# What's the Big Deal About Big Data?

Natasha Balac, Ph.D.

Jun, 2016

# Background

- Over 25 Years of Experience in Data Mining

- Ph.D. in Machine Learning – with emphasis on Big Data and Mobile Robots

- Director of Predictive Analytics center of Excellence at the Supercomputer Center at UCSD

- Lecturer
  - UCSD MAS in Data Science and Engineering
  - UCSD Extension Data Mining Certificate
  - Coursera Big Data Specialization

# University of California, San Diego UCSD

Student-centered, research-focused, service-oriented public institution

Recognized as one of the top 15 research universities worldwide

Culture of collaboration sparks discoveries that advance society and drive economic impact

UC San Diego's rich academic portfolio includes six undergraduate colleges, five academic divisions and five graduate and professional schools

# CalIT2 – Qualcomm Institute

> **Calit2 represents a new mechanism to address large-scale societal issues by bringing together multidisciplinary teams of the best minds.** *Larry Smarr, Director, Calit2*

Spitzer Space Telescope (Infrared)

Hubble Space Telescope (Optical)

Calit2 is taking ideas beyond theory into practice, accelerating innovation and shortening the time to product development and job creation. Where the university traditionally has focused on education and research, Calit2 extends that focus to include development and deployment of prototype infrastructure for testing new solutions in a real-world context.

# Data Insight Discovery

- Founded in January 2014 – San Diego, CA
  - Women Owned

- Service Offerings Include
  - Predictive Analytics Services
  - Business Intelligence and Analytics
  - CRM system development
  - Condition Based Maintenance
  - Systems and IoT Integration Services
  - ETL, RDBMS, No-SQL data store development
  - Big Data Technologies

- **Key Strengths include**
  - Numerous successfully deployed Predictive Analytics Projects
  - Over 25 years of experience and expertise

data.
insight.
discovery.

# Data Insight Discovery Areas of Expertise

**Marketing**
- Customer Segmentation
- Campaign Analysis
- Market Mix Modelling
- Dashboards and Reporting
- Promotion Effectiveness
- Price Optimization

**Sales**
- Sales forecasts
- Incentive Optimization
- Propensity Modelling
- Consumer Insights
- Branding Solutions

**Digital**
- Social and Web Analytics
- Sentiment Analysis
- Cross media interaction
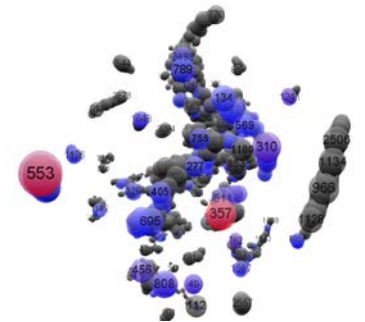- Digital Media ROI
- Anomaly Detection
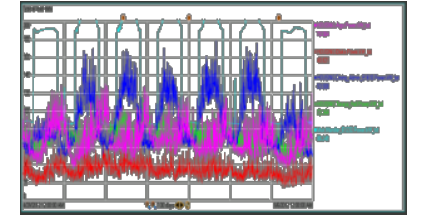
**Condition Based Maintenance**
- Sensor Analytics
- Predictive Maintenance
- Intelligent workflow Solution
- Risk modeling and estimation

# PACE – Predictive Analytics Center of Excellence
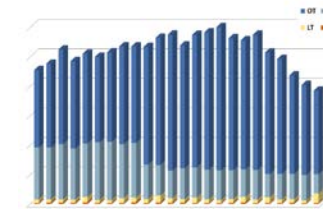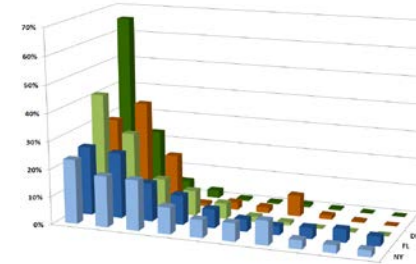# Closing the gap between Government, Industry and Academia
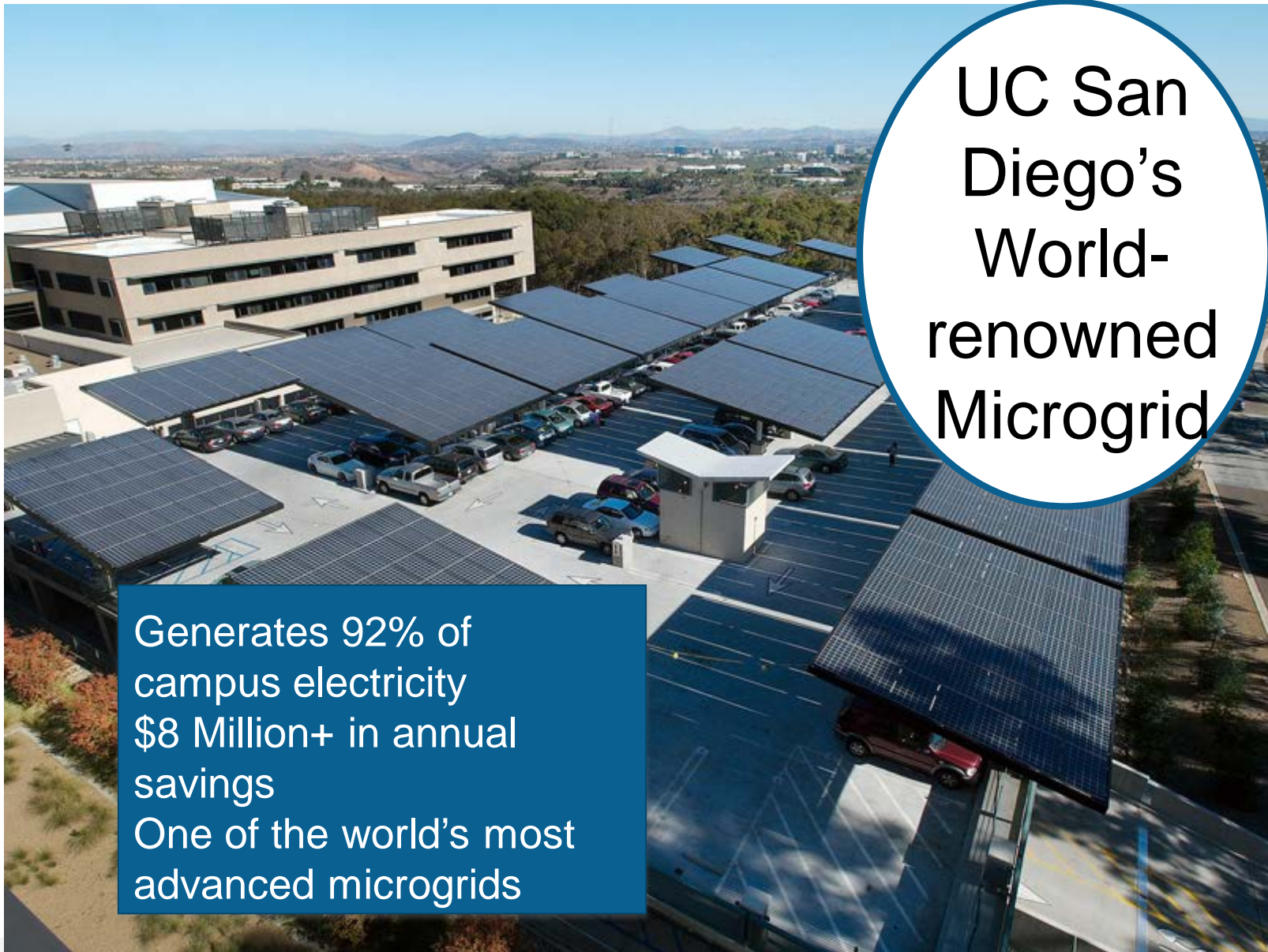


