# Splunk Machine Learning Toolkit In Action

Iman Makaremi, Splunk Principal Data Scientist
Andrew Stein, Splunk Principal PM for Machine Learning
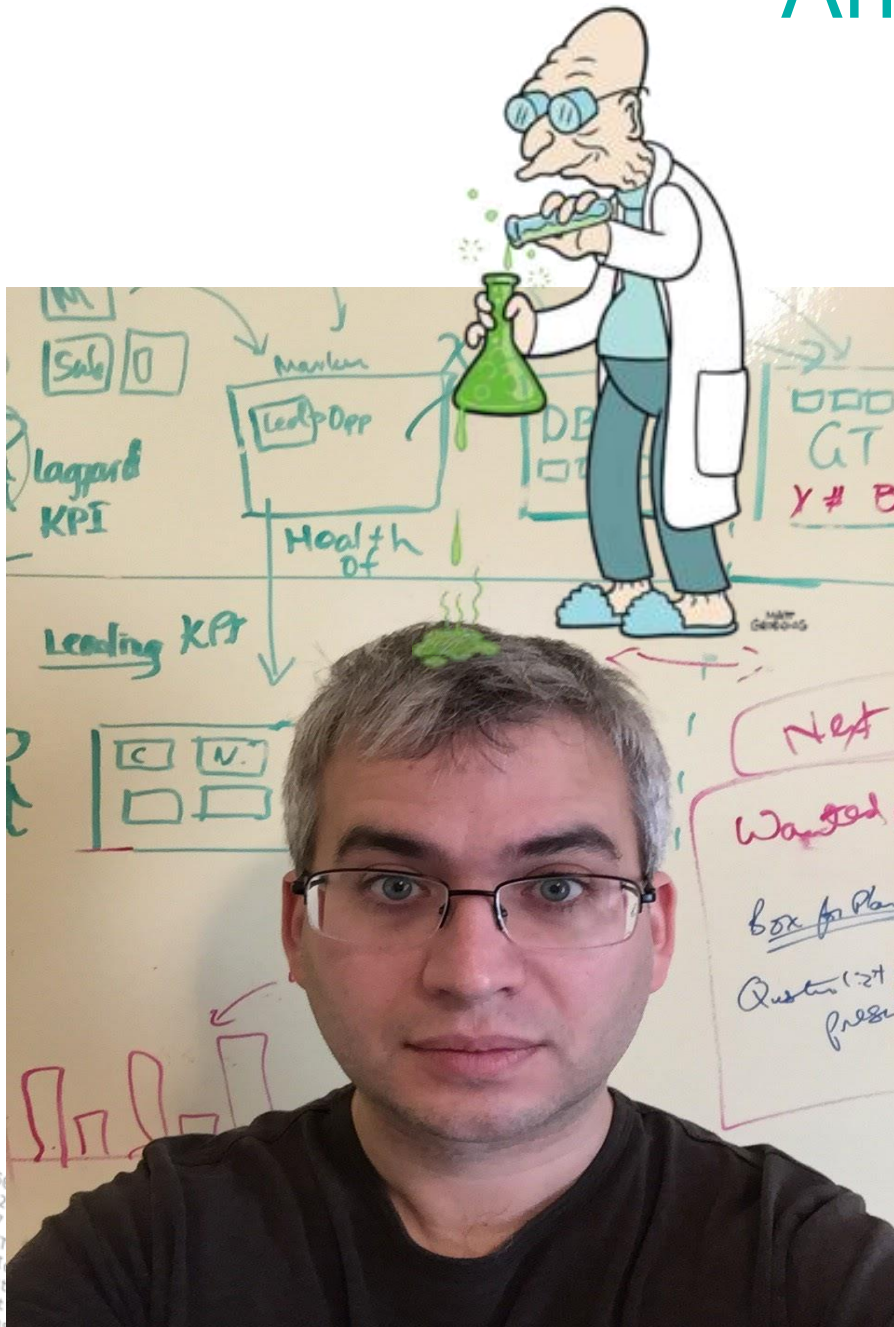
Oct 2018 | Version 1.0

