

## Characteristics of Big Data and Dimensions of Scalability

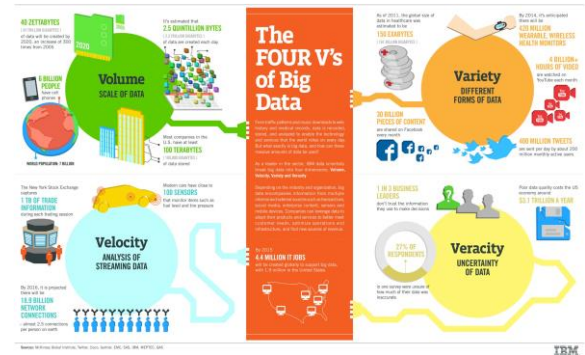
Pejman Rasti

Email: [prasti@esaip.org](mailto:prasti@esaip.org)  
[pejman.rasti@univ-angers.fr](mailto:pejman.rasti@univ-angers.fr)

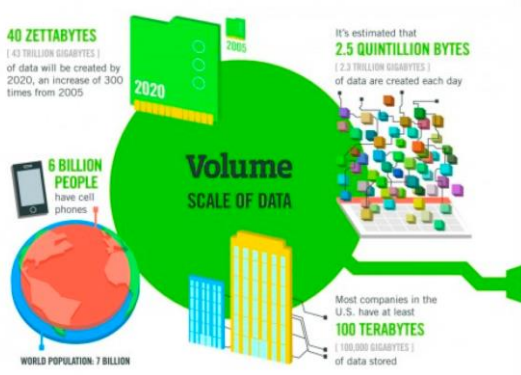
Course Website: Access from your "Moodle" portal

