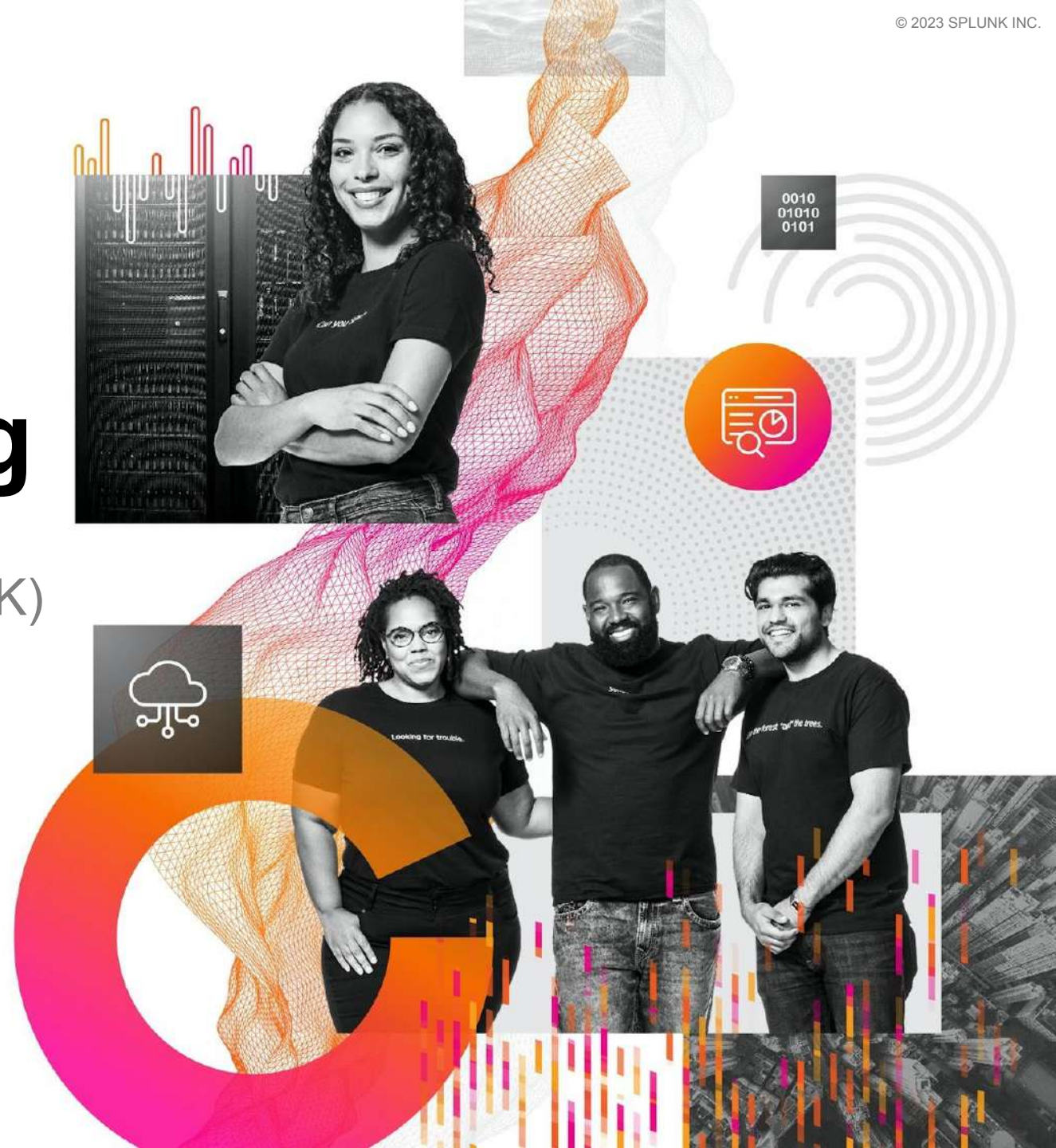# Splunk4Ninjas - Machine Learning

Hands on Introduction to the
Splunk Machine Learning Toolkit (MLTK)

splunk>

# Forward-Looking Statements

This presentation may contain forward-looking statements regarding future events, plans or the expected financial performance of our company, including our expectations regarding our products, technology, strategy, customers, markets, acquisitions and investments. These statements reflect management's current expectations, estimates and assumptions based on the information currently available to us. These forward-looking statements are not guarantees of future performance and involve significant risks, uncertainties and other factors that may cause our actual results, performance or achievements to be materially different from results, performance or achievements expressed or implied by the forward-looking statements contained in this presentation.

For additional information about factors that could cause actual results to differ materially from those described in the forward-looking statements made in this presentation, please refer to our periodic reports and other filings with the SEC, including the risk factors identified in our most recent quarterly reports on Form 10-Q and annual reports on Form 10-K, copies of which may be obtained by visiting the Splunk Investor Relations website at www.investors.splunk.com or the SEC's website at www.sec.gov. The forward-looking statements made in this presentation are made as of the time and date of this presentation. If reviewed after the initial presentation, even if made available by us, on our website or otherwise, it may not contain current or accurate information. We disclaim any obligation to update or revise any forward-looking statement based on new information, future events or otherwise, except as required by applicable law.

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. We undertake no obligation either to develop the features or functionalities described, in beta or in preview (used interchangeably), or to include any such feature or functionality in a future release.

**splunk>**

# Agenda

- Welcome/Introduction

- Intro to Machine Learning at Splunk

- Demo of Machine Learning Toolkit (MLTK) with Q&A

- Intro to the Trackday Dataset

- Four Different Challenges (~30 mins each)

  **Challenge 1** - Explore the track_day.csv Dataset

  **Challenge 2** - Detect Numeric Outliers

  **Challenge 3** - Supervised Learning: Predict Categorical Fields

  **Challenge 4** - Unsupervised Learning: Clustering

- Wrap Up, Discussion and Feedback

splunk>

source

© 2023 SPLUNK INC.

# Disclaimer

What this session is <u>not about</u> and what it <u>is about</u>

- **NO** replacement for a PhD in machine learning, data science or AI

- **NO** replacement for Splunk's Education class for Data Science

- **NO** comprehensive lecture about all possible concepts and algorithms in ML … but,

- **YES** first introduction into Machine Learning @ Splunk

- **YES** getting to know of Splunk's Machine Learning Toolkit

- **YES** guided hands-on challenges to explore a few typical ML tasks
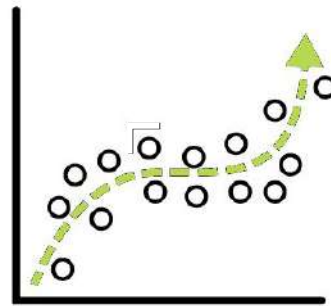
splunk>

# Machine Learning Tour

splunk>

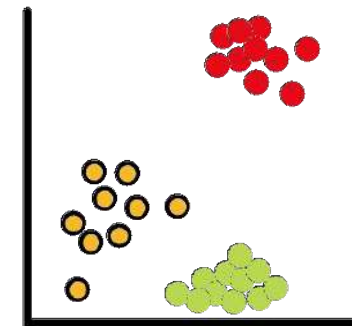# Splunk Customers Want Answers from their Data

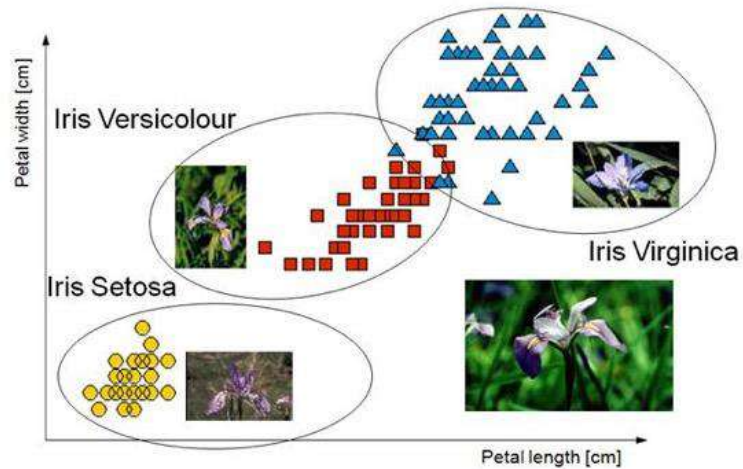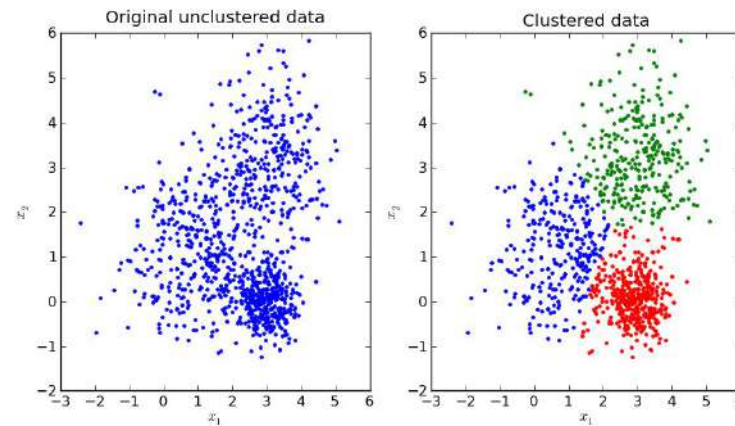| Anomaly detection | Predictive Analytics | Clustering |
|---|---|---|
|  |  |  |
| ► Deviation from past behavior<br>► Deviation from peers<br>► (aka Multivariate AD or Cohesive AD)<br>► Unusual change in features<br>► **ITSI MAD** | ► Predict Service Health Score/Churn<br>► Predicting Events<br>► Trend Forecasting<br>► Detecting influencing entities<br>► Early warning of failure | ► Identify peer groups<br>► Event Correlation<br>► Reduce alert noise<br>► Behavioral Analytics<br>► **ITSI Event Analytics** |

