = TN / (TN + FP)). The **ROC curve** is thus the sensitivity as a function of fall-out.

#### Area under the Curve (AUC)





# **ML Tools**

### **Example Machine Learning Tools**

#### Open source

- R
- Python
- Weka
- Spark
  - ...

#### **Commercial**

- SAS JMP
- SAS Enterprise Miner
- SAS Viya VDMML
- IBM SPSS
- Oracle Data Mining
- Rapid Miner



## **Predictive ML** Demonstration and Hands-on Exercise







## **Predicting Device Types**

This demonstration illustrates how to predict unknown device types given a sample of known devices

#### **Predictive Device Labeling – Enterprise Miner**

#### Using Labeled Devices to Predict Unlabeled Devices



#### **Device Types on the Network**

ID	DEVICE TYPE	COUNT
1	Servers	7498
2	Networks (devices)	1810
3	Audio Video Device	856
4	Storage	181
5	Telecommunications	126
6	Application Server	77
7	Facilities	37
8	Critical Workstations	27
9	Printing	17
10	Blade Chassis	8
11	UPS	5
12	Applications	3
13	Clusters	3
14	Service	1
15	Application Code	1
16	Databases	1
17	WebSites	1
		10,652

Unknown devices: 34,615





## **Exercise Review**


# MODEL MANAGEMENT CHAMPION / CHALLENGER

Need to test multiple methods (hint: deep learning is not ALWAYS the best!)



# Why were most of predictions 2?

ID	DEVICE TYPE	COUNT
1	Servers	7498
2	Networks (devices)	1810
3	Audio Video Device	856
4	Storage	181
5	Telecommunications	126
6	Application Server	77
7	Facilities	37
8	Critical Workstations	27
9	Printing	17
10	Blade Chassis	8
11	UPS	5
12	Applications	3
13	Clusters	3
14	Service	1
15	Application Code	1
16	Databases	1
17	WebSites	1
		10,652

Unknown devices: 34,615





## Do not assume given labels are statistically significant!

. . . .

.

Deterministic (organizational groupings)

Nomina	- Logisti	C FIT T	or peer	Group	.oae			
Whole N	lodel Te	st						
Model	-LogLike	lihood	DF	ChiSqu	iare	Prob>ChiSq		
Difference	33	05.967	440	6611	.934	<.0001*		
Full	242	87.706						
Reduced	275	93.673			Г			
						over random		
RSquare (U)			0.1198			groupings		
AICc			49536.7					
BIC			52845					
Observation	ns (or Sum	Wgts)	10740					
Measure		Irain	ing Detin	ition				
Entropy RSo	quare	0.1	198 1-Log	like(mod	lel)/Lo	oglike(0)		
Generalized	RSquare	0.4	624 (1-(L(	0)/L(mod	del))^	(2/n))/(1-L(0)^(2		
Mean -Log	р	2.2	614 ∑ -Loo	g(ρ[j])/n				
RMSE		0.8	581 √ ∑(y[	JJ-PU				
Mean Abs I	)ev	0.8	507 <u>2 I</u> YUJ	-PUJ/	Classi	5% correct		
Misclassific	ation Kate	1074	438 <u>2</u> (Po		Class	incation rate		
N		1074	Jn					
Lack Of I	Fit							
Source	DF	-LogLi	kelihood	ChiSqua	are			
Lack Of Fit	209740	2	4210.109	48420	).22			
Saturated	210180		77.598	Prob>C	hiSq			
Fitted	440	2	4287.706	1.0	000			

HISTOGRAM It3 P\_VALUES PEERGROUP ADJUSTED VALUES







# Python on Viya CAS

Running Python on CAS via SWAT

# How to program against CAS?



### SAS and Open Source Integration

# Base SAS

 Execute open source code using system commands via the DATA step.

#### Interactive Matrix Language (IML)

- Submit R code within IML from SAS Code Editor.
- Pass data between R and SAS.

#### **Enterprise Miner**

- Execute R code from the open source integration node.
- Execute Python code using a system command.