1

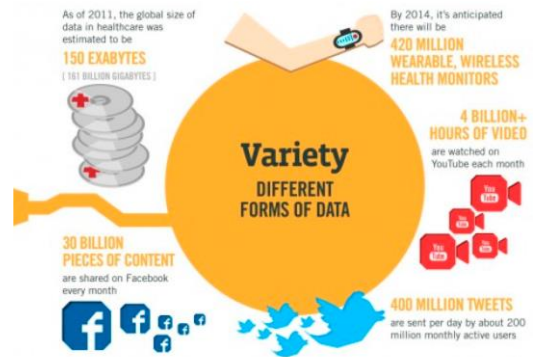
## Characteristics of Big Data



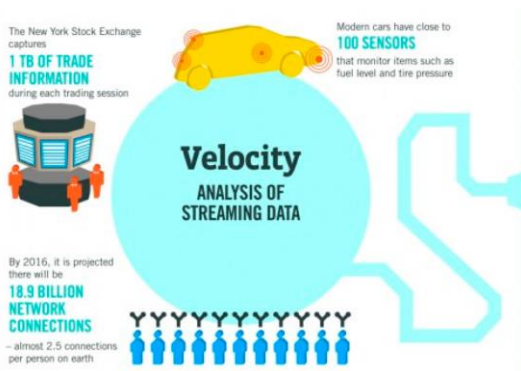
2



3

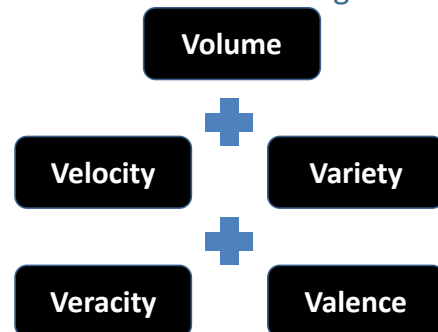


4

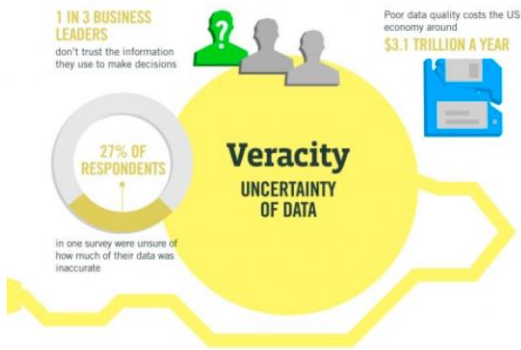


5

## Characteristics of Big Data

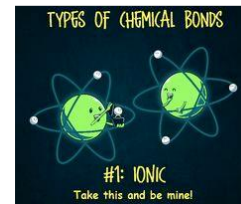


6



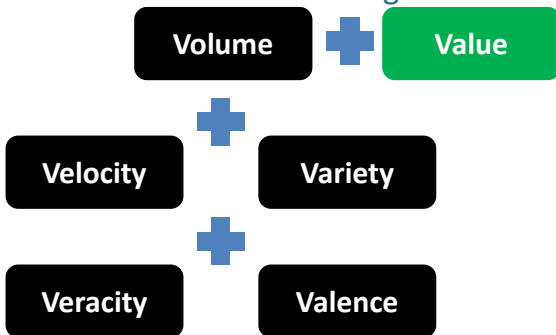
7

## Valence



8

## Characteristics of Big Data



9

## Volume

Volume = Size

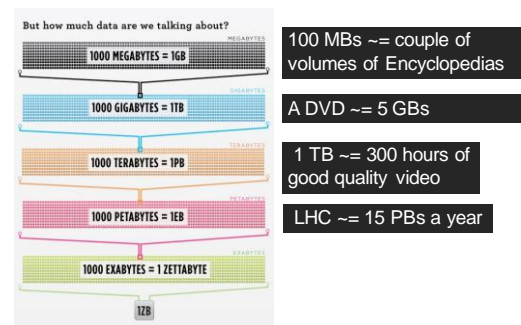


10

## Every minute...

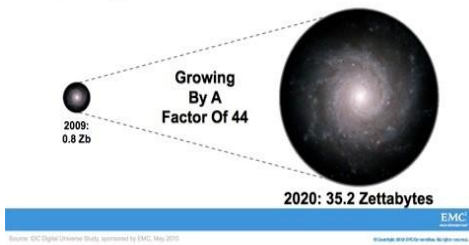


11



12

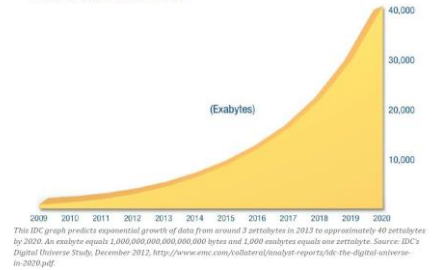
## The Digital Universe 2009-2020



13

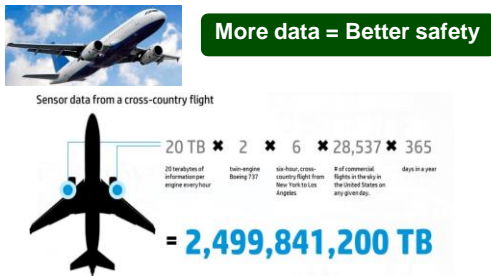
## Exponential data growth!

The Digital Universe: 50-fold Growth from the Beginning of 2010 to the End of 2020

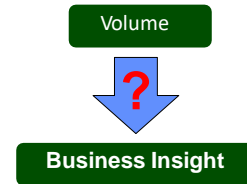


14

## Relevance of Volume for Us?

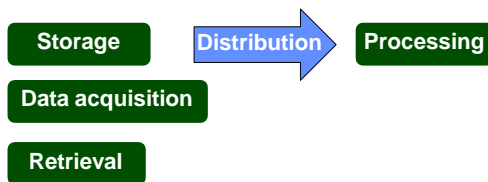


15



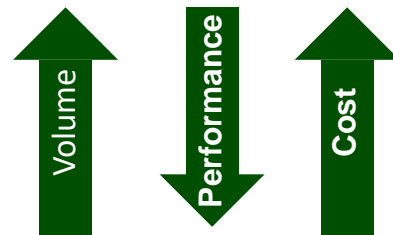
16

## Challenges: Storage and more...



17

## Processing Big Data

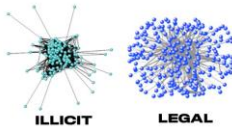


18



## Variety within a Type

- Think of an email collection
  - Sender, receiver, date... **Well-structured**
  - Body of the email **Text**
  - Attachments **Multi-media**
  - Who-sends-to-whom



25

## Variety within a Type

- Think of an email collection
  - Sender, receiver, date... **Well-structured**
  - Body of the email **Text**
  - Attachments **Multi-media**
  - Who-sends-to-whom **Network**
  - A current email cannot reference a past email **Semantics**

26

## Variety within a Type

- Think of an email collection
  - Sender, receiver, date... **Well-structured**
  - Body of the email **Text**
  - Attachments **Multi-media**
  - Who-sends-to-whom **Network**
  - A current email cannot reference a past email **Semantics**
  - Real-Time? **Availability**

27

## Scalability Issues

- Impact of data variety
  - Harder to ingest
  - Difficult to create common storage
  - Difficult compare and match data across variety
  - Difficult to integrate
  - Management and policy challenges



28

## Velocity

Velocity == Speed

$$\bar{v} = \frac{\Delta x}{\Delta t}$$

Speed of creating data  
Speed of storing data  
Speed of analyzing data

29

Big Data → Real-time action



30

Late decisions



Missing opportunities

31



32

How to decide what to pack ?

Use weather information  
of last year at this time?



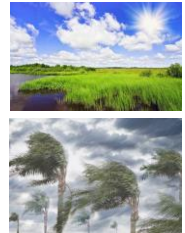
33

How to decide what to pack ?

Use weather information  
of last month ?

OR

Use weather status of  
this week or today ?