splunk>

# Types of Machine Learning

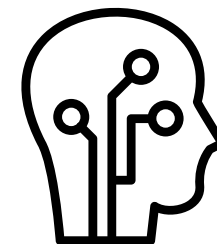| Supervised Learning (labeled data) | Unsupervised Learning (unlabeled data) | Mixed Models (with reinforcement or feedback) |
|---|---|---|
| ▶ Regression<br>▶ Classification | ▶ Clustering<br>▶ Anomaly Detection | ▶ Human in the Loop<br>▶ Autonomous Systems |



splunk>

# Overview of Machine Learning at Splunk
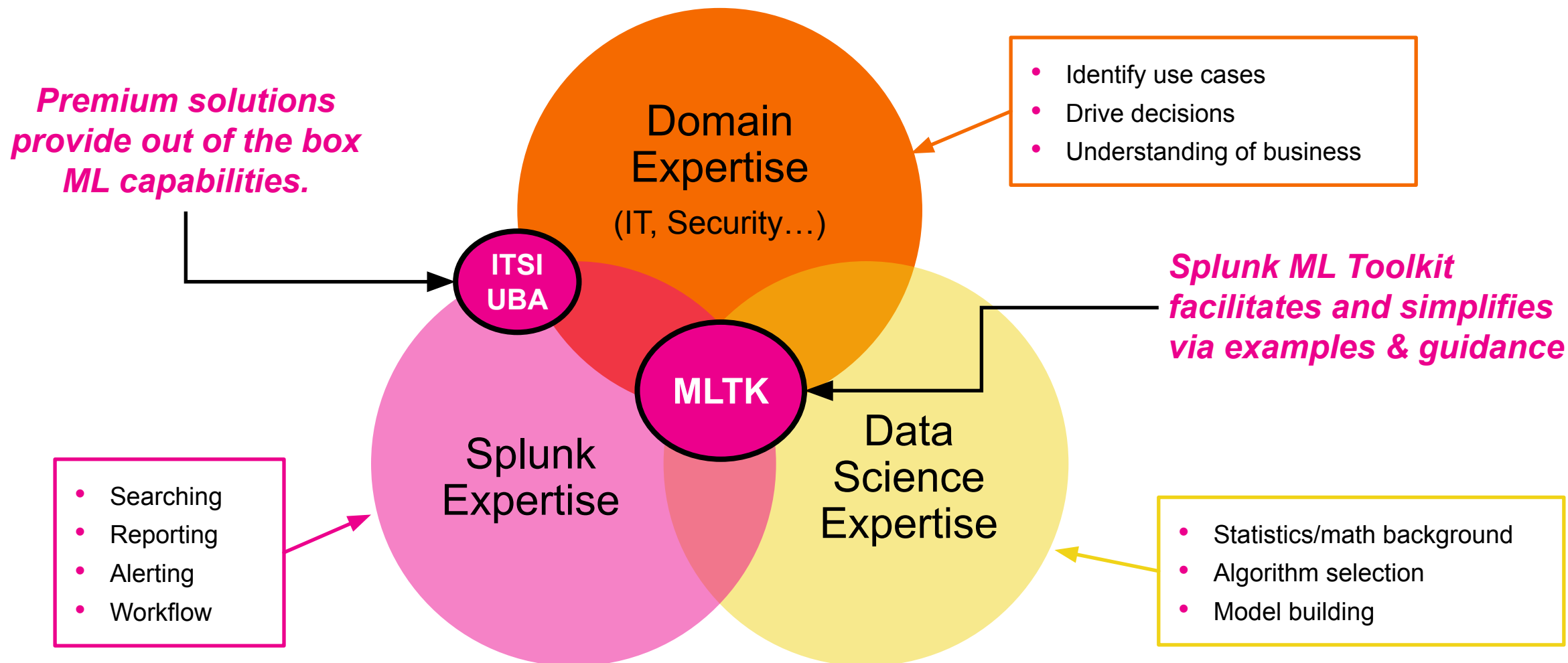
**CORE PLATFORM SEARCH**

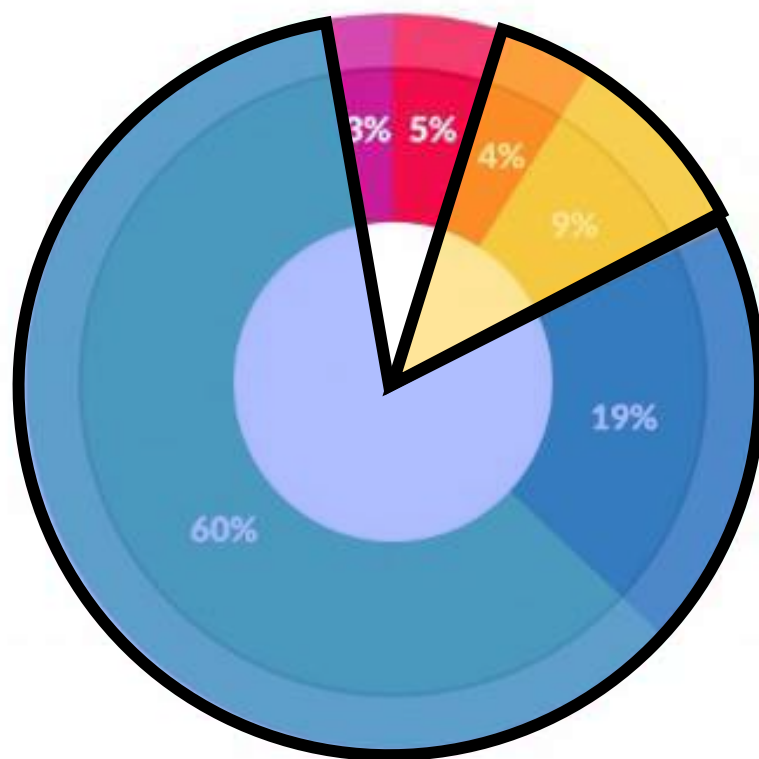**PACKAGED PREMIUM SOLUTIONS**

**MACHINE LEARNING TOOLKIT**

**splunk>** Platform for Operational Intelligence

# Skill Areas for Machine Learning @ Splunk

*Premium solutions provide out of the box ML capabilities.*

Domain Expertise (IT, Security…)

ITSI UBA

- Identify use cases
- Drive decisions
- Understanding of business

*Splunk ML Toolkit facilitates and simplifies via examples & guidance*

MLTK

Splunk Expertise

Data Science Expertise

- Searching
- Reporting
- Alerting
- Workflow

- Statistics/math background
- Algorithm selection
- Model building

splunk>

# What Data Scientists Really Do

**Data Preparation** accounts for about 80% of the work of data scientists



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%

"Cleaning Big Data: Most Time-Consuming, Least Enjoyable Data Science Task, Survey Says", Forbes Mar 23, 2016

