

## Using The Latest Features From The Splunk Machine Learning Toolkit To Create Your Own Custom Models

Adam J. Oliner | Director of Engineering Manish Sainani | Director of Product

September 2017 | Washington, DC

splunk

#### **Forward-Looking Statements**

During the course of this presentation, we may make forward-looking statements regarding future events or the expected performance of the company. We caution you that such statements reflect our current expectations and estimates based on factors currently known to us and that actual events or results could differ materially. For important factors that may cause actual results to differ from those contained in our forward-looking statements, please review our filings with the SEC.

The forward-looking statements made in this presentation are being made as of the time and date of its live presentation. If reviewed after its live presentation, this presentation may not contain current or accurate information. We do not assume any obligation to update any forward looking statements we may make. In addition, any information about our roadmap outlines our general product direction and is subject to change at any time without notice. It is for informational purposes only and shall not be incorporated into any contract or other commitment. Splunk undertakes no obligation either to develop the features or functionality described or to include any such feature or functionality in a future release.

Splunk, Splunk>, Listen to Your Data, The Engine for Machine Data, Splunk Cloud, Splunk Light and SPL are trademarks and registered trademarks of Splunk Inc. in the United States and other countries. All other brand names, product names, or trademarks belong to their respective owners. © 2017 Splunk Inc. All rights reserved.

#### Outline

- Splunk Machine Learning Toolkit
- Platform Extensions: ML-SPL, etc.
- The Assistants: Guided Machine Learning
- What's New
- Demo
  - DIY Anomaly Detector



# Splunk Machine Learning Toolkit

platform extensions and guided modeling dashboards



## **Machine Learning**

- A process for generalizing from examples
- Examples
  - A, B,  $\dots \rightarrow \#$  (regression)
  - A, B,  $\dots \rightarrow a$  (classification)
  - $X_{past} \rightarrow X_{future}$
  - like with like
  - |X<sub>predicted</sub> X<sub>actual</sub>| >> 0

- (forecasting)
- (clustering)
- (anomaly detection)





splunk

## **Data Gathering And Prep**

Source: CrowdFlower



duct.screen?product id=FL-DSH-01&JSESSIONID=SD

#### What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

#### **Machine Learning Process With Splunk**



## **Overview Of Machine Learning At Splunk**



## **Splunk Machine Learning Toolkit**

extends Splunk with new tools and guided modeling

- Assistants: Guide model building, testing, & deployment for common tasks
- Showcase: 25+ interactive examples from IT, security, business, and IoT
- Algorithms: 25+ standard algorithms plus an extensibility API
- SPL ML Commands: New commands to fit, test, and operationalize models
- Python for Scientific Computing Library: 300+ open-source algorithms

screen?product id=FL-DSH-01&JSESS





### Machine Learning Customer Success



## **Machine Learning Toolkit Customer Use Cases**

Reducing customer service disruption with early identification of difficult-to-detect network incidents

Minimizing cell tower degradation and downtime with improved issue detection sensitivity

\_\_\_\_\_

ZIIOW Speeding website problem resolution by automatically ranking actions for support engineers

**docomo** Ensuring mobile device security by detecting anomalies in ID authentication



Predicting and averting potential gaming outage conditions with finer-grained detection Preventing fraud by Identifying malicious accounts and suspicious activities



Jct.screen?product\_1d=FL-DSH-01&JS

TELUS

Improving uptime and lowering costs by predicting/preventing cell tower failures and optimizing repair truck rolls

# Platform Extensions: ML-SPL, etc.

custom search commands for machine learning



## SPL, Macros, & Viz

#### Oh, my!

#### Commands (ML-SPL)

- fit
- apply
- summary
- listmodels
- deletemodel
- sample

#### Macros

- regressionstatistics
- classificationstatistics

.Screen?product id=FL-DSH-01&JSH

- classificationreport
- confusionmatrix

- forecastviz
- histogram
- modvizpredict
- splitby(1-5)
- Viz
  - Outliers Chart
  - Forecast Chart
  - Scatter Line Chart
  - Histogram Chart
  - Downsampled Line Chart
  - Scatterplot Matrix