34

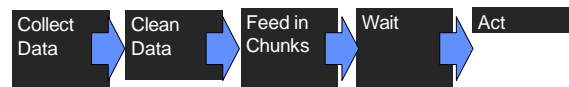


Action

35

Real-time Processing

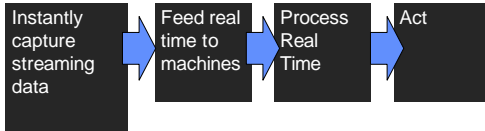
Batch Processing



36

## Real-time Processing

### Real-Time Processing



37

Batch Processing → Incomplete

Real-Time Processing → Fast

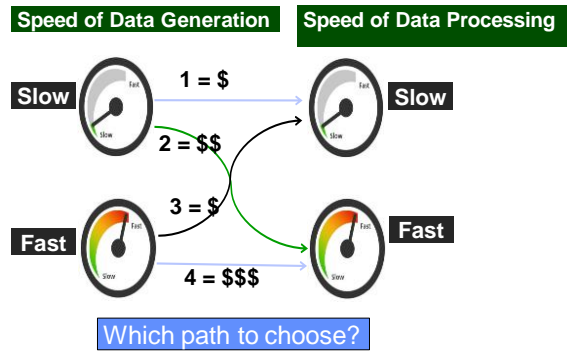
38

Rate needed for data-driven actions



Rate of generation and processing of data

39



40

Veracity == Quality

Validity  
Volatility



42

Veracity == Quality

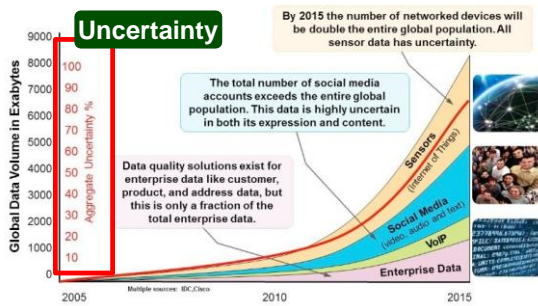
Accuracy of data

Reliability of the data source

Context within analysis

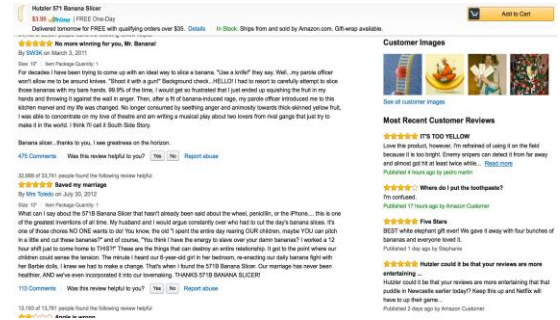
43



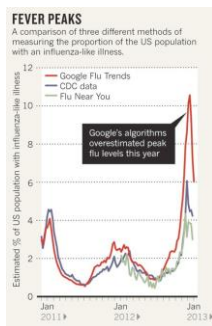


44

## When sentiment analysis doesn't work?



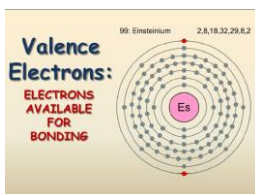
45



46

Valence == Connectedness

Valence – a Concept from Chemistry



48

Veracity == Quality

Accuracy of data

Reliability of the data source

Context within analysis



Uncertainty

Provenance

47

Valence – Measure of Connectivity



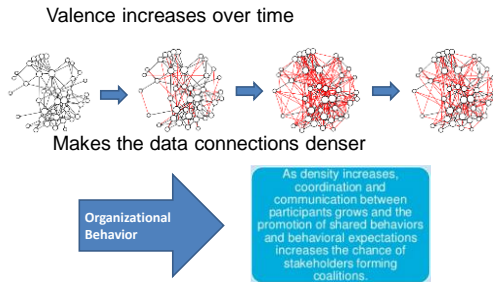
Data Connectivity

- Two data items are connected when they are related to each other
- Valence
  - Fraction of data items that are connected out of total number possible connections

49



## Why worry about Valence?



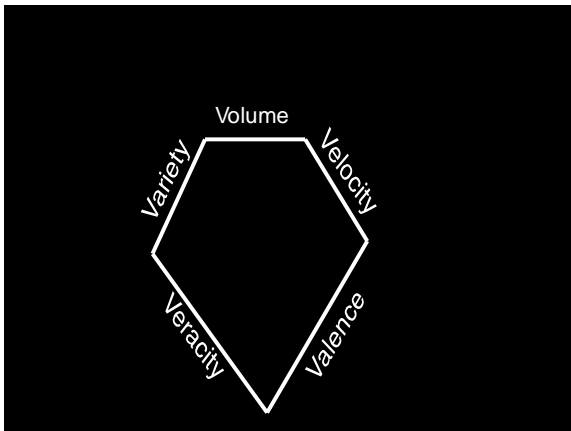
50

## Valence: Challenges

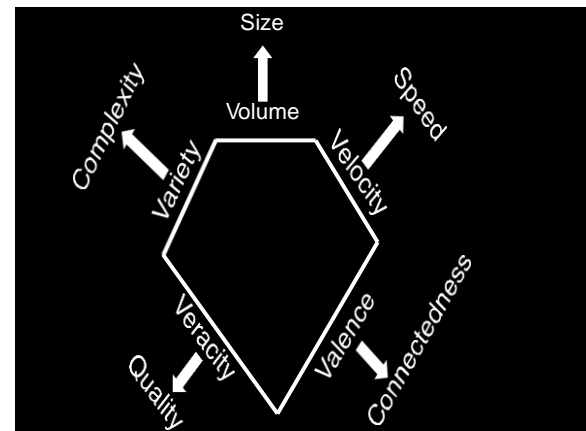
- More complex data exploration algorithms
- Modeling and prediction of valence changes
- Group event detection
- Emergent behavior analysis



