

Steps in the Data Science Process

ACQUIRE

PREPARE

ANALYZE

REPORT

ACT

ACQUIRE

PREPARE

ANALYZE

REPORT

ACT

Step 1: Acquire Data



Identify data sets

Retrieve data

Query data

ACQUIRE

PREPARE

ANALYZE

REPORT

ACT

Step 2: Prepare Data

Step 2-A: Explore

Step 2-B: Pre-process



ACQUIRE

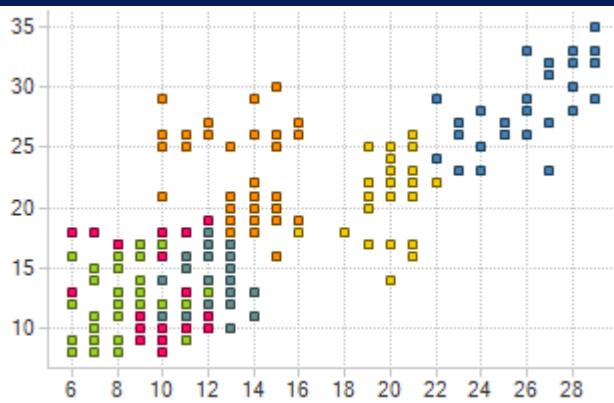
PREPARE

ANALYZE

REPORT

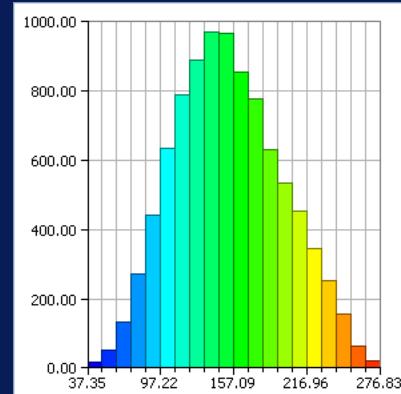
ACT

Step 2-A: Explore Data



Understand
nature of data

Preliminary
analysis



ACQUIRE

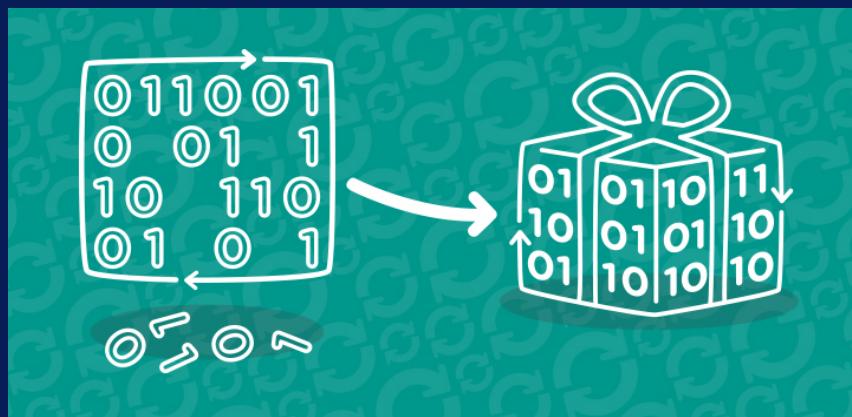
PREPARE

ANALYZE

REPORT

ACT

Step 2-B: Pre-process Data



Clean

Integrate

Package

ACQUIRE

PREPARE

ANALYZE

REPORT

ACT

Step 3: Analyze Data



Select analytical techniques

Build models

Image courtesy of Krom Krating / FreeDigitalPhotos.net

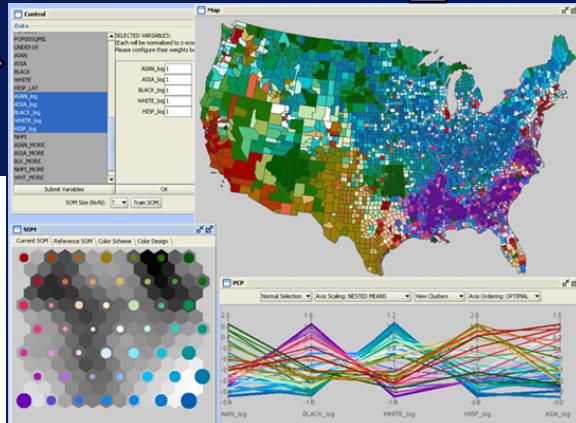
ACQUIRE

PREPARE

REPORT

ACT

Step 4: Communicate Results



ACQUIRE

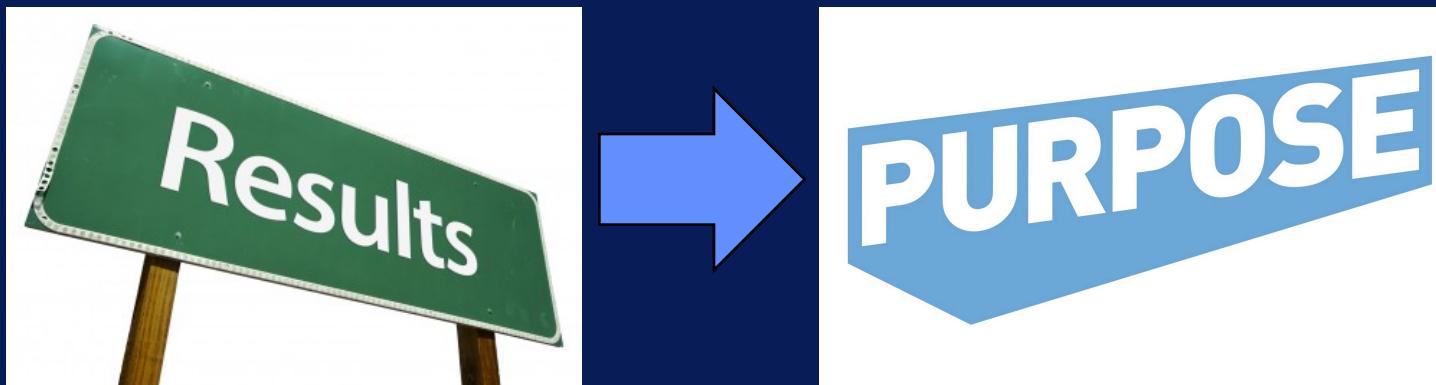
PREPARE

ANALYZE

REPORT

ACT

Step 5: Apply Results



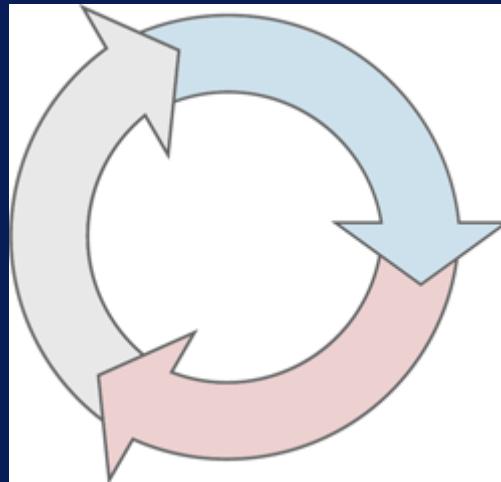
ACQUIRE

PREPARE

ANALYZE

REPORT

ACT



Iterative process

Step 1:

Acquiring Data

Big Data Engineering

Computational Big Data Science

ACQUIRE

PREPARE

ANALYZE

REPORT

ACT

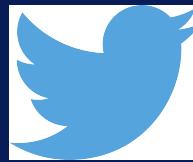
Where's the data?



Identify suitable data

Acquire all available data

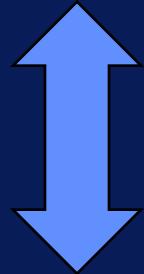
Data comes from many places...