- Fraud Detection
- Modeling behaviors
- Biomedical Informatics
- Smart Grid Analytics
- Solar powered system modeling
- Microgrid anomaly detection
- Battery Storage Analytics
- Sport Analytics
- Genomics
- Population Health
- Nano-engineering

# CMS Fraud, Waste and Abuse Detection and Prediction

- Descriptive Statistics
  - Claims summary information
  - History and trends
  - Distributions across periods, transactions, etc.

- Exploratory Analysis
  - Profiles of provider transactions
  - Provider similarity according to profiles
  - Visual summaries of large amounts of data
  - Eligibility data link to provider billing

- Predictive analytics
  - Adjustments
  - Equipment, Service Codes
  - Long term vs. short term hospital stay
  - Provider profiles

**UC San Diego's World-renowned Microgrid**

Generates 92% of campus electricity
$8 Million+ in annual savings
One of the world's most advanced microgrids

# UCSD Smart Grid

- **UCSD Smart Grid sensor network data set**
  - **45MW peak micro grid; daily population of over 54,000 people**

- Smart Grid data – over 100,000 measurements/sec
  - **Sensor and environmental/weather data**
    - Large amount of complex data streaming from sensor networks
  - **Predictive Analytics throughout the Microgrid**

Collaboration between public, private, and academic organizations working to develop and implement initiatives that will improve the San Diego region's energy independence, empower consumers to embrace clean technologies, reduce greenhouse gas emissions, and drive economic growth
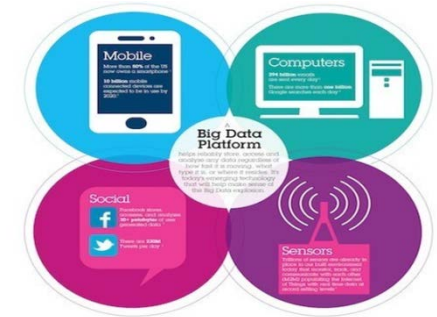
**White House Big Data Event:** "Data to Knowledge to Action" – Launch Partners Award

# What is "Big Data"?

- Big data

  "Extremely large data sets that may be analyzed computationally to reveal patterns, trends, and associations, especially relating to human behavior and interactions"
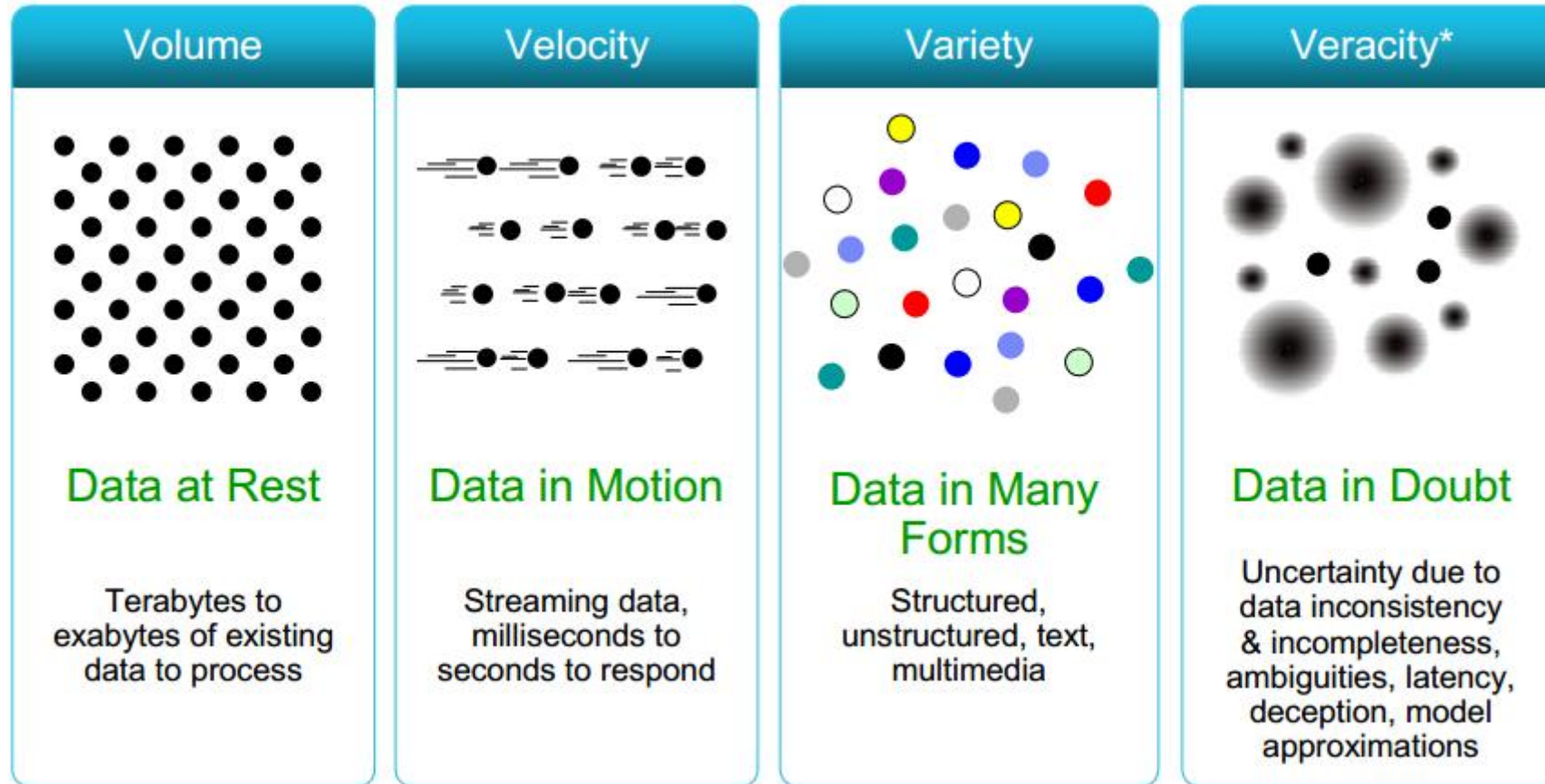
- Big data is in the eye of the beholder