#### ML-SPL What is it?

- A suite of SPL commands specifically for machine learning
  - modeling
  - sampling
- Most are implemented using modules from the Python for Scientific Computing add-on for Splunk
  - scikit-learn
  - numpy
  - pandas
  - statsmodels
  - scipy



splunk

### **ML-SPL Commands**

► Fit (i.e., train) a model from search results

- Apply a model to obtain predictions from (new) search results
   ... | apply <MODEL>
- Inspect a model (e.g., display coefficients)
  summary <MODEL>

#### **ML-SPL Commands:** fit

optional

Examples:

- ... | fit LinearRegression
   system\_temp from cpu\_load fan\_rpm
   into temp\_model
- ... | fit KMeans k=10
  - downloads purchases posts days\_active visits\_per\_day
    into user\_behavior\_clusters



## **ML-SPL: Algorithms**

#### ► 25+ algorithms OotB

- prediction, clustering, forecasting, feature engineering
- Extensibility API for 300+ more
- Pipeline for advanced use cases
- ... | fit TFIDF message

...

- fit StandardScaler files bytes
- fit KMeans message\_tfidf\_\* SS\_\* k=5
- fit PCA message\_tfidf\_\* k=2



#### ML-SPL Commands: apply

... | apply <MODEL>

Examples:

- ... | apply temp\_model
- ... | apply user\_behavior\_clusters



#### **ML-SPL Commands:** summary

... | summary <MODEL>

Examples:

- ... | summary temp\_model
- ... | summary user\_behavior\_clusters



#### **ML-SPL Commands**

listmodels

deletemodel <MODEL>

[] "GET /Category.screen?category\_id=GIFTS&ISESSIONID=SOISLAFF10ADFF10 HTTP 1.1" 404 720 "http://buttercup-shopping.com/category.id=category\_id=GIFTS&ISESSIONID=SOISLAFF10ADFF10 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/category.id=GIFTS&ISESSIONID=SOISLAFF10ADFF10 HTTP 1.1" 406 3323 #ttp://buttercup-shopping.com/category.id=GIFTS&ISESSIONID=SOISLAFF10ADFF10 HTTP 1.1" 406 3323 #ttp://buttercup-shopping.com/category.id=GIFTS&ISESSIONID=SOISLAFF10ADFF10 HTTP 1.1" 406 3323 #ttp://buttercup-shopp



## ML-SPL Commands: sample

#### Randomly sample or partition events

... | sample <PARAMETERS>

#### ► Four modes

•	Ratio		sample	0.01
•	Count	•••	sample	20
•	Proportional		sample	<pre>proportional="some_field"</pre>
			-	

...

Partition





#### Plug: Machine Learning In SPL With The Machine Learning Toolkit

Tuesday, September 26<sup>th</sup> @ 2:15pm in Ballroom B

- Custom search commands!
- Tabular munging!
- Jacob Leverich!
- Shang Cai!







# The Assistants

guided machine learning



#### **Guided ML With The Assistants**

#### Guides you through an analytic

- Prepare, fit, validate, and deploy
- Automatically generates all the relevant SPL

Fit a model on all your data in search 🛽		×
inputlookup server_power.csv		
<pre>  fit StandardScaler "total-cpu-utilization", "total-disk- accesses", "total-disk-blocks", "total-disk-utilization", "total-instructions_retired", "total- last_level_cache_references", "total- memory_bus_transactions", "total-unhalted_core_cycles" with_mean=true with_std=true into example_server_power_StandardScaler_0</pre>	<pre>// apply preprocessing steps</pre>	
fit LinearRegression fit_intercept=true "ac_power" from "SS_*" into "example_server_power"	<pre>// fit and save a model using the entire dataset and provided parameters</pre>	



