

# Data Science - IoT Data

Provincial Electricity Authority co.,



## Data Science for IoT - Agenda

- 09.00 - 10.30 **Introduction to Data Science**
- **Coffee Break** 
- 10.45 - 12.00 **Data Modeling**
- **Lunch**  
- 13.00 - 14.30 **Model Validation**
- **Coffee Break** 
- 14.45 - 16.30 **Case Study**



# Paul

Wichit Sombat, PhD  
[wichit.s@ubu.ac.th](mailto:wichit.s@ubu.ac.th)



# K

Kriengsak Treeprapin, PhD  
[kriengsak.t@ubu.ac.th](mailto:kriengsak.t@ubu.ac.th)



# Tom

Phaichayon Kongchai, PhD  
[phaichayon.k@ubu.ac.th](mailto:phaichayon.k@ubu.ac.th)

# 1

## Intro to Data Science

แนะนำวิทยาการข้อมูล

“Unify **statistics**<sup>1</sup>, **data analysis**<sup>2</sup>,  
**machine learning**<sup>3</sup> to scientifically  
understand data.

Wikipedia, 2018

Maths + Stats + Info Sci + Com Sci

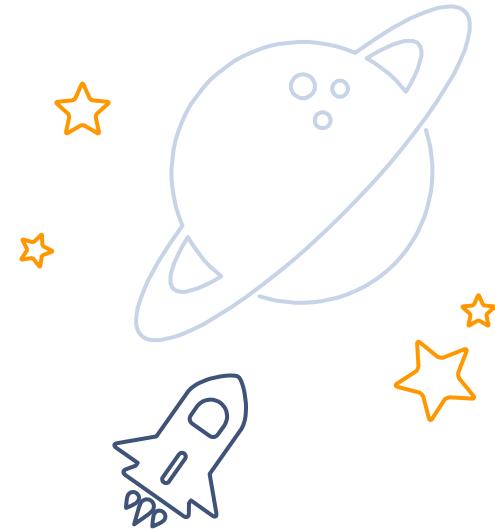
# 1.1

## Introduction

ที่มา หลักการและเหตุผล

# Why??? Understanding Data

ทำไมต้องเรียนข้อมูลมาคีกษา?

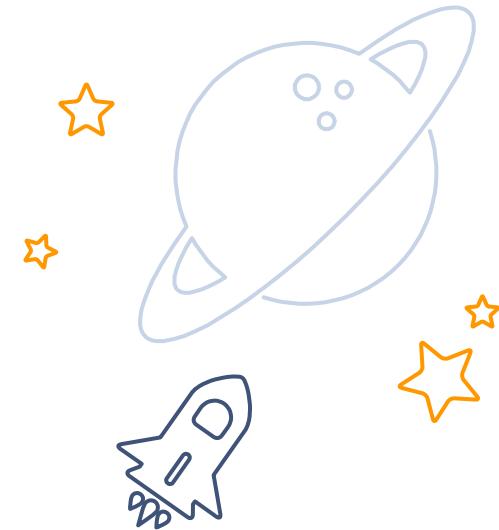




เรียนรู้สาเหตุและป้องกันภัย



คาดการและปรับปรุงระบบการทำงานตามสถานะการ



# Data from <40,000<sup>1</sup> planes

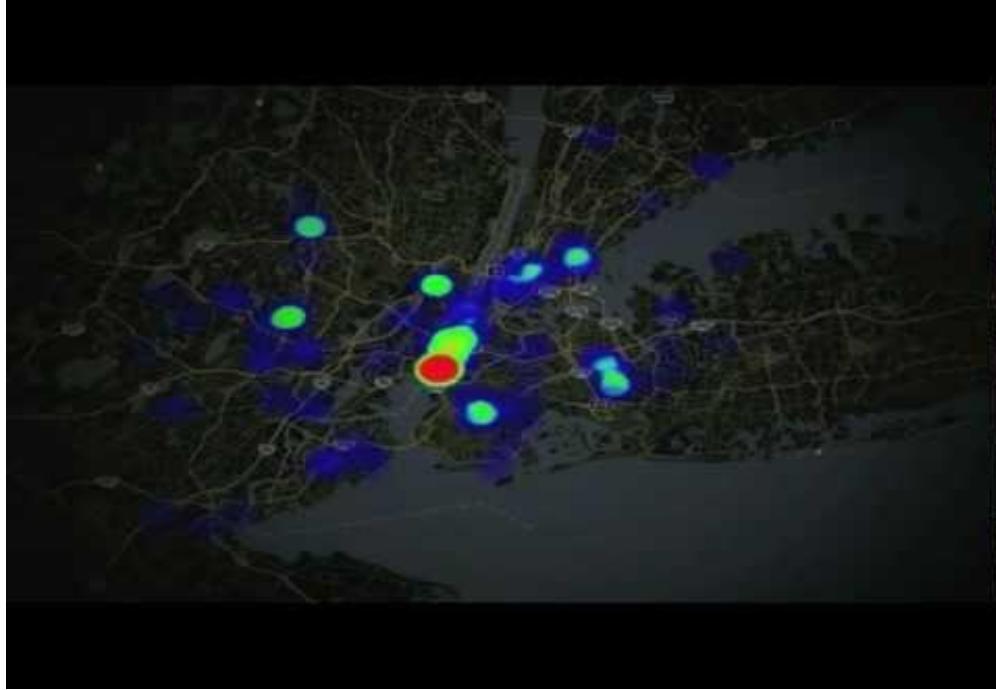
ข้อมูลการบินจาก เครื่องบิน สีหมีน ลำ(ทั่วโลก)

<sup>1</sup>[telegraph.co.uk](http://telegraph.co.uk)



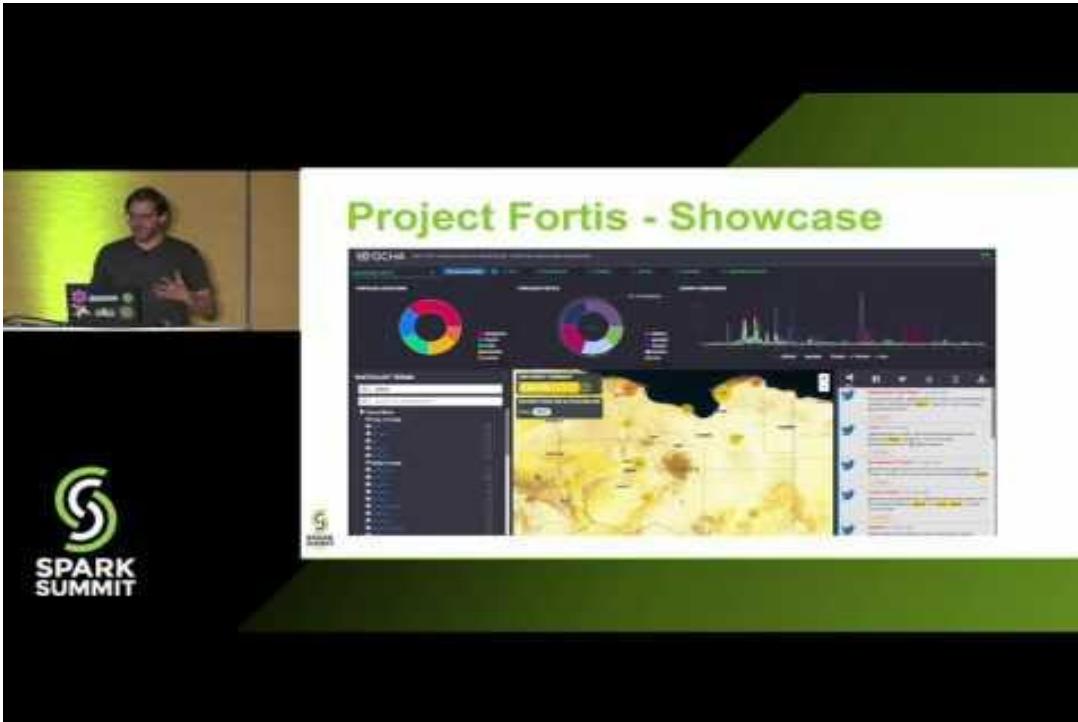
<sup>1</sup>statista.com

Air Traffic from social network posts (twitter)  
ข้อมูลมาได้จากหลากหลายแหล่ง twitter ~ 336M<sup>1</sup>



Tracking Influenza Via Twitter - 2013  
โปรแกรม GermTracker<sup>1</sup> จากข้อมูล twitter

<sup>1</sup>University of Rochester

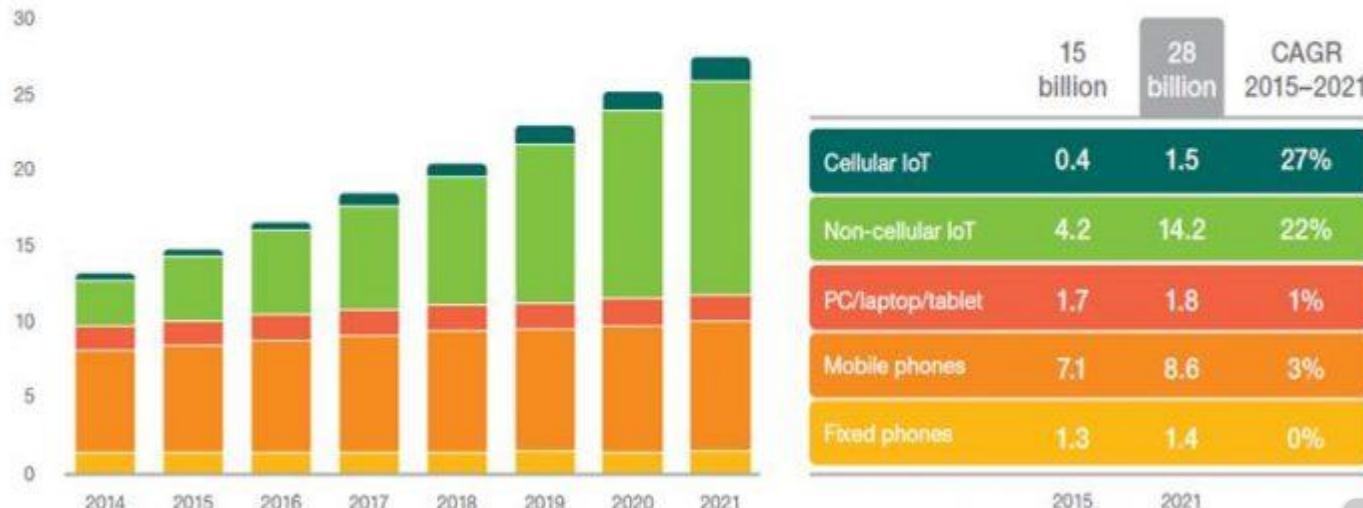


Social Media for UN Project Fortis for humanitarian aid plan

<sup>1</sup>Spark Summit 2017

# THE INTERNET OF THINGS

Connected devices (billions)



