

# Introduction to Data Science and Analytics

*Summer School 2015*

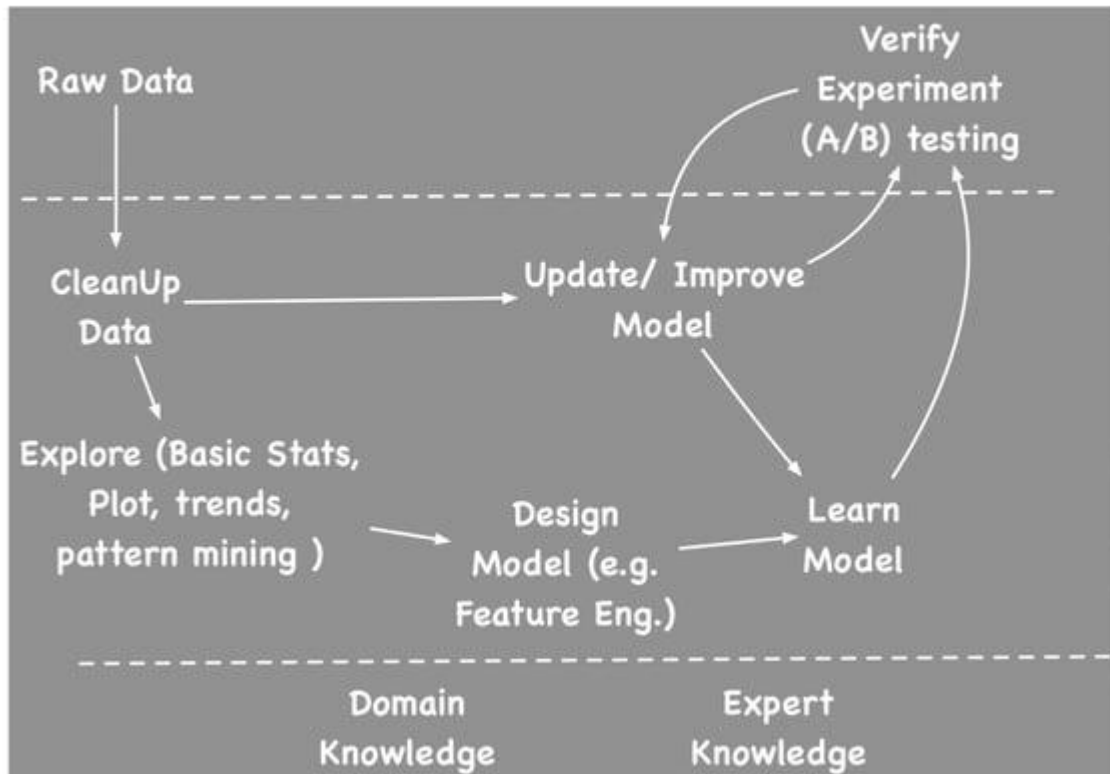
# What is Data Science?

Extraction of knowledge from large volumes of data that are structured or unstructured.

It is a continuation of the fields **data mining** and **predictive analytics**

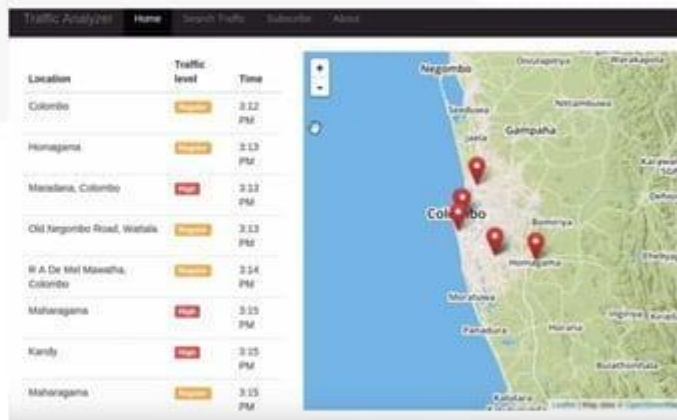


# Data Science Pipeline



# Example ( Road.lk) traffic Feed

1. Data as tweets
2. Extract time, location, and traffic level using NLP
3. Explore data
4. Model based on time, and it is a holiday
5. Predict traffic given a time and location.