# 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> .conf18

# Andrew Stein

- Splunk Principal Product Manager – Machine Learning
- 18 years creating mathematical modeled solutions as a data scientist.
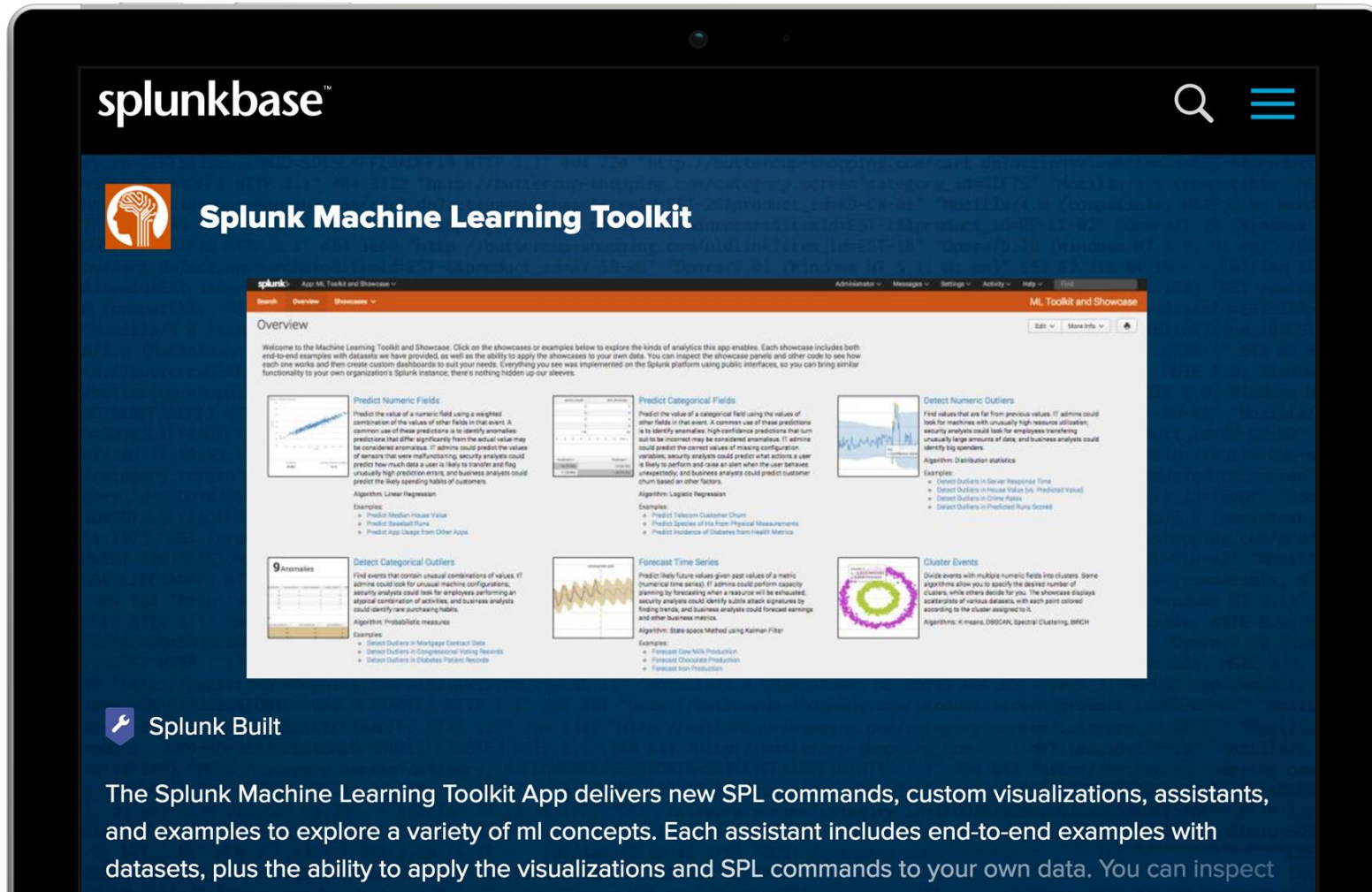- I spend 80 percent of time preparing data and 20 percent of time complaining about the need to prepare data.

splunk> .conf18

# Iman Makaremi

- Splunk Principal Data Scientist
- I like math and physics a lot and know a few things about them.
- I spend 80 percent of time preparing to complain and 20 percent of time complaining about the need to complain.

# Agenda

**Am I in the right room?**

▸ What?