#### SAS Viya

 Use open source software to take control of analytical tools.

# What is SWAT?

- SAS Scripting Wrapper for Analytics Transfer
- Enables Python, Lua, and R to interface with CAS





# Scripting Wrapper for Analytics Transfer (SWAT)

- The SWAT package enables you to interface with CAS from R or Python.
- You can write an R or Python program that connects to a CAS server, load data into CAS, analyze large in-memory data sets quickly and efficiently using CAS actions, and work with results of your analyses using familiar data wrangling techniques in the open source language.



# How do you submit code?

**CAS Session** Used to enable clients to communicate with the server to request actions.









# **Predicting Security Violations**

This demonstration illustrates how to predict when a user requesting secure access should be denied based on past behavior and profile

💭 Jupyter		OPEN JUPY	TER NOTEBOOK			
Files Running	Clusters					
Select items to perform a	ctions on them.					
0 - 1						
Contacts						
🔲 🗅 Desktop						
Documents						
Downloads		yter ython_Factorization_Machine_Demo.ipynb ython_Factorization_Machine_Exercise.ipynb ython_Image_Classification_Demo.ipynb ython_Machine_Learning_Demo.ipynb ython_Machine_Learning_Exercise.ipynb	Python Machine Lear ✓ Overview of the Predictiv A secure organization has a system for ap- incidents where improper access was requ- being recorded. Denied access requests a wishes to semi-automate access denial by outstanding requests, and job role variable administrator. Data After analyzing past security data (5966 cc	re Modeling proving access to uested. Denied ac rer relatively high ( predicting wheth as, the company w	no - Sec g Case b highly secure of ccess involves a (~20%) and tak vants to build a n wants to build a n	tata. Each user has a security permission level and access is tracked, including n assessment when improper access is requested, resulted in a security incident e up a great deal of security administrative time and resources. The organization is tresembles past denied requests (incidents). By using permissions, past activity, model to predict whether an applicant would be denied secure access by an
-			data. These variables, along with their mo Name I BAD CLAGE CINO	del role, measurer Model Role Measur Target Input	irement Level Binary Interval	Consider the second of th
			SECPEM DELINQ DEROG	Input Input Input	Interval Ra Interval	atio of outstanding secure requests to security permission level Number of delinquent requests for clarification on access Number of past access denials

JOB

Input

Nominal

Occupational categories



# Python API on SAS Viya CAS

- Anaconda 3 (distribution of tools): jupyter notebook, Spyder, Pscikit-learn, matplot lib, NumPy, pandas
- >> DOCUMENTS
  - Python\_Machine\_Learning\_Demo.ipynb
  - Python\_Machine\_Learning\_Exercise.ipynb (completed)
- Three ways to run code
  - Run icon (top bar)
  - Ctrl + Enter (selected segment)
  - Shift + Enter (run & move to next cell)