## Data Cleanup

Real data is messy, often needs to be cleaned up before useful.

- Bad formats - ignore or treat like missing data
- Missing Data - extrapolate or remove data line
- Useless variables - remove
- Wrong data - e.g. aaa, bbb, joe, some might be deliberate lie, or 99 may be a code for N/A

## Data Cleanup (Contd.)

- Transform variables ( date formats, String to int)
- Create derived variables
  - Derive country from IP
  - age from ID card number
- Normalize strings
  - e.g. stemm or use phonetic sounds
  - different spellings and nicknames ( William->Bill)
- Feature value rescaling (e.g. most ML algorithms needs value to rescaled to 0-1 range).
- Enrich (e.g. lookup and add age from profile)

## Data Exploration

Understand, and get a feel for what is **expected** (models => densities, constraints) and **unexpected/ residuals** (errors, outliers)

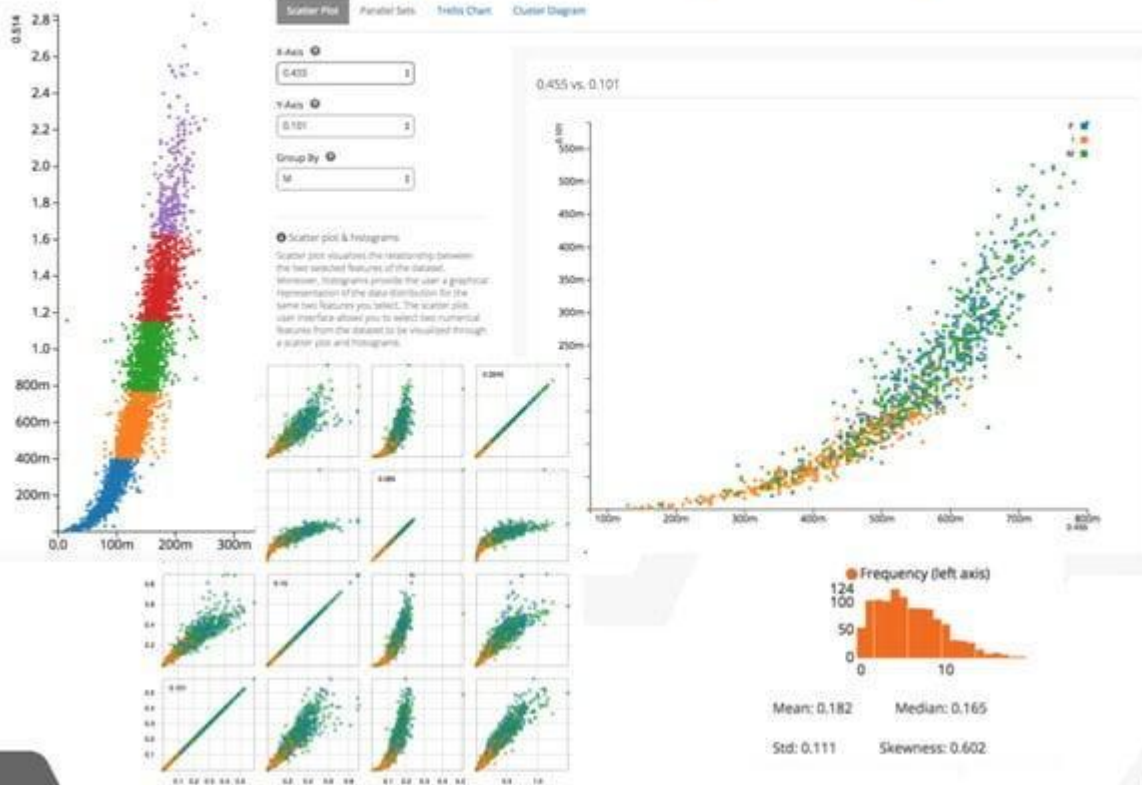
- o think what this is data about? domain, background, how it is collected, what each fields mean and range of values.
- o head, tail, count, all descriptives (Mean, Max, median, percentiles .. ) - Five number Summary. Min. 1st Qu. Median Mean 3rd Qu. Max.
- o run a bunch of count/group-by statements to gauge if I think it's corrupt.