▸ Assistant : Detect Numeric Outliers

- IOT & IT Example

▸ Assistant : Predict Numeric Fields

- Forecasts verse Prediction
- IT Example

▸ Assistant : Clustering

- Security Example

# What?

**Splunk has a Machine Learning Toolkit App!**



▸ What is Splunkbase

▸ What is the App

▸ Where can I go to learn more

splunk> .conf18

# Workflow. Workflow. Workflow. Bees. Or Math.

**Problem**: <Stuff in the world> causes big time and money expense.
**Solution:** Build ML model to learn the behaviors at scale and take action.

Get and **explore** data

> Select and **fit** an algorithm, generating a model

> **Apply** and **validate** models

> **Surface** model to consumers to solve problems

**Operationalize**

splunk>live!

# What is the ML Advisory Program?

## Partners a Splunk Data Science Resource to Help Operationalize a ML Use Case

### Machine Learning Customer Advisory Program FAQs

- What is the Machine Learning Customer Advisory Program? ⊕
- Are there examples from the advisory program? ⊕
- This program is free...what's the catch? ⊕
- This sounds interesting! How do I know if I qualify to apply? ⊕
- Anything else I should know? ⊕
- I meet the criteria and am interested in applying! What's next? ⊕
- I don't meet the criteria for the advisory program, but am interested in leveraging Splunk for machine learning. What options do I have? ⊕

▸ Early Access to new and enhanced MLTK features

▸ Opportunity to shape the development of the product

▸ Assistance in operationalizing a production quality ML model

splunk> .conf18

# Assistant :Detect Numeric Outliers

**Most Popular Assistant**

splunk> .conf18

☰ **About UNLV**

# Detecting and Resolving Data Outages

▸ *"Everything is Data, All of it is important."*

▸ Data

- Splunk index logs, Enrollment of Students through time
- Database logs, Any data where outages matter to you

▸ Action Taken with the Detect Numeric Outliers

- *Splunk Admin after taking the Splunk EDU Data Science course.*
- Detecting data source outages that are critical for supported research and operational centers, using custom seasonality
- Automatically Impute and replace missing information in summary index, send alert to administrator for further action.

splunk > listen to your data

# Mobility 3GPP Core KPI



‣ *Telus cell towers create a valuable and complicated network to maintain customer phone connectivity. The usability of our network by our customers is critical for our business.*

‣ Data

- The 3GPP Core receives transactions from each subscriber **to maintain connection**.
- The Telus custom KPI captures **the behavior of TELUS's network.**

‣ Action Taken *with a Custom Machine Learning Model in Splunk*

- Monitor this **dynamic** KPI and alert on **contextual performance degradation**.
- Radio engineers informed about deviations from expectations immediately, creating the opportunity for in-the-field technician corrective actions.

# Numerical Outlier Detection in MLTK

# Iterating over different threshold methods



Standard Deviation

Absolute Median Deviation

# Missed Outlier!

# Customized Outlier Detection

$x[t]$     Field to Monitor

$S$     Window Size

$H$     No. of Historical References

$T$     History Step Size

$c$     Confidence Interval Tuner

$P$     Vote Percentage

$$d_h[t] = \sum_{S=0}^{S} x[t-s] \cdot x[t-hT-s]$$

$$m_h = \underset{t}{\text{median}}\left( d_h[t] \right)$$

$$M_h = \underset{t}{\text{median}}\left( \left| d_h[t] - m_h \right| \right)$$

$$o_h = \begin{cases} 0 & m_h - c M_h < d_h[t] < m_h + c M_h \\ 1 & o.w. \end{cases}$$

$$\text{is-outlier} = \begin{cases} 0 & \frac{1}{H}\sum_{h=1}^{H} o_h < P \\ 1 & o.w. \end{cases}$$

130.60.4 - [07/Jan 18:10:57:153] "GET /category.screen?category_id=GIFTS&JSESSIONID=SD1SL4FF10ADFF10 HTTP 1.1" 404 720 "http://buttercup-shopping.com/cart.do?action=view&itemId=EST-6&product...
128.241.220.82 - [07/Jan 18:10:57:123] "GET /product.screen?product_id=FL-DSH-01&JSESSIONID=SD5SL7FF6ADFF9 HTTP 1.1" 404 3322 "http://butte.cup-shopping.com/category.screen?category_id=GIFTS...
317 27.160.0.0 - [07/Jan 18:10:56:156] "GET /oldlink?item_id=EST-26&JSESSIONID=SD5SL9FF1ADFF3 HTTP 1.1" 200 1318 "http://JSESSIONID=SD9SL4FF4ADFF7 HTTP 1.1" 200 2423 "http://buttercup-shopping...
ows NT 5.1: SV1: .NET CLR 1.1.4322)" 468 125.17 14...

