Augment your Security Monitoring Use Cases with MLTK's Machine Learning

Oliver Kollenberg, Philipp Drieger Splunk .conf19, Las Vegas, Oct 24th 2019

Unrestricted © Siemens 2019

siemens.com/datacenters

SIEMENS

Ingenuity for life

EAGLE DataCenter





Since 12 years In 4 locations 8 Petabytes storage Jet-zero carbon footprint

> Splunk at EAGLE: IT Ops and Security 70 indexes 3,200 data sources 210 data types

Unrestricted © Siemens 2019 Page 2 2019-10-24

Agenda





- Intro
- Machine Learning Concepts
- Use Cases
 - Find Changes in Proxy Communication
 - Identify System Compromises
 - Baseline Privileged Account Behavior
- Wrap up, Q&A

Machine Learning for Security Monitoring: The Project



Project setup

- Started Q1 2019
- Project staff of four
- Intention:
 - Can we extend SecMon with ML, to:
 - Baseline the data center, analyze behaviors?
 - Find security relevant patterns?



Project Team

SIEMENS Ingenuity for life

Markus

IT Sec deepwaters



Philipp ML project & science coach



Michael

ML concepts and fundamentals



Oliver Splunk> that!



Unrestricted © Siemens 2019

Page 5 2019-10-24



Machine Learning Concepts

Typical Machine Learning Scenarios



Anomaly detection Predictive Analytics









Deviation from past behavior

Proxy Communication

Classify malware communication

Identify System Compromises

Behavioral Analytics

Protect Privileged Accounts

Machine Learning Project Timeline







Use Case 1: Proxy Communication

Use Case One: Proxy Communication Idea, Concept, Results



Idea: Can we **profile system communications** outbound via webproxy and **detect changes** in behavior to ultimately **find unwanted outliers**?

Concept

- Collect eight features (each by Source IP) onto summary index.
- One DensityFunction per feature. | fit every night.
- | apply density models in intervals. Historize by again collecting onto summary index.
- Dashboard visualizes findings, guides investigation workflow.

Results

- · Changes are reliably reported.
- Findings are consolidated into meta alarms (aka "change of change").
- Change caused by unwanted software (HTTP tunneling of SSH traffic) correctly detected.

Use Case One: Proxy Communication Detected Malware Communication



Detect the Unwanted: HTTP tunneling of SSH traffic

Meta Alert:

Triggers on first change, then change of change.



Alerted on outlying characteristics:

Anomaly Details for (displaving all events having outliers)										
_time \$	src_ip \$	outlierValue_newUserAgents	anomaly_score	outlier_newUserAgent \$	outlier_dcDstNetloc \$	outlier_countConns \$	outlier_sumBytesIn \$	outlier_sumBytesOut \$		
2019-07-22 14:00	a sure	cur1/7.60.0	6.0	1	1.0	1.0	1.0	1.0		

Unrestricted © Siemens 2019

Page 11 2019-10-24

Use Case One: Proxy Communication Concept Building Blocks: SPL-Foo



Three scheduled Searches:

Search #1: Base Collector and SummaryFiller

- Create eight feature fields (dcDstIp, dcDstNetloc, dcCategories, categories, dcUserAgent, UserAgents, sumBytesIn, sumBytesOut), splitted by _time (1h) and src_ip
- Collect onto summary index.

Search #2: Training / Model Fitter

- Train one DensityFunction per feature.
- Full | fit every night.

Search #3: Alert Generator, Delta Fields Generator and SummaryFiller

- AlertGen: Find outliers by | apply 'ing density models over latest timeslice (every few hours).
- DeltaFieldsGen: Also cover last month of data to generate additional feature "newUserAgent", via streamstats implementing a historybrain.
- Collecting again onto summary index, to historize and make available for dashboard.