## Data Exploration (Contd.)

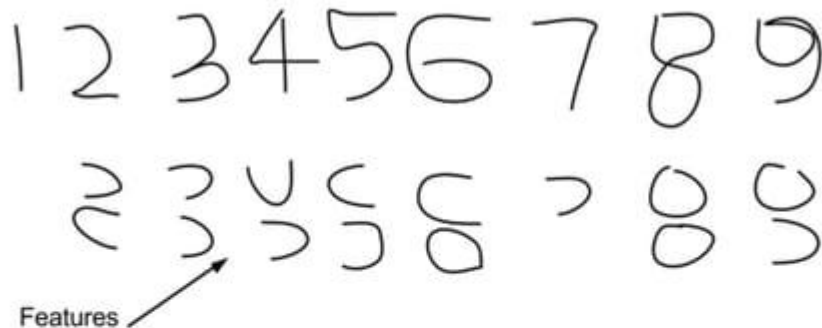
- o Plot - take random sample and explore ( scatter plot)
  - o e.g. Draw scatter plot or Trellis Plot
- o Find Dependencies between fields
  - o Calculate Correlation
  - o Dimensionality reduction
  - o Cluster and look visualize clusters
- o Look at frequency distribution of each field and try to find a known distribution if possible.



# Data Exploration (Contd.)



# Feature Engineering



- o Feature engineering is the art of finding feature that leads simplest decision algorithm. ( Good features allow a simple model to beat a complex model.)
- o Best features may be a subset, or a combination, or transformed version of the features.

# How to do Feature Engineering?

- Manually pick by domain experts and trial and error.
- Search the possible combinations by training and combining subsets (e.g. Random Forest)
- Use statistical concepts like correlation and information criteria
- Reduce the features to a low dimension space using techniques like PCA.
- Automatic Feature Learning though Deep Learning
- ...