...with many ways to access it



Traditional databases

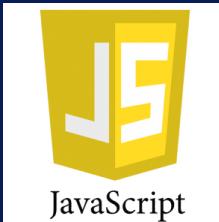


SQL and query browsers

Text files



Ruby



JavaScript



Perl

Scripting languages

SOAP



Remote data



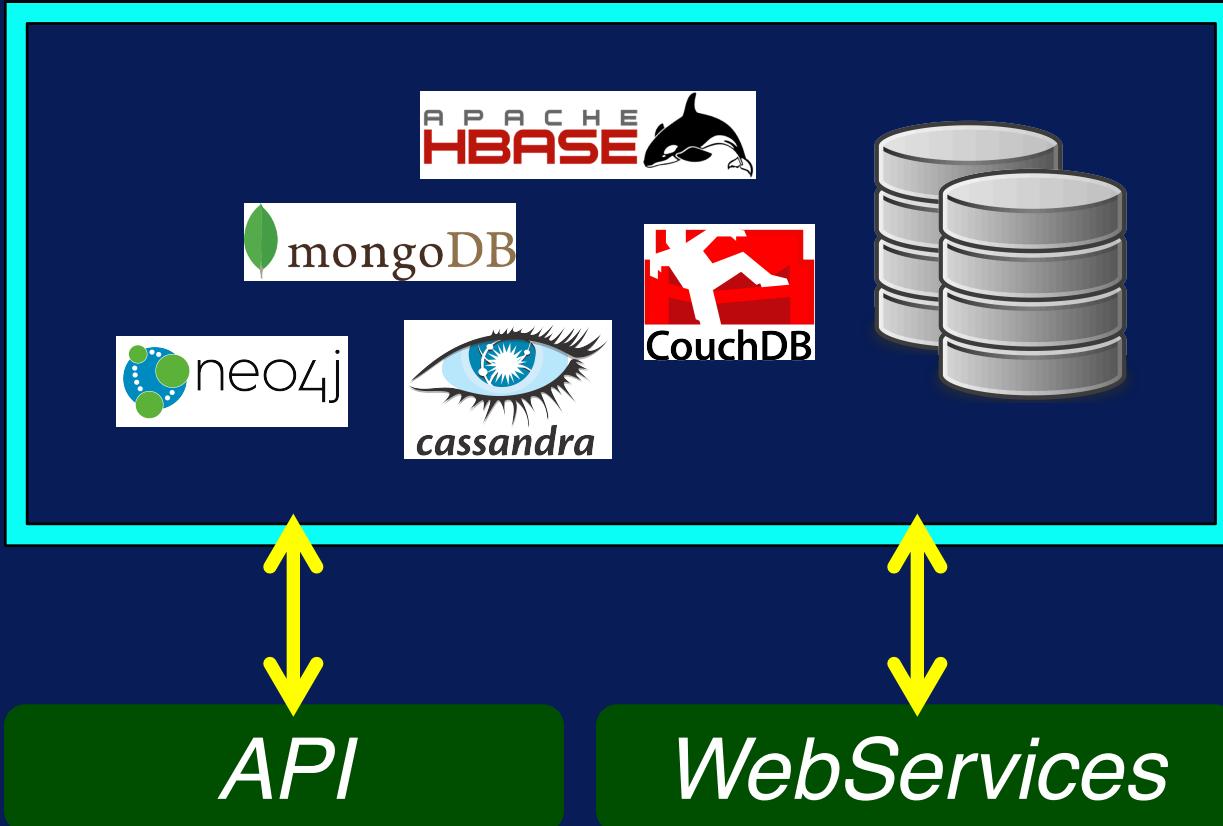
REST



WebSocket

Web Services

NoSQL storage



Acquiring Data From WIFIRE

Historical weather

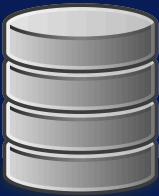


Current weather



Real-time tweets
near fires





Traditional databases

SQL and query browsers



Text files

Scripting languages



Remote data

Web Services



NoSQL storage

Web Services

Programming Interfaces

Step 2-A:

Exploring Data

ACQUIRE

PREPARE

ANALYZE

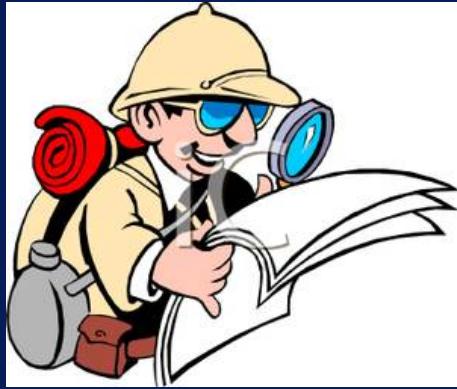
REPORT

ACT

Step 2-A: Explore

Step 2-B: Pre-process





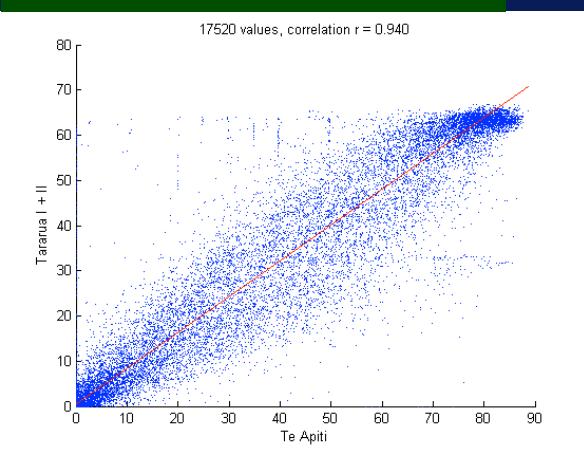
Why Explore?

Goal: Understand your data



Why Explore?

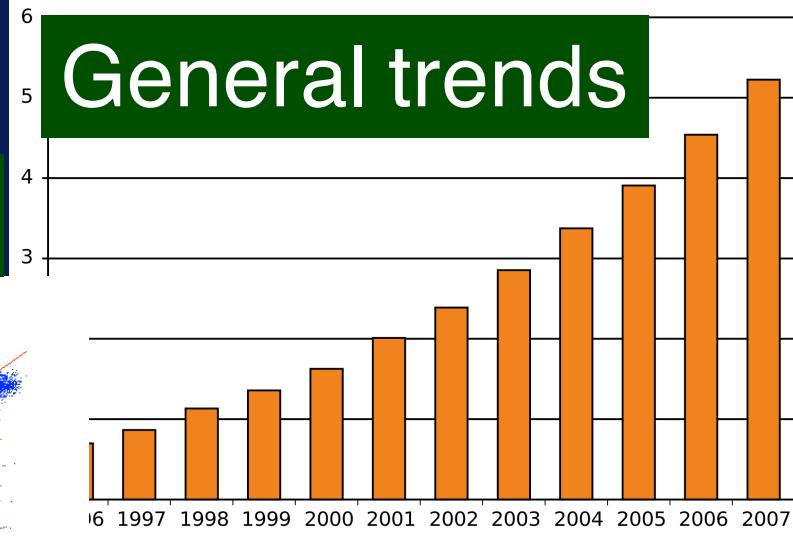
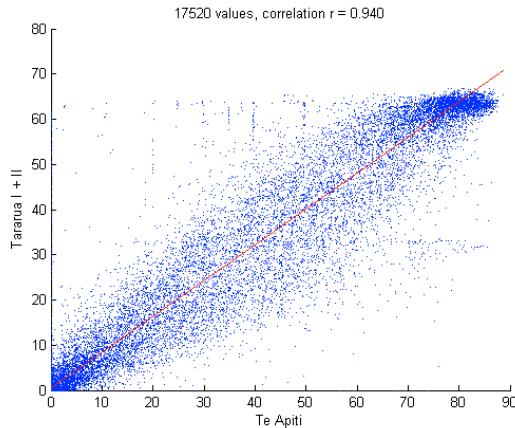
Correlations





Why Explore?

Correlations

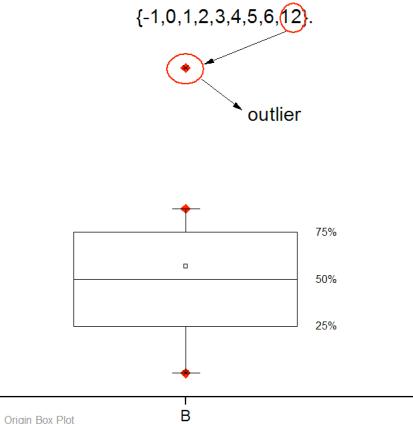
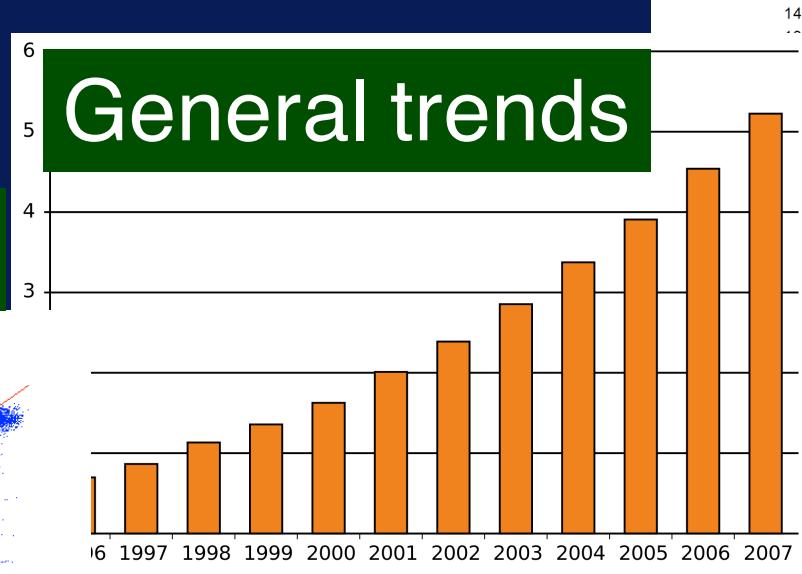
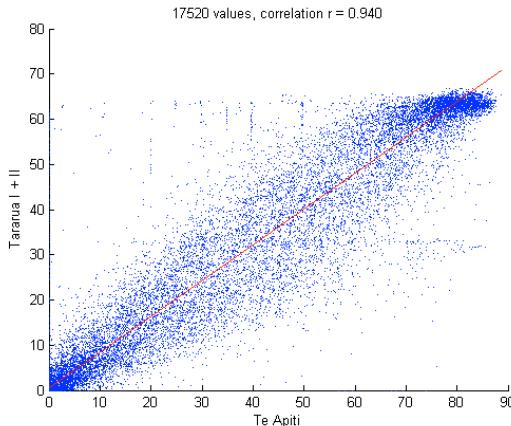