splunk>

# Custom ML with the Splunk Platform

**Ecosystem** — Splunk's App Ecosystem contains 1000's of free add-ons for getting data in, applying structure and visualizing your data giving you faster time to value.
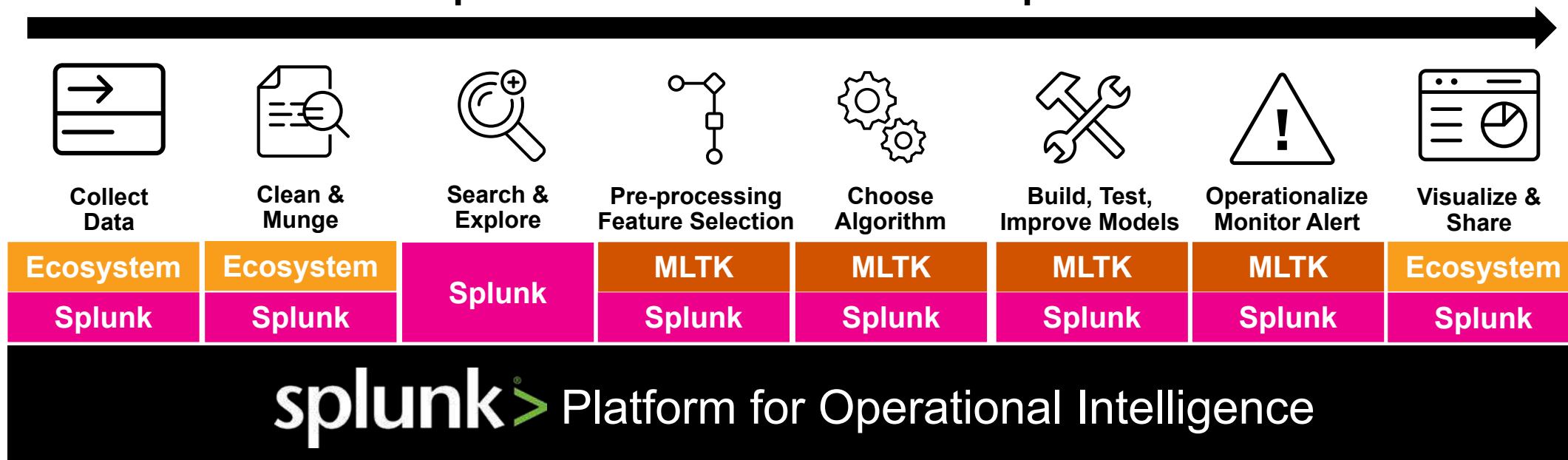
**MLTK** — The Machine Learning Toolkit delivers new SPL commands, custom visualizations, assistants, and examples to explore a variety of ml concepts.
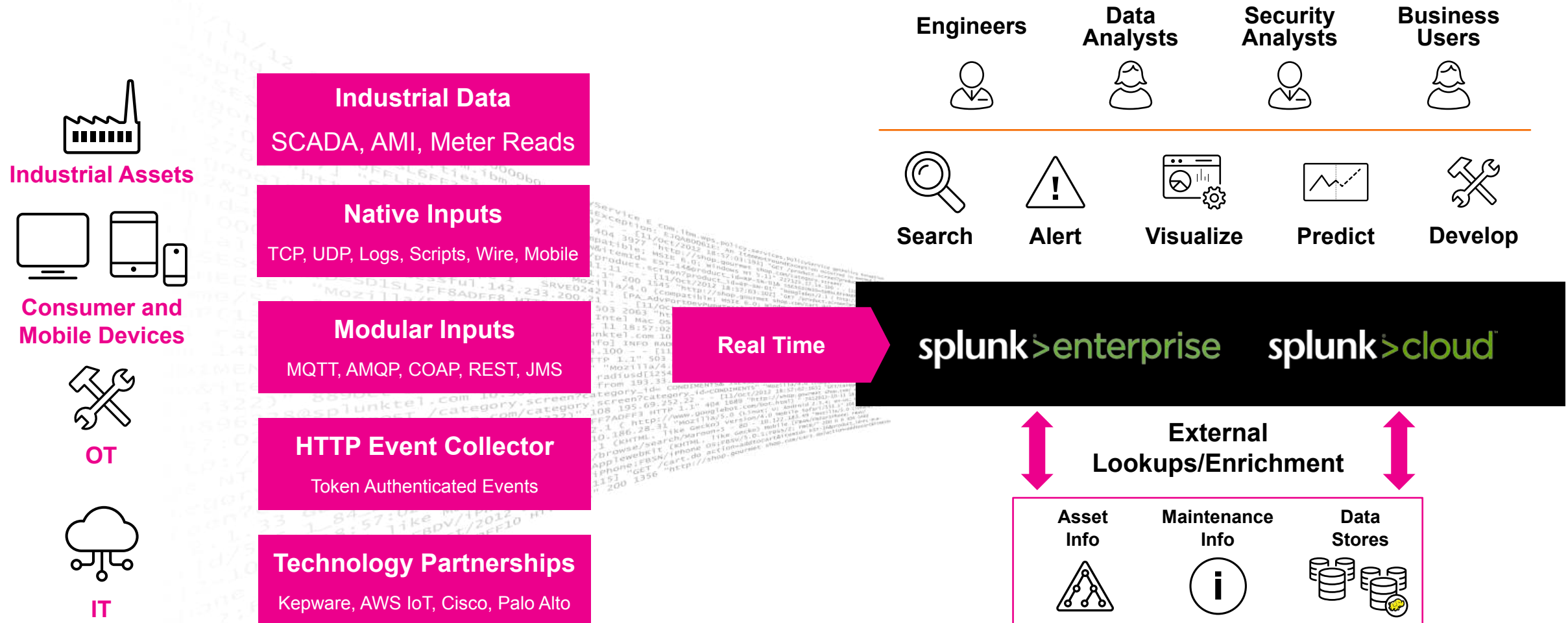
**Splunk** — Splunk Enterprise is the mission-critical platform for indexing, searching, analyzing, alerting and visualizing machine data.
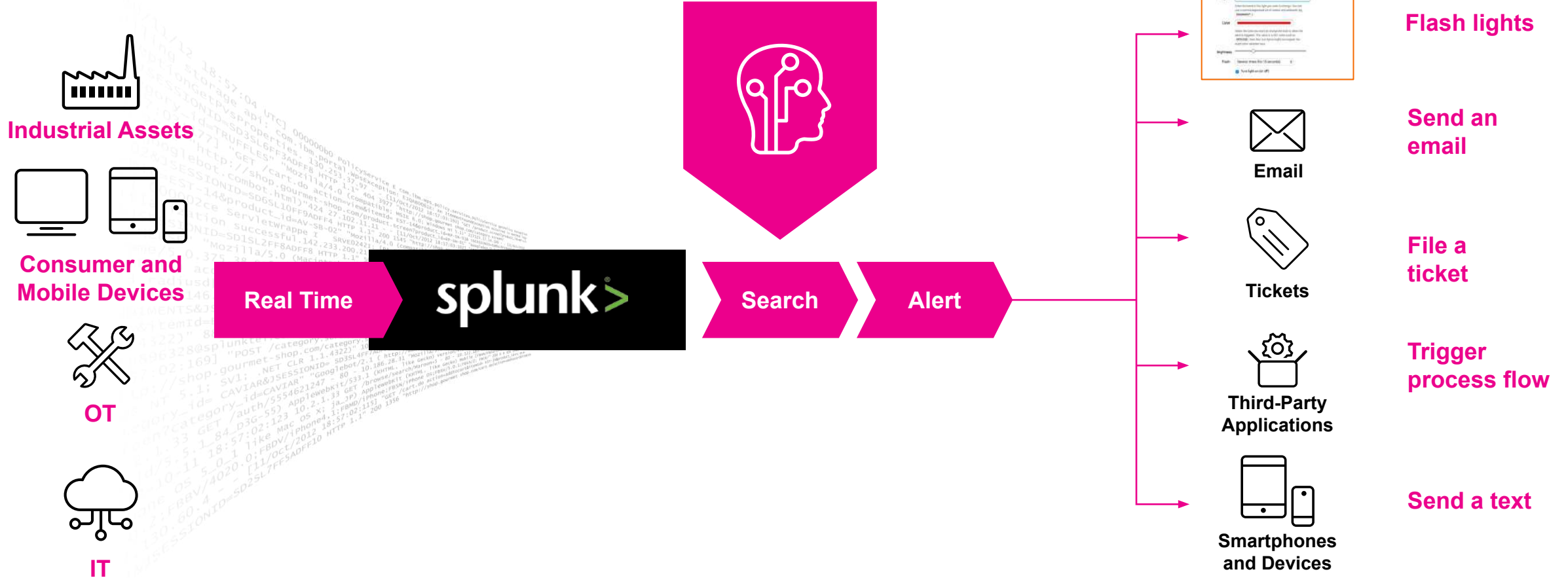
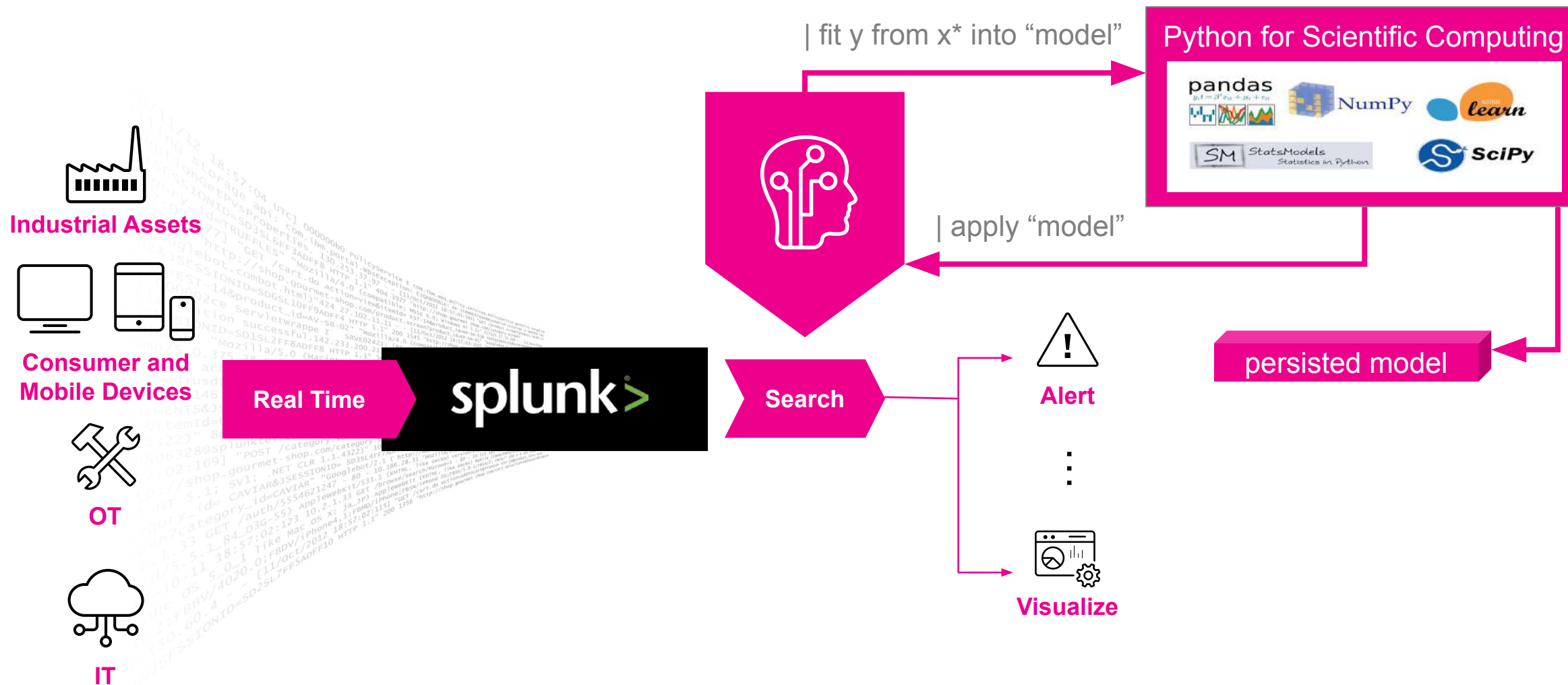## Operationalized Data Science Pipeline