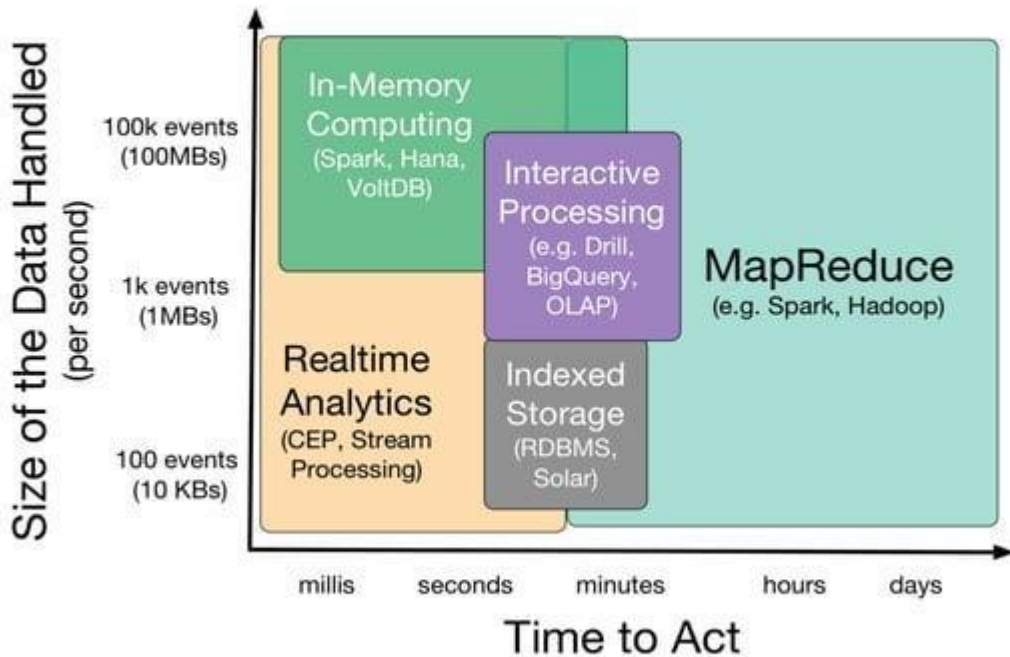
# Analysis

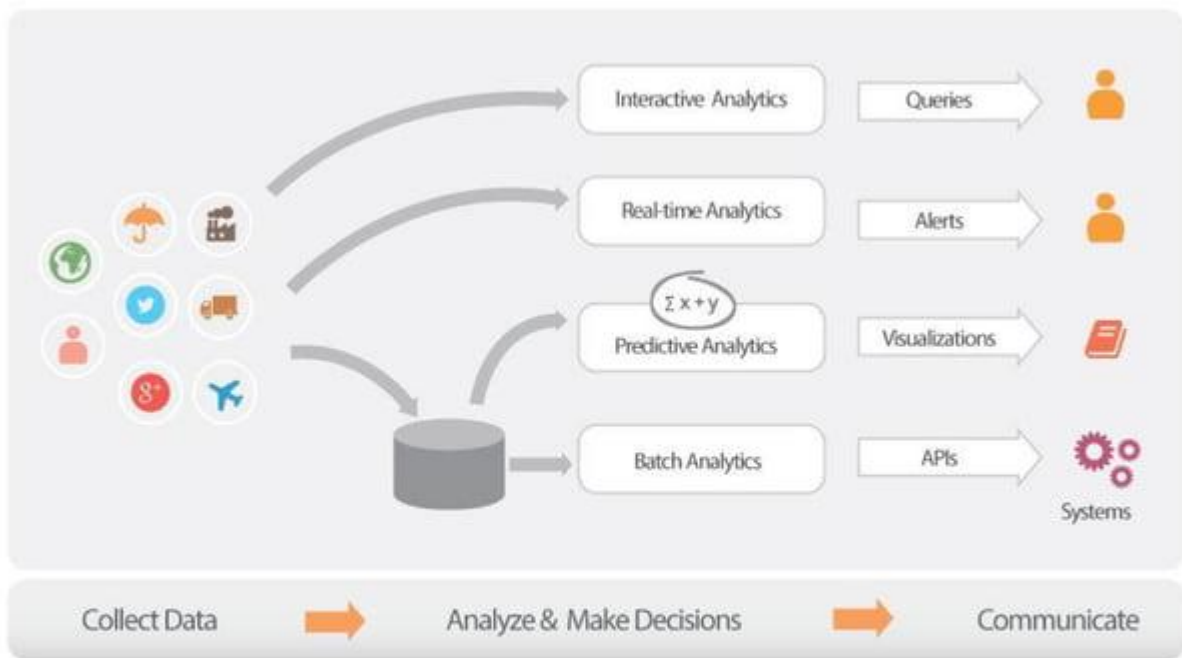
- o Goal of analysis is to extract knowledge
- o This knowledge usually come in one of the two forms
  - o KPI (Key Performance Indicators)
    - Describe key measurement for what is being measured. (e.g. revenue per year, profit margin, revenue for sqft in retail, revenue per employer)
  - o Models to describe or predict the data
    - e.g. Machine Learning models or Statistical models