# Jupyter Notebook



# Jupyter Notebook Results



TargetBinary1 = security incident recorded, 0 = no incidentsInputIntervalAge of oldest outstanding request for clarification on accessInputIntervalNumber of sensitive secure access permissionsInputIntervalRatio of outstanding secure requests to security permission levelInputIntervalNumber of delinquent requests for clarification on accessInputIntervalNumber of delinquent requests for clarification on accessInputIntervalNumber of delinquent requests for clarification on accessInputIntervalSecurity sensitivity of past access denialsInputIntervalSecurity sensitivity of requested data accessInputIntervalNumber of recent access to secure dataInputIntervalNumber of recent access clarification inquiriesInputBinaryReason for access request	Measurement Level	Model Role	Name
InputIntervalAge of oldest outstanding request for clarification on accessInputIntervalNumber of sensitive secure access permissionsInputIntervalRatio of outstanding secure requests to security permission levelInputIntervalNumber of delinquent requests for clarification on accessInputIntervalNumber of delinquent requests for clarification on access denialsInputIntervalNumber of past access denialsInputIntervalSecurity sensitivity of requested data accessInputIntervalNumber of recent access to secure dataInputIntervalNumber of recent access clarification inquiriesInputBinaryReason for access request	Binary 1 = security incident reco	Target	BAD
InputIntervalNumber of sensitive secure access permissionsInputIntervalRatio of outstanding secure requests to security permission levelInputIntervalNumber of delinquent requests for clarification on accessInputIntervalNumber of delinquent requests for clarification on accessInputIntervalNumber of delinquent requests for clarification on accessInputIntervalOccupational categoriesInputIntervalSecurity sensitivity of requested data accessInputIntervalNumber of recent access to secure dataInputIntervalNumber of recent access clarification inquiriesInputBinaryReason for access request	Interval Age of oldest outstanding request for	Input	CLAGE
InputIntervalRatio of outstanding secure requests to security permission levelInputIntervalNumber of delinquent requests for clarification on accessInputIntervalNumber of delinquent requests for clarification on accessInputIntervalOccupational categoriesInputIntervalSecurity sensitivity of requested data accessInputIntervalRecent access to secure dataInputIntervalNumber of recent access clarification inquiriesInputBinaryReason for access request	Interval Number of sensitive secu	Input	CLNO
InputIntervalNumber of delinquent requests for clarification on accessInputIntervalNumber of past access denialsInputNominalOccupational categoriesInputIntervalSecurity sensitivity of requested data accessInputIntervalRecent access to secure dataInputIntervalNumber of recent access clarification inquiriesInputBinaryReason for access request	Interval Ratio of outstanding secure requests to se	Input	SECPEM
InputIntervalNumber of past access denialsInputNominalOccupational categoriesInputIntervalSecurity sensitivity of requested data accessInputIntervalRecent access to secure dataInputIntervalNumber of recent access clarification inquiriesInputBinaryReason for access request	Interval Number of delinquent requests for	Input	DELINQ
InputNominalOccupational categoriesInputIntervalSecurity sensitivity of requested data accessInputIntervalRecent access to secure dataInputIntervalNumber of recent access clarification inquiriesInputBinaryReason for access request	Interval Number	Input	DEROG
InputIntervalSecurity sensitivity of requested data accessInputIntervalRecent access to secure dataInputIntervalNumber of recent access clarification inquiriesInputBinaryReason for access request	Nominal O	Input	JOB
InputIntervalRecent access to secure dataInputIntervalNumber of recent access clarification inquiriesInputBinaryReason for access request	Interval Security sensitivity of a	Input	SECS
InputIntervalNumber of recent access clarification inquiriesInputBinaryReason for access request	Interval Recent	Input	RECA
Input Binary Reason for access request	Interval Number of recent acces	Input	NINQ
	Binary Reas	Input	REASON
Input Interval Security permission level	Interval Se	Input	SECL
Input Interval Years at present job	Interval	Input	YOJ



# **Semi-Supervised Learning**

## **Machine Learning Segmentation and Classification**



https://medium.com/datadriveninvestor/differences-between-ai-and-machine-learning-and-why-it-matters-1255b182fc6

### Human-in-the-Loop: Self-Improving Cybersecurity Analytics Cycle





# Semi-Supervised Approach



#### Machine learning tasks

Regression Decision trees Gradient boosting Random forests Autoencoders Text Processing Image Processing

### Semi-supervised Learning

- Helpful when volume or variety of data is too high to allow labelling
- Utilized labeled and unlabeled examples
- Classification, regression, prediction

#### Considerations

- Best of both worlds?
- More automated



# Moving from Anomalies to Focused Incident Detection





# **Reinforcement Learning**

#### EXAMPLES

Monte Carlo

- Neural Networks
- Deep Forward Neural Networks
- Convolutional Neural Networks

Recurrent Neural Networks

## **Reinforcement Learning**

- Like teaching someone a game
- The machine takes actions and learns from results
- Maximize an expected future reward Considerations
- Data intensive Data drives insight
- Hard to understand what was learned

Random Walk





- Walls block the agent's path
- Agent's action has randomness
- Big reward comes at the end