<sup>1</sup>ericsson.com



## End Users IoT dashboards

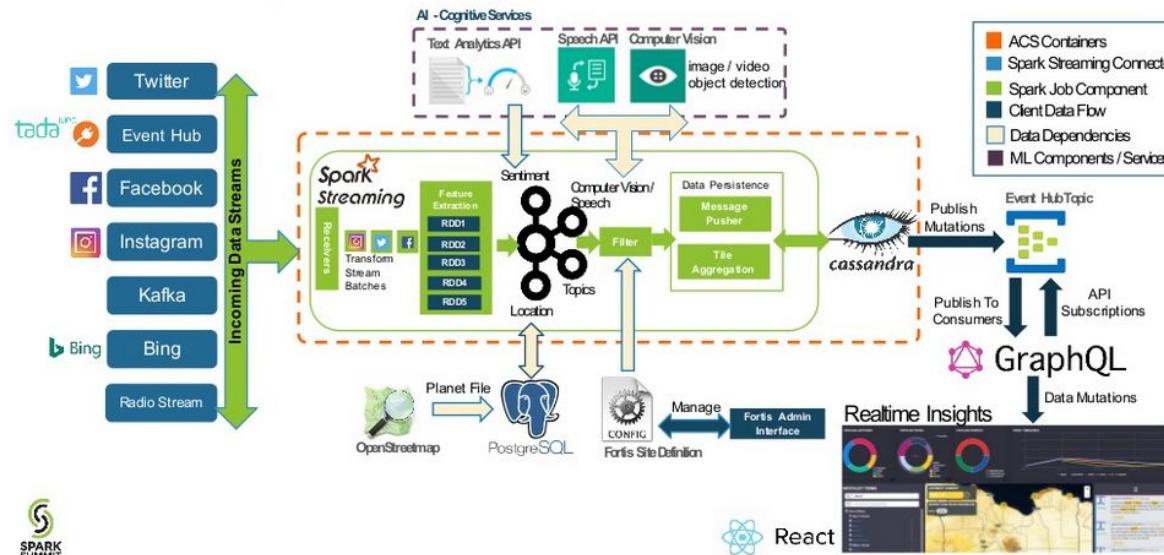
Data per IoT Device<sup>1</sup>

<sup>1</sup>Thingsboard



# Architecture for (Social) Data Analysis<sup>fortis</sup>

## Data Pipeline: Functional Architecture

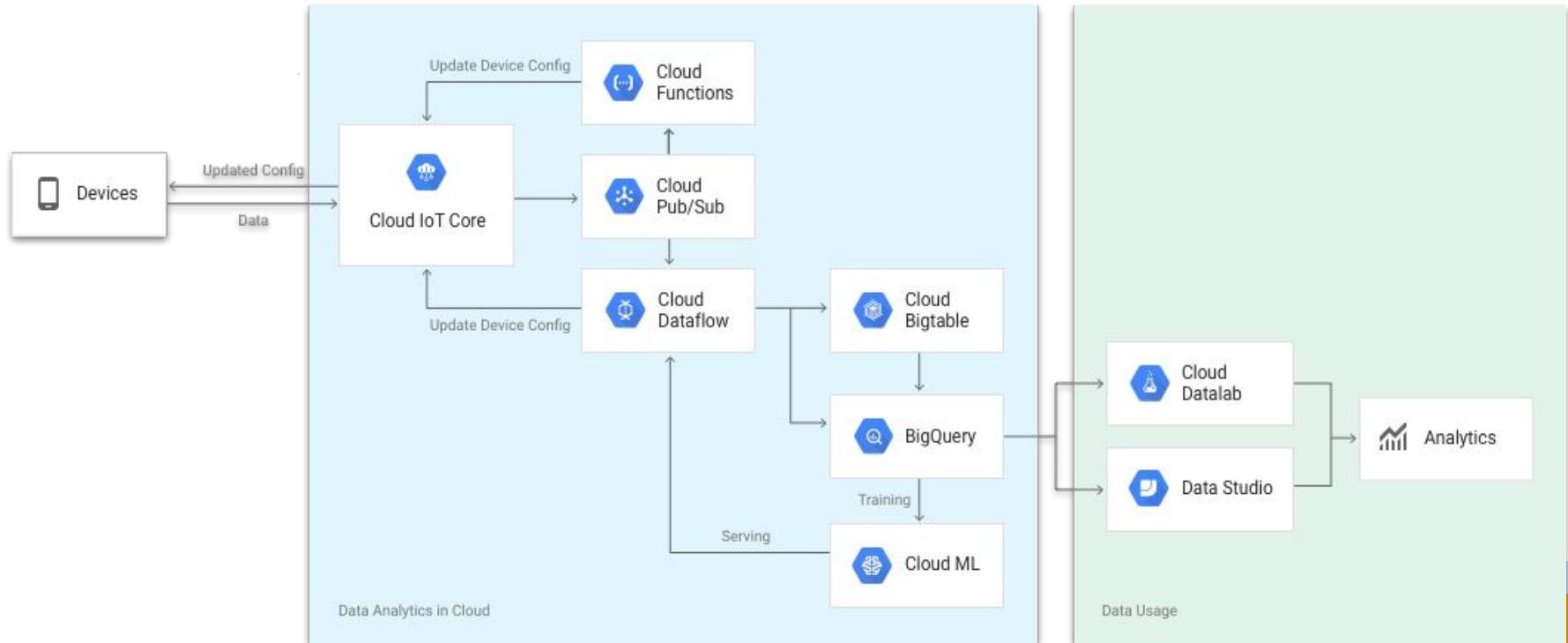


<sup>1</sup>Spark Summit 2017



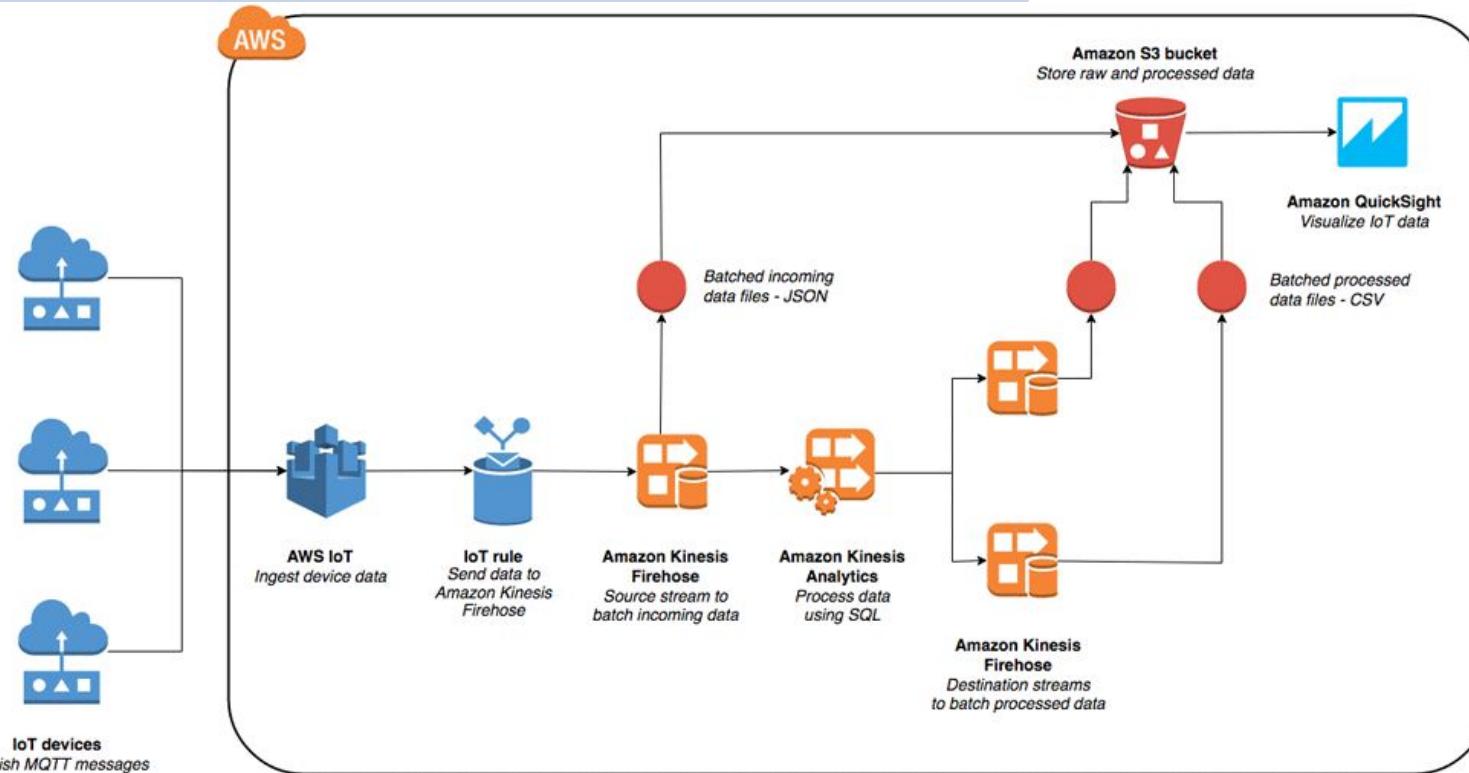


# Architecture for Data Analysis Google Cloud



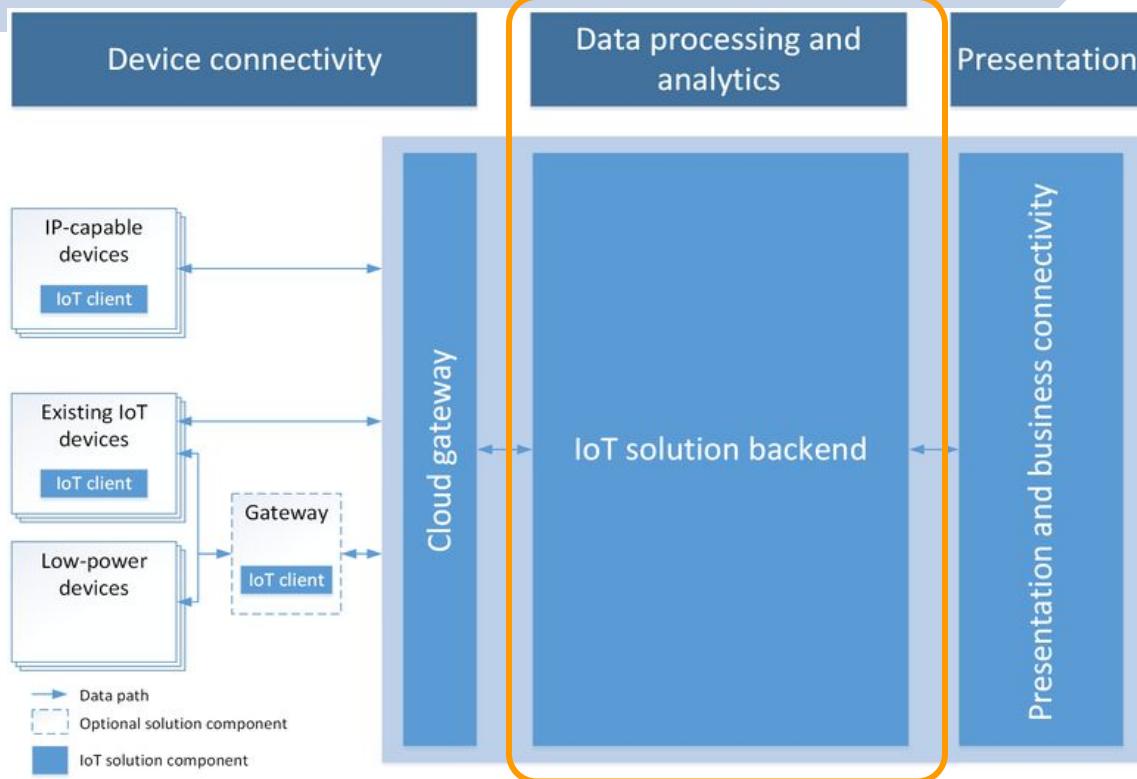


# Architecture for Data Analysis





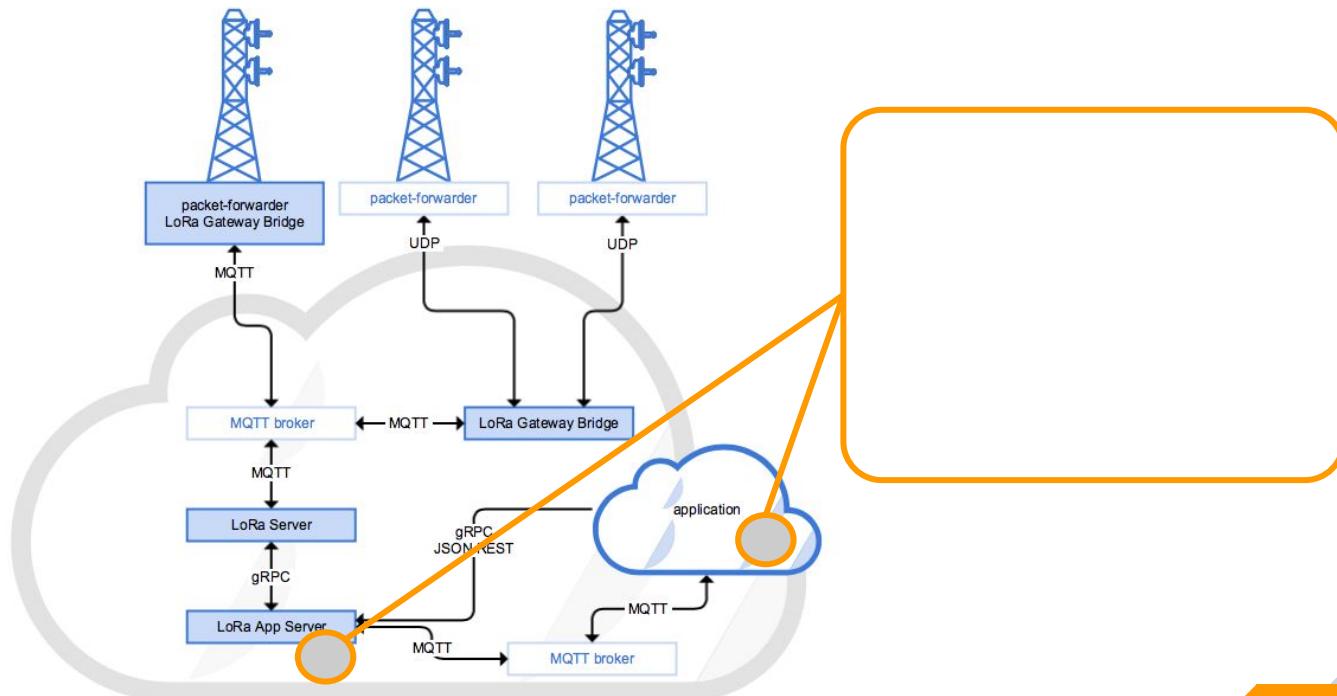
# Architecture for Data Analysis<sup>Ms Azure</sup>





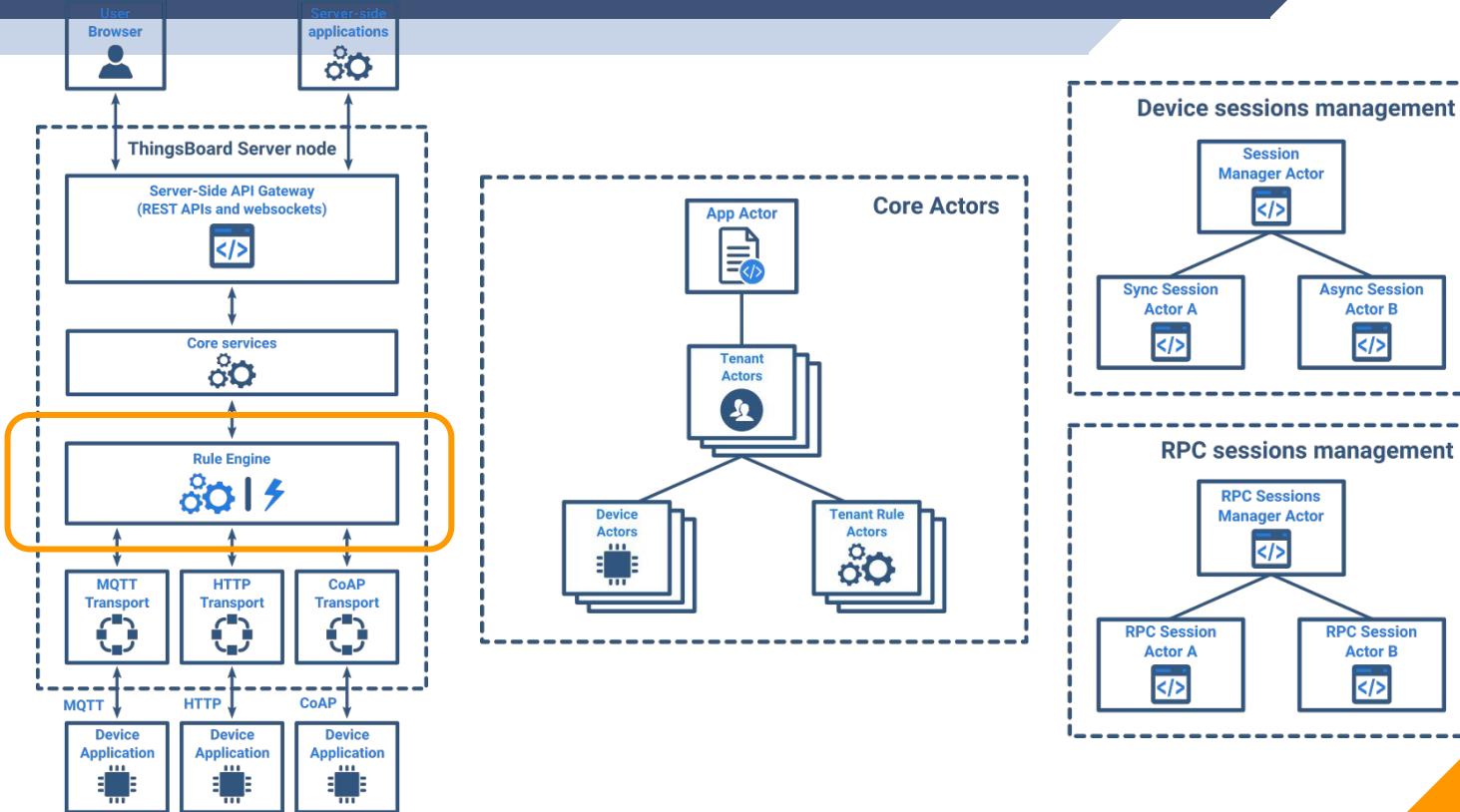
# IoT Architecture for Data Analysis

Loraserver





# Architecture for Data Analysis<sup>thingsboard</sup>





# Architecture for Data Analysis thingsboard

Districts

Districts Map

District A US, San Francisco, CA, Sunset District 35,405 kWh 29,535 gal

District B US, San Francisco, CA, Mission District 43,131 kWh 26,272 gal

All Districts: Alarms

Realtime - last 7 days

<input type="checkbox"/>	Created time	Originator	Type	Severity	Status
<input type="checkbox"/>	2017-10-11 19:53:54	Living room Thermostat	HighTemperature	Critical	Active Unacknowledged
<input type="checkbox"/>	2017-10-11 19:22:14	Kitchen Thermostat	LowTemperature	Critical	Active Unacknowledged
<input type="checkbox"/>	2017-10-11 19:22:14	Bedroom Thermostat	HighTemperature	Critical	Active Unacknowledged