  Data that exceeds the ability to handle with current resources

# 4 V's of Big Data



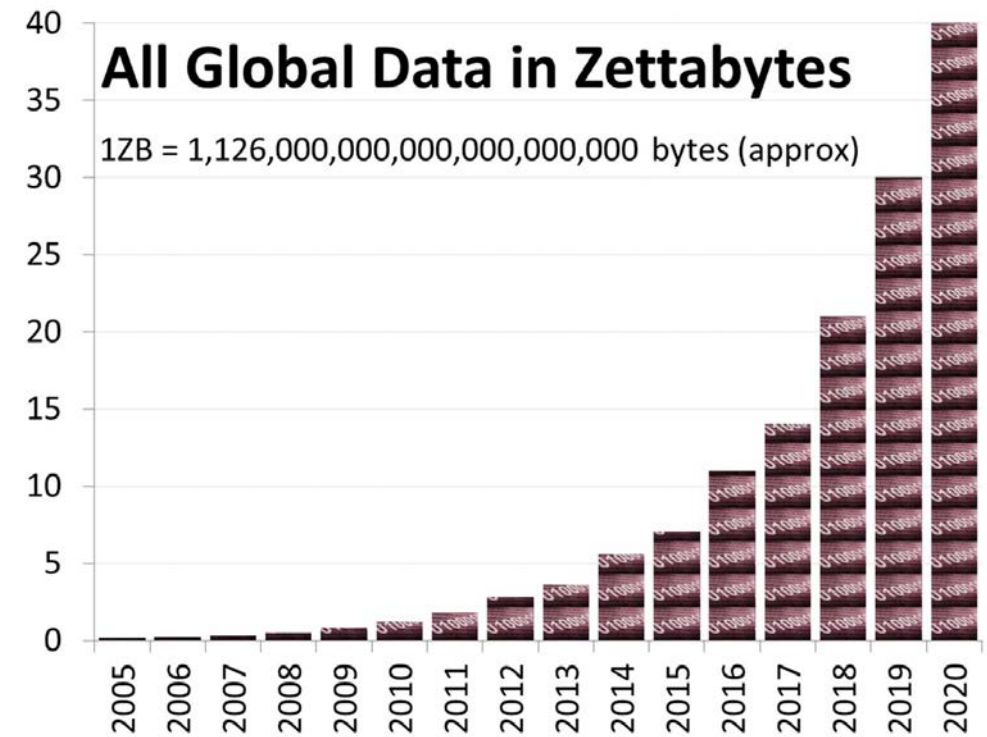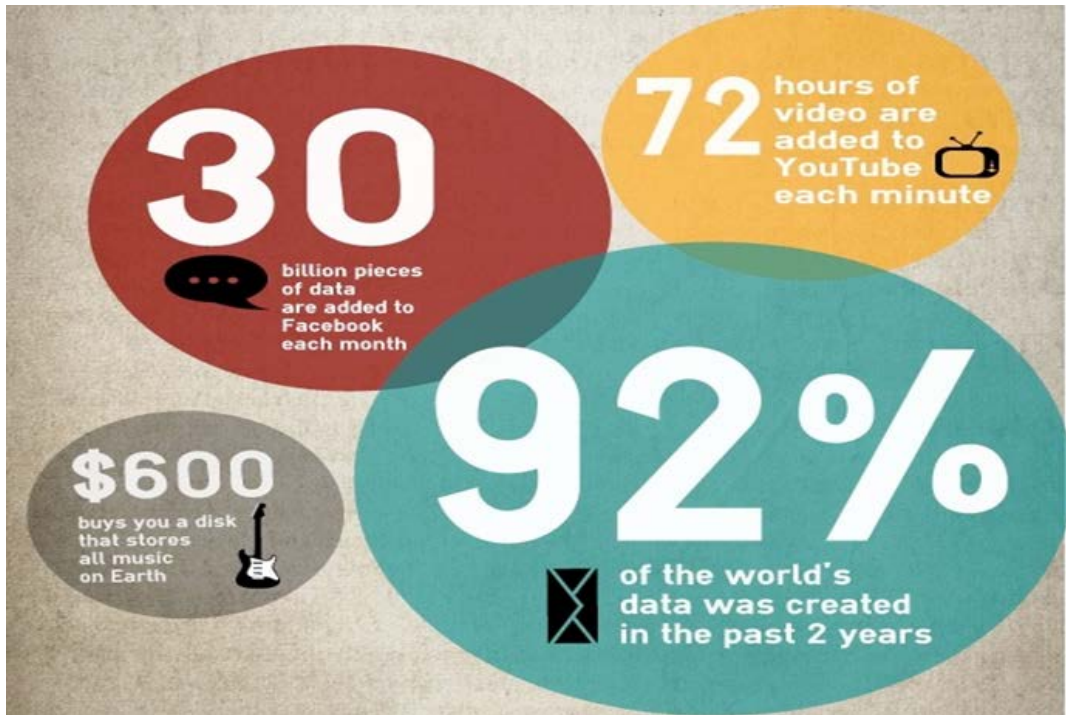| Volume | Velocity | Variety | Veracity* |
|---|---|---|---|
| **Data at Rest** | **Data in Motion** | **Data in Many Forms** | **Data in Doubt** |
| Terabytes to exabytes of existing data to process | Streaming data, milliseconds to seconds to respond | Structured, unstructured, text, multimedia | Uncertainty due to data inconsistency & incompleteness, ambiguities, latency, deception, model approximations |