## 4 Analysis types by time to decision

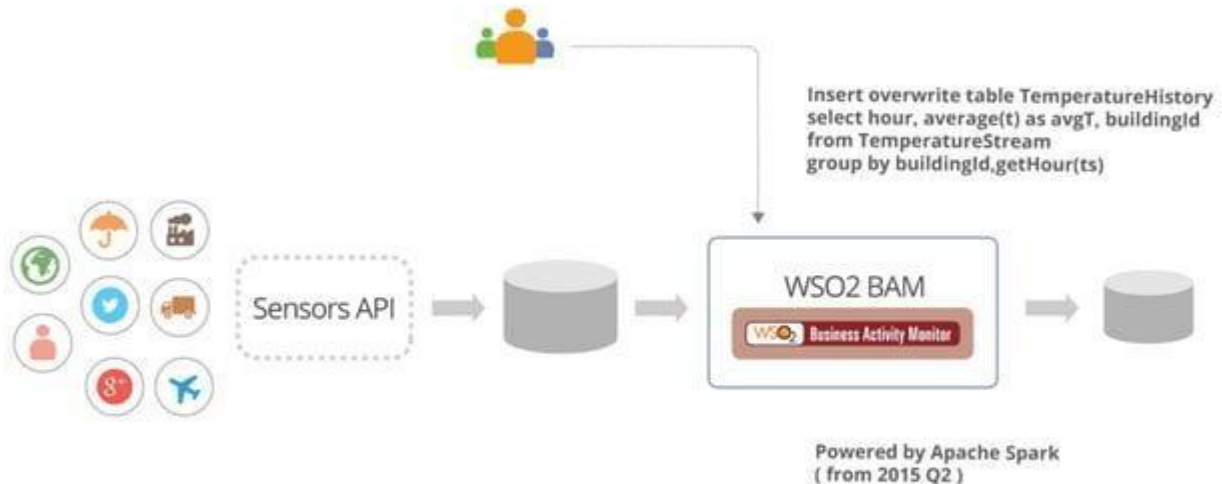
- Hindsight ( what happened?)
  - Done using Batch Analytics like MapReduce
- Oversight ( what is happening?)
  - Done using Realtime Analytics technologies like CEP
- Insight ( why things happening?)
  - Done with Data Mining and Unsupervised learning algorithms like Clustering
- Foresight ( what will happen?)
  - Done by building models using Machine learning or one of other techniques

# Data Analytics Tools Landscape





# Batch Analytics: SparkSQL



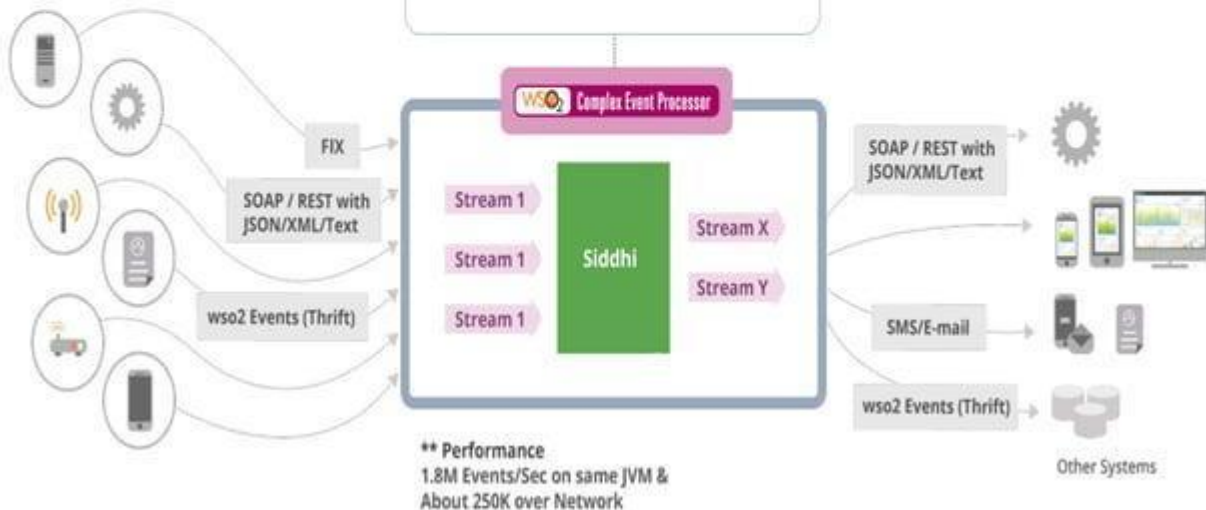


# Realtime Analytics: Complex Event Processing



```
from stockQuotes#window.time(1 min)
join Tweets#window.time(1 min)
  on StockQuotes.symbol==Tweets.company
select *
insert into PredictedstockQuotes;
```

Filter Transformation Window +  
{ Aggregation, group by}  
Join Event Sequence Event Table