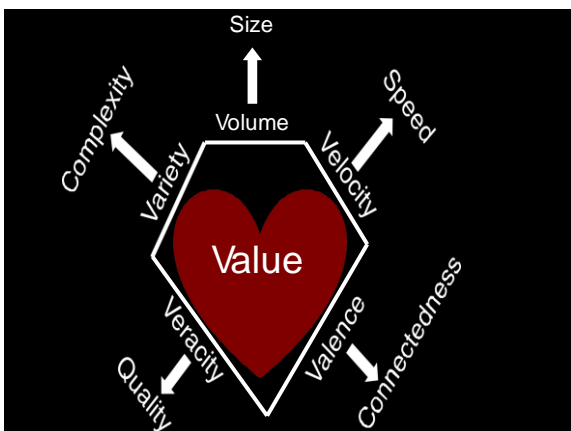
51



52



53



54



## Eglen Inc. Big Data Case:

Catch The Pink Flamingo

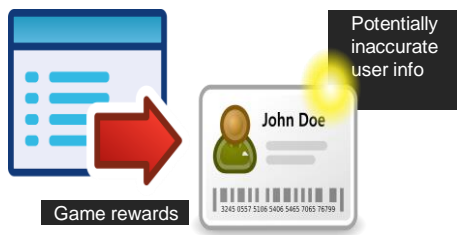
55



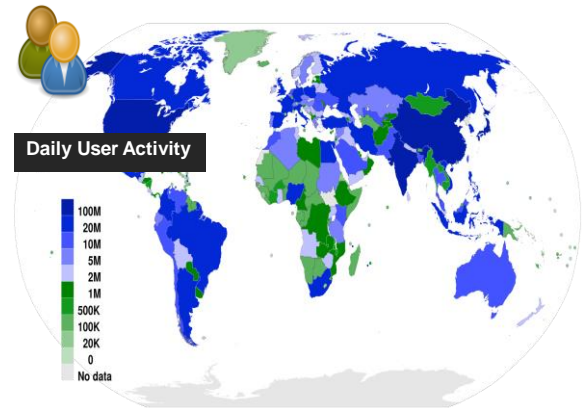
56



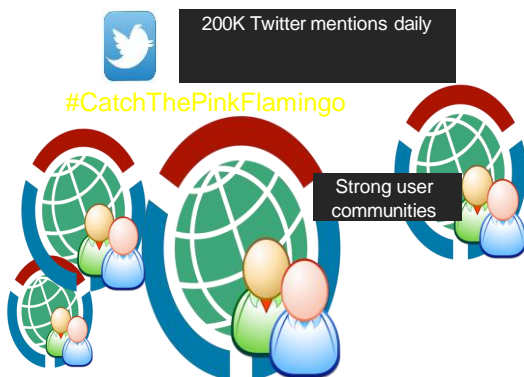
57



58



59



60



61

**Data Source**

<b>Machine</b>	<ul style="list-style-type: none"> <li>User activity logs</li> </ul>
<b>People</b>	<ul style="list-style-type: none"> <li>Twitter conversations</li> </ul>
<b>Organization</b>	<ul style="list-style-type: none"> <li>User demographic info</li> <li>Game stats</li> </ul>

62

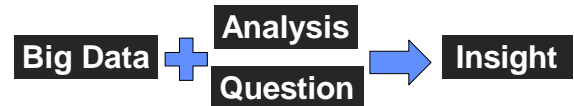
Dimension	
<b>Volume</b>	<ul style="list-style-type: none"> <li>Big daily workload and associated data on players and game stats</li> </ul>
<b>Variety</b>	<ul style="list-style-type: none"> <li>Multiple types of data</li> </ul>
<b>Velocity</b>	<ul style="list-style-type: none"> <li>Real-time analysis of usage activity</li> </ul>
<b>Veracity</b>	<ul style="list-style-type: none"> <li>Demographic info not accurate</li> </ul>
<b>Valence</b>	<ul style="list-style-type: none"> <li>Connections between players</li> </ul>