Follow a good policy



# Wrap-Up





# **Section Review**



# **Review: Machine Learning**

#### Unsupervised



- No targets required
- Finds similar 'groups' within the data
- Not a direct indicator of
  incident
- Will generally patterns in all data (need to have right data)

#### Supervised



- Easy to build
- Easy to understand

0

- Only known cases!
- Data skew

## Semi-Supervised



- Useful with rare events
- Utilized both labeled & unlabeled examples
- Requires careful selection
  of cases
- Requires careful 'tuning'

#### Reinforcement



- Self learning
- Improves over time
- Pattern analysis

Black box

٠

- Large training set required
- Overfitting a problem

# **Machine Learning as a Process**





http://www.oreilly.com<sup>6</sup>/data/free/archive.html



What do we do when there are very few examples of known incidents?

# **Cybersecurity Analytics Maturity**



# **Cybersecurity Analytics Maturity**



#### Cybersecurity Data Science (CSDS) Lifecycle **Monitor** -0-Frame DECIDE **Explore** Deploy DETECTION **DISCOVERY**

En

Model

Validate

Engineer



# **APPENDIX**

# **Deep Learning**



# What about AI?


## **Deep Learning**

O'REILLY'

W1.3  $w_{2,3}^{(2)}$   $w_{3,3}^{(2)}$ 

 $w_{1,3}^{(1)}$ 

 $w_{2,3}^{(1)}$   $w_{3,3}^{(1)}$ 



## Deep Learning A specialization of machine learning



- Neural networks with many layers and different types of ...
  - Activation functions
  - Network architectures
  - Sophisticated optimization routines

Each layer represents an optimally weighted, non-linear combination of the inputs

- Automatic feature generation
- Extremely accurate results if well-trained; use for classification, prediction, or pattern recognition, especially in unstructured data



#### **Practical Applications:**

- Fraud (& cyber)
- Marketing analytics
- Text pre-processing

### Deep Forward Neural Networks

- Simplest type of Deep Learning algorithm, where information is only fed forward, from input to output
- While often combined with other networks to from new networks, small or shallow DFNNs have been shown to be effective with certain tasks, making them more desirable for resource constrained environments (such as cell phones) – versus other, more resource-hungry algorithms
- Can also be applied to NLP tasks including: language identification, part-of-speech tagging, word segmentation, and preordering for statistical machine translation

#### Components





Fully-connected layers

### **Basic Neurons and Fully-Connected Layers**



#### In a Deep Forward Neural Net:

- Each neuron is connected to all other neurons (fully-connected layers)
- All connected previous values, multiplied by their weights, are added to bias and then fed through an activation function



#### Practical Applications:

- Anomaly detection
- Feature extraction
- Dimension reduction (VAE)

# Auto Encoders

- Unsupervised learning models used to encode information via a fully-connected, hour-glass shaped architecture which is always symmetrical around the middle layer(s), where the info is most compressed
- In its simplest form, the output layer is a copy of the input layer and has the task of reconstructing its own inputs (rather than predicting target values y given inputs x)
- The aim of an autoencoder is to learn a representation (encoding) for a set of data, typically for the purpose of dimensionality reduction and/or noise reduction. Recently, the autoencoder concept has become more widely used for learning generative models of data



#### Practical Applications:

- Image/object recognition
- Video analysis
- Text Classification

### **Convolutional Neural Networks**

- Versatile, though best-known for their prowess at image analysis tasks. CNNs can detect components in images such as edges and curves, and learn how these features are combined in order to detect larger structures such as faces, numbers, etc.
- Because CNNs preserve spatial structure, they are well-suited to array data where nearby values are correlated (images, sound, video, speech, etc.).
- CNNs are relatively fast because convolutions are critical in computer graphics and can be implemented on GPUs
- For text analysis, word-level or character-level CNNs can be applied for tasks such as sentiment analysis, categorization, or spam detection. Note: because CNNs require fixed input and output sizes, padding is often required for inputs.