| Collect Data | Clean & Munge | Search & Explore | Pre-processing Feature Selection | Choose Algorithm | Build, Test, Improve Models | Operationalize Monitor Alert | Visualize & Share |
|---|---|---|---|---|---|---|---|
| Ecosystem | Ecosystem | Splunk | MLTK | MLTK | MLTK | MLTK | Ecosystem |
| Splunk | Splunk | | Splunk | Splunk | Splunk | Splunk | Splunk |

**splunk> Platform for Operational Intelligence**

splunk>

# Continuous Data Ingest at Scale

© 2023 SPLUNK INC.

**Industrial Assets**

**Consumer and Mobile Devices**

**OT**

**IT**

**Industrial Data**
SCADA, AMI, Meter Reads

**Native Inputs**
TCP, UDP, Logs, Scripts, Wire, Mobile

**Modular Inputs**
MQTT, AMQP, COAP, REST, JMS

**HTTP Event Collector**
Token Authenticated Events

**Technology Partnerships**
Kepware, AWS IoT, Cisco, Palo Alto

**Real Time**

splunk>enterprise   splunk>cloud

Engineers   Data Analysts   Security Analysts   Business Users

Search   Alert   Visualize   Predict   Develop

**External Lookups/Enrichment**

Asset Info   Maintenance Info   Data Stores

splunk>

# Sense and Respond

Every Search Can
Use Machine Learning

Industrial Assets

Consumer and
Mobile Devices

OT

IT

Real Time

splunk>

Search

Alert

Flash lights

Email
Send an
email

Tickets
File a
ticket

Third-Party
Applications
Trigger
process flow

Smartphones
and Devices
Send a text

splunk>

# MLTK + Python for Scientific Computing

| fit y from x* into "model"

Python for Scientific Computing

pandas  NumPy  learn

SM StatsModels Statistics in Python  SciPy

| apply "model"

Industrial Assets

Consumer and Mobile Devices

Real Time  **splunk>**  Search

Alert

⋮

Visualize

persisted model

OT

IT

splunk>

# Splunk Machine Learning Toolkit (MLTK)

Extends Splunk platform functions and provides a guided modeling environment

## Built for the Citizen Data Scientist

- **Experiments and Assistants**: Guided model building, testing, and deployment for common objectives
- **Algorithms**: 80+ standard algorithms (supervised & unsupervised)

## Extensible to operationalize any use case

- **Python for Scientific Computing Library**: Access to 300+ open source algorithms
- **Deep Learning Toolkit** : Supports NN and GPU accelerated machine learning
- **ML-SPL API**: Import any open-source or proprietary algorithm



splunk>

# Example: MLTK powered DGA App for Splunk

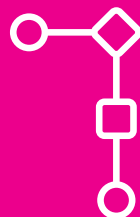Detect Malicious Domain Names using Machine Learning

# Overview of ML including DL at Splunk

(not covered in this workshop)

**CORE PLATFORM SEARCH**

**PACKAGED PREMIUM SOLUTIONS**

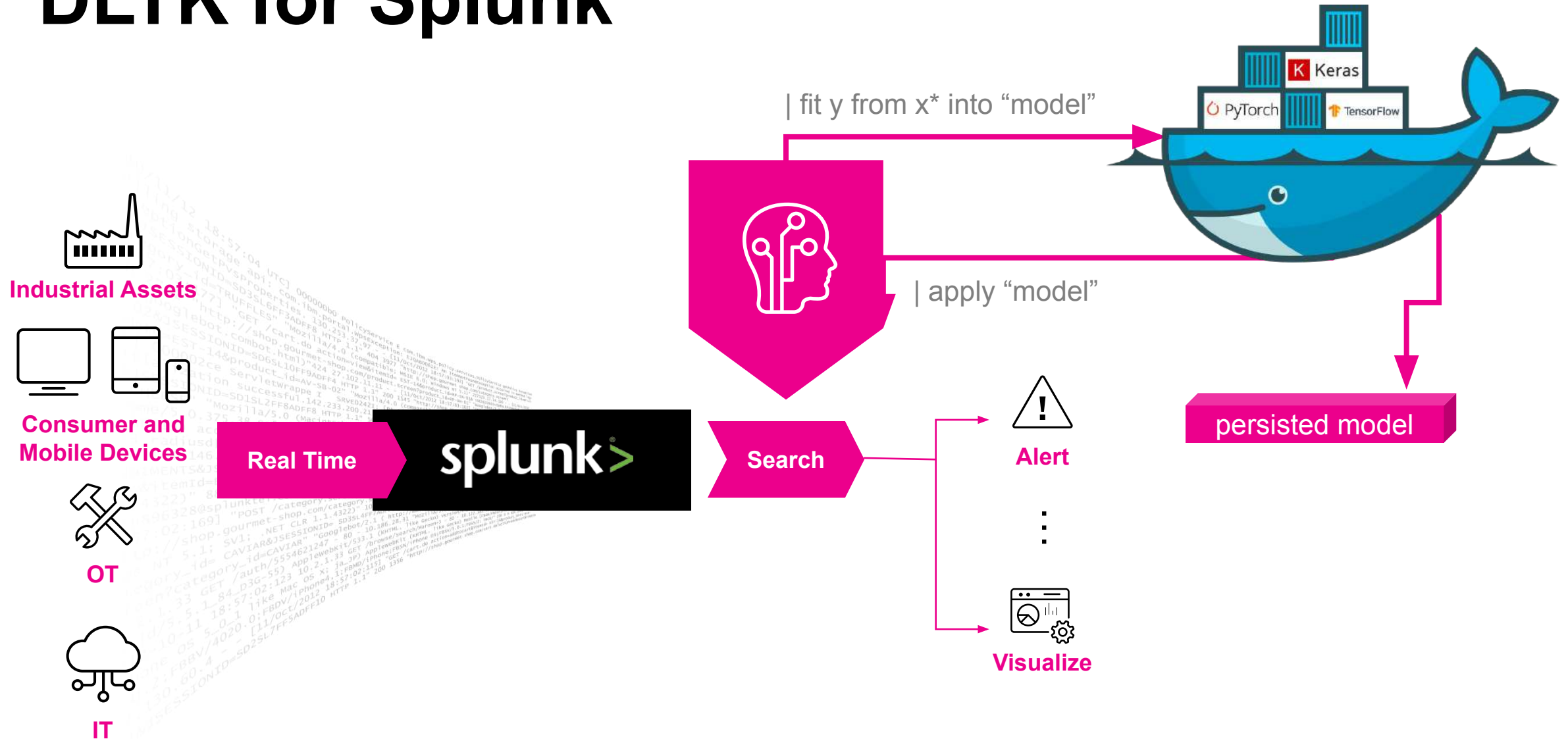**MACHINE LEARNING TOOLKIT**

**DEEP LEARNING TOOLKIT**

**splunk>** Platform for Operational Intelligence

splunk>

# DLTK for Splunk

Keras

PyTorch  TensorFlow

| fit y from x* into "model"

Industrial Assets

| apply "model"

Consumer and
Mobile Devices

Real Time  **splunk>**

persisted model

Search

Alert

OT

Visualize

IT

splunk>

# Hyatt ensured <u>every</u> customer experience with hotel wifi across the globe

1. All the data showing every customer sign in to the Hotel Wifi programs along with other relevant data (provider and lookup file of provider info, hotel ID and lookup file of hotel info, indexes for real time data of customers, local and global holidays).