# Interactive Analytics

- Define Indexes on Collected data ( Streams)
- Issue, dynamic queries and get results right away. ( Powered by Apache Lucene)
- Shows multiples events from same activity together using custom defined activity IDs
- Useful for data exploration
- Powered by Apache Lucene, with support for Index Sharding

```
Welcome to interactive analytics SQL shell
This interactive shell lets you execute Spark SQL commands against a Spark cluster
Initializing Spark client...
SparkSQL> CREATE TEMPORARY TABLE testTable USING CarbonAnalytics OPTIONS (tableName
SparkSQL> SELECT * FROM testTable;
Show 10 entries
```

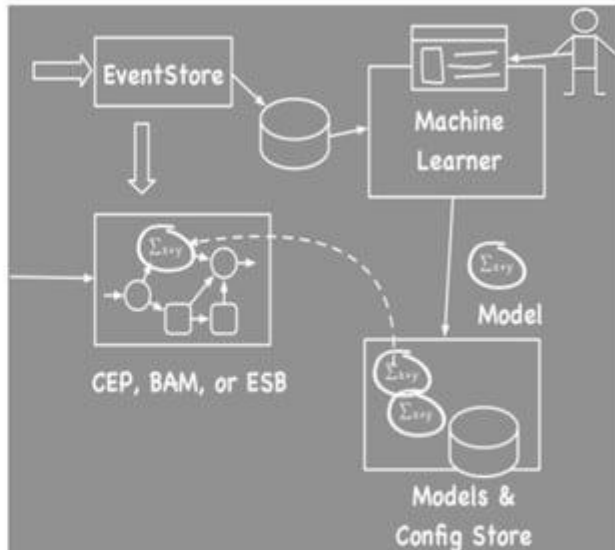
x	y
1	100
2	200
3	300
4	400
5	500
6	600
7	700
8	800
9	900
10	1000

```
Showing 1 to 10 of 30 entries

SparkSQL> CREATE TEMPORARY TABLE destTable USING CarbonAnalytics OPTIONS (tableName
SparkSQL> INSERT INTO TABLE destTable SELECT * FROM testTable;
```

# Predictive Analytics

- o Build models and use them with WSO2 CEP, BAM and ESB using WSO2 Machine Learner Product ( 2015 Q3)
- o Build model using R, export them as PMML, and use within WSO2 CEP



# WSO2 Machine Learner

- o Sample, explore, and understand data through visualizations
- o A wizard to configure, train machine learning models, and select the best model
- o Find and use those models with WSO2 CEP, BAM and ESB
- o Powered by Apache Spark MLlib



# Building Decision Models

A model describe how a system behave when inputs changes. There are many ways to build models.

- o Regression models and ML Models
  - Time series models
- o Statistical models
- o Physical Models - based on physical phenomena. They include 6-DoF flight models, space flight models Weather models.
- o Mathematical Models

*the signal and the  
and the noise and  
the noise and the  
noise and the no  
why so many and  
predictions fail –  
but some don't t  
and the noise and  
the noise and the  
nate silver noise  
noise and the no*

see <https://icrunchdatanews.com/what-are-predictive-models/>

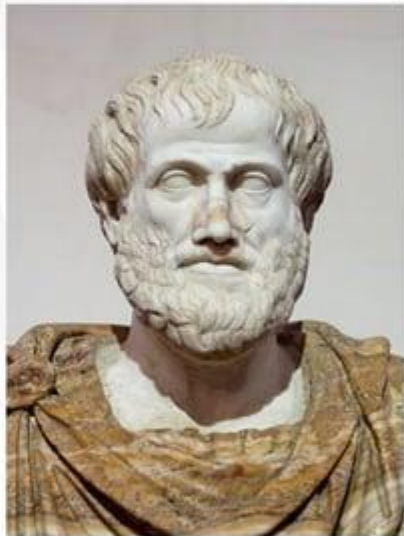
# Verification

- o All is good, now you have a model. You must verify that it is correct before using it in the real world.
- o Prediction can be verified by waiting for events to occur
- o Relationships like causality (e.g. having free shipping leads a customer to buy more) must be verified with A/B testing
- o Let's look at few of pitfalls



# Pitfalls: Experiment vs Observation

- If you follow scientific method, you would do experiments, and they have control sets ( A/B) tests.
- Bigdata does not have a control set, it is rather observations. ( we observe the world as it happens)
- So what we can tell are limited.
- Correlation does not imply Causality!!
  - Send a book home example [1]
  - All big buyers have free shipping