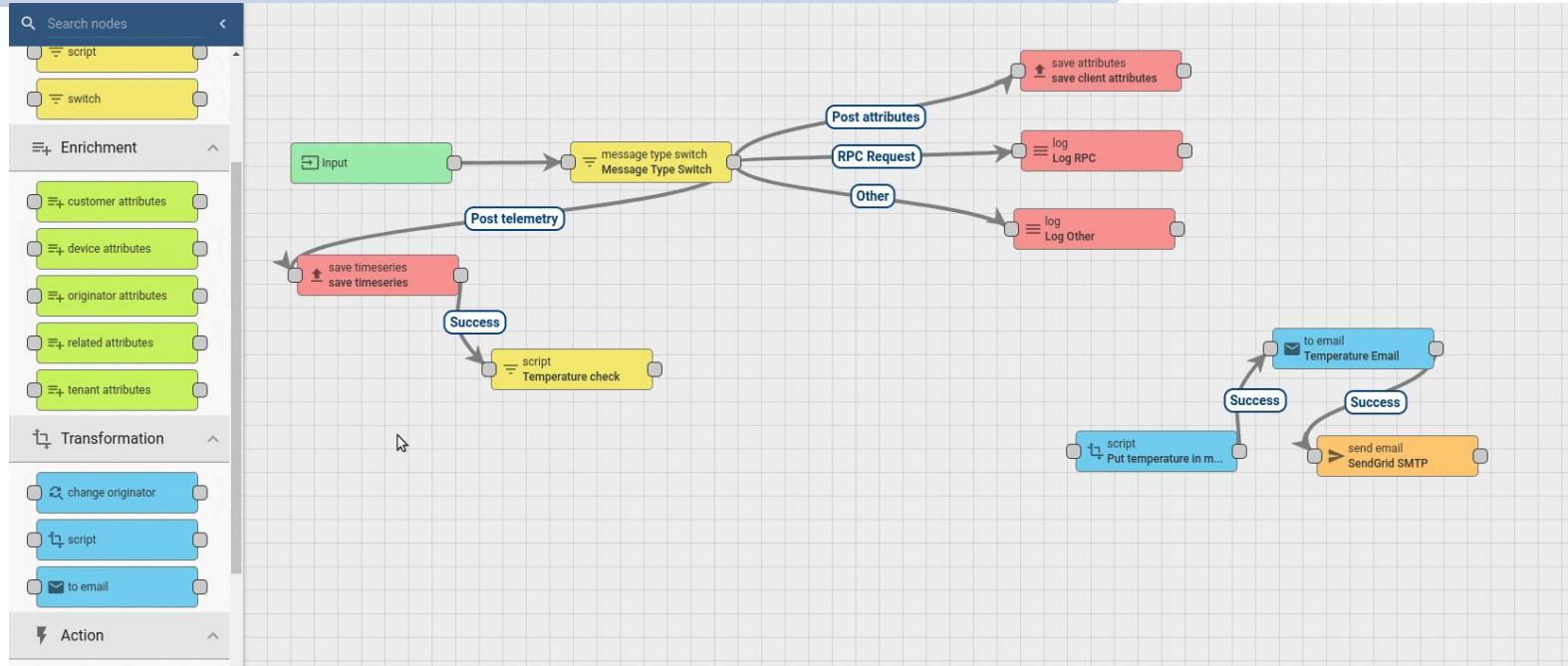
All Districts 78,536 kWh All Districts 55,807 gal

Energy consumed Water consumed

Powered by Thingsboard v1.3.1



# Architecture for Data Analysis<sup>thingsboard</sup>





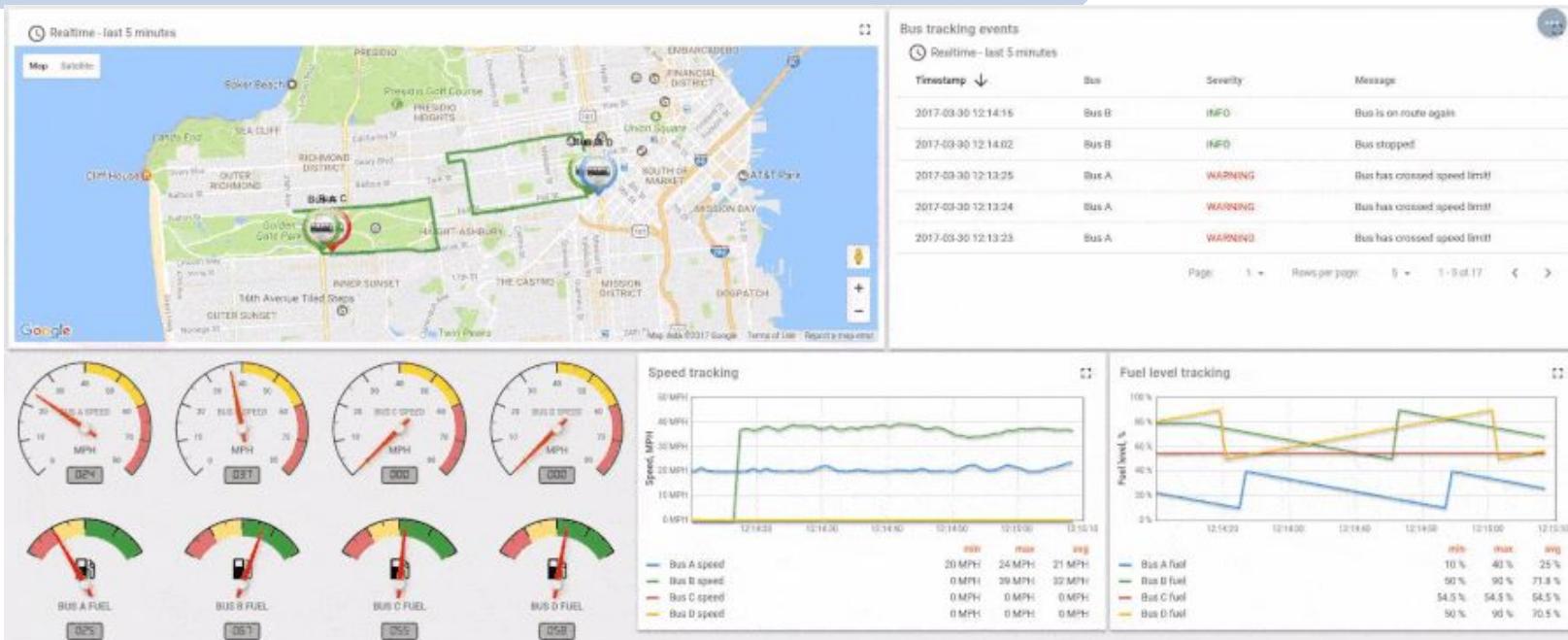
# Architecture for Data Analysis thingsboard





# Architecture for Data Analysis

thingsboard



# 1.2

## Data Science Process

ขั้นตอนหลักของวิทยาการข้อมูล



## How Fortis UN generate heatmap?

The image is a composite of two parts. On the left, a man is giving a presentation at a podium during the Spark Summit 2017. He is wearing a dark t-shirt and gesturing with his hands. The background is a yellow wall. On the right, there is a screenshot of a software interface titled "Project Fortis - Showcase". The interface features a map with a heatmap overlay, several circular charts, and various data tables and graphs. The overall theme is data visualization and humanitarian aid planning.

<sup>1</sup>Spark Summit 2017

Social Media for UN Project Fortis for humanitarian aid plan



## Humanitarian aid plans are manually composed

<sup>1</sup>Spark Summit 2017





## Fortis UN project goals

- **Accelerate** the construction of aid **planning**
- **Improve** its data **accuracy**
- Provide **deeper insights** and **trends**
- **Real-time** analytics
- More intelligence and insight to enable better  
**forecasting**

<sup>1</sup>Spark Summit 2017



## Fortis UN project tasks

### Geospatial Requirements



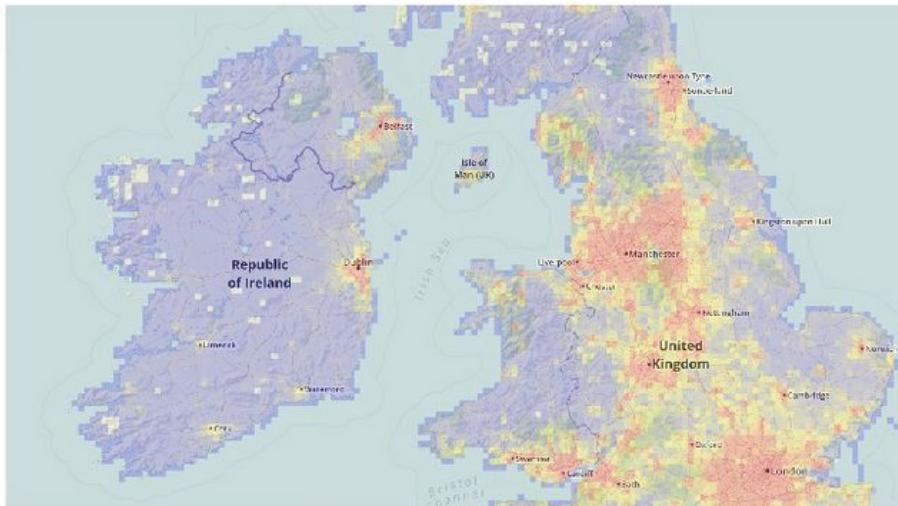
- Geotag mentioned places from conversations
- To query activities within a bounding box
- Filter activities based on topic(s)



## Fortis UN project tasks

### Heatmap Generation

Aggregate detected places across geographic regions



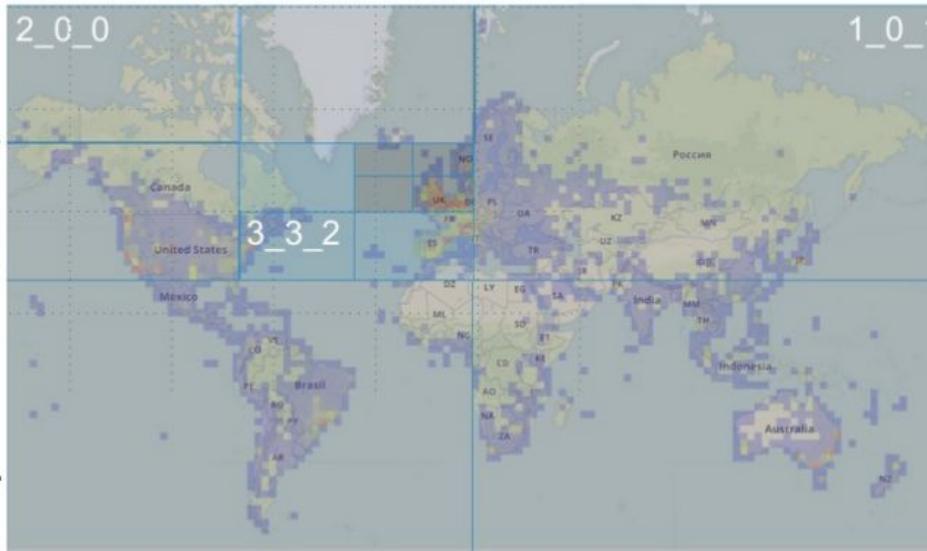
<sup>1</sup>Spark Summit 2017



## Fortis UN project tasks

### XYZ Tiles for summarization

- Divides world up into tiles.
- Each tile has four children at the next higher zoom level.
- Maps multiple layers to a single space dimension.