#### **Assistants: Prepare**

Prenrocessing Stens

reprocessing steps					
✓ StandardScaler				q	×
Preprocess method	Fields to preprocess	Standardize Fields			
StandardScaler	business_acres property_tax_rate	✓ with respect to mean ✓ with respect to standard deviation			
	distance_to_employment_center				
Apply ~ PCA	Fields to proprocess	K (# of Components)	0	đ	×
Preprocess method		R (# or components)			
PCA	<pre>* highway_accessibility_index * distance_to_employment_center</pre>	2			
	× pupil_teacher_ratio × crime_rate				
Apply					

404 3322

Y\_id=GIFTS&JSESSIONID=SDISL4FF10ADFF10 HTTP

/product.screen?product id=FL-DSH-01&JSESSIONID=SD15L4FF10ADFF10 /old1.screen?product id=FL-DSH-01&JSESSIONID=SD5SL7FF6ADFF9 / /old1.screen?product id=FL-DSH-01&JSESSIONID=SD5SL7FF6ADFF9



#### **Assistants: Fit**

Create New Model	Load Exist	ting Settings						
1: Enter a search								
inputlookup server_power.csv All time ~ (							Q	
→ → → → → → → → → → → → → → → → → → →						∕lode ∽		
2: Field to predict		3: Fields to use fo	or predicting		4: Split for training / test: 50 / 50	5: Save the model as		
ac_power	v	× total-cpu-utili	zation × total-disk-accesses	× total-disk-blocks	·	example_server_power		
		× total-disk-util	ization × total-instructions_re	tired				
Fit Model Oper	n in Search	Show SPL						



#### **Assistants: Validate**





1 //ategory.screen?category\_id=GIFTS&ISESSIONID=SDISL4FF10ADFF10 HTTP 1.1" 404 720 "http://buttercup-shopping.com/category.screen?category.screen?category.screen?category\_id=GIFTS&ISESSIONID=SDISL4FF10ADFF10 HTTP 1.1" 404 322 "http://buttercup-shopping.com/cattinopping.com/category.screen?category.screen?category\_id=GIFTS&ISESSIONID=SDISL4FF10ADFF10 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/cattinopping.com/category.screen?category\_id=GIFTS&ISESSIONID=SDISL4FF10ADFF10 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/cattinopping.com/category.screen?category\_id=GIFTS&ISESSIONID=SDISL4FF10ADFF10 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/category.screen?category\_id=GIFTS&ISESSIONID=SDISL4FF10ADFF10 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/category.screen?category.id=GIFTS&ISESSIONID=SDISL4FF10ADFF10 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/category.screen?category.id=GIFTS&ISESSIONID=SDISL4FF10ADFF10 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/category.id=GIFTS&ISESTIONID=SDISL4FF10ADFF10 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/category.screen?category.id=GIFTS&ISESTIONID=SDISL4FF10ADFF10 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/category.screen?category.id=GIFTS&ISESTIONID=SDISL4FF10ADFF10 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/category.screen?category.id=GIFTS&ISESTIONID=SDISL4FF10ADF

#### **Assistants: Deploy**



Creen?category\_id=GIFTS&JSESSIONID=SD1SL4FF10ADFF10 HTTP /product.screen?product\_id=GIFTS&lSESSIONID=SD1SL4FF10ADFF10HITF 1. T /olditate:

/oldlink?item

?6&JSESSIONID=SDSSI 9FFIADFF3 HTTP 1.1"





404 3322

200 1318



#### **The Assistants**

- Predict Numeric Fields
- Predict Categorical Fields
- Detect Numeric Outliers
- Detect Categorical Outliers
- Forecast Time Series
- Cluster Numeric Events





# What's New

since last .conf



© 2017 SPLUNK INC

## What's New

(since .conf 2016)

- Detect Numeric Outliers improvements
- Preprocessing
- Model Management
- ML-SPL extensibility API
- Spark Support (private beta)
- ► New algorithms:
  - ACF & PACF
  - ARIMA
- Load Existing Settings is per-user
- Downsampled Line Chart supports drilldown