# Evaluate

# Current vs Historic and Delta Calculation

# DETECTED!

# I want Anomalies…

## What is the minimum requirements for each workflow

▸ Detect Numeric Outliers Workflow

I have one number I care about, with possible seasonality (time) effects or some combination of identities (the "by" clause from stats for example).

**I want to find anomalies in one number moving through time.**

**How to Blog:**

*Statistical Anomalies and Forecasts (parts 1 ,2,3)*

▸ Or use many workflows…

I understand the statistics workflow, but I have many fields describing my problem or a measurable ground truth.

**I want to find complex anomalies and my data is organized.**

**How to Blog:**

*Anomalies like Neapolitan Ice Cream*

# Assistant: Predict Numeric Fields

**Most Misunderstood Assistant**

splunk> .conf18

# Predict vs Forecast

## English kind of sucks….

**predict** | prɪˈdɪkt |

verb *[with object]*

say or estimate that (a specified thing) will happen in the future or will be a consequence of something: *it is too early to predict a result* | *[with clause] : he predicts that the trend will continue* | *(as adjective **predicted**) : the predicted growth in road traffic.*

**forecast** | ˈfɔːkɑːst |

verb (past and past participle **forecast** or **forecasted**) *[with object]*

predict or estimate (a future event or trend): *rain is forecast for Scotland* | *[with object and infinitive] : coal consumption in Europe is forecast to increase.*

**Source: Mac's Dictionary**

splunk> .conf18

# Predict vs Forecast

**predict** | prɪˈdɪkt |

verb *[with object]*

say or estimate that (a specified thing) will happen in the future or will be a consequence of something: *it is too early to predict a result* | *[with clause] : he predicts that the trend will continue* | (as adjective **predicted**) : *the predicted growth in road traffic.*

**forecast** | ˈfɔːkɑːst |

verb (past and past participle **forecast** or **forecasted**) *[with object]*

predict or estimate (a future event or trend): *rain is forecast for Scotland* | *[with object and infinitive] : coal consumption in Europe is forecast to increase.*

The Splunk stock is influenced by interest rates, global economic conditions, road map, CFO's blood pressure, density of CEO's beard…

**≠**

The Splunk stock is cyclical, and every July stock price in the future will look like the July stock in the past +/- trending.

splunk> .conf18

# Predict vs Forecast

**predict** | prɪˈdɪkt |

**verb** *[with object]*

say or estimate that (a specified thing) will happen in the future or will be a consequence of something: *it is too early to predict a result* | *[with clause] : he predicts that the trend will continue* | (as adjective **predicted**) : *the predicted growth in road traffic.*

**forecast** | ˈfɔːkɑːst |

**verb** (past and past participle **forecast** or **forecasted**) *[with object]*

predict or estimate (a future event or trend): *rain is forecast for Scotland* | *[with object and infinitive] : coal consumption in Europe is forecast to increase.*
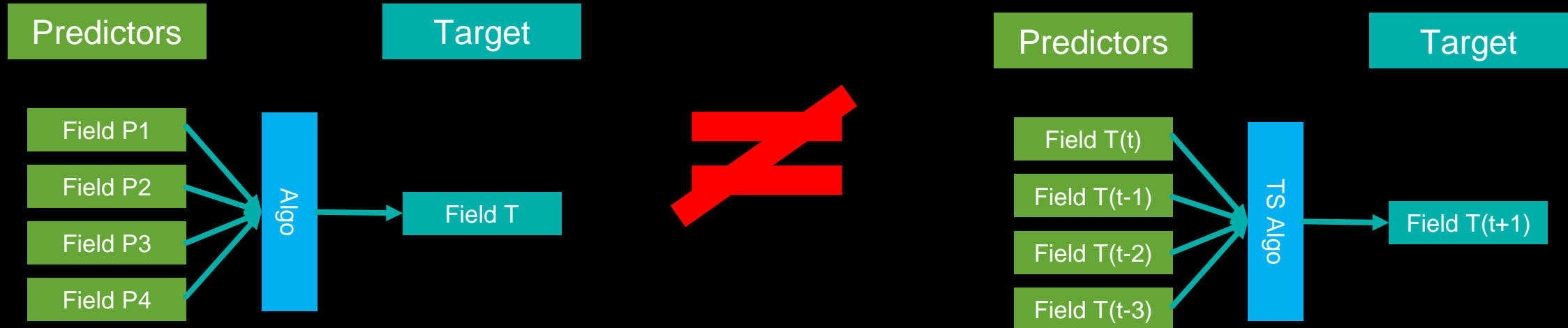
# Predict vs Forecast