2. A workflow (dashboard(s)) to show customers at every hotel signing in, normalize _time so a hotel in one time zone at 8 am to can be compared to another hotel at 8am.

3. Forecasting the likely wifi logins based on each property, day of week, and local and global holiday out two days into the future to show our expectations to executives.

4. Detecting meaningful anomalies as real time data comes in and is compared to the forecast . I must be able to insert business rules into the anomalies based on analyst feedback in a quick and nimble way. I should update my learning every night.
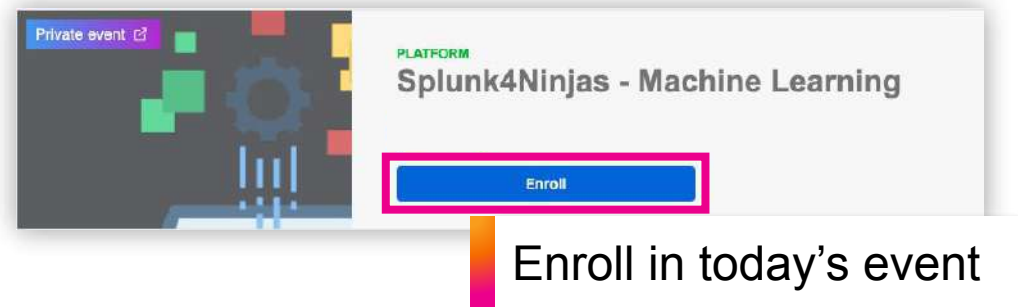
splunk>

# Demo: Machine Learning Toolkit

splunk>

# Hands-on Challenges

splunk>

# Enroll in Today's Workshop

## Tasks

1. Get a splunk.com account if you don't have one yet:
https://splk.it/SignUp

2. Enroll in the Splunk Show workshop event:
https://show.splunk.com/event/<eventID>

3. Download a copy of today's slide deck:
https://splk.it/S4N-ML-Attendee

## Goal



Enroll in today's event

# Fun Facts about the Track Day dataset

A popular private event of racing and sportscar affine Splunkers in the early days.

**Simple concept**

Go on a race track, have fun and collect some car data to get insights about driving behavior etc.

A subset of this data is available in MLTK!



Image Source: https://www.youtube.com/watch?v=meBjI-ay9-U

splunk>

# Today's Challenges

We are going to create four dashboards:

**1** **Explore the Dataset:** Create a sample dataset and explore it using different types of visualizations such as SPL

**2** **Detect Numeric Outliers:** Explore the MLTK showcase and adapt it to start a new experiment with your own dataset

**3** **Use a Classification Model:** Create a classification model and use it to predict vehicle types from your sensor data
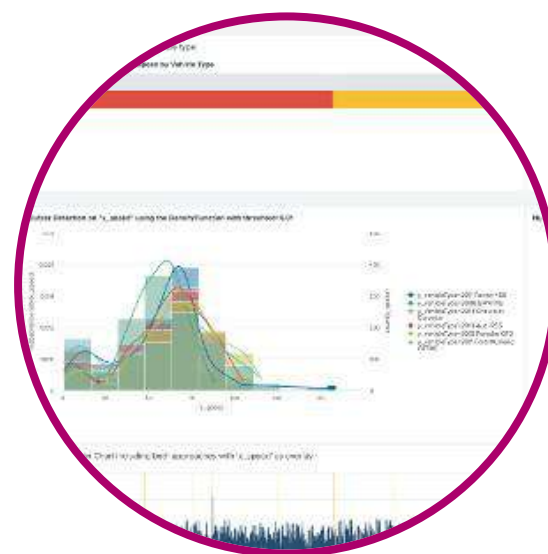
**4** **Use a Clustering Model:** Create a clustering model and and use it to analyze your dataset



We're aiming for a dashboard like this!

splunk>

# Workshop Goals

- Getting to know Splunk in the context of Machine Learning

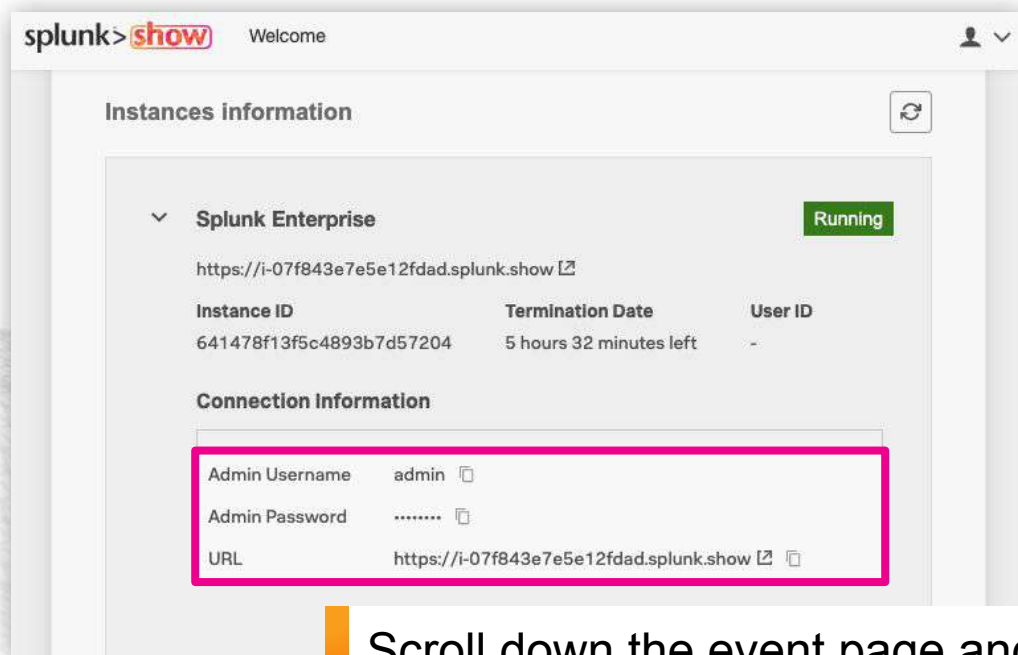- Prepare and analyze a dataset and summarize results on 4 dashboards



splunk>

# Login to Splunk

Locate your instance URL and credentials
in the Splunk Show event
https://show.splunk.com

Log in to your Splunk instance



Scroll down the event page and
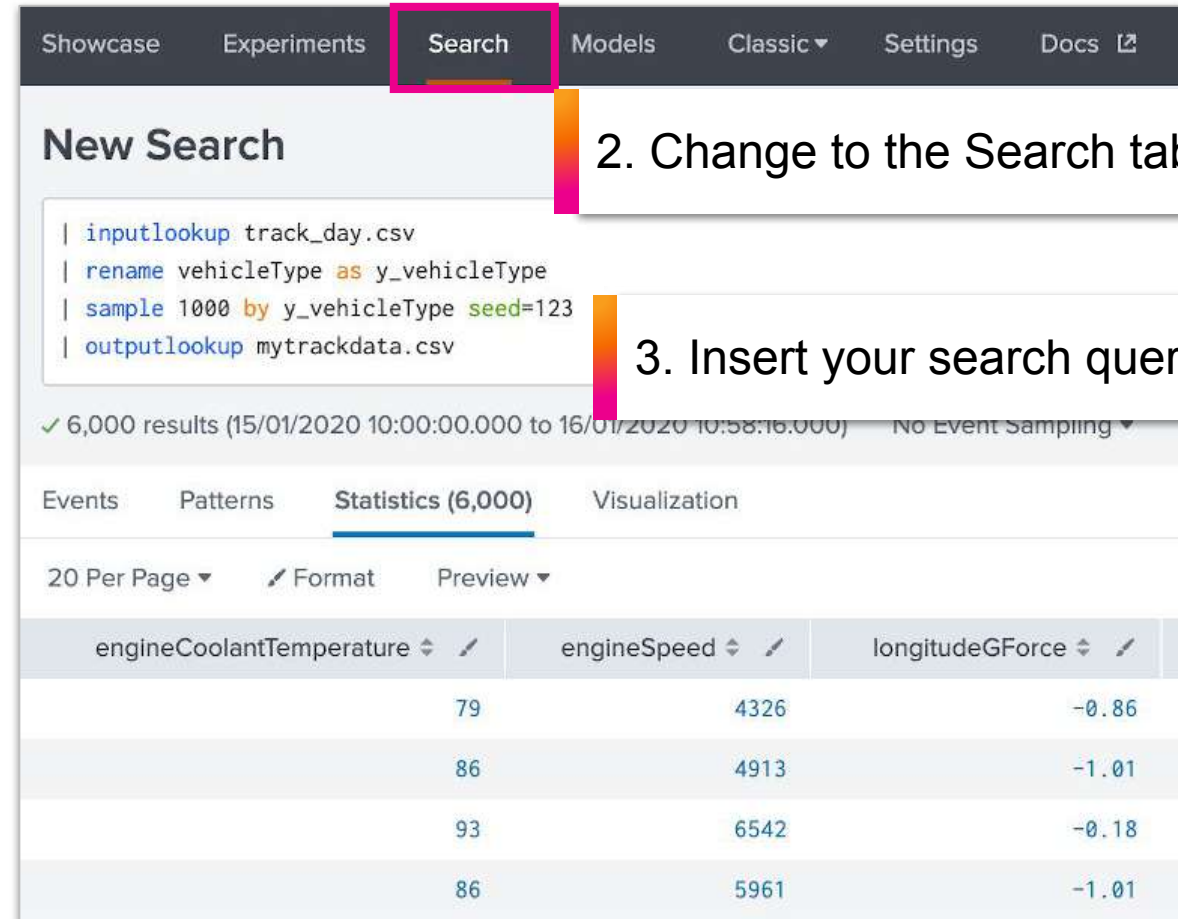expand the **Splunk Enterprise**
section to view your login
details

Login using the credentials
from Splunk Show

# Create a Sample Dataset

**1**

1. Access the Splunk Machine Learning Toolkit

**?** What's the benefit of renaming variables?

2. Change to the Search tab

3. Insert your search query

# ① Use Fieldsummary to Explore your Dataset



Eliminate unwanted Fields with
`| fields - values`

**?**

What's going on with the engine coolant temperature?