## **Detect Numeric Outliers**

split-by support

Detect Nume Find values that differ sign	ric Outliers hificantly from previous values.				?	
Detect Outliers	Load Existing Settings					
Enter a search						
inputlookup sup	ermarket.csv   head 1000				All time ∽	Q
✓ 1,000 results (12/31/6)	9 4:00:00.000 PM to 8/11/17 4:26:57.000 PM)			Job∨ II I	Smart M	lode ∽
Field to analyze	Threshold method	Threshold multiplier	Sliding window (# of values)	Fields to split by		
quantity	Standard Deviation     The second secon	5	0 Include current point	× shop_id		
Detect Outliers	pen in Search Show SPL					



### **Detect Numeric Outliers**

#### data distribution viz

#### Data Distribution



#### Preprocessing

#### build a pipeline of data prep

#### ▶ In Predict Numeric, Predict Categorical, and Cluster Numeric assistants



✓ StandardScaler				q	×
Preprocess method		Fields to preprocess	Standardize Fields		
StandardScaler	Ŧ	business_acres property_tax_rate	with respect to mean  with respect to standard deviation		
		distance_to_employment_center			
Apply					
∼ PCA			0	0	×
Preprocess method		Fields to preprocess	K (# of Components)		
PCA	v	× highway_accessibility_index × distance_to_employment_center	2		
		× pupil_teacher_ratio × crime_rate			
Apply					
7/Jan 18:10:57:1531 "GET / Category.scr 0 [07/Jan 18:10:57:1231 "GET / Category.scr 18:10:57:1231 "GET / Produ 0004 [10:14.437:1256:156] [10:17]	een?category	/_id=GIFTS&JSESSIONID=SD1SL4FF10ADFF10 HTTP 1.1" 404 720 "http://buttercup-shopping.com/cat http://buttercup-shopping.com/cat microsciencescomercescome	int.dofaction=view&itemId=Est-g&product g.com/category.screenid=Est-g&product g.com/category.screenid=Est-g&product Tot.int 200 243 "http://dwshordwat.gav.gav.gav. Tot.int 200 243 "http://dwshordwat.gav.gav.gav.gav.gav.gav.gav.gav.gav.gav	.co	nf20

## **Model Management**

#### RBAC & more

- Assign permissions to models to control access
- Manage models via the UI
- Experiments
  - Assistant configurations
  - May produce 1+ models

Edit Permissions			×	
Model Title Model ID Owner App Allow access for	itle Buttercup Store Purchases ID Buttercup_ Store_Purchases ner admin Machine Learning for Owner App All Apps			
Role		Read	Write	
Everyone				
admin			$\checkmark$	
can_delete				
poweruser				
Splunk-System-Role				
user		$\checkmark$		



## **ML-SPL Extensibility API**

featuring: primo documentation

- Make more algos available to fit / apply
  - 300+ in PSC
  - Custom algorithms
- Expose new or different parameters
- Docs include examples
  - Correlation Matrix
  - Agglomerative Clustering
  - Support Vector Regressor
  - Savitzky-Golay Filter
- Use in your apps / dashboards / etc.!





#### Plug: Advanced Machine Learning using the Extensible ML API

Wednesday, September 27<sup>th</sup> @ 4:35pm in Ballroom B

- Implementation details!
- Extensibility API!
- Alexander Johnson!
- Zidong Yang!





splunk

## **Spark Support**

private beta open now

- Use your existing Spark cluster with MLTK
  - Distributed fit on massive datasets
  - Apply MLlib models for supported algos
- sfit / sapply
- Contact <u>mlprogram@splunk.com</u>
  - What is your use case (e.g., predicting server downtime)?
  - Why do you want / need Spark (i.e., why isn't MLTK sufficient)?



# Demo





# Q&A

Adam J. Oliner | Director of Engineering Manish Sainani | Director of Product



# Thank You

# Don't forget to rate this session in the .conf2017 mobile app