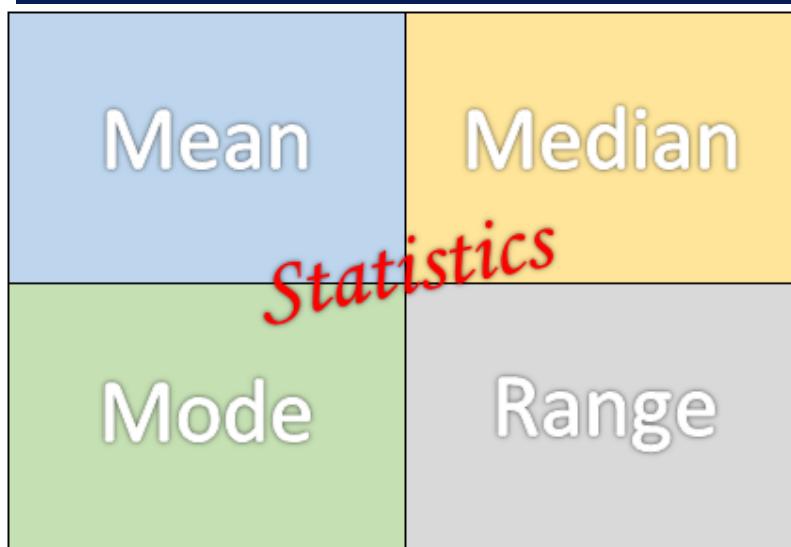
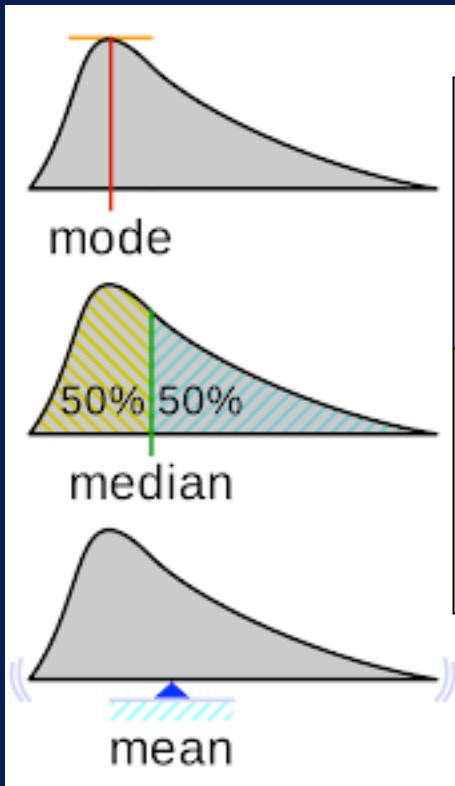
Why Explore?

Outliers

Correlations



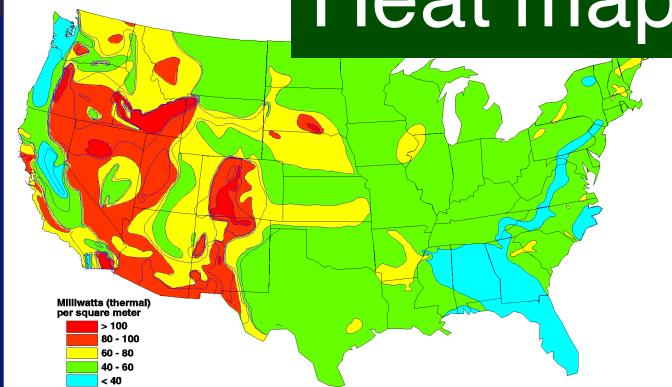
Describe Your Data



Visualize Your Data

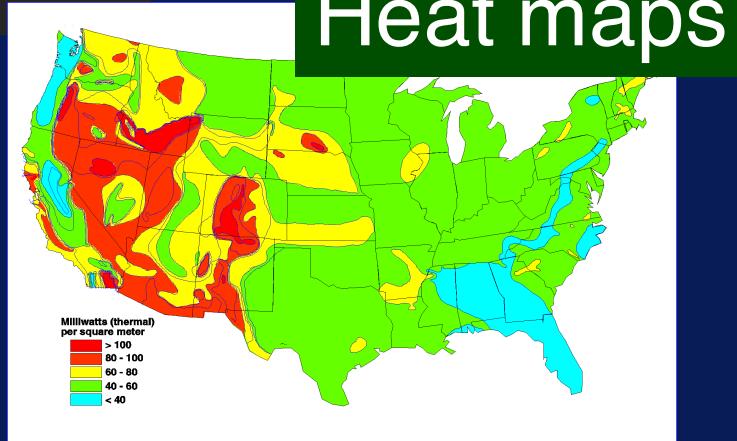
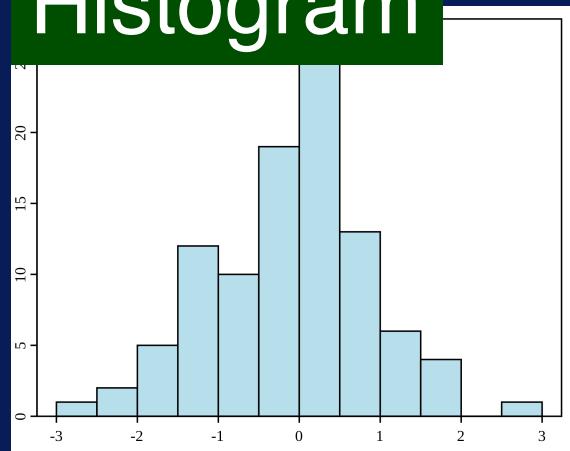
Visualize Your Data