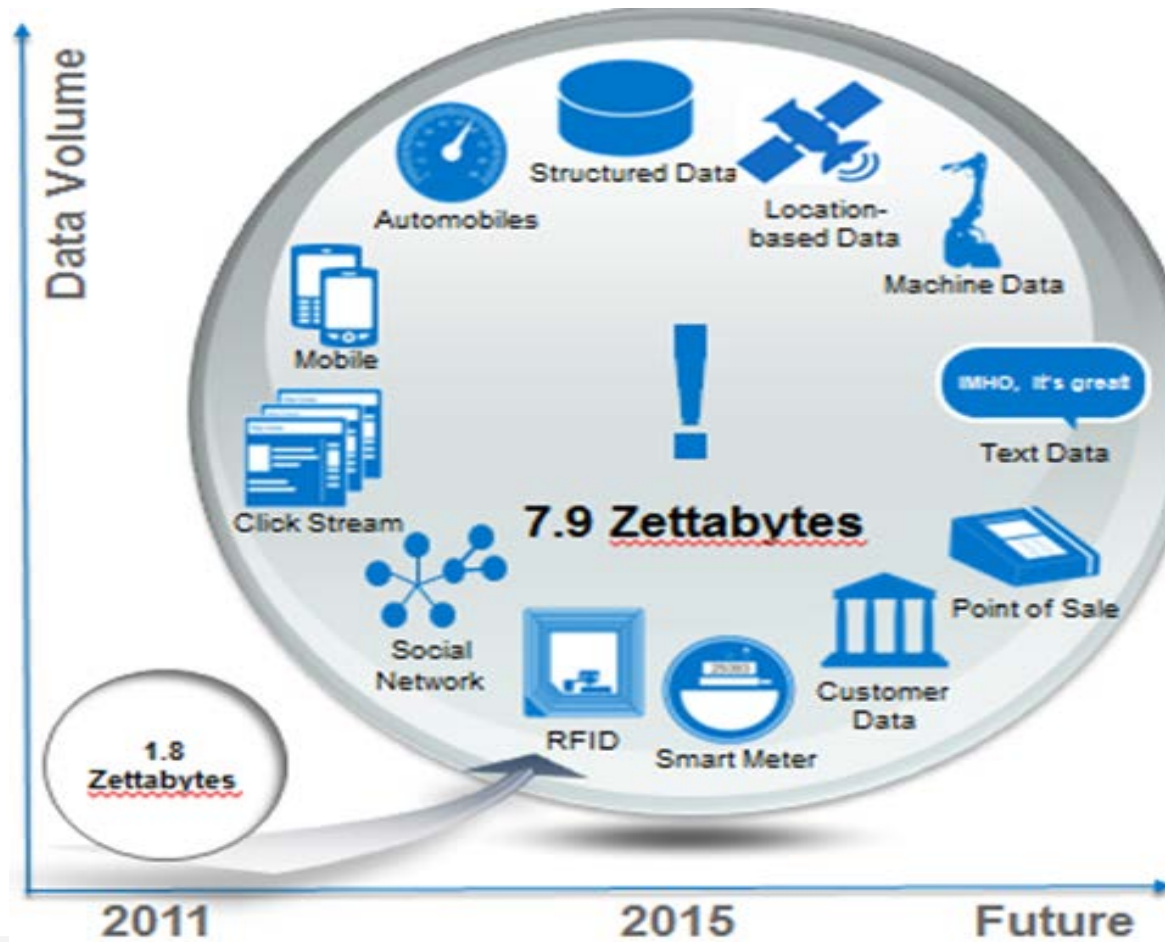IBM, 2012

# Big Data Definition

- *Wikipedia*: an all-encompassing term for any collection of data sets so large and complex that it becomes difficult to process using on-hand data management tools or traditional data processing applications

- *Oxford English Dictionary:* data of a very large size, typically to the extent that its manipulation and management present significant logistical challenges
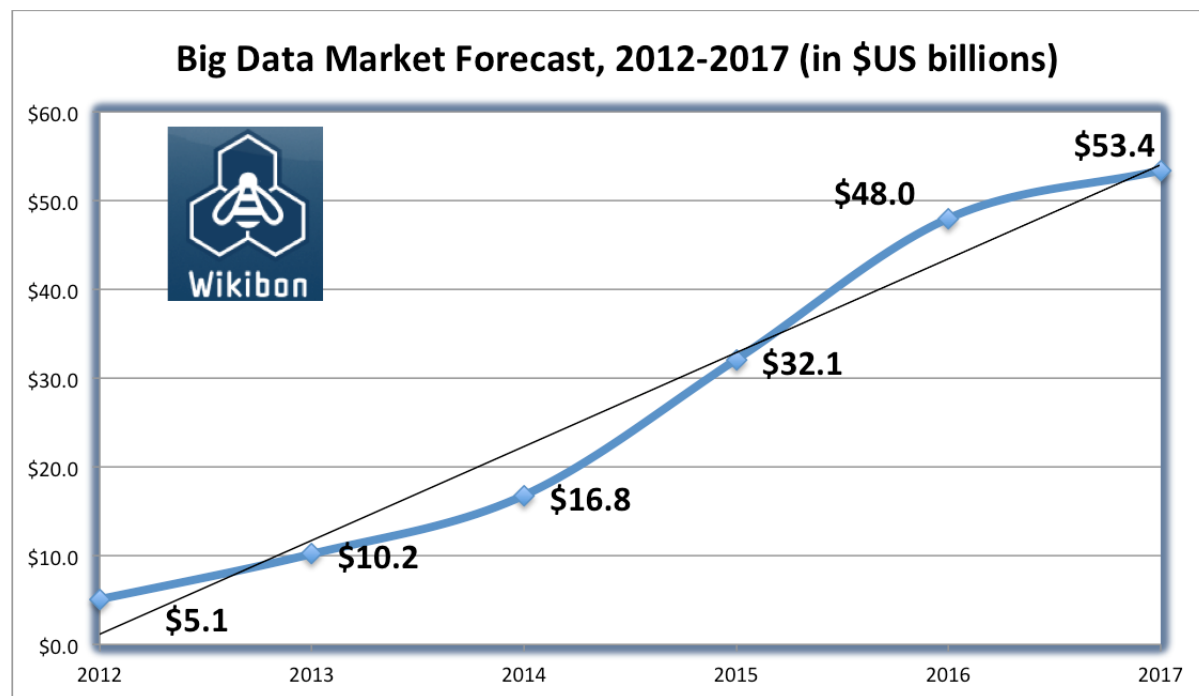
# How Big is Big Data?