# Causality: What can we do?

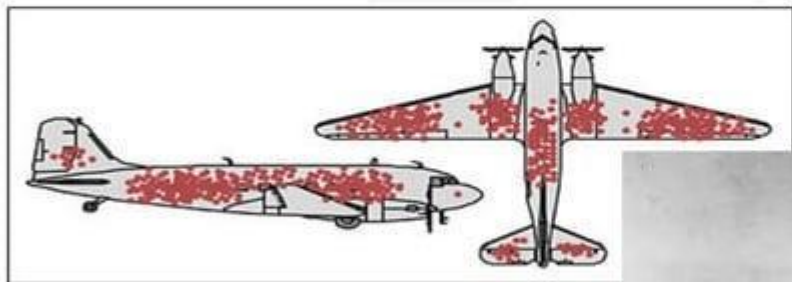
- o Option 1: We can act on correlation if we can verify the guess or if correctness is not critical (Start Investigation, Check for a disease, Marketing )
- o Option 2: We verify correlations using A/B testing or propensity analysis



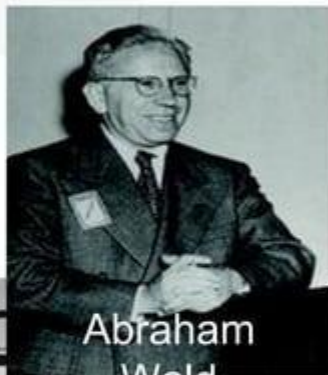


## Pitfalls: Think about the Missing Data

- o WW II, Returned Aircrafts and data on where they were hit?
- o How would you add Armour?



Cre



Abraham  
Wald



<http://www.fastcodesign.com/1671172/how-a-story-from-world-war-ii-shapes-facebook-today>, Pic from <http://www.phibetaiota.net/2011/09/defdog-the-importance-of-selection-bias-in-statistics/>

# Communicate: Dashboards

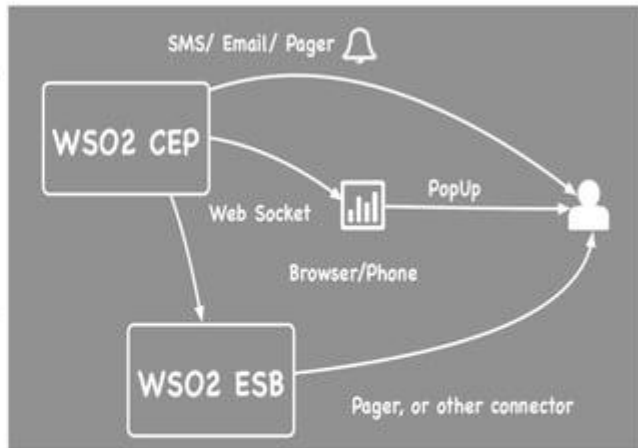


- o Dashboard give an “Overall idea” in a glance (e.g. car dashboard)
- o Support for personalization, you can build your own dashboard.
- o Also the entry point for Drill down
- o How to build?
  - o WSO2 DAS supports a gadget generation Wizard
  - o Or you can write your own Gadgets using D3 and Javascript.

## Communicate: Alerts

- o Detecting conditions can be done via CEP Queries. Key is the “Last Mile”.

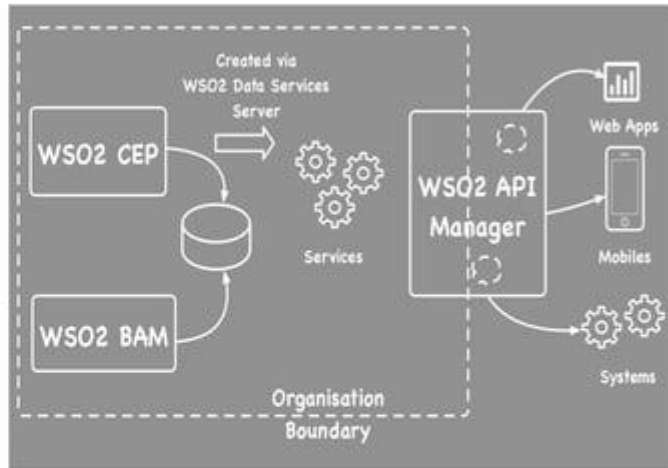
- o Email
- o SMS
- o Push notifications to a UI
- o Pager
- o Trigger physical Alarm