Heat maps



Visualize Your Data

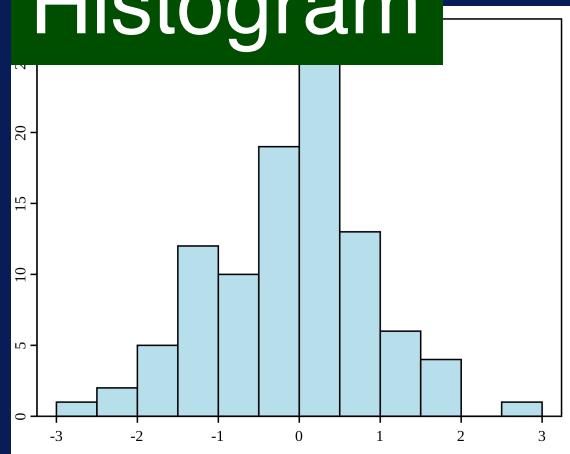
Histogram



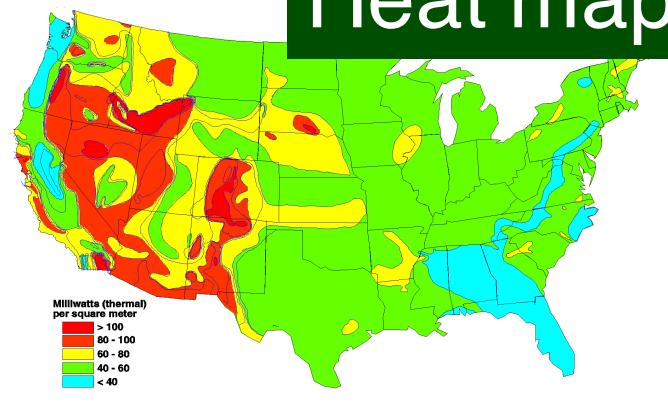
Heat maps

Visualize Your Data

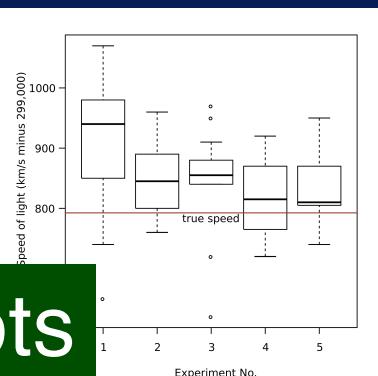
Histogram



Heat maps

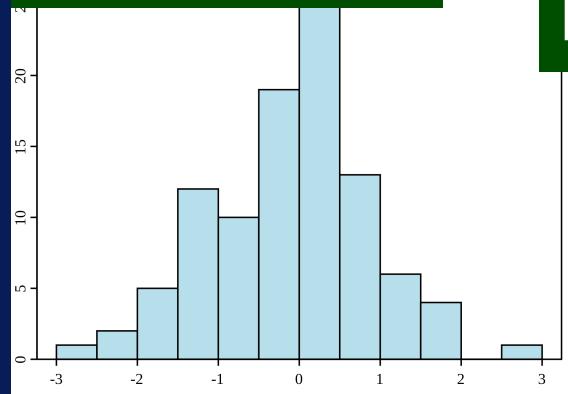


Boxplots

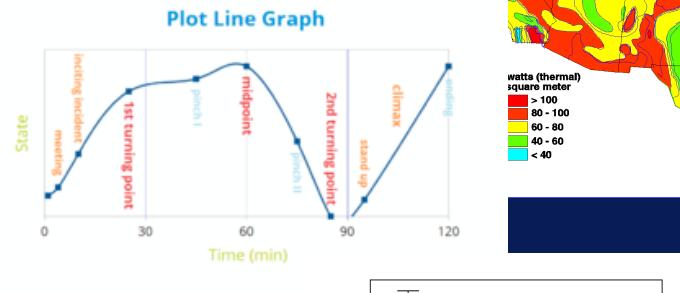


Visualize Your Data

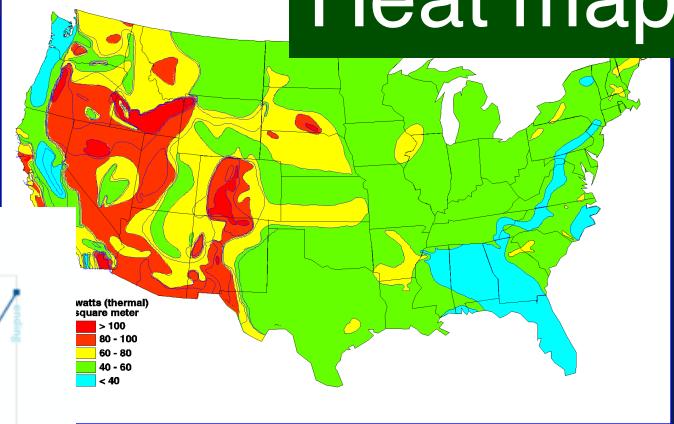
Histogram



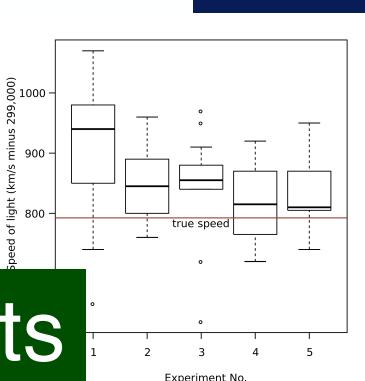
Line graphs



Heat maps

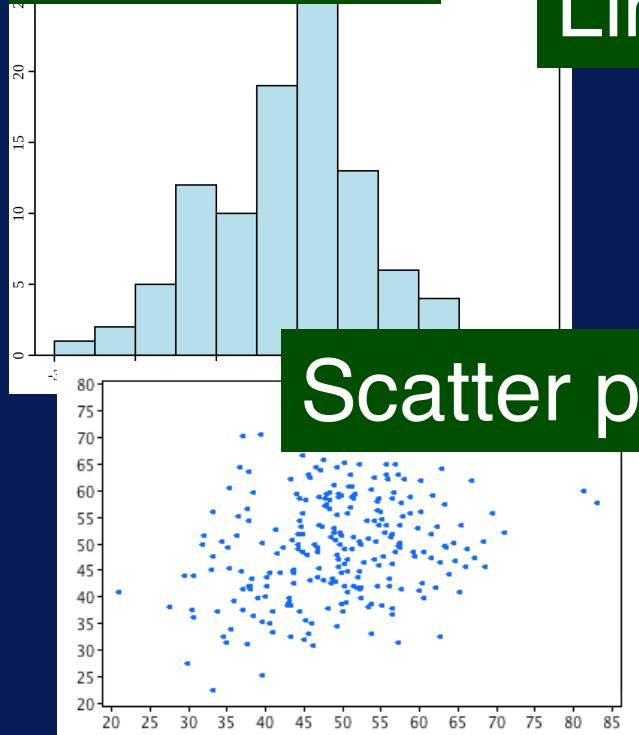


Boxplots



Visualize Your Data

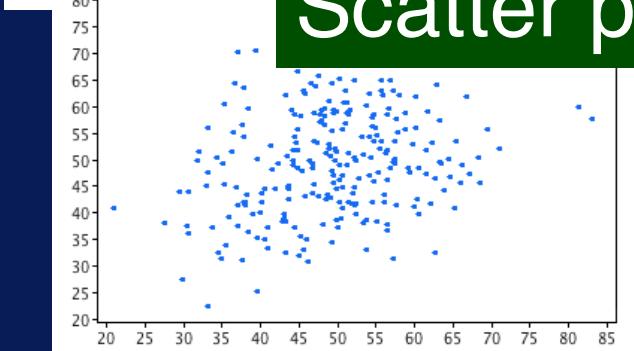
Histogram



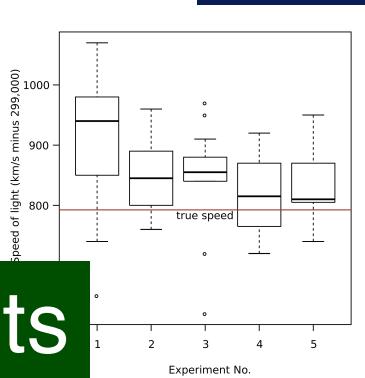
Line graphs



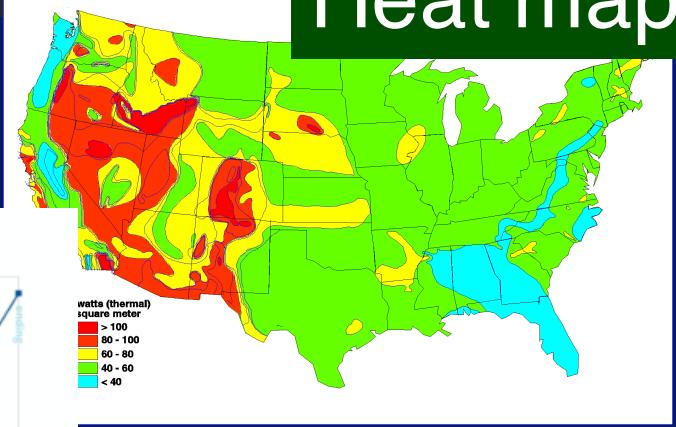
Scatter plots



Boxplots



Heat maps



Data
Exploration

Data
Understanding

Informed
Analysis



Step 2-B:

Pre-processing Data

ACQUIRE

PREPARE

ANALYZE

REPORT

ACT

Step 2-A: Explore

Step 2-B: Pre-process

Clean



Transform



Real-world data is messy!