63



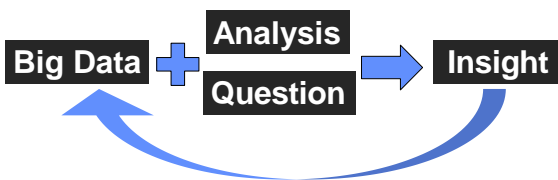
64

**Insight** → **Data Product**



65

**Insight** → **Data Product**



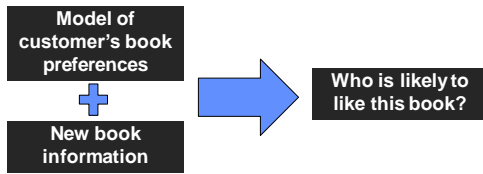
66

**Book Recommendations**



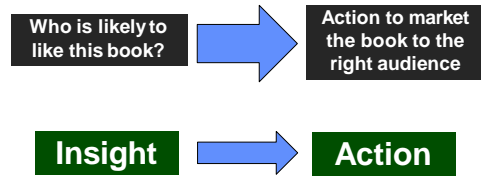
67

## Find Potential Audience for a Book



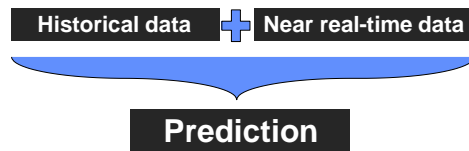
68

## Market a New Book



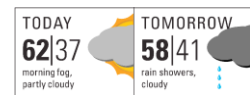
69

## Actionable Information



70

## Prediction

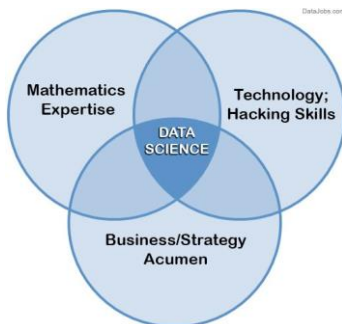


## Action

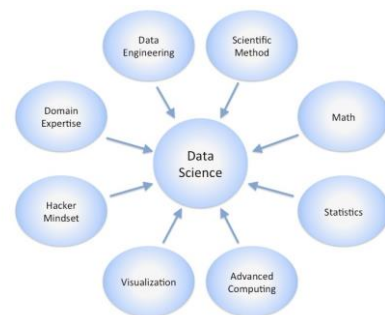


71

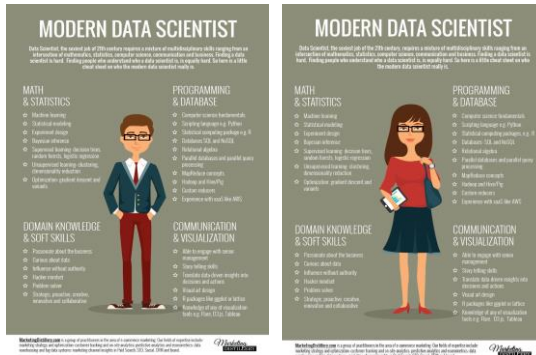
## Data Science is Team Work!



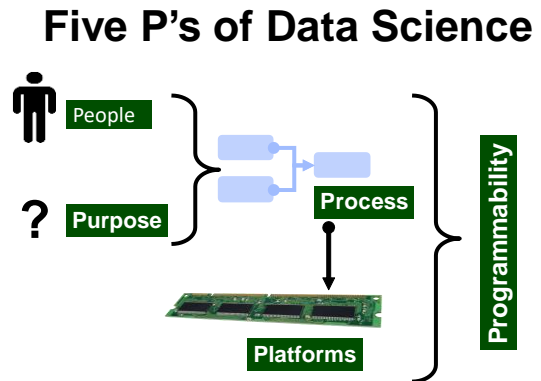
72



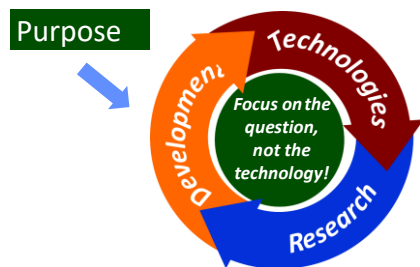
73



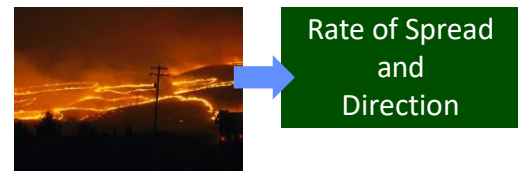
74



75



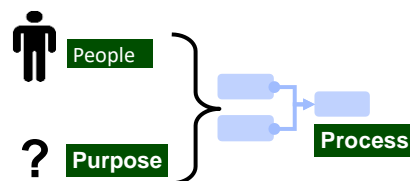
76



77

Let's not dive into the techniques yet! What is the problem at large? How do you see yourself solving it?