**predict** | prɪ'dɪkt |

**verb** *[with object]*

say or estimate that (a specified thing) will happen in the future or will be a consequence of something: *it is too early to predict a result* | *[with clause] : he predicts that the trend will continue* | *(as adjective* **predicted***) : the predicted growth in road traffic.*

**forecast** | 'fɔːkɑːst |

**verb** (past and past participle **forecast** or **forecasted**) *[with object]*

predict or estimate (a future event or trend): *rain is forecast for Scotland* | *[with object and infinitive] : coal consumption in Europe is forecast to increase.*

MLTK Assistants
- Predict a Numeric Field
- Predict a Categorical Field

≠

Splunk
- *predict* Command

MLTK Assistant
- Forecast Time Series

*+ some optional time Travel SPL.*

How to Blog:

*ITSI and Sophisticated Machine Learning*

How to Blog:

*Statistical Anomalies and Forecasts (parts 1,2,3)*

splunk> .conf18

# TransUnion Invests in Splunk Solutions for Enterprise Monitoring, Machine Learning

" Understanding customer volume patterns is important for the business. If traffic falls outside of a certain range, an alert is created. Splunk machine learning allows us to investigate early to ensure a seamless customer experience."

– *Lead Splunk Developer, TransUnion*

▸ With the Splunk Machine Learning Toolkit and Splunk Machine Learning Customer Advisory Program Hyatt:

▸ Helping to meet customer SLAs

▸ Discovering incident root causes in minutes instead of hours

▸ Reducing the number of false alerts

▸ Increasing revenue by improving transaction processing

splunk> .conf18

# I want a prediction…

## What are the minimum requirements for each workflow

▶ Predict a Numeric Field

I want to generalize the relationship between one target **numeric** field and a series of descriptive fields.

I need a | table with the target **numeric** field and the descriptive fields , with _time if I am going to predict the future. The use of  | table is not required, this is just  good formatting step

How to Blog:

*Custom Anomaly Detection with Splunk IT Service Intelligence and Machine Learning Toolkit*