# Big Data Growth



1 Terabyte = 1024 Gigabytes
1 Petabyte = 1024 Terabytes
1 Exabyte = 1024 Petabytes
1 Zettabyte = 1024 Exabytes

# Tip of the Iceberg

**Big Data Market Forecast, 2012-2017 (in $US billions)**



A line chart showing: $5.1 (2012), $10.2 (2013), $16.8 (2014), $32.1 (2015), $48.0 (2016), $53.4 (2017). Source: Wikibon.
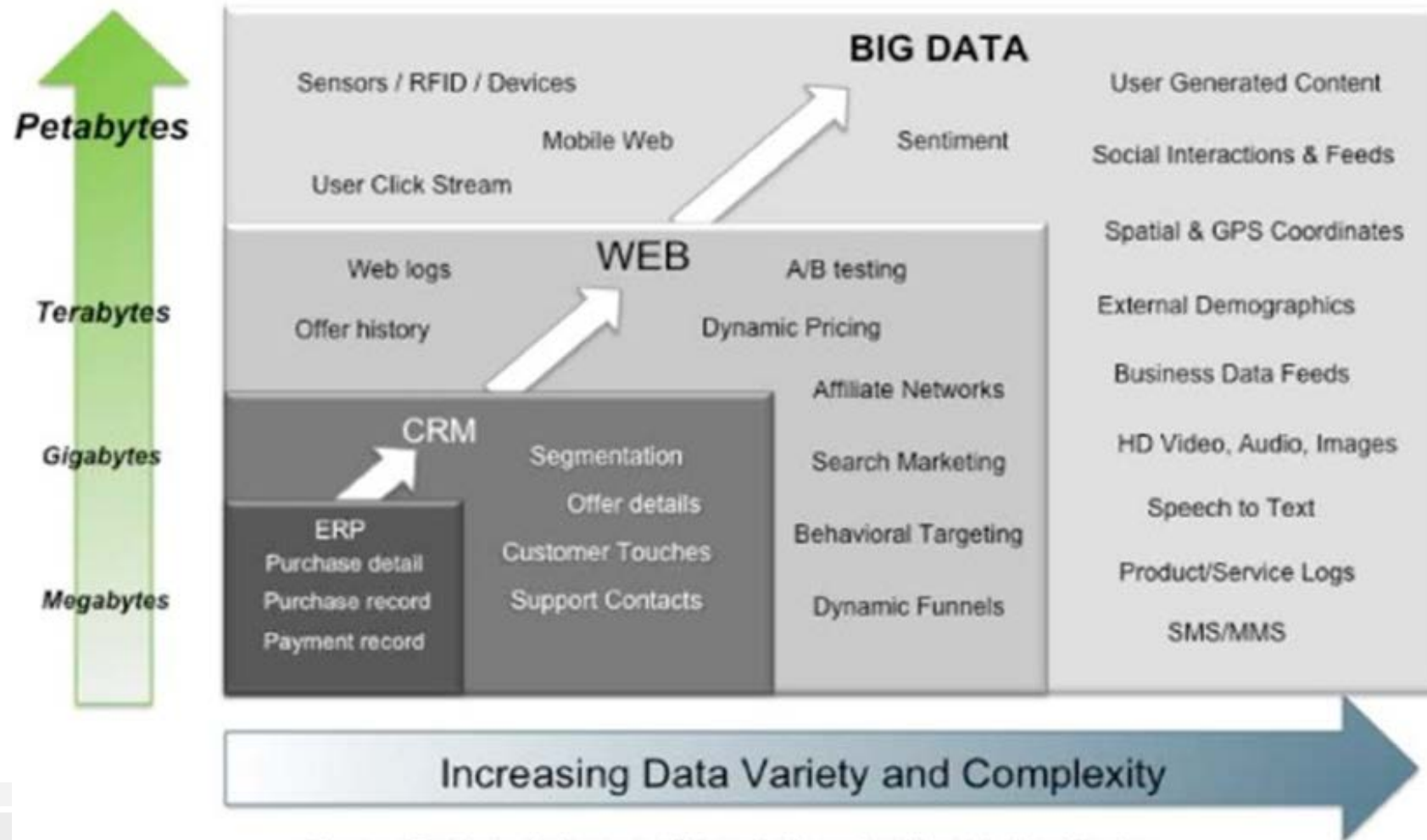
# What to do with big data?