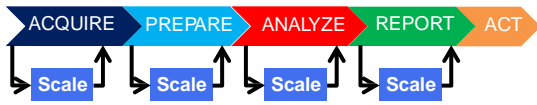
78



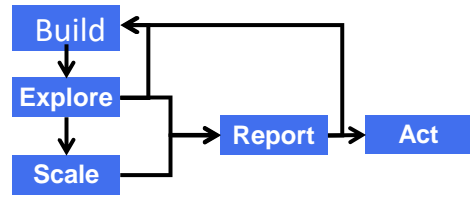
79

Big Data Engineering

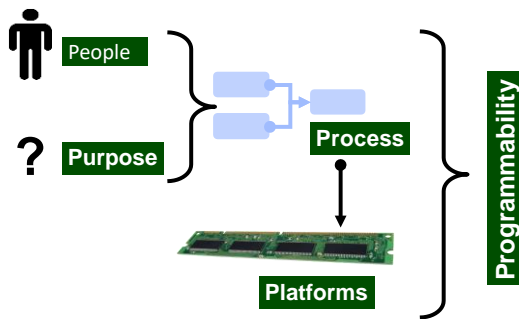
Computational Big Data Science



80



81

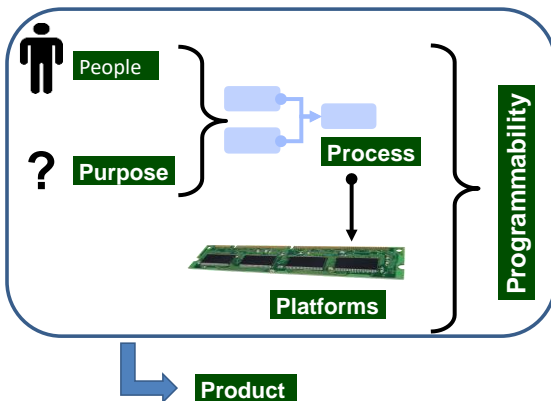


82

Process: Build metrics for accountability

Cost  
Timeline  
Planning of deliverables  
Expectations  
Purpose

83



84

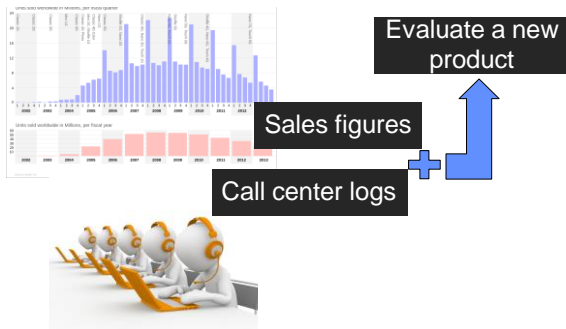
Define the Problem



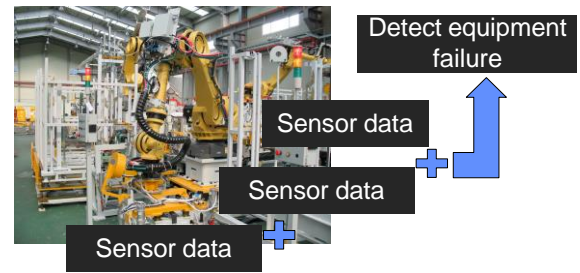
"A problem well defined  
is a problem half  
solved."