Unrestricted © Siemens 2019

Use Case One: Proxy Communication Concept Building Blocks: Scheduling the ML System



schedules

174

timechant when was in

NVESTIGATE

relea = Summary-alornin

1.91.7.123 753



Unrestricted © Siemens 2019

offestricted @ Sterriens 2019

Page 13 2019-10-24

Use Case One: Proxy Communication SPL-Foo: Search #1 Base Collector and SummaryFiller





Unrestricted © Siemens 2019

Page 14 2019-10-24

Oliver Kollenberg

-

1697

583

877

4465611

pcr 🌲 🦨

-64398571

-6297

-1448

-2188

Use Case One: Proxy Communication SPL-Foo: Search #2 Training / Model Fitter





Runs once per day, nightly.

Use Case One: Proxy Communication SPL-Foo: Search #3 AlertGen, DeltaFieldsGen and SmryFillr



	Earliest time - d@d
Run longer for	Time specif
historybrain 4 table _time, src_ip, sc_filter_result, countConns, dcDstIp, dcDstNetloc, dcCategories, categories, dcViruses, dcUserAgent, UserAgents, sumBytesI	Latest time -2h@h
6 rename UserAgents as UserAgent 7 sort 0 _time	
Historybrain:	14 addinfo 15 eval secsInHour=60*60
Find new values 9 eval outlier_newUserAgent = if(countUserAgent == 1, 1, 0) 10 eval outlierValue_newUserAgent = if(outlier_newUserAgent == 1, UserAgent, pull()	<pre>16 eval timeRemainder=info_max_time % secsInHour 17 eval timeMaxSnappedHour=info_max_time - timeRemainder 18 eval myEarliest=timeMaxSnappedHour - 8*secsInHour</pre>
null())	<pre>19 eval myLatest=timeMaxSnappedHour 20 d where time >= myEarliest AND time < myLatest</pre>
Throw away events which were	21 fields - secsInHour, timeRemainder, timeMaxSnappedHour
only needed for historybrain 34 apply Mdl sumBytesIn threshold=0.	.005
<pre>26 rename IsOutlier(sumBytesIn) as outlier_sumBytesIn</pre>	0.005
Apply to get outliers. Versionize. Collect to historize. 42 eval version=11 43 eval categoriesZipped=mvjoin(categories, `tok`), UserAger (UserAgents, `tok`), outlierValue_newUserAgentZipped=mv (outlierValue_newUserAgent, `tok`) 44 fields - categories, UserAgents, outlierValue_newUserAger 45 collect testmode=false addtime=true index=	ntsZipped=mvjoin vjoin nt
Unrestricted © Siemens 2019 marker="search_name	
Page 16 2019-10-24	Oliver Kollenberg