ERIC SALL's list

- Big Data Exploration

  – To get an overall understanding of what is there

- 360 degree view of the customer

  – Combine both internally available and external information to gain a deeper understanding of the customer

- Monitoring Cyber-security and fraud in real time

- Operational Analysis

  – Leveraging machine generated data to improve business effectiveness

- Data Warehouse Augmentation

  – Enhancing warehouse solution with new information models and architecture

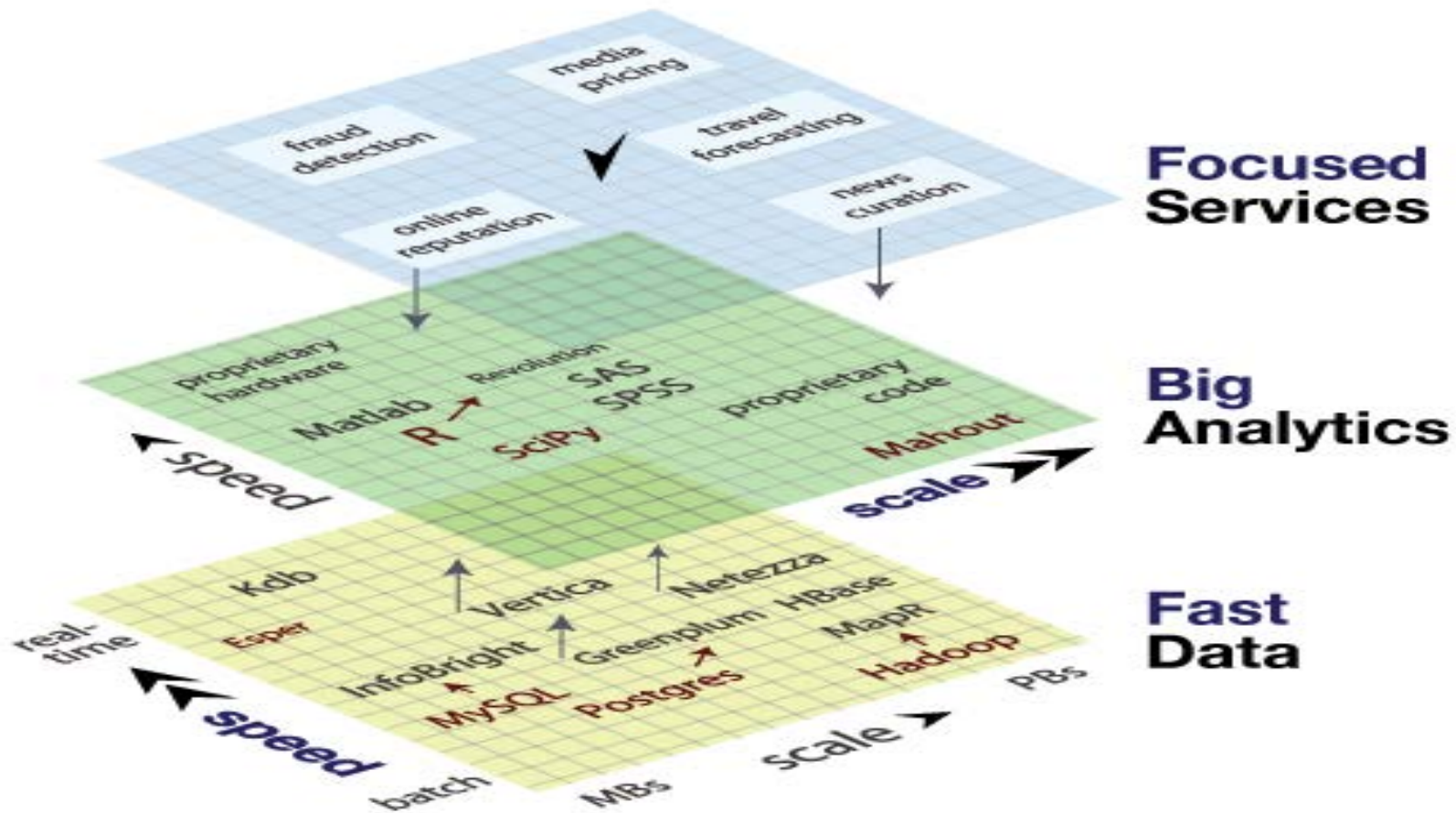# Transactions, Interactions & Observations



Source: Contents of above graphic created in partnership with Teradata, Inc.

# The Emerging Big Data Stack

# Big Data Landscape 2016 (Version 2.0)

Last Updated 2/12/2016 — © Matt Turck (@mattturck), Jim Hao (@jimrhao), & FirstMark Capital (@firstmarkcap) — FIRSTMARK

# Gartner Hype Cycle for Emerging Technologies 7/14

# Gartner Hype Cycle for Emerging Technologies 7/15

# What is IoT?

# IoT and Big Data



INSTRUMENTED INDUSTRIAL MACHINE

Intelligence flows back into machines

Extraction and storage of proprietary machine data stream

PHYSICAL AND HUMAN NETWORKS

data
Storage database

INDUSTRIAL DATA SYSTEMS

SECURE, CLOUD-BASED NETWORK

Data sharing with the right people and machines

Machine-based algorithms and data analysis

REMOTE AND CENTRALIZED DATA VISUALIZATION

BIG DATA ANALYTICS

# Growing Internet of Things Data

Today data scientist uses **Yottabytes** to describe how much government data the NSA or FBI have on people altogether.

In the near future, **Brontobyte** will be the measurement to describe the type of sensor data that will be generated from the IoT (Internet of Things)

$10^{27}$

**Brontobyte**
This will be our digital universe tomorrow...

**Yottabyte**
This is our digital universe today
= 250 trillion of DVDs

$10^{24}$

$10^{21}$

Zettabyte
1.3 ZB of network traffic by 2016

$10^{18}$

Exabyte
1 EB of data is created on the internet each day = 250 million DVDs

$10^{15}$

Petabyte
The CERN Large Hadron Collider generates 1PB per second

$10^{12}$

Terabyte
500TB of new data per day are ingested in Facebook databases

$10^{9}$

Gigabyte

$10^{6}$

Megabyte

# INTERNET OF THINGS LANDSCAPE

## Platforms & Enablement (Horizontals)

**Connectivity:** IFTTT · Symple · ioBridge · ARRAYENT · Twitter · electric imp · haystack Beta · sensinode · ThingWorx · NODE · EVRYTHNG · bugswarm