# Explore your Dataset with Visualizations

# 1

# Using Splunk's Histogram Macro

**OR**

Check Macro in settings

Check Macro with
Cmd + Shift + E (Mac) or
Ctrl + Shift + E (Windows)

splunk>

# Adjusting the Histogram Macro

**1**



**?** How can we get from the top to the bottom histogram?



splunk>

# Adjust the Macro to Split by Vehicle Type

### 1

`| chart count by y_vehicleType`

| vehicleType | count |
|---|---|
| Ferrari | 641 |
| Audi | 42 |
| BMW | 51 |
| Chevrolet | 44 |
| Ford | 95 |

`| chart count over x_engineSpeed by y_vehicleType`

| batteryVoltage | Ferrari | Audi | BMW | Chevrolet | Ford |
|---|---|---|---|---|---|
| 13 | 0 | 0 | 0 | 0 | 1 |
| 14 | 0 | 0 | 0 | 1 | 1 |
| 15 | 1 | 0 | 1 | 0 | 0 |
| 16 | 1 | 1 | 1 | 0 | 0 |
| 17 | 1 | 0 | 0 | 0 | 1 |

splunk>

# Working with the Boxplot Macro

**1**

**?** How can this query be improved?

```
| inputlookup "mytrackdata.csv"
| `boxplot`
```

Visualization

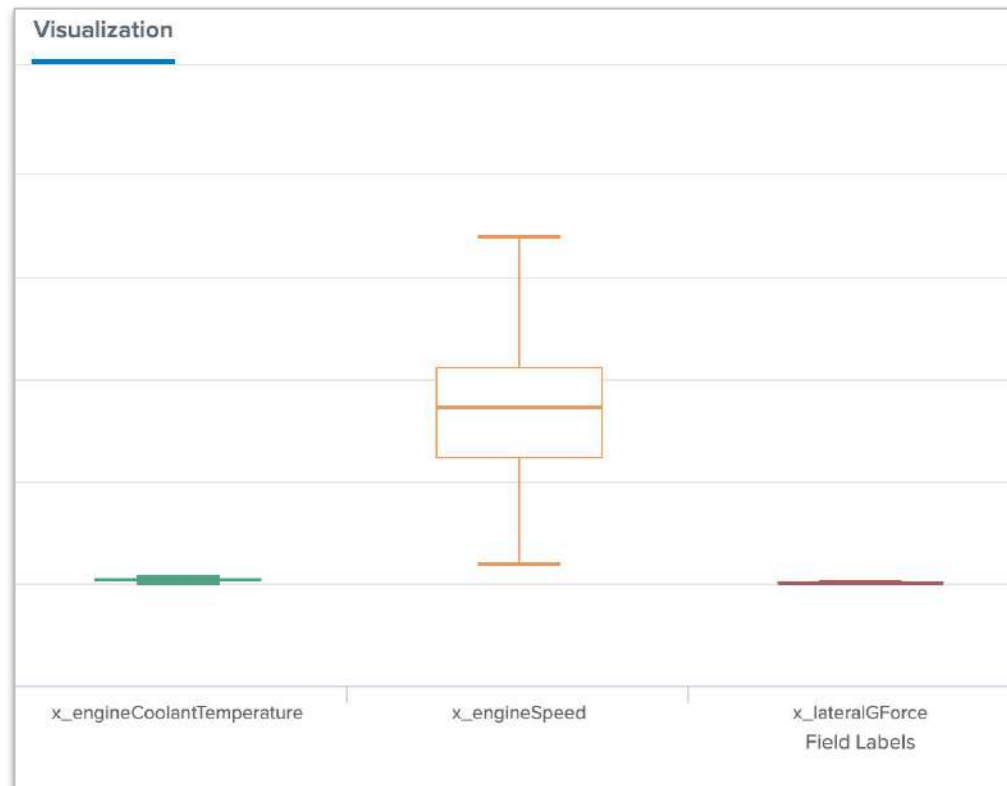x_engineCoolantTemperature    x_engineSpeed    x_lateralGForce
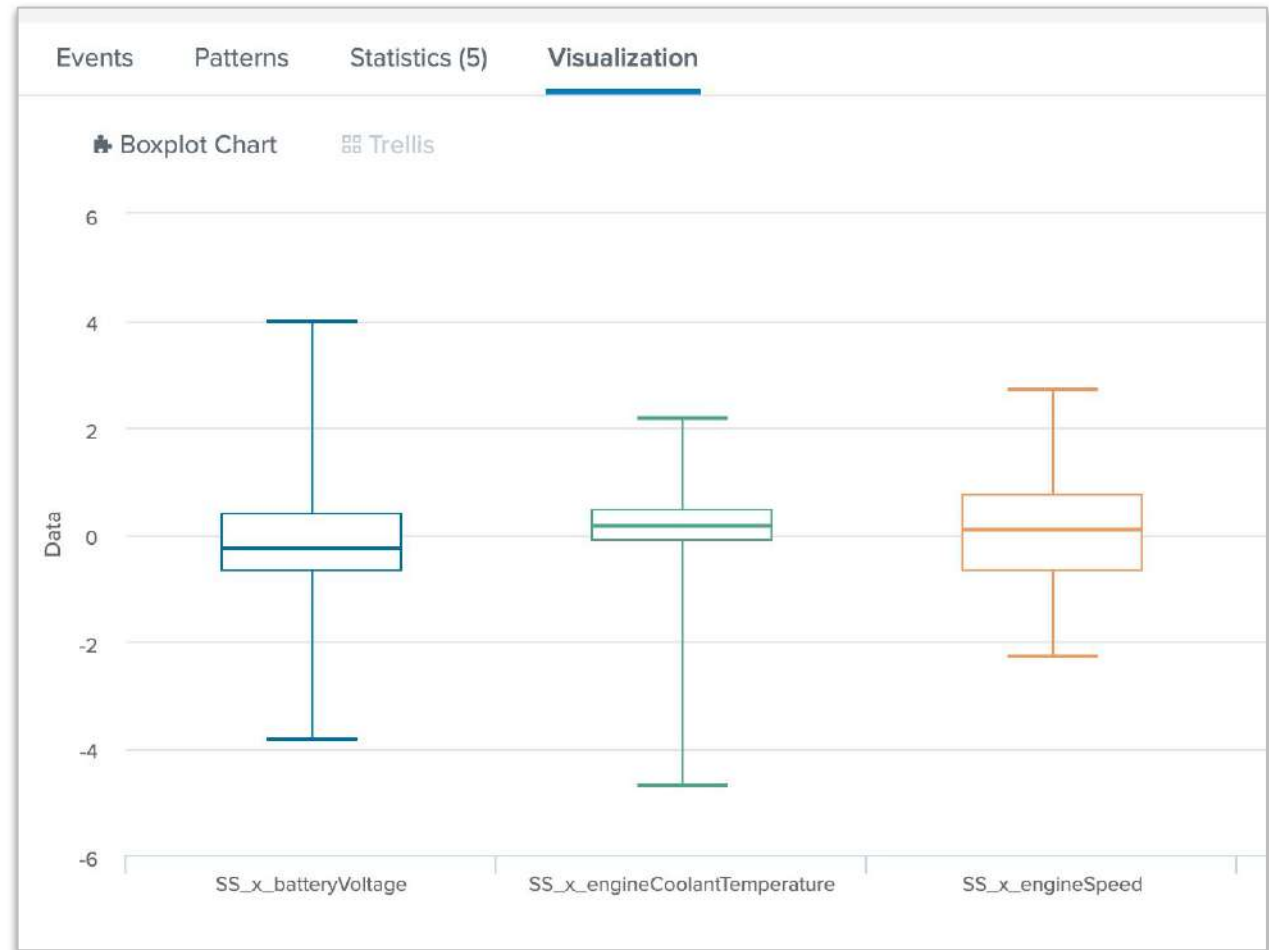Field Labels

**Hints:**
Scale numeric values using the `fit` command with the StandardScaler

splunk>

# Explore the Dataset with Box Plots

**1**

```
| inputlookup "mytrackdata.csv"
| fit StandardScaler x_*
| table SS_*
| `boxplot`
```

- Standardized data fields have a mean of 0 and a standard deviation of 1
- The box plots are less stretched and can be analyzed more easily



splunk>

# **Detect Numeric Outliers:**
## Explore the MLTK showcase and adapt it
## to start a new experiment with your own dataset

**2**



Explore the Outlier
Detection
Showcases

Start your own
Outlier Detection
Experiment

Optionally try to
compare different
outlier detection
approaches

splunk>

# 2 **Explore the Outlier Detection Showcases**



- Switch to the Showcase tab of the MLTK and explore the assistant to detect outliers in server response time
- We are now going to use statistics to detect the outliers



**Detect Outliers**

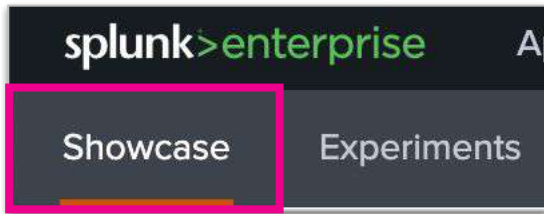View examples that detect numeric and categorical values that differ significantly from values in the rest of the data. Identified outliers are indicative of interesting, unusual, and possibly dangerous events.