Use Case One: Proxy Communication Investigation Workflow

1. Derivate Alerts: Act on these



2. Deriving from sum(outliers) (=scoring) per src_ip



3. See anomaly_score trend (per src_ip)









4. Drilldown on specific src_ip: Check outliers by visually presenting base values and its delta



5. Drilldown on specific src_ip: Check raw values

in i	11.12.3	methos 4	4	Disstances:	$\frac{1}{2} \frac{1}{2} \frac{1}$	ance of the second	seller_continuition	outer_Activation i	out equilibrium a	out an a de Drais. B	autice_orCompanys	out an constituent de	reflec_sackposit
115 67 22 85 96			1 100	Home (1) (1) (1) (1) (1)	(947) 14 (Stan go)	. A &	1	3.1	11	2.0	8.8	# A	14
112-07-22 14:00			1973	csr1/7.50.0 core	L217.61.8	6.8		1.8	14	8.0	9.8	-1.0	
19-67-12, 59:80			- 1977	Rect/1 14 (linux-gra) rane	1216					4.1	9.8		
019-07-38 33-04			1.01	ceil/i.se.e		5.0	4	9.8	1.8	8.0	9.8		
15-07-58 22:80			1011	car1/1.90.0 note		5.0				8.0	9.8		
46-47-38.35-86			0.02	car1/1.40.0		5.8		0.0	0.1	2.0	2.2		
119-67-38 28:86			. 1912	cert/i.se.e nose		5.4				8.0	9.2		
10.67 19 19.00			1273	eer1/7 59.0		5.8	9	2.8	2.8	2.0	3.2		
14-67-38 13-86			0.00	reri/T 30.0 NEW						2.6	5.2		
112-07-59 17:00			3071	cs=1/7.50.0		3.0	9	3.8	0.8	2.0	0.2		
19 87 39 35 86			322.8	41/1.50.0 10/0		5.8				8.0	3.8		
112-07-39 15:80			1911	car1/7.50.0		3.0	9	9.8	9.8	2.0	2.2		
112-07-33 14:00			1071	cur1/7.50.0 nove		5.8				2.0	9.8		
11-62-54 TS-RE			8.82	card/r.90.0		: //	1	9.8	9.8	8.0	9.8		
12-67-59 12:80			1.01	car1/1.50.0 note		5.8				8.0	9.8		
14-47-38 11-86			122.0	curl/1-lin in		25.4			1.6	2.5	2.2		
112-67-58 58:86			- 191	Cardini, SR, B. Rohe		>1				1.0	9.8		
15 67 18 89.96			222.9	cur1/7.30.0		(5.9		1.0	3.8	2.0	3.2		
114-47-38 83-84			2.4.4	rani/i 44.6 Iok		2.8	1			3.6	8 a		
15-67-39.67:80			1071	cur1/7.50.0		3.0		3.0	5.0	2.0	2.2		

Oliver Kollenberg



Use Case 2: Identifying System Compromises

Use Case Two: Webserver Monitoring Identify compromises by their behavior

Idea

- Can we use MLTK to ease alert fatigue?
- Instead of further evolving attack signatures:
- Let's find signals that identify successful compromises
 - A hacked webserver will change its behavior. How?
 - Files written where they should not, e.g. Webshell
 - Webpage or DB table modified
 - OS Command executed, e.g. cmd.exe

Concept

• ...

- Use Windows logs (standard and sysmon) to generate features:
 - [Webserver: File Written into Webroot by Webserver Worker Process
 Classic SPL]
 - [All Systems: Lsass CrossProcs
 Classic SPL]
 - Webserver: Suspicious commands started by webserver worker process

 ML

• ...



Use Case Two: Webserver Monitoring Find Suspicious Commands started by W3Worker Process



Observations

- How can we find unusual command from w3wp.exe?
 - □ "unsual" can be tricky, CommandLines have a high variance (e.g. date/time parameter).
- Instead, how about suspicious commands? But not every process binary is bad
 - □ Need to consider whole CommandLine.
- A single command may not be significant enough
 - □ Need to consider multiple commands, their sequence.

Concept

- Define a list of suspicious words to extract matches from the CommandLine
- Collect extracted matches into bag of words pattern of suspicious sequences
- Classify malicious sequences:
 - | fit TFIDF from bag of words
 - Classify with | fit LogisticRegression from a set of positive examples
- Identify malicious sequences with a supervised learning approach
- Maintain lookup of positive examples to train and refine the model over time and adapt to new situations

Use Case Two: Webserver Monitoring Find Suspicious Commands



Apache webshell usage





SQL Server CmdInjection





Unrestricted © Siemens 2019



Use Case 3: Baseline Privileged Account Behavior

Use Case Three: Baseline Privileged Account Behavior Idea, Concept, Gain



Idea: Can we **profile Windows admin logon sessions** and **detect changes** in behavior to ultimately **find suspicious outliers**?

Concept

- Collect features describing Windows admin logon sessions onto summary index
- Cluster training data (last [30;1[days) using | fit XMeans models every night. One model each for focusing on source (ip address), session (behavior) and target (host)
- | apply models each day to testing data (last day). Historize by collecting onto summary index.
- Detect anomalies based on changes in cluster_distances (| fit DensityFunctions and | outputlookup percXX statistics) and cluster memberships (| outputlookup account to cluster relations)
- Dashboard visualizes findings and guides investigation workflow.