**Open Source Platforms:** sense · spark · Nimbits · ThingSpeak

**Software Platforms:** sense · SmartThings · Withings · NINJABLOCKS · xively · TWINE · OSITO · zonoff

**Sensor Networks:** MESH SYSTEMS · SAFECAST

**Enabling Networks:** FreedomPop · SocialSign.in · Open Garden · SIGFOX

**Corporates:** IBM · GE · LG · CISCO · Honeywell

## Applications (Verticals)

### Quantified Self

**Wearable Computing:** GLASS · Pebble

**Fitness:** FUEL · amiigo · M · Withings · fitbit · JAWBONE

**Health:** BASIS · LUMOback · HAPIfork · wahoo FITNESS · M signal · NuMetrex by textronics

**Family:** REST · Live!y · Good Night Lamp · Withings · EVADO FILIP

### Lifestyle

**Leisure:** blossom · ICA kitchen · Thimble · remee · iGrill · HEXBRIGHT · sobi

**Pets:** gibi · FITBARK

**Toys:** sifteo · MakieLab · KAROTZ · greenGOOSE!

**Music:** gtar

**Gardening:** BITPONICS · plantlink · Koubachi

**Home Improv.:** Radiator Labs · netatmo

### Connected Home

**Home Automation:** SmartThings · stick AND · NINJABLOCKS · NODE · revolv · Ubi · lapka · electric imp · Wovyn

**Energy Efficiency:** knut · nest · wemo · tado° · LIFX · ecobee · belkin echo · micasaverde

**Security:** Kwikset · ALARM.COM · BOSCH · Lockitron · CANARY · HomeMonitor · iSmartAlarm

**New Interfaces:** NeuroSky · sphero · EQUISO · emotiv · gestigon · PrimeSense · InteraXon · LEAP

### Industries

**Retail:** Nomi · euclid · CP placemeter

**Healthcare:** ViSi MOBILE · AdhereTech · AliveCor · TELCARE · intelligentM

**Automotive:** A · mojio · Dashlabs · SYNC Powered by Microsoft · OpenXC · Toyota · entune

**Smart Buildings:** APOGEE ANYWHERE · Johnson Controls · Schneider Electric

### Industrial Internet

**Robotics:** KIVA Systems · ROMOTIVE · Double Robotics · Airware · liquid robotics · ROBOTEX · 3DRobotics UAV TECHNOLOGY · MOMENTUM

**Greentech:** BigBelly SOLAR · Axeda · enlighted · GRIDMOBILITY

**3D Printing:** 3DSYSTEMS · MezzoMill · Stratasys FOR A 3D WORLD · formlabs · shapeways · MakerBot INDUSTRIES · RepRap

## Building Blocks

**Connection Protocols:** neul · ZigBee · macheen · RFID · NFC · WiFi · Bluetooth · M-Bus · MQTT · 2G 3G 4G

**Telecom:** at&t · verizon · T··Mobile · Virgin mobile · boost mobile

**M2M:** CROSSBRIDGE SOLUTIONS · gemalto · Jasper wireless · Numerex · Telit · ERICSSON

**Software:** amazon web services · hadoop

**Mobile:** iOS · Android · Parse

**Hardware:** opengate · ARDUINO · Raspberry Pi · beagleboard.org · spark

**Parts /Kits:** MAKEY MAKEY · TinkerForge · littleBits · readydiymate · MOSORO

**Services:** TCH International · DRAGON innovation · makexyz · TINKERCAD · adafruit · CIRCUIT LAB

**Incubators:** BOLT · LEMNOS Labs · HAXLR8R · springboard();

**Funding:** KICKSTARTER · indiegogo

**Distribution:** GRAND St. · Anvil

© Matt Turck (@mattturck), Sutian Dong (@sutiandong) & FirstMark Capital (@firstmarkcap)

# What Is The Value In Big Data?

# Transforming Data Into Insight For Making Better Decisions
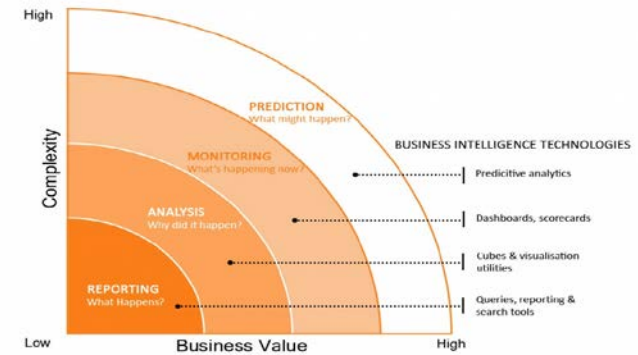


Gartner, 2013

# Why Data Mining?

- Automated data collection tools and mature database technology lead to tremendous amounts of data stored in databases, data warehouses and other information repositories.

*"We are drowning in data, but starving for knowledge!"*

(John Naisbitt, 1982)

- Data Mining is derived from the necessity to address the Data Explosion

# What is Data Mining?



- A set of technologies that uncovers relationships and patterns within large volumes of data that can be used to predict future behavior and events

- Predictive Analytics is technology that learns form experience to predict the future outcomes in order to drive better business decisions
  - Extracting / "Mining"
    - Information/Meaning from data
    - Interesting knowledge (rules, regularities, patterns, constraints) from raw data
    - Implicit, previously unknown and unexpected, potentially extremely useful information from data

# Data Mining is Multidisciplinary Field

- Database technology
- Artificial Intelligence
- Machine Learning including Neural Networks
- Statistics
- Pattern recognition
- Knowledge-based systems/acquisition
- High-performance computing
- Data visualization
- Other Disciplines



34

# Data Mining is NOT...

- Data Warehousing

- (Deductive) query processing
  - SQL/ Reporting

- Software Agents

- Expert Systems

- Online Analytical Processing (OLAP)

- Statistical Analysis Tool

- Data visualization

- BI – Business Intelligence

- Workflows

# Terminology

# What's Involved?

- **Learn the application domain**
  - Relevant prior knowledge and goals of application
- **Create a target data set:**
  - Data selection
  - Clean and preprocess the data (may take > 60% of effort!)
    - Data reduction and transformation
    - Find useful features, dimensionality/variable reduction, representation
- **Choose data mining functions**
  - Summarization, classification, regression, association, clustering
- **Apply data mining methods/algorithm(s)**
  - Data mining: search for patterns of interest
- **Evaluation, validation, and knowledge presentation**
  - Visualization, transformation, removing redundant patterns, etc.
- **Use and integrate discovered knowledge**

# Data Mining Functions

- Exploratory Data Analysis
  - Visualization, Statistical analysis

- Descriptive Modeling/Discovering patterns rules
  - Cluster analysis/segmentation
  - Association/Dependency  rules
  - Sequential patterns

- Predictive Modeling
  - Classification and Regression
  - Temporal sequences
  - Deviation detection

# CRISP-DM Iterative Process

# Predictive Analytics Process

# Data Mining Process

# Big Data and Predictive Analytics Processing

# Analytics Maturity Levels

# Numerous Applications

Improve ability to classify and treat cancer, tumors, diseases

Adjust credit scores as transactions are occurring to account for risk fluctuations

Apply inferred customer social relationships to prevent churn

Increase revenue and customer satisfaction by discovering passengers who are likely to miss their flight

Hospital

Loan officer

Call Center

Airline

# Condition Based Maintenance (CBM)

❑ Consumes online data - conditions are constantly monitored and their signatures evaluated

❑ Focuses on sensors and communications

❑ Machines are networked collecting large volumes of data

❑ What to do with all this data?