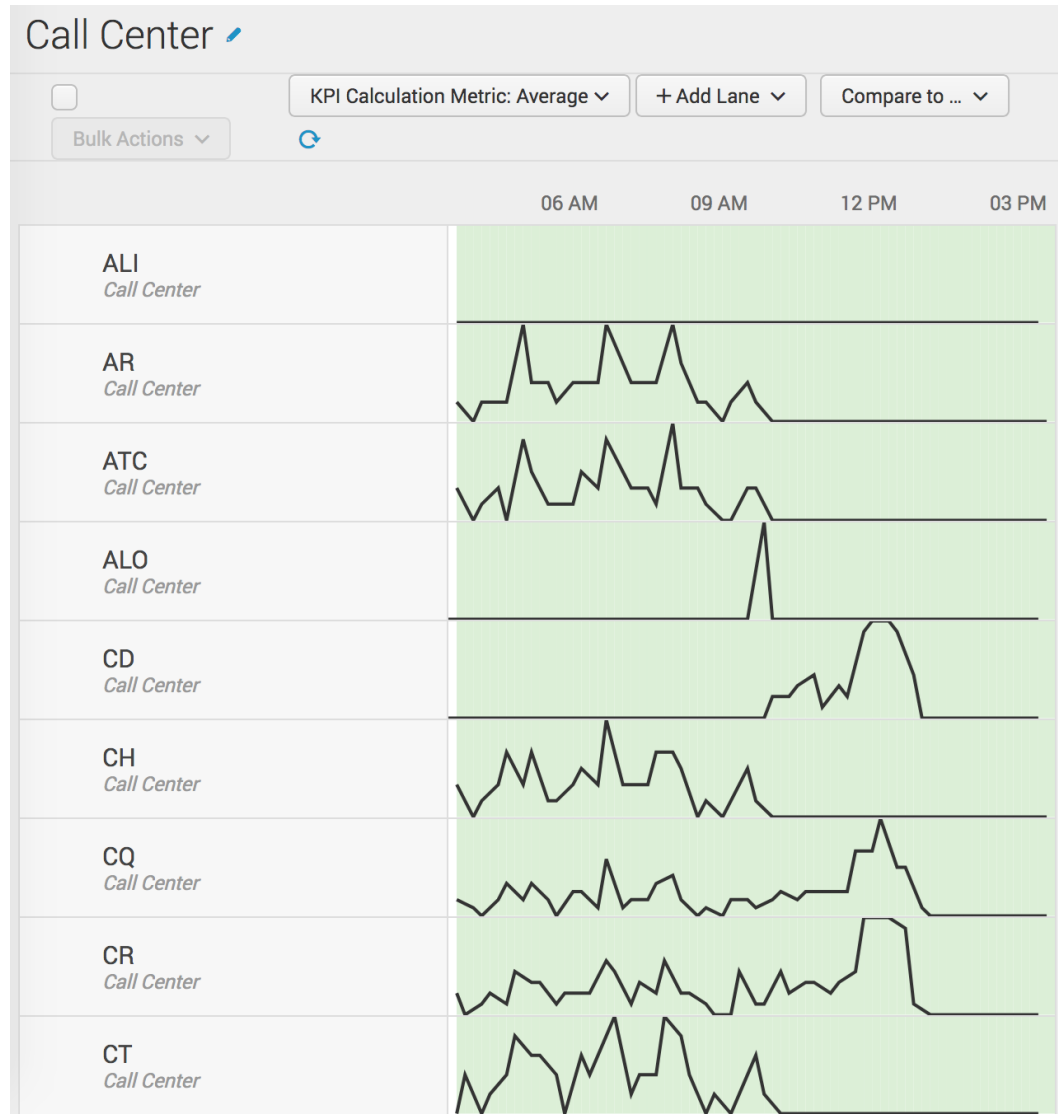▶ Predict a Future Value of a Field

I want to generalize the relationship between one target **numeric** field and a series of descriptive fields, but I want to have that relationship be explicitly in the future