14 Examples Available



**Detect Outliers in Server Response Time**

This example uses the Detect Numeric Outliers Assistant and threshold method of Median Absolute Deviation to look for outliers in server response time.

*IT*

splunk>

# **2** **Explore the Outlier Detection Showcases**

Enter a search

| `| inputlookup mytrackdata.csv`

✓ 6,000 results (01/01/1970 01:00:00.000 to 23/01/2020 12:05:17.000)

| Field to analyze | Threshold method | Threshold multiplier | Sliding window (# of values) |
|---|---|---|---|
| x_speed ▾ | Standard Deviation ▾ | 3 | (optional) |

**Detect Outliers**    Open in Search    Show SPL

Pick an appropriate threshold method
(E.g. Standard deviation +/- 3)

View the corresponding SPL query to
the assistant's settings

splunk>

## **2** **Detecting Outliers with the Density Function**

- Switch to the Experiments tab of the MLTK and create a new experiment
- Instead of an approach based on statistics
  we are now going to use the density function to detect outliers

splunk>

# 2 Create Your Own Smart Outlier Experiment

Click here to get to the next step

Smart Outlier Detection: Find Anomalies in Hard Drive Metrics

Cancel     Next >

**Define**

**Define Data Source**

Q Search     Datasets

**Learn**

| inputlookup mytrackdata.csv

All time ▾     Q

✓ 6,000 results (01/01/1970 01:00:0

Job ▾   ‖   ▪   ♥ Smart Mode ▾

Look up the dataset you want to work with

Data Preview     Visualization

Review

splunk>

# SPL for MLTK:
# The `fit` and `apply` Commands

**(3)**

## Examples:

```
<your search> | fit <model name>

<your search> | apply <model name>
```

```
| inputlookup mytrackdata.csv
| apply car_outlier_df_speed
```

- The `fit` command produces a machine learning model based on the behaviour of a set of events. It applies the model to the current search results in the search pipeline

- The `apply` command applies the machine learning model that was learned using the fit command

splunk>

# SPL for MLTK:
# The `fit` and `apply` Commands

**Examples:**

```
<your search> | fit StandardScaler <fields> into <model name>

<your search> | apply <model name> | `<macro name>`

<your search> | fit SVM "X X X" from "XXX" "XXX" kfold_cv=3
```
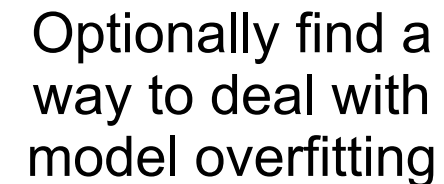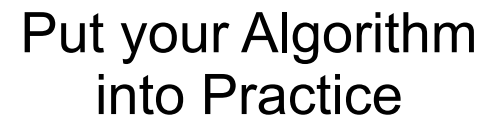
Check out the confusion matrix
and classification statistics
macros!

- The StandardScaler algorithm uses the scikit-learn StandardScaler algorithm to standardize data fields

- Splunk's MLTK allows you to cross-validate your models right from the search queries that train them. Simply specify the number of cross-validation folds you want by setting the `fit` command's parameter `kfold_cv`

splunk>

**(3)**

# Use a Classification Model:
Create a classification model and use it
to predict vehicle types from your sensor data



Explore the
Classification
Assistant

Put your Algorithm
into Practice

Optionally find a
way to deal with
model overfitting

splunk>

# **3** **Explore the Classification Assistant**

**Prediction Results** ↗

| y_vehicleType ⇕ | predicted(y_vehicleType) ⇕ |
|---|---|
| 2015 Porsche GT3 | 2011 Ford Mustang GT500 |
| 2011 Ferrari 458 | 2011 Ford Mustang GT500 |
| 2014 Chevrolet Corvette | 2014 Chevrolet Corvette |
| 2011 Ferrari 458 | 2011 Ford Mustang GT500 |
| 2011 Ford Mustang GT500 | 2011 Ford Mustang GT500 |
| 2015 Porsche GT3 | 2011 Ford Mustang GT500 |
| 2008 BMW M3 | 2011 Ford Mustang GT500 |
| 2011 Ford Mustang GT500 | 2011 Ford Mustang GT500 |
| 2011 Ferrari 458 | 2011 Ferrari 458 |
| 2011 Ferrari 458 | 2011 Ferrari 458 |

**Experiment Settings**     **Experiment History**

Enter a search

| inputlookup mytrackdata.csv

Algorithm

SVM ▼

Field to predict

y_vehicleType ▼

**?** Why is SVM doing so bad?

splunk>

# **3** Save your Classification Model



Publish your model in the app of your choice

# 3

# Apply your Classification Model

Models have been published                                    ✕

You can find the models under Settings > Lookups > Lookup Table Files ⬏. You may now
use the published models to create alerts or schedule model training.

To create an alert, you can use the below SPL snippet in your search:

```
... | apply car_classifier_StandardScaler_1 | apply car_classifier
```

OK

```
| inputlookup mytrackdata.csv
| apply car_classifier_StandardScaler_0
| apply car_classifier
| table y_vehicleType "predicted(y_vehicleType)" *
|  `confusionmatrix("y_vehicleType","predicted(y_vehicleType)")`
```

splunk>

# Which Car Gets Classified Worst?
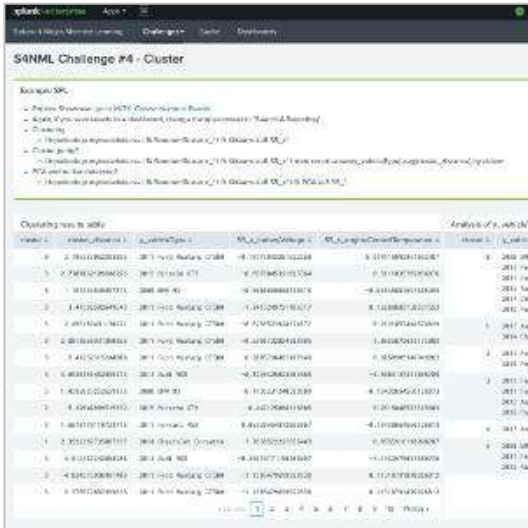
**3**

**?** How can you find out where your model is off?

Further analysis of wrong classifications

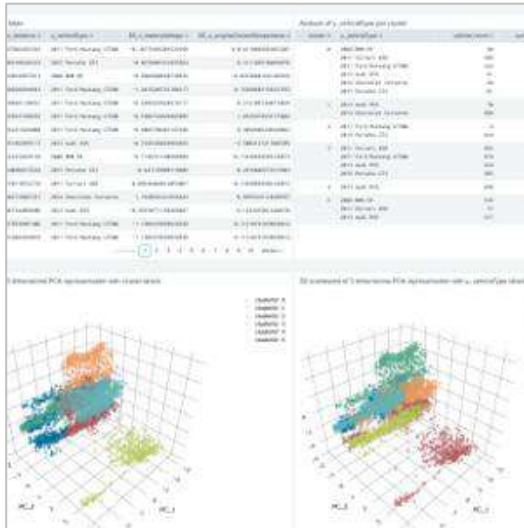| Predicted actual ⇕ | Predicted 2011 Ferrari 458 ⇕ | Predicted 2011 Ford Mustang GT500 ⇕ | Predicted 2013 Audi RS5 ⇕ | Predicted 2015 Porsche GT3 ⇕ |
|---|---|---|---|---|
| 2011 Ferrari 458 | 0 | 0 | 6 | 0 |
| 2011 Ford Mustang GT500 | 0 | 0 | 0 | 15 |
| 2013 Audi RS5 | 1 | 0 | 0 | 0 |
| 2015 Porsche GT3 | 0 | 8 | 0 | 0 |

splunk>

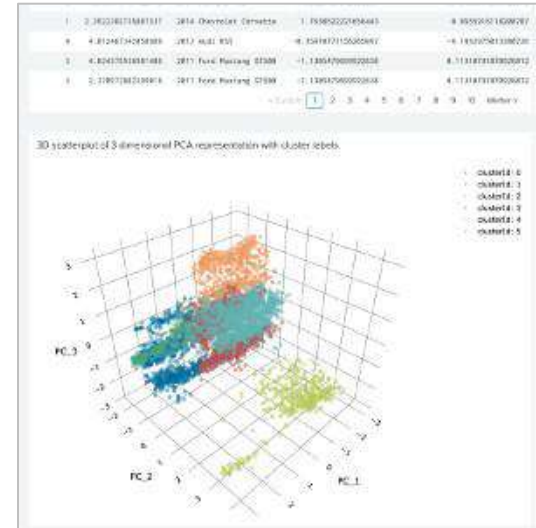# Use a Clustering Model:
## Create a clustering model and use it to analyze your dataset
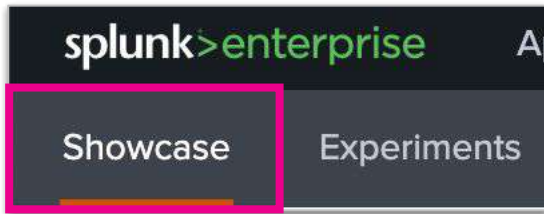
**3**