Charles F. Kettering

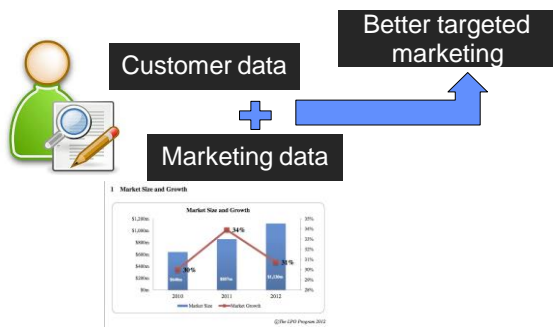
85



86



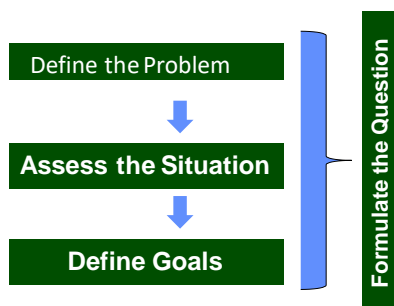
87



88



89



90

**Steps in the Data Science Process**

91





### Step 1: Acquire Data



Identify data sets  
Retrieve data  
Query data

92



### Step 2: Prepare Data

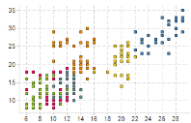
Step 2-A: Explore

Step 2-B: Pre-process

93

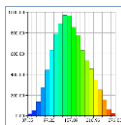


### Step 2-A: Explore Data



Preliminary analysis

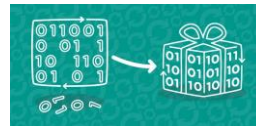
Understand nature of data



94



### Step 2-B: Pre-process Data



Clean Integrate Package

95



### Step 3: Analyze Data



Select analytical techniques  
Build models

96



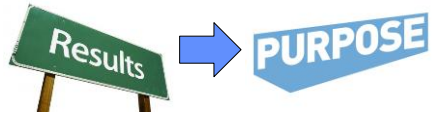
### Step 4: Communicate Results



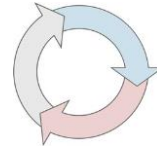
97



### Step 5: Apply Results



98



**Iterative process**

99



100

**Data comes from many places...**



**...with many ways to access it**

101

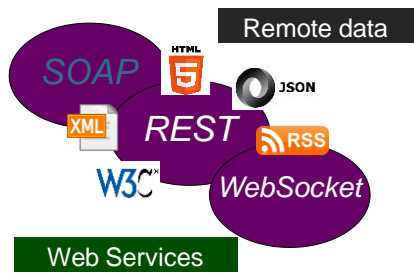


**SQL and query browsers**

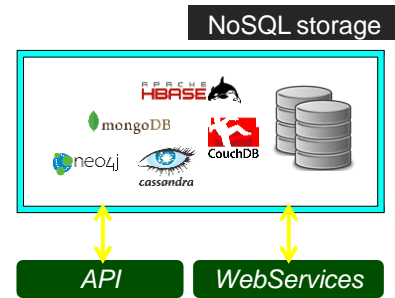
102



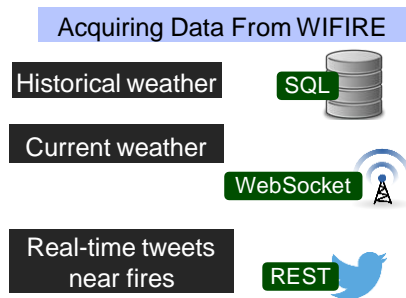
103



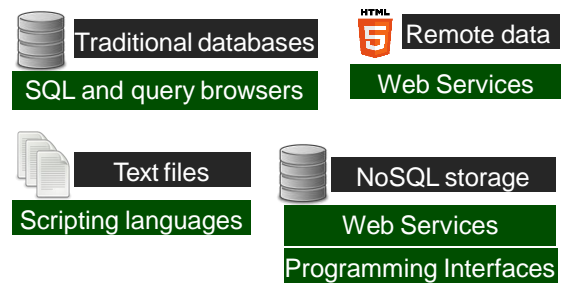
104



105

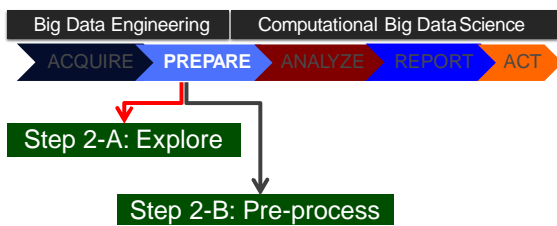


106



107

## Exploring Data



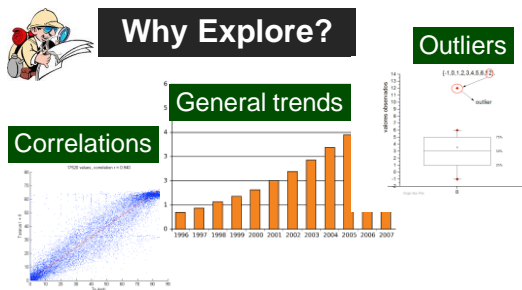
108



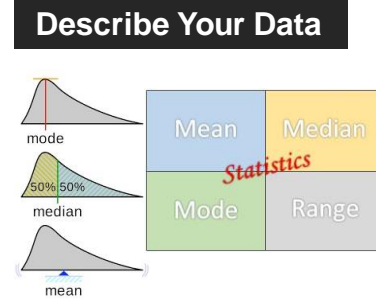
## Why Explore?

Goal: Understand your data

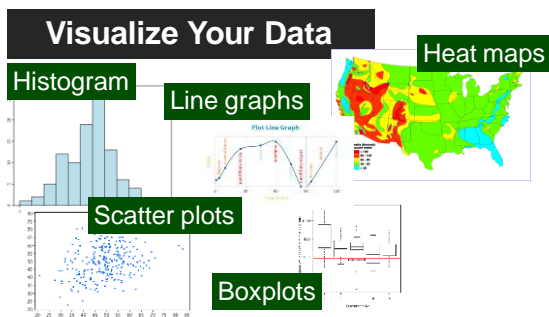
109



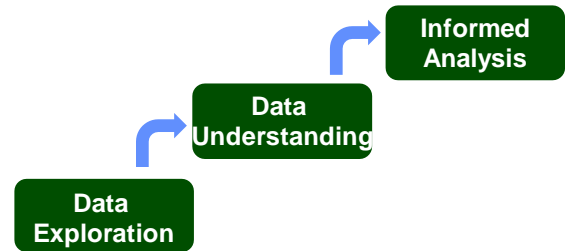
110



111

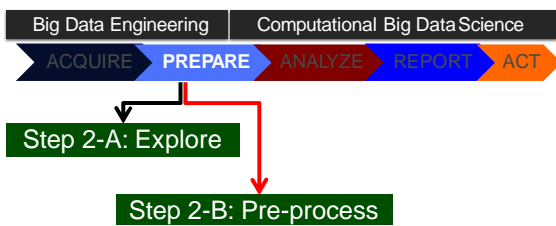


112



113

### Pre-processing Data

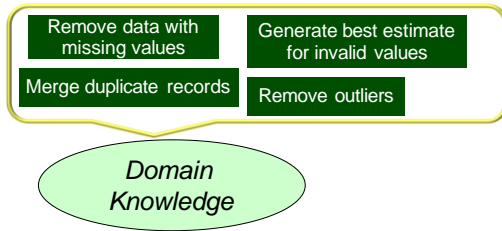


114



115

## Addressing Data Quality Issues



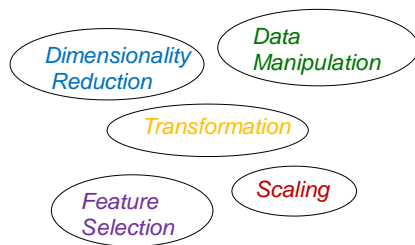
116

## Getting Data in Shape



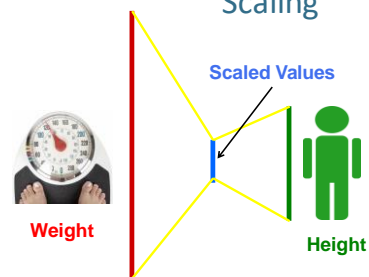
117

## Data Munging



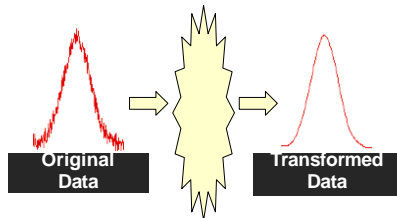
118

## Scaling



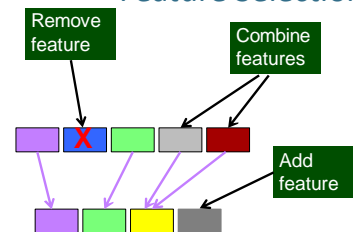
119

## Transformation



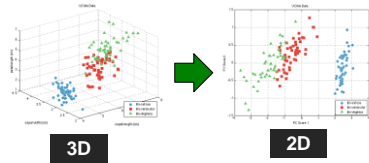
120

## Feature Selection



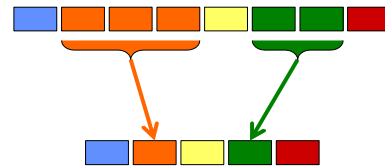
121

## Dimensionality Reduction



122

## Data Manipulation



123

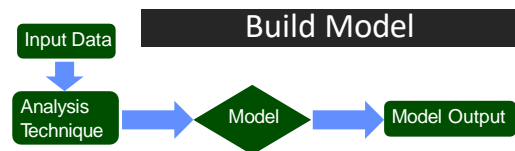
## Always Remember!