Data Quality Issues

Inconsistent values

Duplicate records

Missing values

Invalid data

Outliers

Addressing Data Quality Issues

Addressing Data Quality Issues

Remove data with
missing values

Addressing Data Quality Issues

Remove data with
missing values

Merge duplicate records

Addressing Data Quality Issues

Remove data with
missing values

Generate best estimate
for invalid values

Merge duplicate records

Addressing Data Quality Issues

Remove data with
missing values

Generate best estimate
for invalid values

Merge duplicate records

Remove outliers

Addressing Data Quality Issues

Remove data with
missing values

Generate best estimate
for invalid values

Merge duplicate records

Remove outliers

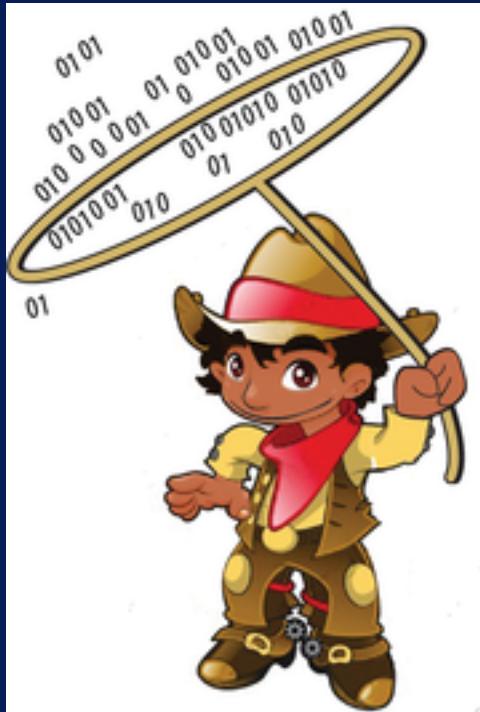