- o How?
  - o Select Email sender “Output Adaptor” from CEP, or send from CEP to ESB, and ESB has lot of connectors

# Communicate: APIs

- With mobile Apps, most data are exposed and shared as APIs (REST/Json ) to end users.
- Need to expose analytics results as API
- Following are some challenges
  - Security and Permissions
  - API Discovery, Billing, throttling, quotas & SLA
- How?
  - Write data to a database from CEP event tables
  - Build Services via WSO2 Data Service
  - Expose them as APIs via API Manager

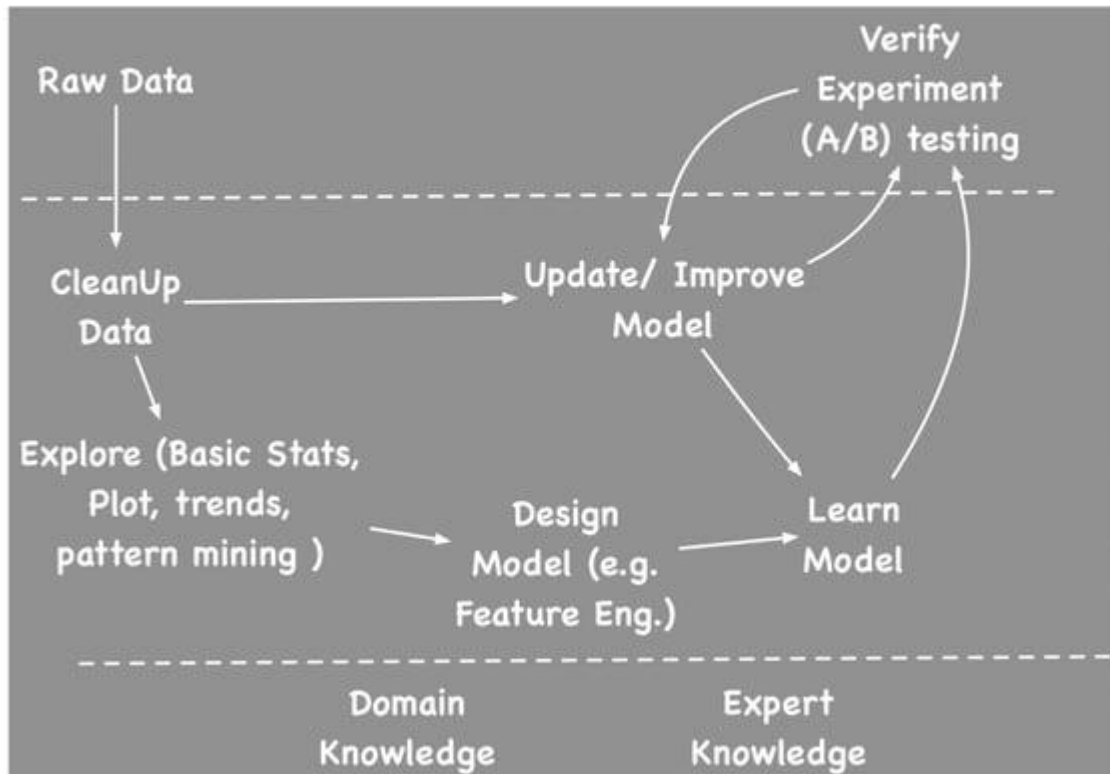


# Communicate: Realtime Soccer Analytics



<https://www.youtube.com/watch?v=nRI6buQONOM>

# Data Science Pipeline



# Conclusion

- o Data Science is extracting knowledge by analyzing data
- o Discussed the pipeline and tools you can use to do that
- o Rest of summer school will look at different aspects in detail.
- o All tools discussed are available free under Apache Licence.