Garbage in = Garbage out



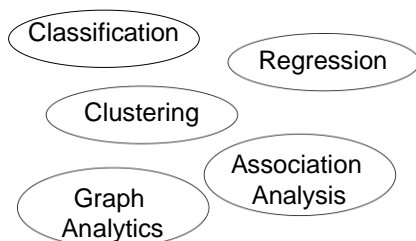
Data preparation is very important for meaningful analysis!

124



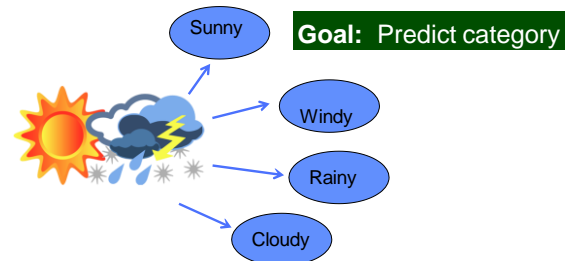
125

## Categories of Analysis Techniques



126

## Classification



127

## Regression

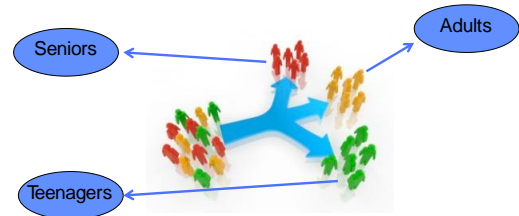
**Goal:** Predict numeric value



128

## Clustering

**Goal:** Organize similar items into groups



129

## Association Analysis

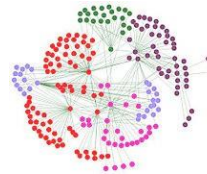
**Goal:** Find rules to capture associations between items



130

## Graph Analytics

**Goal:** Use graph structures to find connections between entities



131

## Modeling

Select technique



Build model



Validate model

132

## Evaluation of Results

133



### Classification & Regression

Predicted  
Value



Correct  
Value

134

### Clustering



135

### Association Analysis & Graph Analytics



Investigate



Validate

136

### Determine Next Steps



Repeat analysis?

Take deeper dive?

Act on results?

137

Select technique

Build model

Evaluate

Classification  
Regression  
Clustering  
Association  
Analysis  
Graph Analytics



138

Big Data Engineering

Computational Big Data Science

ACQUIRE

PREPARE

ANALYZE

REPORT

ACT

### What to Present



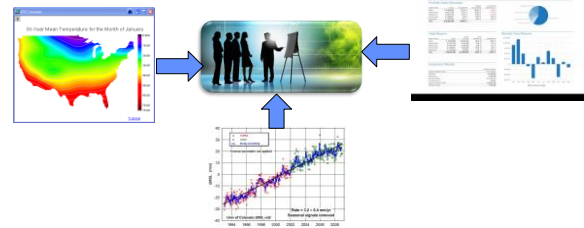
139

## What to Present



140

## How to Present



141

## Visualization Tools



142

## Present



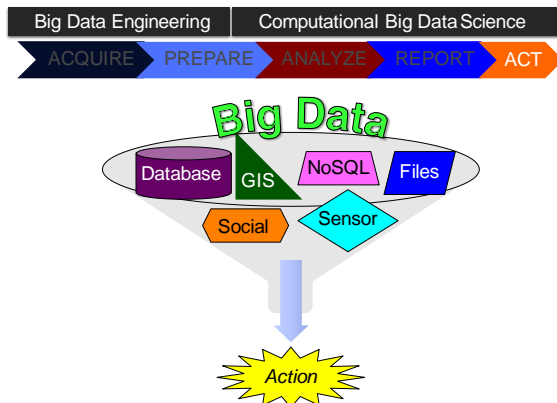
with



using

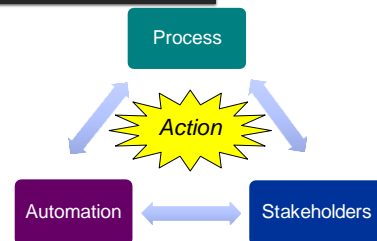


143

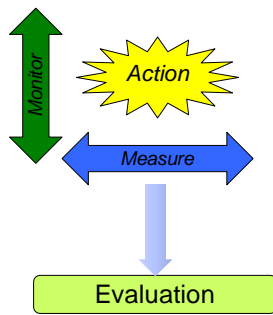


144

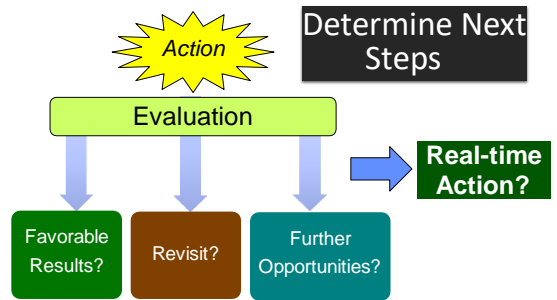
## Implementation



145



146



147

Copyright: University of California San Diego