Results

- Three types of anomalies detected: Source, Target and Session Metadata
- Enables detection of credential theft, malicious account usage (e.g. lateral movement), ...
- Differentiation between local and global outliers: New accounts trigger only global alerts

Unrestricted © Siemens 2019

Use Case Three: Baseline Privileged Account Behavior Investigation Workflow

SIEMENS Ingenuity for life

> 0 20k - 40i

						4. Compare sources	
1. Ale	rts overv	iew (gra	phical)			Sec. secid day Strsecid Strsecid	
talaga kalaga kalaga kalaga	Nu-lag Sur-lag Nu-lag	4 The Bar B	instant instant instant	The face if the face of the fa		Image: state	-
2. Ale	rts overv	iew (tab	ular)			src_secid day h	ost 3
day \$ src_secid \$	lastname ¢	firstname ¢	department \$	alert_reasons ≎	alert_types ‡		
2019-09- 05	-	Training .		too many global outliers too many outlier types too many outliers	session target		
2019-09-		Received.		too many outlier types	session target		
2019-09- 05	10/10/0	tara -		too many outlier types	session target	training 2019-08-21	
2019-09-				ton many outliers	session	6. Compare session features	
3. Dril	lldown on	Alert: V	/iew Detai	ls		4	-
alert_type time	9 0	src_secid	¢ IsOutl	ier_DF_local \$	IsO	utlier_DF_global \$ 3.5	
session 2019-	-09-05 23:00	and the second	unkno	vn_secid_cluster_rela	tion	1.0	
session 2019-	-09-05 23:00		unkno	vn_s	Souority *	IsOutlier LiserClusterBol *	
session 2019-	-09-05 23:00	-	unkno	vn_s	_ Sevency ¢		
session 2019-	-09-05 23:00		unkno	vn s			
Unrestri	cted © Sieme	ens 2019					
Page 24		2019-10-24	1		3		



Graphs Special

Graphs Special: Consider Relations Between Use Cases



Scenario: Compromise, C2. Gives these Anomalies:



Relate e.g. via Accounts' Connections. Find ConnectedComponents in the Graph:



comment("relate via Connections using ConnectedComponents")`

- eval FromSys = FromSystem + ":" + ConnectionId
- eval ToSys = ToSystem + ":" + ConnectionId

graph algo=WeaklyConnectedComponents linkSrc=FromSys linkTgt=ToSys outfield=componentId=

Calculate AnomalyScore per System, Connection and per ConnectedComponents. Map Score to Color:

eventstats sum(objAnomalyScore) as componentAnomalyScore by componentId `bilinearinterpolatecolorgradient(objAnomalyScore, componentAnomalyScore, ..., outcolor)`_

Unrestricted © Siemens 2019

Page 26 2019-10-24







Wrap Up

Summary, Key Benefits and Takeaways



Key Benefits

- New SecMon capabilities:
 - 20 Mio Console Cmds classified every day by machine
 - Pattern analytics of core infrastructure components: Communications, compromises, privileged accounts
- Efficiency gained: Smoother day-to-day ops, more time for use case development.

Takeaways

- Classic SPL is here to stay
- ML vastly expands SecMon capabilities:
 - Allows to look at the "grey" between the black&white
 - Enables feature-rich use cases
 - Easy tuning
- New data insights may impact project plan. Stay flexible.





Unrestricted © Siemens 2019

Page 28 2019-10-24

Thank you! Q&A

Oliver Kollenberg Security Consultant EAGLE DataCenter SIEMENS AG oliver.kollenberg (a t) siemens.com

Philipp Drieger Staff Machine Learning Architect Splunk Inc. philipp (a t) splunk.com

ttp://siemens.com/datacenters