<sup>1</sup>Spark Summit 2017



## Fortis UN project tasks - Heatmap Spark

<sup>1</sup>Spark Summit 2017

```
final val MAX_ZOOM_LEVEL = 16
final val MIN_ZOOM_LEVEL = 5

def tile_id_mapper(location: (Double, Double)): List[(String, Int)] = {
  (for (zoom <- MIN_ZOOM_LEVEL to MAX_ZOOM_LEVEL)
    yield (TileUtils().tile_id_from_lat_long(location._1, location._2, zoom).tileId, 1)).toList
}
```

```
val locationListSamples = List[(Double, Double)]((30.294221,-97.7760937),
(30.294221,-97.7760937), (30.4007241,-97.7368647),(30.4007241,-97.7368647),....)
val locationRdd = sc.parallelize(locationListSamples)

locationRdd.flatMap(tile_id_mapper)
.reduceByKey(_+_)
.saveToCassandra("fortis", "tilesExample", SomeColumns("tileId", "count" append))
```



## Fortis UN project tasks

<sup>1</sup>Spark Summit 2017

### Aggregated Tile Schema

All inbound events are aggregated by  
tile\_x, tile\_y, tile\_z, period, source, publisher, topic, lang

#### Aggregated Tile Data - Cassandra: fortis\_tiles\_tbl

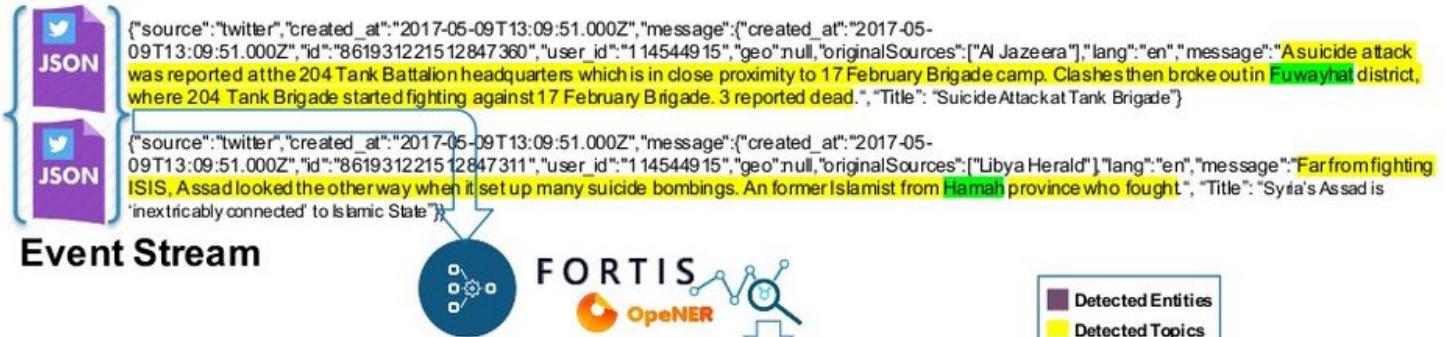
tile_x	tile_y	tile_z	period	Source	publisher	Topic	Lang	feature_collection
15	13346	15	month-2017-06	Twitter	Al Jazeera	isis	en	{"mention_count": 1213, "sentiments_avg": ["neg", ".78", "pos": .02], "entities": [{"name": "bashar assad", "mention_count": 627, "ref_id": "3256"}]}}
231	3345	14	month-2017-06	Facebook	Times of Libya	Isis	en	{"mention_count": 453, "sentiments_avg": ["neg", ".78", "pos": .02], "entities": [{"name": "bashar assad", "mention_count": 124, "ref_id": "3256"}]}}
76	98242	13	month-2017-06	Facebook	Times of Libya	isis	en	{"mention_count": 453, "sentiments_avg": ["neg", ".78", "pos": .02], "entities": [{"name": "bashar assad", "mention_count": 124, "ref_id": "3256"}]}}



# Fortis UN project tasks

<sup>1</sup>Spark Summit 2017

## NLP Feature Extraction



### Event Stream

### Event Details – Cassandra: fortis\_events\_tbl

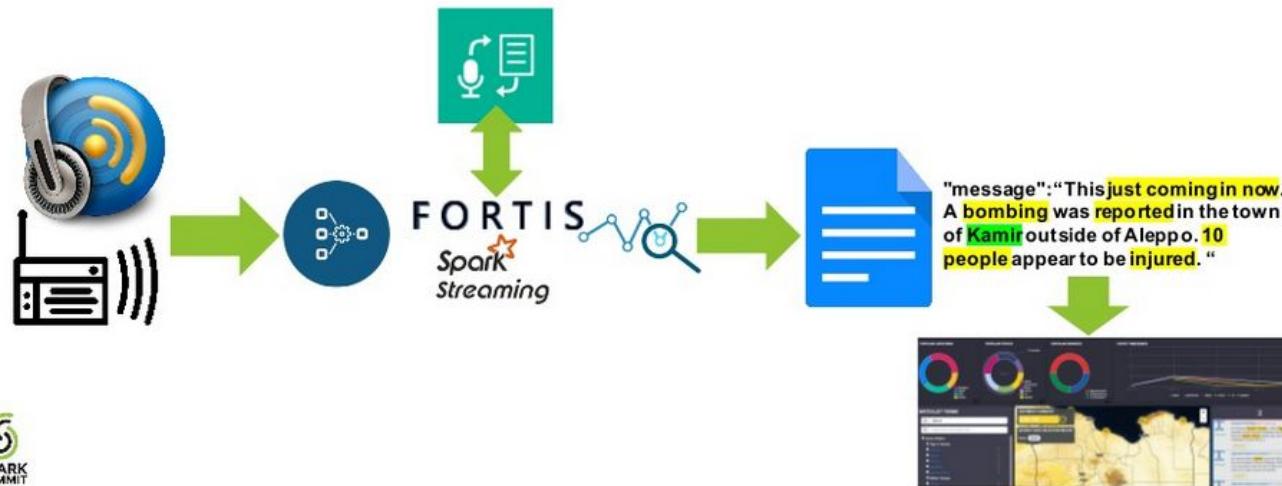
event_id	source	Title	src_url	detected_features	Publisher	Topics	lang	message_body	Feature_collection
851144530 439090180	twitter	Syria's Assad is 'inextricably connected' to Islamic State	<a href="https://twitter.com/AJEenglish/status/85114453043909018">https://twitter.com/AJEenglish/status/85114453043909018</a>	["wof85678363"]	Al Jazeera	["isis", "bomber"]	en	Far from fighting ISIS, Assad looked the other way when it set up many suicide bombings. An former Islamist from Hamah province who fought	{ "sentiments_avg": { "neg": ".78", "pos": ".02"}, "entities": [ { "name": "basher assad", "ref_id": "3256" } ] }
8619312 2151284 7360	twitter	Suicide Attack at Tank Brigade	<a href="https://twitter.com/AJEenglish/status/85114453043909018">https://twitter.com/AJEenglish/status/85114453043909018</a>	["wof85678323"]	Libya Herald	["suicide", "attack", "clashes"]	en	A suicide attack was reported at the 204 Tank Battalion headquarters which is in close proximity to 17 February Brigade camp. Clashes then broke out in Fuwayhat district, where 204 Tank Brigade started fighting against 17 February Brigade. 3 reported dead.	{ "sentiments_avg": { "neg": ".78", "pos": ".02"}, "ref_id": "3256" } ] }



## Fortis UN project tasks

### Realtime Audio Streaming with Spark

- Speech Recognition APIs through radio broadcasts



<sup>1</sup>Spark Summit 2017

# 1.2

## Data Science Skills

ทักษะที่สำคัญของวิทยาการข้อมูล



## ทักษะที่จำเป็น ของวิทยาการข้อมูล

- **การวิเคราะห์ข้อมูล**

- ▷ Data Mining, Machine Learning, NLP, NN, etc.

- **การเขียนคำสั่งโปรแกรม**

- ▷ Python, R, Scala, JavaScript, SAS, etc.

- **การจัดการข้อมูลและนำเสนอ**

- ▷ Data Visualization, Infographics, Plots, etc.