Explore the Clustering Assistant
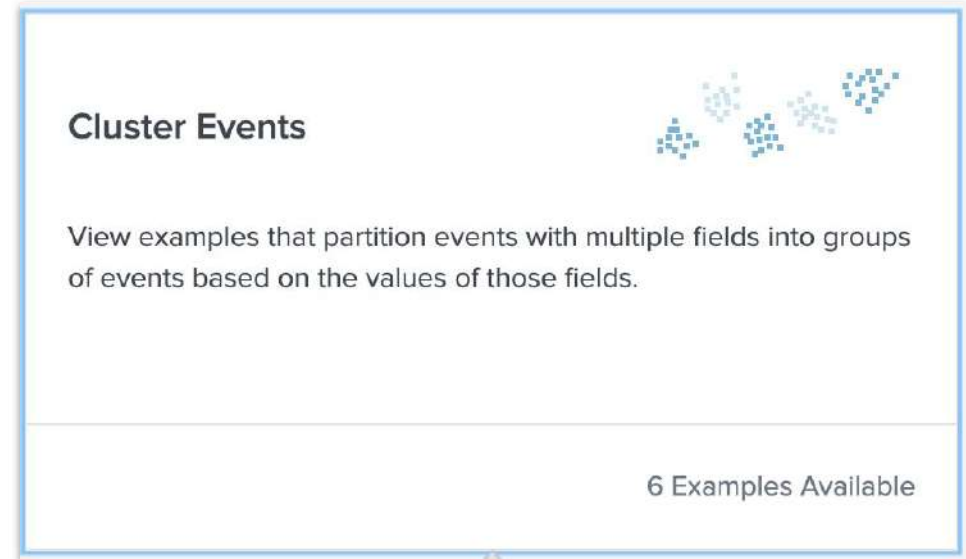


Cluster Analysis of the mytrackdata-Dataset



Optionally try and detect outliers

splunk>

# 4

# Explore the Cluster Showcases

splunk>enterprise A

**Showcase** | Experiments

**Cluster Events**

View examples that partition events with multiple fields into groups of events based on the values of those fields.

6 Examples Available

- Switch to the Showcase tab of the MLTK and explore the assistant to identify clusters of events

**Cluster Vehicles by Onboard Metrics**

This example uses the Cluster Numeric Events Assistant, a preprocessing step using the PCA method, and the Birch algorithm to cluster data on seven fields including battery voltage, engine speed, and vertical G-force.
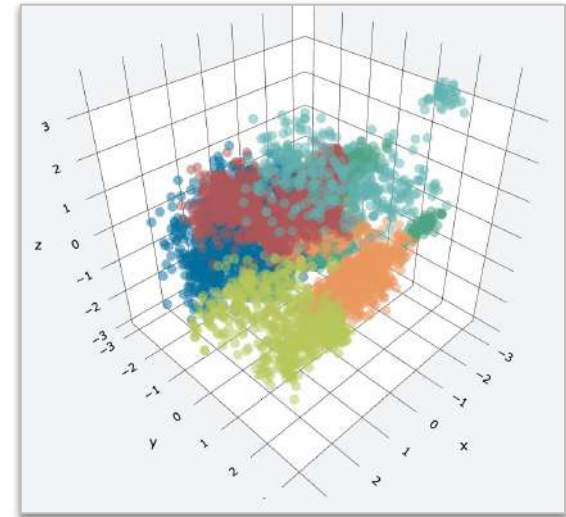
*IoT*

splunk>

# The MLTK Comes with Many Different Algorithms

**(4)**

**Example:**

```
<your search> | fit PCA k=<int> <fields>
```

> Factor analysis with an algorithm such as PCA can reduce the number of variables one must deal with

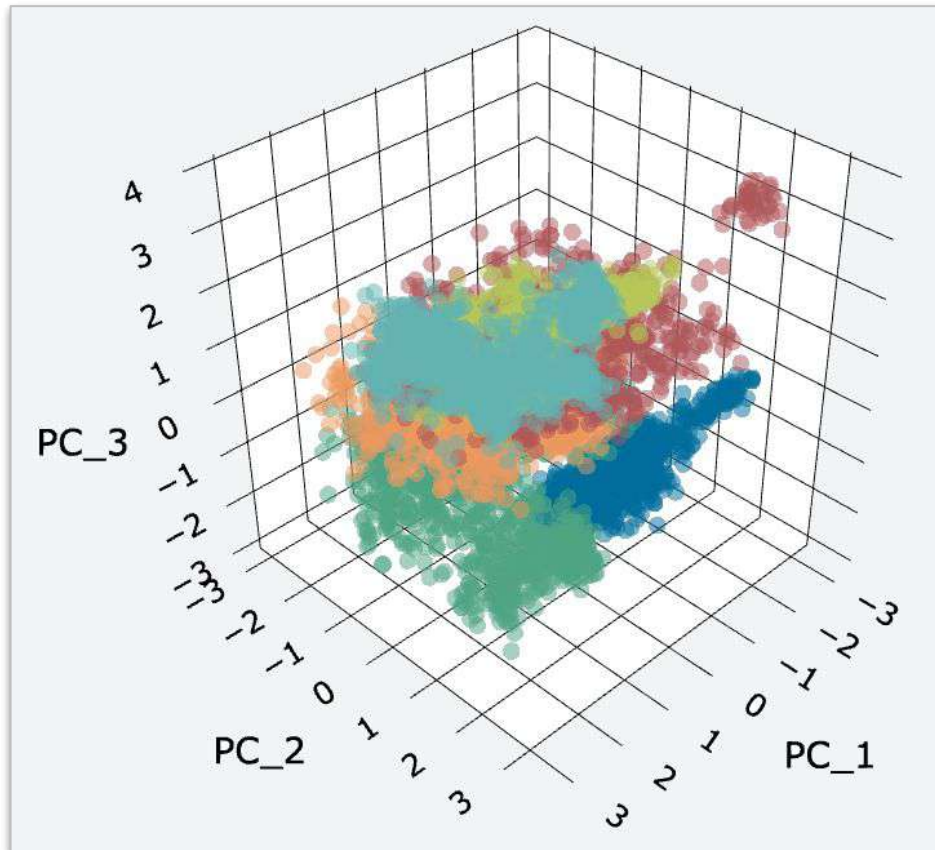> The k parameter specifies the number of features to be extracted from the data

**?** Why is there a cluster with "clusterId: null" ?



- clusterId: 2.0
- clusterId: 0.0
- clusterId: 5.0
- clusterId: 1.0
- clusterId: 3.0
- clusterId: 4.0
- clusterId: null

splunk>

# The MLTK Comes with Many Different Algorithms

**4**



We have missing values in "x_engineCoolantTemperature" that we didn't fix/impute in mytrackdata.csv

- clusterId: 2008 BMW M3
- clusterId: 2015 Porsche GT3
- clusterId: 2011 Ford Mustang GT500
- clusterId: 2013 Audi RS5
- clusterId: 2011 Ferrari 458
- clusterId: 2014 Chevrolet Corvette

splunk>

# Wrap Up

**splunk>**

# Wrap Up

- **Feedback:** How was your experience, what worked well, what did not?

- **Discussion Brainstorming:** How could you transfer the topics learned today to other use cases or departments?

- **You want to learn more about Splunk's Machine Learning?**

  - Check our latest [Splunk Blogs around Machine Learning](#)
  - Watch videos from [Splunk Machine Learning YouTube Channel](#)
  - Take the [Splunk Education Class for Data Science and Advanced Analytics](#)
  - Learn more about [Splunk's Machine Learning Advisory Program](#)

splunk>

# Thank You!

splunk>

# Additional Information

**Login:**

Username: **admin**

Password:  **<See Splunk Show>**

**Challenge Solution Examples:**

We created a dashboard for each challenge with example solutions in the hidden app "Splunk 4 Ninjas Machine Learning". Use this app for preparation, debriefing after the challenges or as assistance for unexperienced attendees.

► https://{your-host}/en-GB/app/s4n_ml/splunk_4_ninjas_ml

   or click the button next to "Splunk 4 Ninjas Machine Learning" on top of the Home dashboard

splunk>