#### Components

#### **Convolutional Layers**

- Each neuron is only connected to cells within a certain proximity and NOT connected to all other neurons in the previous layer
- Also called "filters"



#### Pooling Layers

- Used for down-sampling, with maxpooling being most popular
- Once a feature is known to be present, its exact location is not as important as its relative location to the other features

### **Convolutional and Pooling Layers**







#### Practical Applications:

- Language modeling (e.g., statistical machine translation, word prediction)
- Time-series forecasting
- Image/video captioning

### **Recurrent Neural Networks**

- Versatile, but especially well suited for sequence data like time series or text because an RNN will recognize patterns across time or the order of words in a given language.
- Unlike CNNs, RNNs can handle variable-length inputs and produce variable-length outputs, which is useful for applications like NLG and machine translation.
- Stacked RNNs can form a net capable of more complex output (e.g., used for autonomous vehicles)
- Note: even if your data is not sequential, you can still process it sequentially and reap the benefits of RNNs<sup>80</sup>

#### Components

 $\bigcirc$ 

Recurrent cell



Time-Delayed Layers

### Recurrent (Time-Delayed) Layers



#### Time-Delayed Layers:

• They are connected to all neurons in the previous layer and updated just like basic cells, but with extra weights: connected to the previous values of the cells and most of the time also to all the cells in the same layer. These weights between the current value and the stored previous value work much like a volatile memory (like RAM).

#### **Algorithm Variants**



Long Short Term Memory



Input Cell

Output Cell

Different Memory Cell

Gated Recurrent Unit

### Long Short Term Memory (LSTM) RNNs and Gated Recurrent Unit (GRU) RNNs

- RNN variants that help with the problem of the "vanishing" or "exploding" gradient by introducing gates and an explicitly defined memory cell
- Can learn complex sequences (e.g., write like Shakespeare, compose primitive music) but typically require more resources to run. In cases where extra expressiveness is not needed, GRUs can outperform LSTMs.
- Each neuron has a memory cell and 2-3 gates which safeguard the information by stopping or allowing the flow of it.

#### LSTM RNN

• Input gate

• Output gate

• Forget gate

- **GRU RNN** 
  - Update gate
    - Reset gate

#### Autoencoders





## What are they?

The basic idea behind autoencoders is to <u>encode</u> information <u>auto</u>matically.

It requires at least 3 layers with perfect symmetry on both sides.

Since the output layer is a copy of the input layer there is no target.

The output layer is trying to represent the input layer so anything that stands out from the input layer is anomalous.



Output layer is copy of input layer, so no target

Fully-connected, symmetric



Odd number of hidden layers

Hidden layers symmetric in number of units, typically decreasing to middle hidden layer then increasing

Middle hidden layer ("code" layer) used as new features







Units from the middle hidden layer ("code" layer) are our new features, like PCs



How to use NNET for autoencoders Programmatic Approach using SAS Studio / Code Node

Model Information	
Model	Neural Net
Number of Observations Used	1987
Number of Observations Read	1987
Number of Nodes	76
Number of Input Nodes	25
Number of Output Nodes	25
Number of Hidden Nodes	20
Number of Hidden Layers	3
Number of Weight Parameters	648
Number of Bias Parameters	51
Architecture	MLP (AUTOENCODER)
Seed for Initial Weight	802977670
Optimization Technique	LBFGS
Number of Neural Nets	1
Objective Value	0.7453511545

### **PROC NNET**

Example:

proc nnet data=casuser.train; input cycle X1 - X24 / level=interval; hidden 12; hidden 2; hidden 12; train outmodel=casuser.enginennet; optimization regL2=0.1; score out=casuser.autoencoder; code file="/casuser/nnetscore.sas"; run;

**Best Practices:** 

- Input layers = Output layers
- Use odd number of hidden layers (ex. 3, 5, 7, etc...)
- Center hidden layer should have 2 nodes
- Create a bowtie / hourglass between input & output layers:
  - #nodes in layers between input and center should be decreasing
  - # nodes in layers between center and output should be increasing and match input to center pattern
- Deeper the network  $\rightarrow$  more time is needed to train