<sup>1</sup>Spark Summit 2017

# 1.3

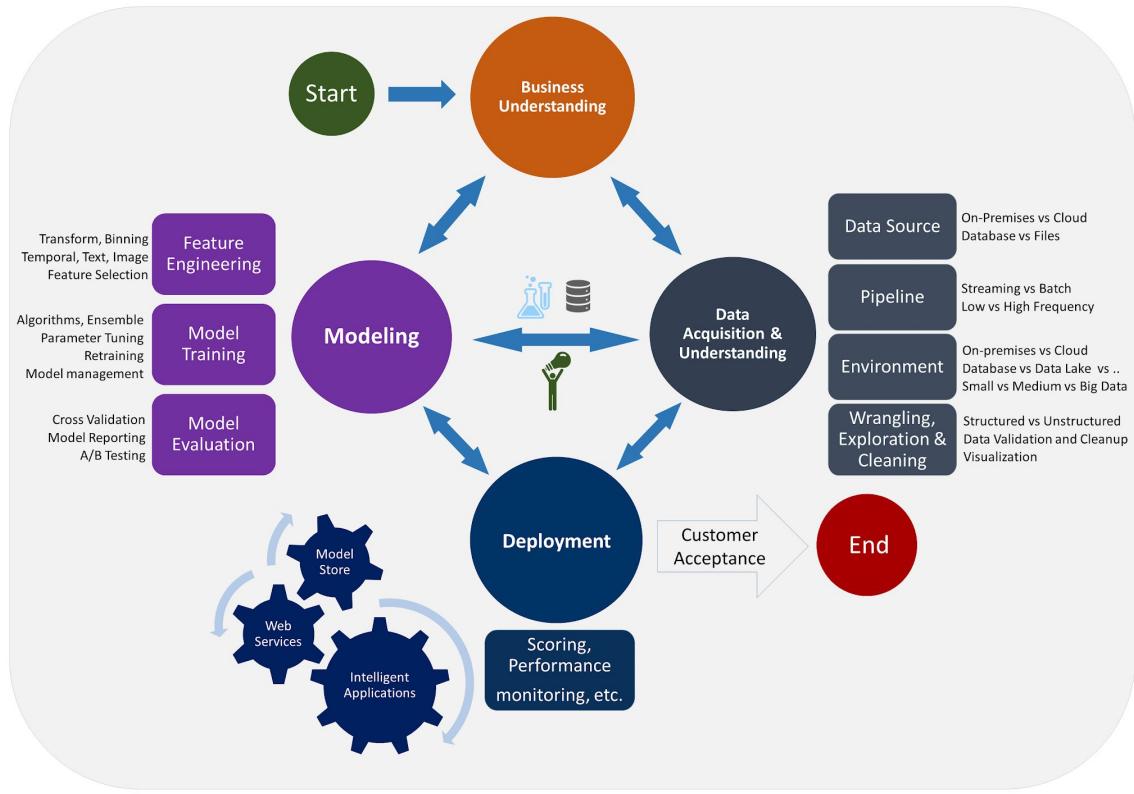
## Data Science Process

ขั้นตอนของวิทยาการข้อมูล



# Data Science Lifecycle<sup>1</sup>

## Data Science Lifecycle

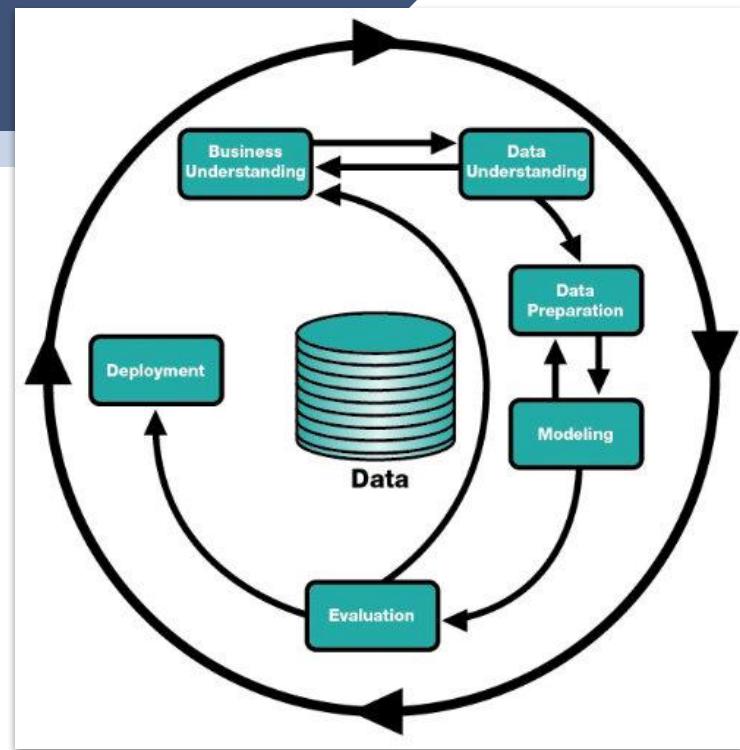
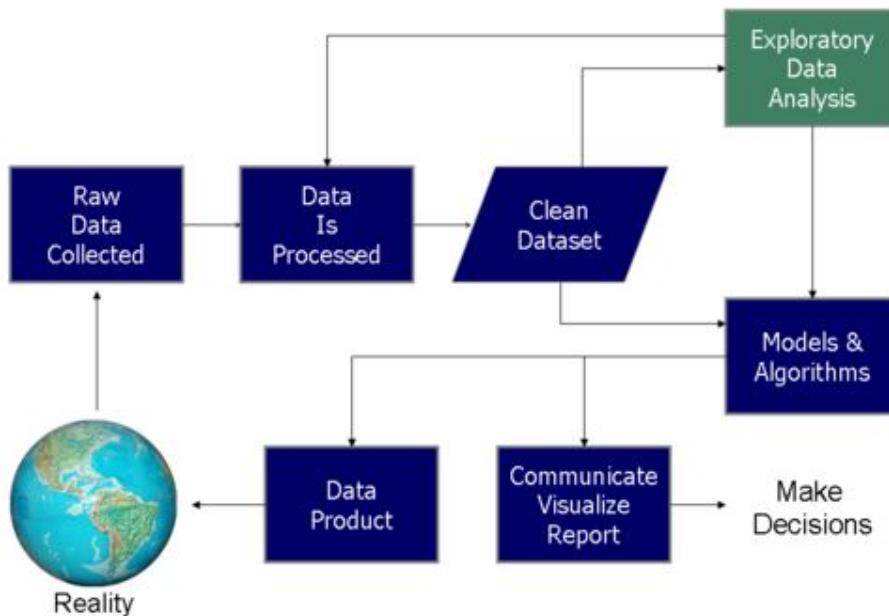


<sup>1</sup>Microsoft



# Data Science Process<sup>1,2</sup>

## Data Science Process



<sup>1</sup>kdnugget



## ขั้นตอนหลักของวิทยาการข้อมูล

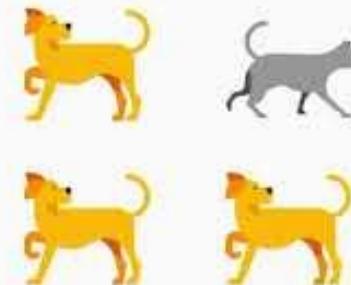
1. **ตั้งคำถาม** - goal? What to predict or estimate?
2. **รวบรวมข้อมูล** - นำเข้าข้อมูล, ทำความสอดคล้องกัน
3. **วิเคราะห์ข้อมูล** - plot, anomalies, patterns
4. **พัฒนาต้นแบบ** สร้าง, ทดสอบ - build, fit, validate
5. **นำเสนอผลลัพธ์** นำต้นแบบไปใช้งาน



## วิทยาการข้อมูล ตอบคำถามอะไรได้บ้าง?

Is this weird?

Anomaly detection algorithms



<sup>1</sup>Daniel Yarmoluk 2018



## วิทยาการข้อมูล ตอบคำถามอะไรได้บ้าง?

1. Is this **A** or **B**? (classification)
2. Is this **weird**? (anomaly detection)
3. How much or **how many**? (regression)
4. How is this **organized**? (clustering)
5. **What** should I do **next**? (reinforcement learning)

# 1.4

## Machine Learning

หลักการเรียนรู้ของเครื่อง



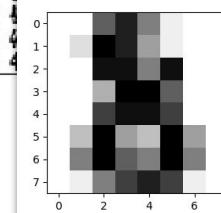
# Classifications



ຂ້ອຍກໍຮັກເລື່ອງນີ້.

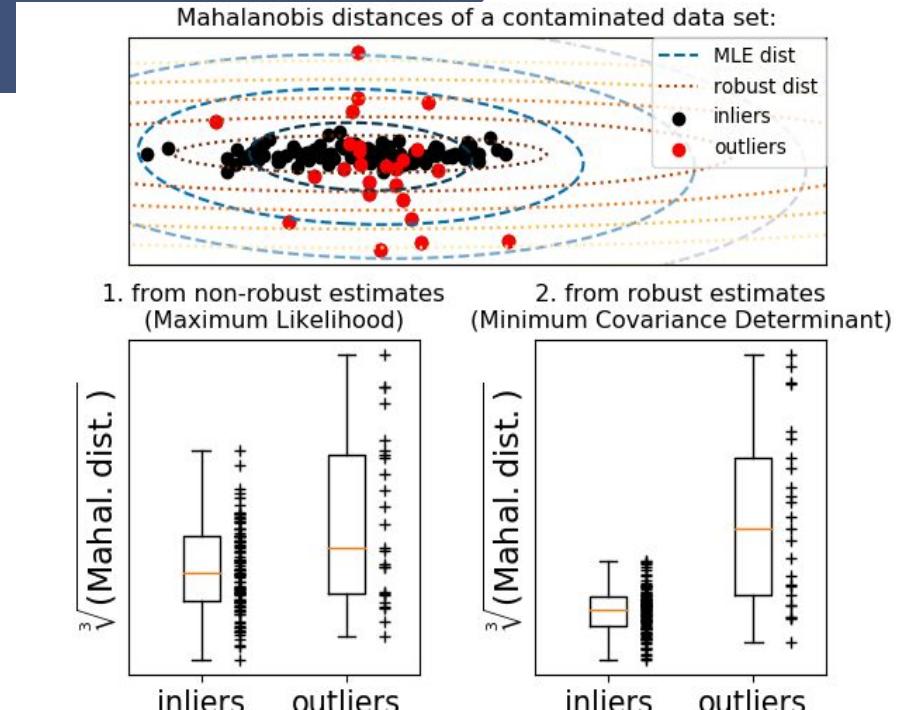
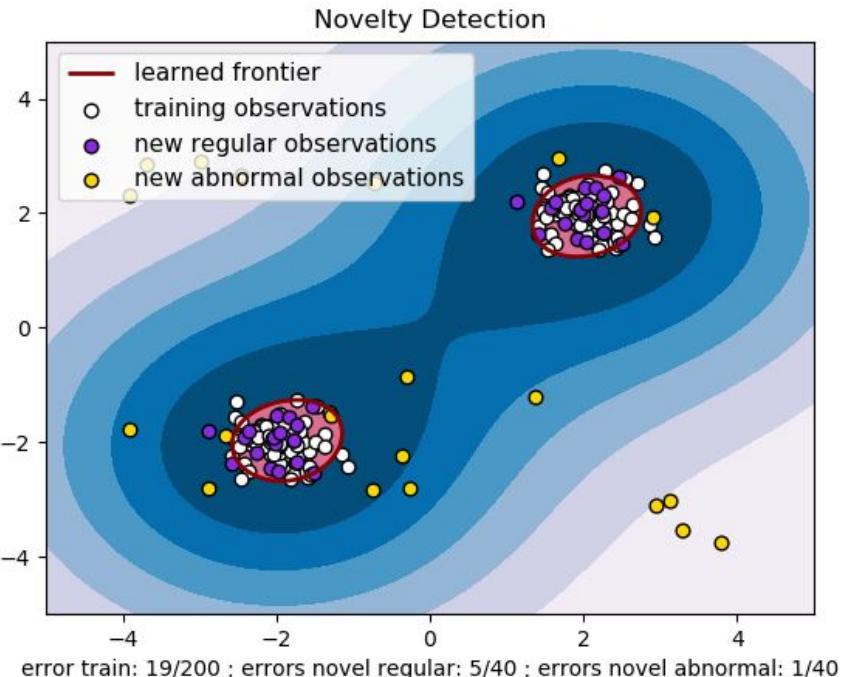
A selection from the 64-dimensional digits dataset