*Domain
Knowledge*

Getting Data in Shape

*Data
Munging*

*Data
Preprocessing*

*Data
Wrangling*



Data Munging

*Dimensionality
Reduction*

*Data
Manipulation*

Transformation

*Feature
Selection*

Scaling

Scaling

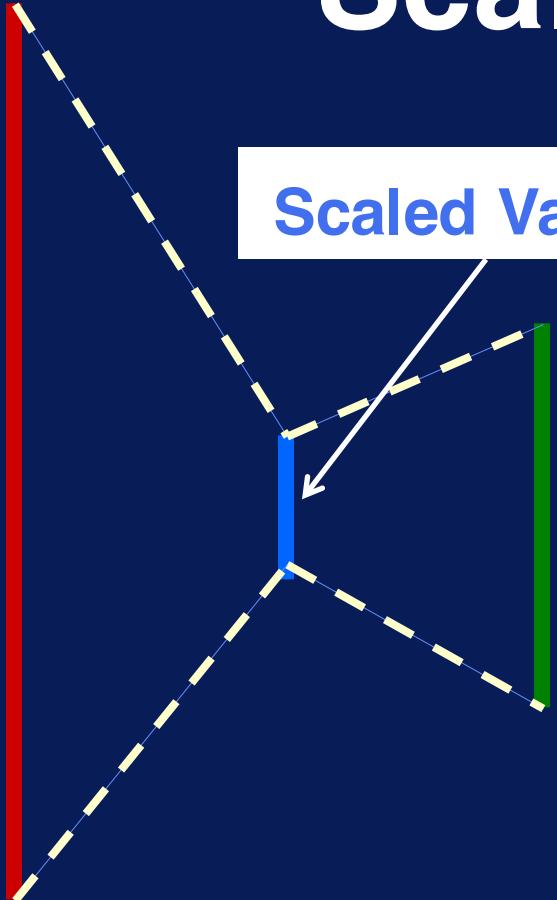


Weight

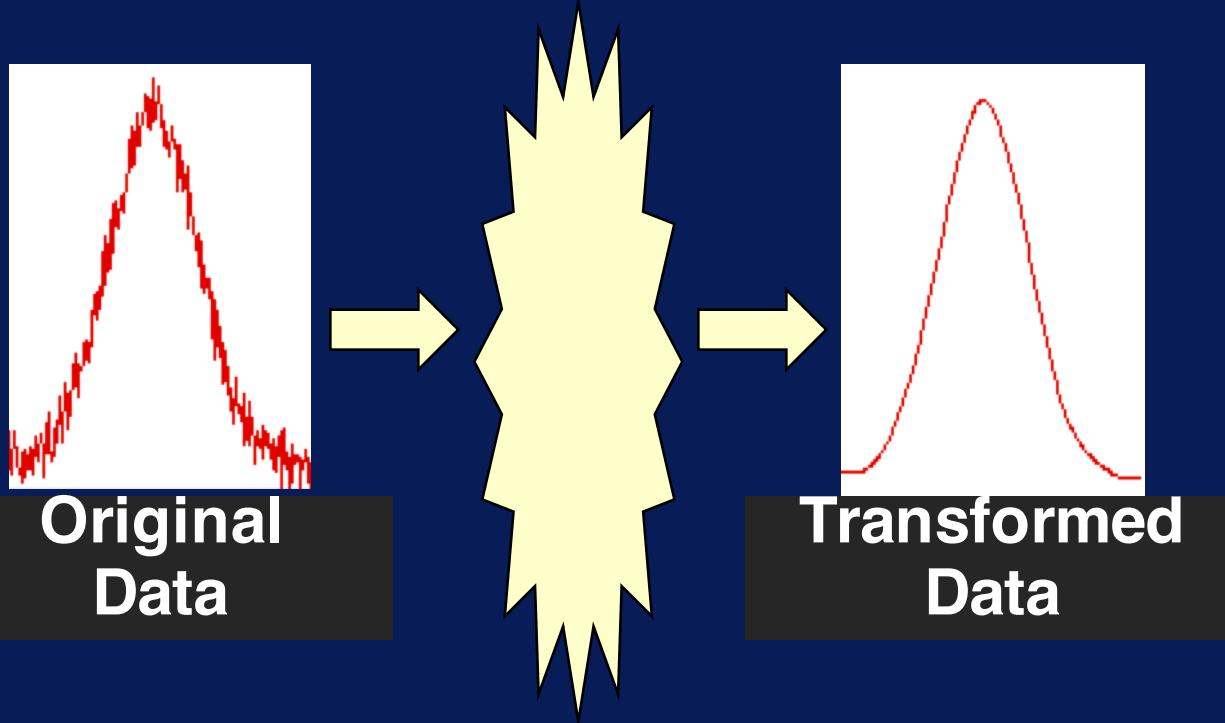
Scaled Values



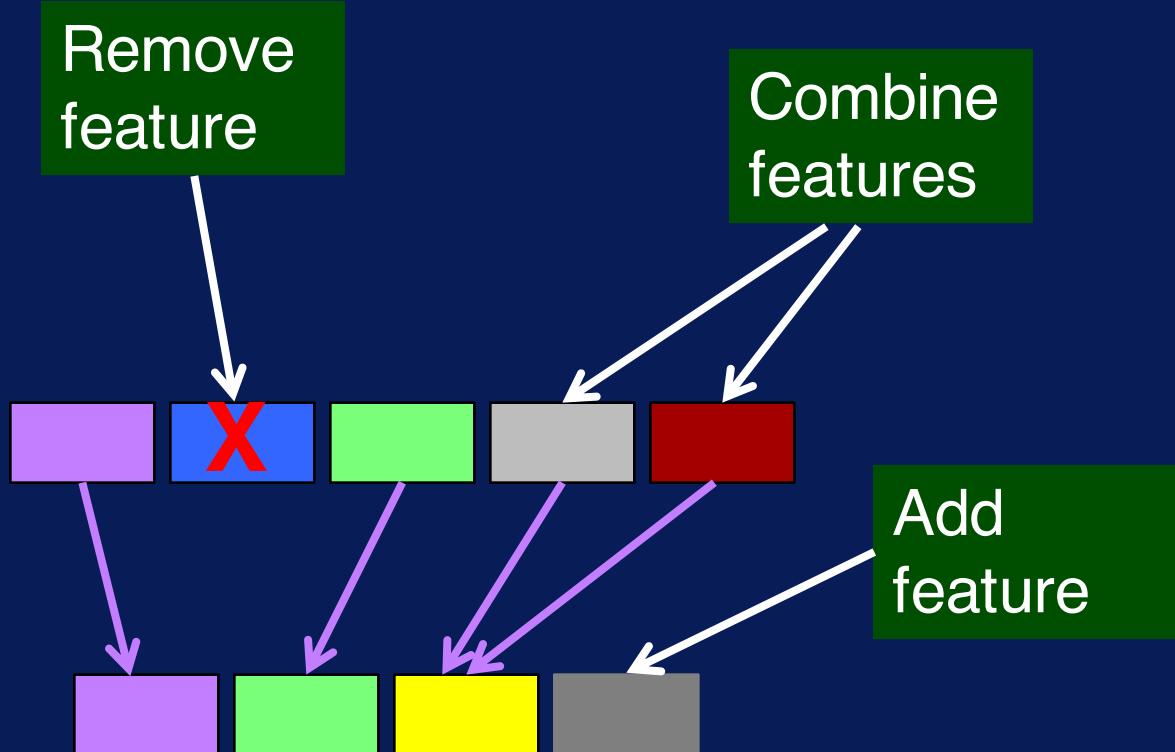
Height



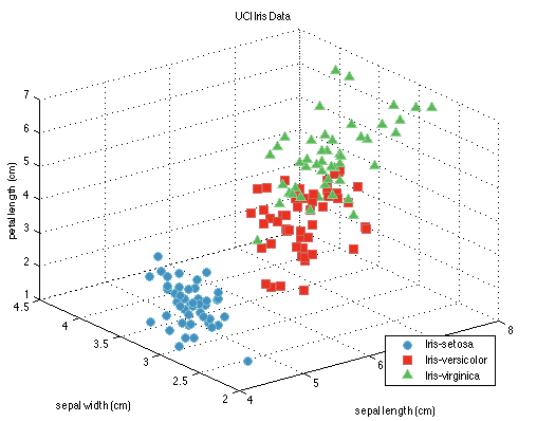
Transformation



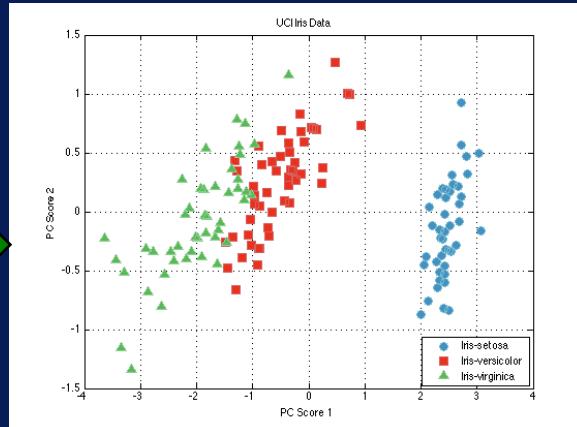
Feature Selection



Dimensionality Reduction

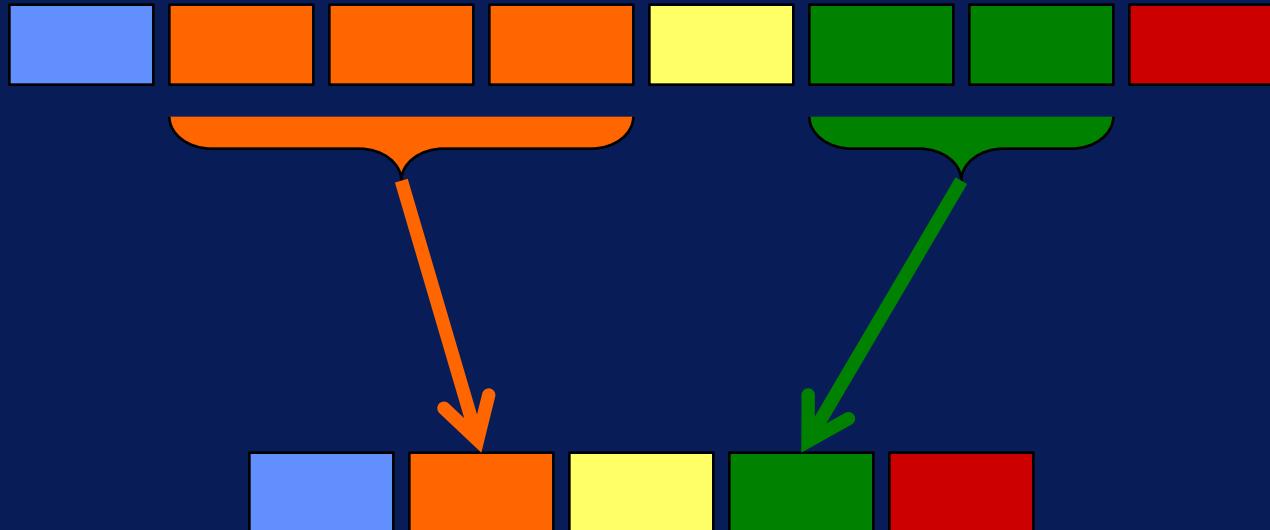


3D



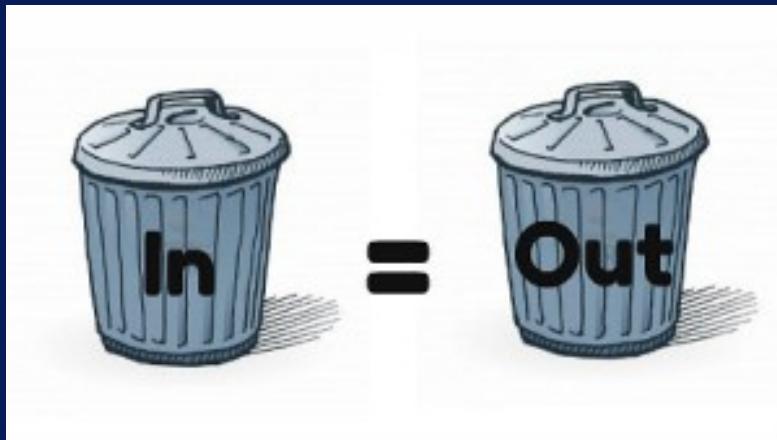
2D

Data Manipulation



Always Remember!

Garbage in = Garbage out



Data preparation is
very important for
meaningful analysis!

Step 3:

Analyze Data

Big Data Engineering

Computational Big Data Science

ACQUIRE

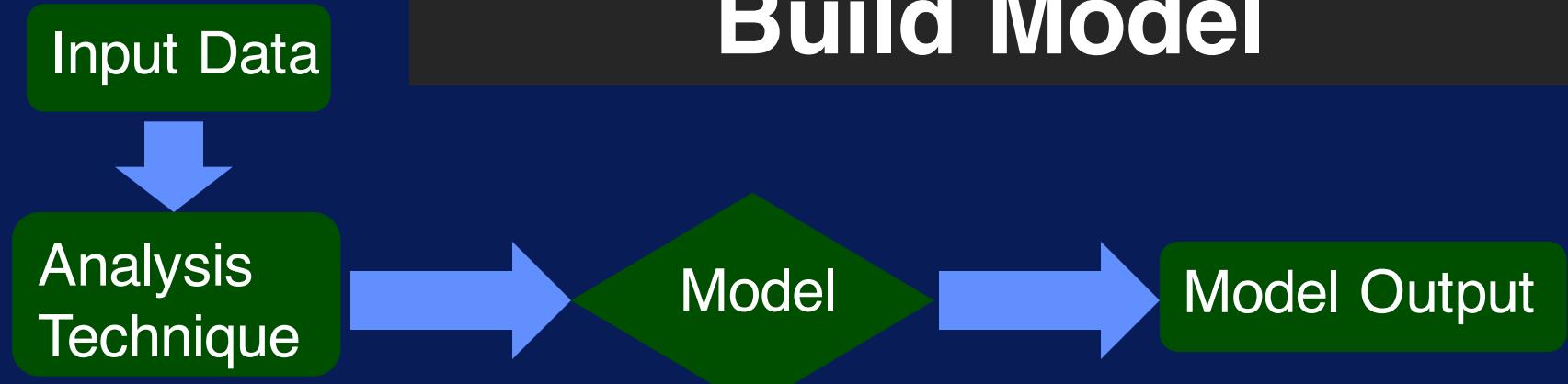
PREPARE

ANALYZE

REPORT

ACT

Build Model



Categories of Analysis Techniques

Classification

Regression

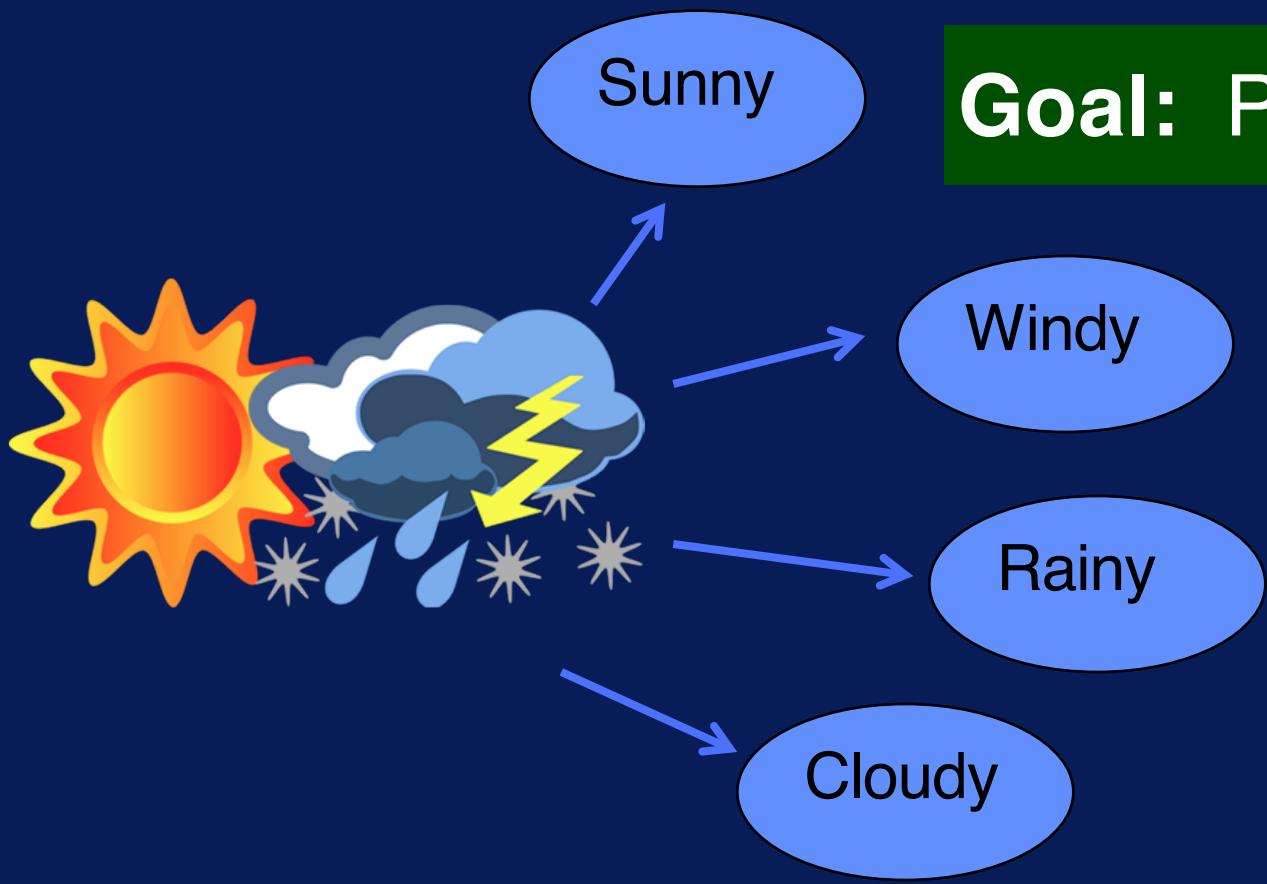
Clustering

Association
Analysis

Graph
Analytics

Classification

Goal: Predict category



Regression

Goal: Predict numeric value

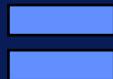


Clustering

Goal: Organize similar items into groups



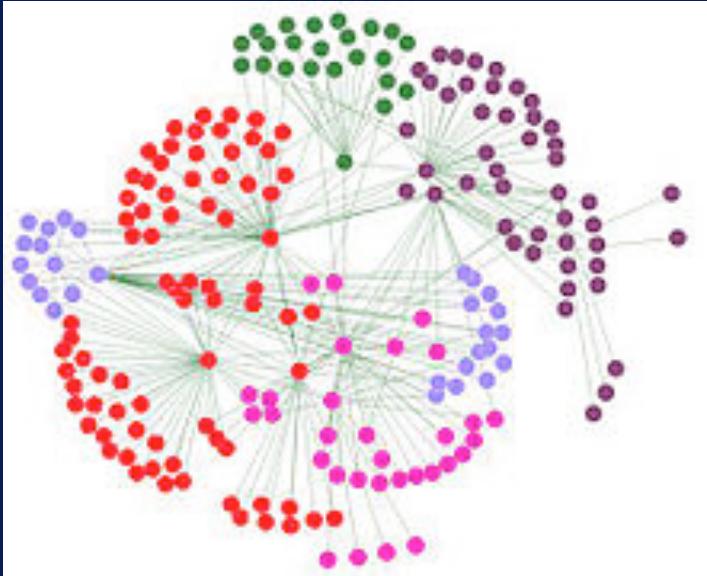
Association Analysis



Goal: Find rules to capture associations between items

Graph Analytics

Goal: Use graph structures to find connections between entities



Modeling

Select technique



Build model



Validate model

Evaluation of Results

Classification & Regression

Predicted
Value



Correct
Value

Clustering



Association Analysis & Graph Analytics



Investigate



Validate

Determine Next Steps



Repeat analysis?

Take deeper dive?

Act on results?

Select technique

Classification
Regression
Clustering
Association
Analysis
Graph Analytics

Build model



Evaluate



Step 4:

Reporting Insights

Big Data Engineering

Computational Big Data Science

ACQUIRE

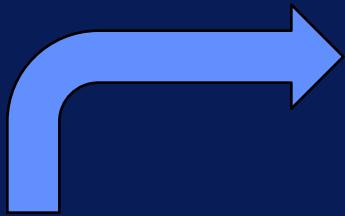
PREPARE

ANALYZE

REPORT

ACT

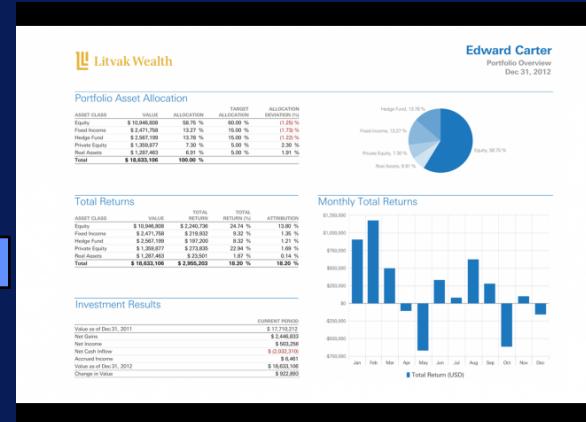
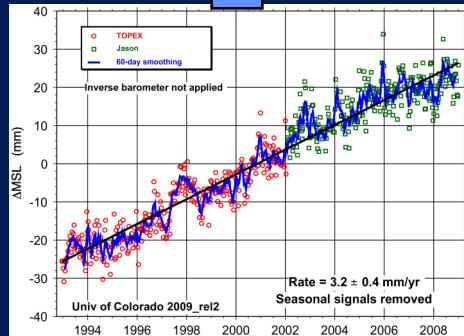
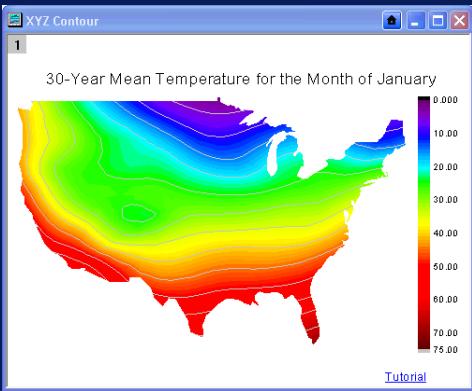
What to Present



What to Present



How to Present



Visualization Tools



Timeline JS

Beautifully crafted timelines that are easy and intuitive to use.

Present



with



using

tabular

Step 5:

Insights into Action

Big Data Engineering

Computational Big Data Science

ACQUIRE

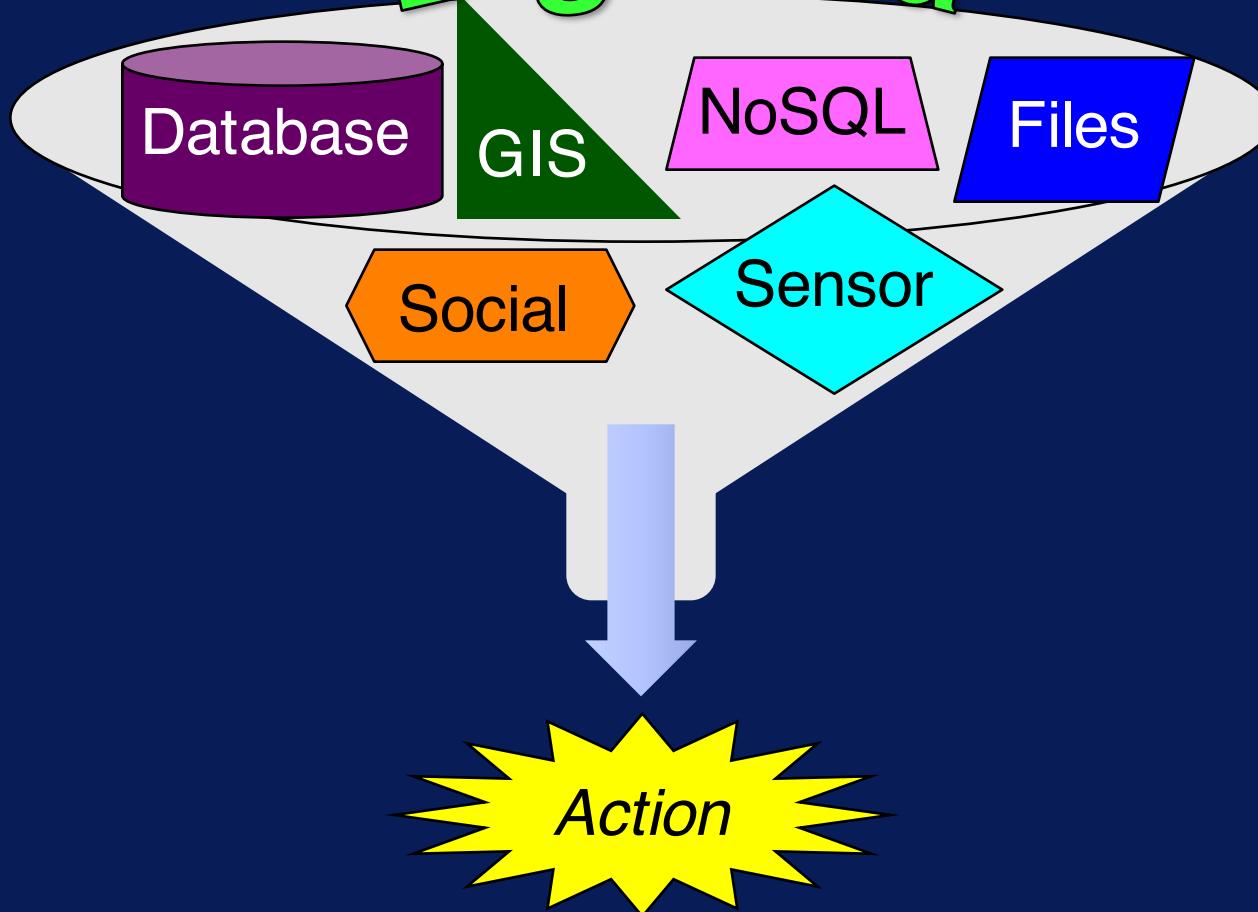
PREPARE

ANALYZE

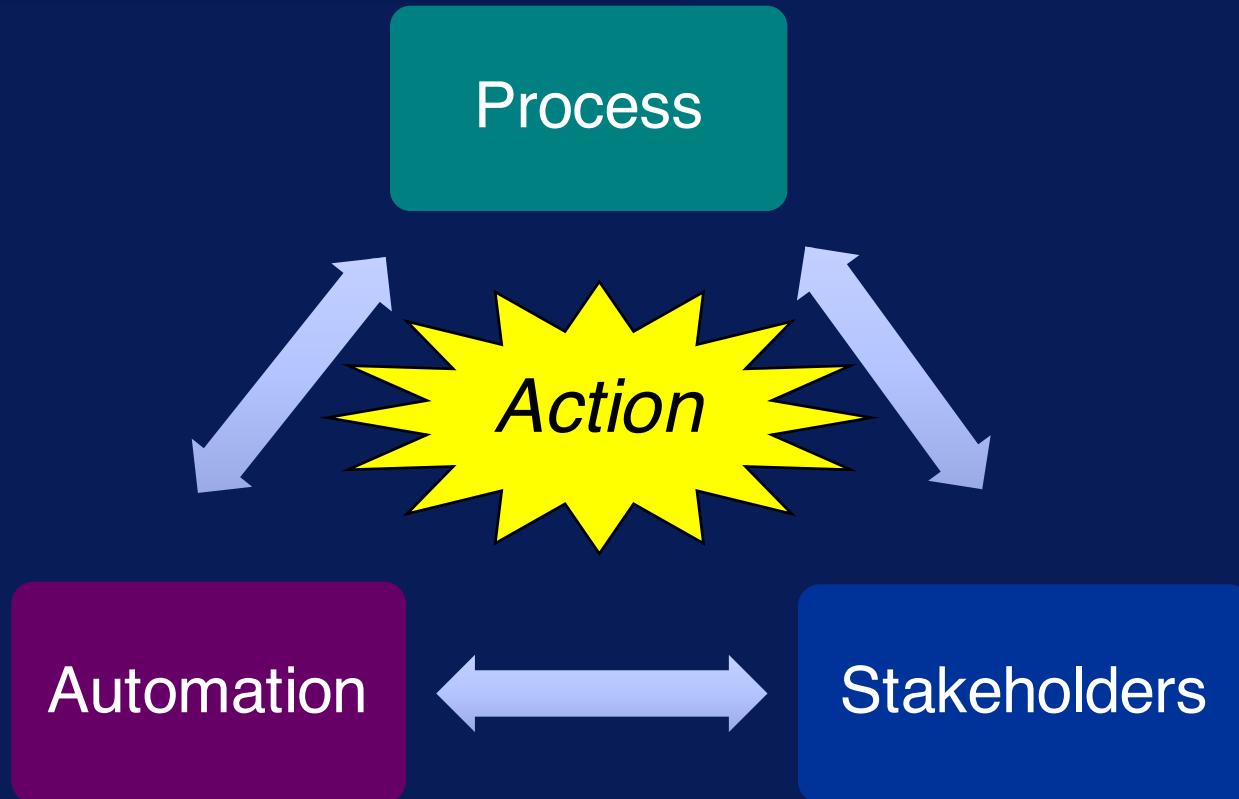
REPORT

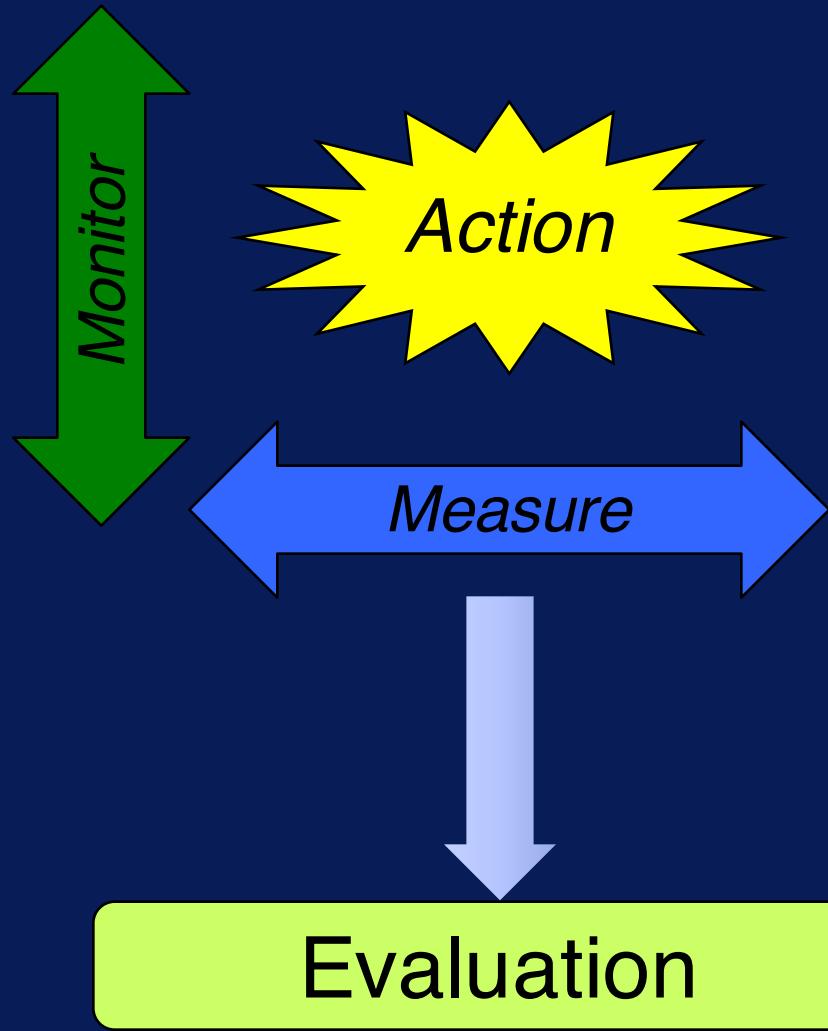
ACT

Big Data



Implementation



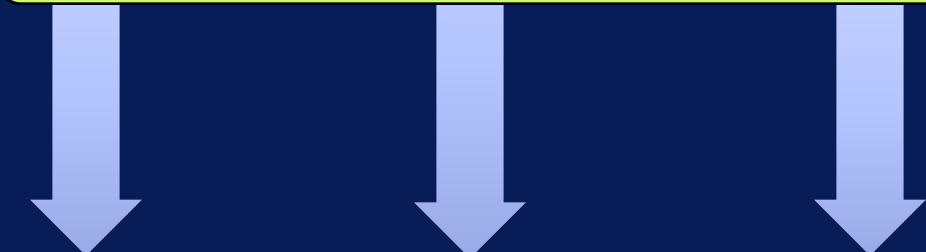


Assess Impact



Determine Next
Steps

Evaluation



Favorable
Results?

Revisit?

Further
Opportunities?