I need a | table with the target **numeric** field and the descriptive fields , with _time if I am going to predict the future. The use of | table is not required, this is just  good formatting step. *I need to move the target field through time.*
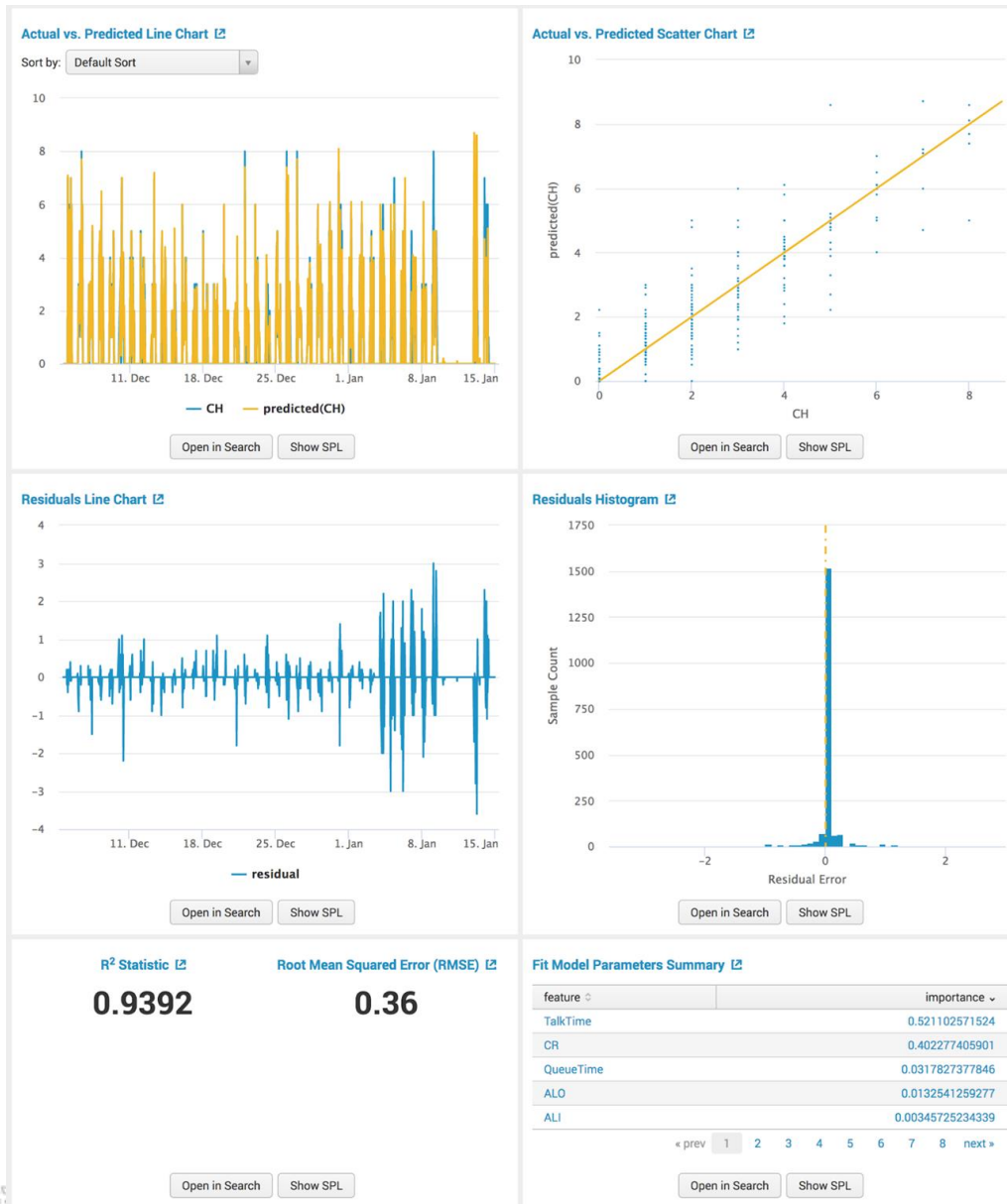
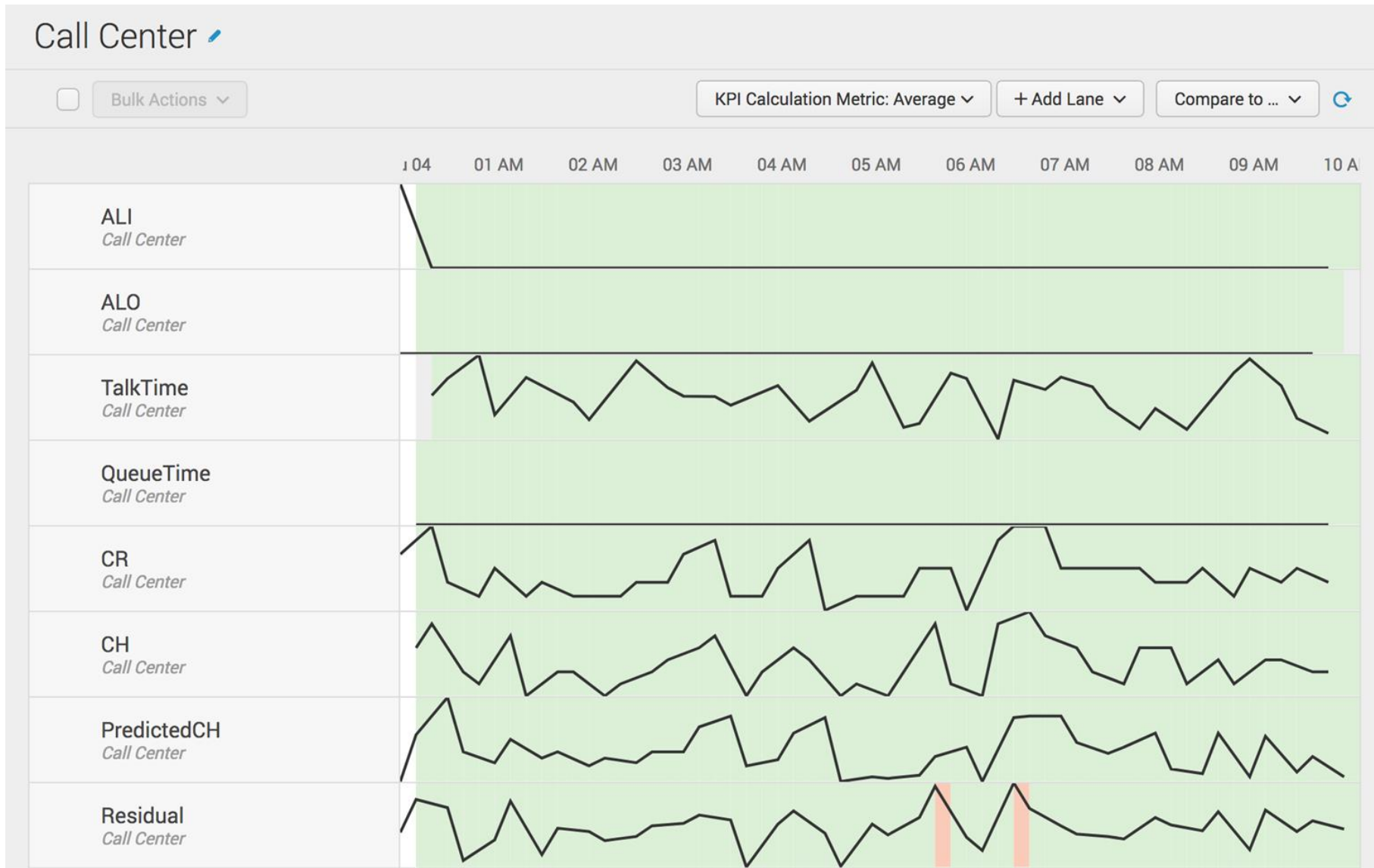How to Blog:

*ITSI and Sophisticated Machine Learning*

splunk> .conf18

# Customer Call Center

# Customer Call Center

# Customer Call Center

# Assistant: Clustering

**Sometimes you need to listen to your data!**

splunk> .conf18

# I have data but… ?

## What is the minimum requirements for each workflow

▸ ## Clustering

I have fields of data but I don't have a target field to generalize a relationship. I just want to know what rows are similar or dissimilar and by how much?

I need to create a | table with the fields I plan on using. I **should** really scale the fields first with StandardScaler or RobustScaler in the preprocessing step, and I **should consider** using PCA to reduce the dimensions pre clustering, and I **should** convert text fields to meaningful numeric values.

The use of |table is not required, this is just good formatting step

App, Videos, and How to Blog:

*DGA app on Splunkbase*

*DGA videos on Splunk Videos*

*Anomalies like Neapolitan Ice Cream*

130.60.4 - - [07/Jan 18:10:57:153] "GET /category.screen?category_id=GIFTS&JSESSIONID=SD1SL4FF10ADFF10 HTTP 1.1" 404 720 "http://buttercup-shopping.com/cart.do?action=view&itemId=EST-6&product_id=FL-SW-01 128.241.220.82 - - [07/Jan 18:10:57:123] "GET /product.screen?product_id=FL-DSH-01&JSESSIONID=SD5SL7FF6ADFF9 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-26&product_id=GIFTS 317 27.160.0.0 - - [07/Jan 18:10:56:156] "GET /oldlink?item_id=EST-18&JSESSIONID=SD9SL4FF4ADFF7 HTTP 1.1" 200 2423 "http://buttercup-shopping.com ows NT 5.1; SV1; .NET CLR 1.1.4322) "GET /oldlink?item_id=FL-DSH-01&JSESSIONID=SD5SL9FF1ADFF3 HTTP 1.1" 200 1318 "http://JSESSIONID=changequantity&itemId=EST-6&JSESSIONID=SD10SLBFF2ADFF9

# Now What?

So you want some ML ?

splunk> .conf18

# How do you replicate at your company?

**Problem**: <Stuff in the world> causes big time and money expense.
**Solution:** Build ML model to learn the behaviors at scale and take action

Get and **explore** data

Select and **fit** an algorithm, generating a model

**Apply** and **validate** models

**Surface** model to consumers to solve problems

**Operationalize**

splunk>live!

© 2018 SPLUNK INC.

# END of LINE

**Q&A**

splunk> .conf18