0	1	2	3	4	5	0	1	3	9	5	0	1	2	3	4	5	0	5
5	0	4	1	3	5	4	1	0	0	2	2	0	4	3	3	3	3	3
4	4	1	5	0	5	4	2	4	0	0	1	3	2	1	4	3	4	1
3	4	4	1	6	5	0	5	1	2	3	4	2	2	5	3	4	6	0
2	3	4	5	0	5	1	2	0	3	3	4	1	2	3	3	5	5	4
1	0	4	4	1	6	5	0	1	0	1	2	3	4	3	3	4	4	4
0	4	4	1	7	0	5	1	2	0	0	1	3	2	1	4	3	4	4
4	0	5	0	7	2	4	1	4	0	1	2	3	5	4	0	5	5	4
5	0	6	5	0	7	2	4	1	4	0	1	2	3	5	4	0	5	5
6	0	5	0	6	5	0	7	2	4	1	4	0	1	2	3	5	5	4
5	0	6	5	0	6	5	0	7	2	4	1	4	0	1	2	3	5	5
4	0	5	0	5	0	6	5	0	7	2	4	1	4	0	1	2	3	5
3	0	4	0	4	0	5	0	6	5	0	7	2	4	1	4	0	1	2
2	1	0	0	4	0	5	0	6	5	0	7	2	4	1	4	0	1	2
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1



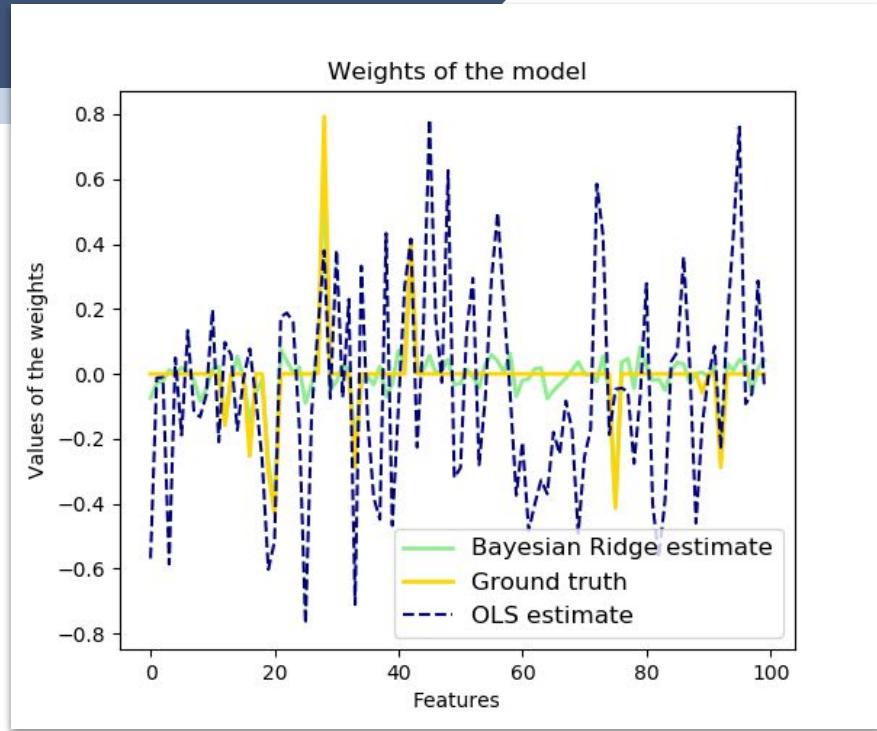
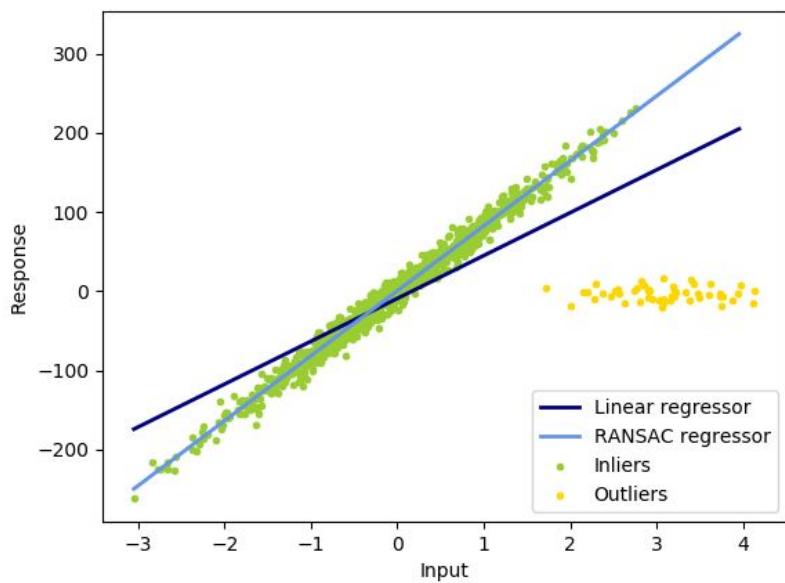


# Anomaly Detection



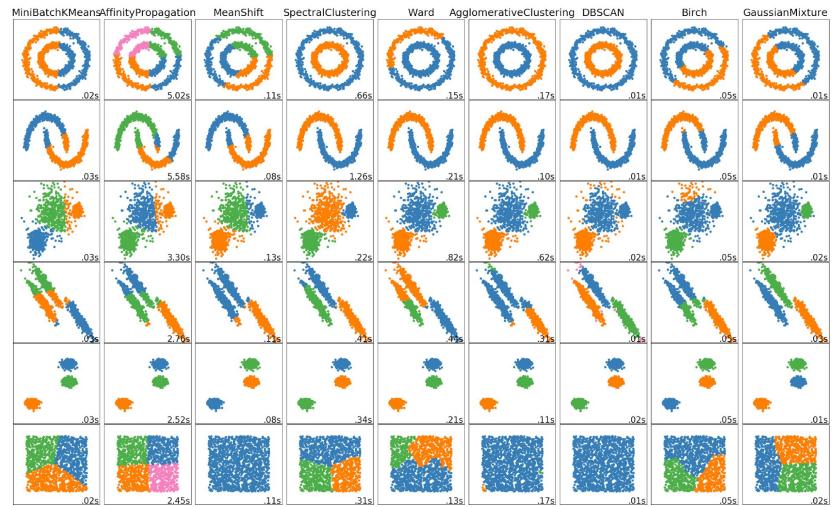
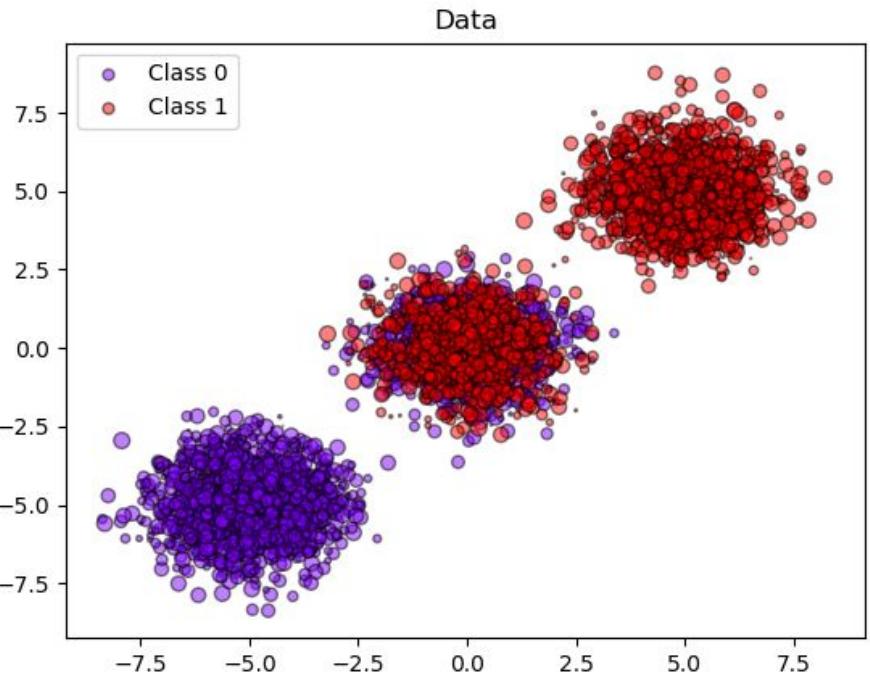


# Regression



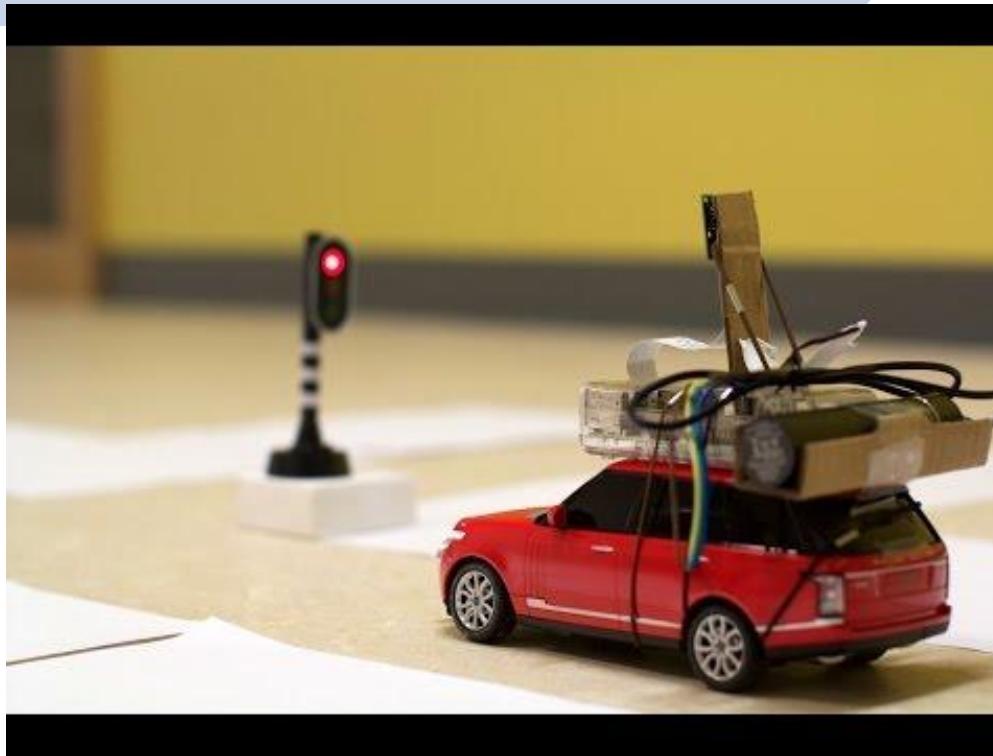


# Clustering



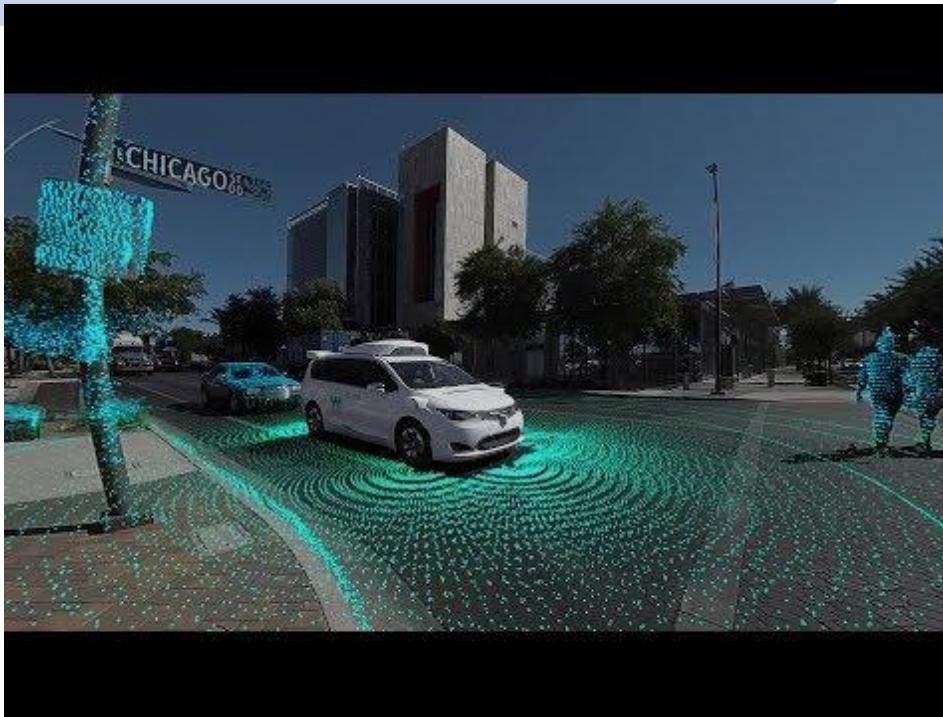


# Reinforcement Learning





# Reinforcement Learning

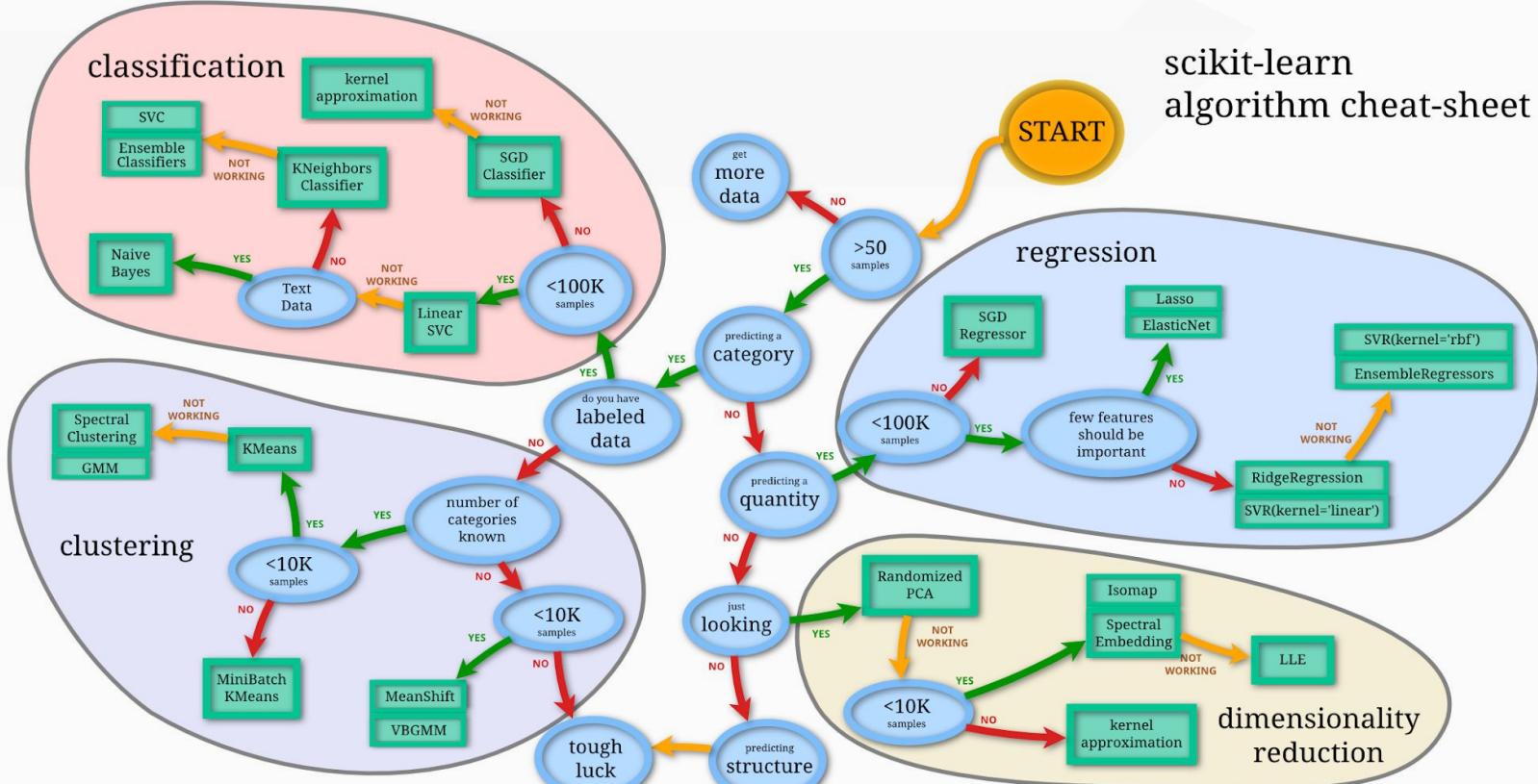


# 1.4

## Algorithms

ขั้นตอนวิธีของการเรียนรู้ของเครื่อง

# scikit-learn algorithm cheat-sheet



Back

scikit  
learn



## Machine Learning Frameworks



**Apache SINGA**

A General Distributed  
Deep Learning Platform

# Caffe



**Amazon** Machine Learning

[aws.amazon.com/machine-learning](http://aws.amazon.com/machine-learning)



Torch  
Accord.NET  
Apache Mahout  
Theano  
Brainstorm  
etc.,



# Python Libraries for Data Science - ชุดคำสั่ง

## Core Libraries

- Numpy
- Scipy
- Pandas
- Statsmodels

## Data Acquisition

- Scrapy
- Quadl
- Requests
- Selenium
- lxml
- BeautifulSoup

## Visualization

- Matplotlib
- Seaborn
- Bokeh
- Plotly

## Machine Learning

- SciKit-Learn
- Theano
- Tensorflow
- TFLearn
- Keras
- NLTK
- Gensim



## Software & Downloads



The screenshot shows the Anaconda Distribution download page. At the top, there's a green header with the Anaconda logo and a menu icon. Below it, a large green banner features the text "Download Anaconda Distribution" and "Version 5.2 | Release Date: May 30, 2018". It also includes icons for Windows, Mac, and Linux operating systems. The main content area has a white background with the text "High-Performance Distribution" and "Easily install 1,000+ [data science packages](#)".

<https://www.anaconda.com/download/>

### Includes...

- Numpy
- Pandas
- Tensorflow
- Scrapy
- Seaborn
- ...



Hand-on Lab <http://gg.gg/b1jmx>

<https://mybinder.org/v2/gh/baldwint/PythonDataScienceHandbook/py35>

- 10.45 - 12.00 **Data Modeling**
- **Lunch**
- 13.00 - 14.30 **Model Validation**
- **Coffee Break**
- 14.45 - 16.30 **Case Study**