❑ We need to turn data into insight

❑ Patterns of past behavior can be "learned" to provide deeper insight

    ❑ How much change or degradation has occurred since the last round

    ❑ What is the chance of failure for each machine?

    ❑ What is the current most likely failure to happen?

    ❑ What is the impact of the failure?

# Reliability Value Chain
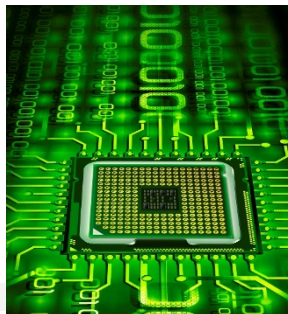
**Unexpected Failures**

Predictive Maintenance

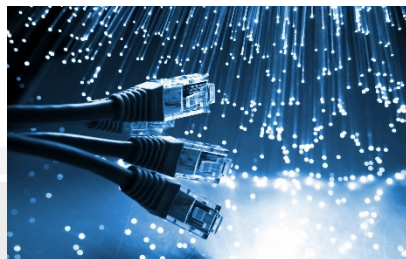# Intelligent CBM Across Industries

## Semiconductor industry

- Capital process equipment
- Anomaly and event detection
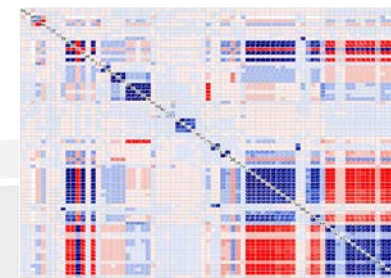- Predictive Maintenance
- Proactive decisions

## Network Data

- Improved flow
- Asset optimization
- Proactive decisions
- Resource allocation
- Data driven risk prediction

## Manufacturing

- Improved process
- Failure prediction
- Quality control
- Prevents damages
- Lowers production cost

# Intelligent CBM for Utilities Projects

## Smart Grid



- Challenge
  - Rising energy demand
  - Aging infrastructure
  - Grid Stability

- Goals
  - Situational awareness
  - Event Detection
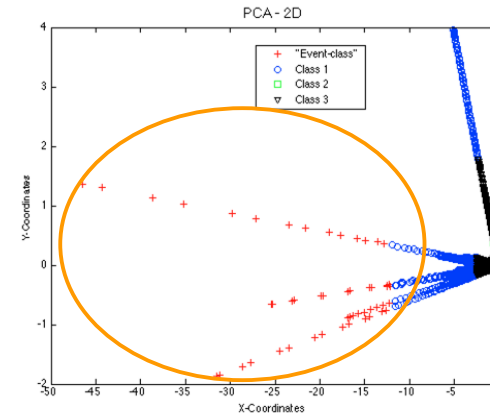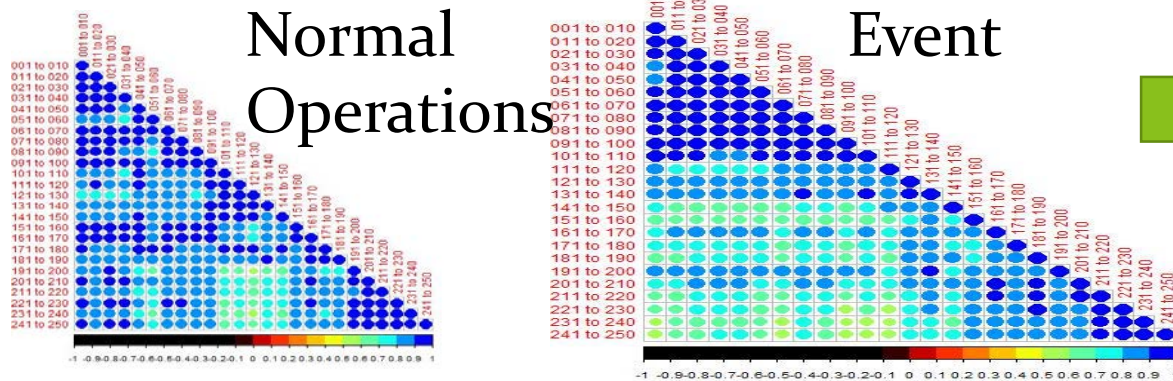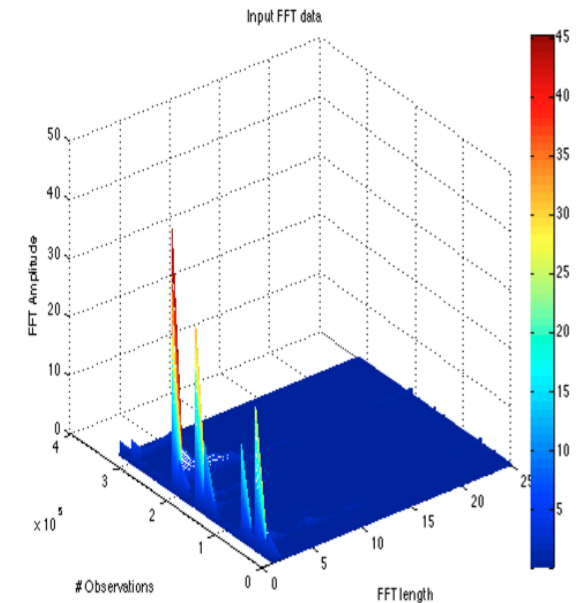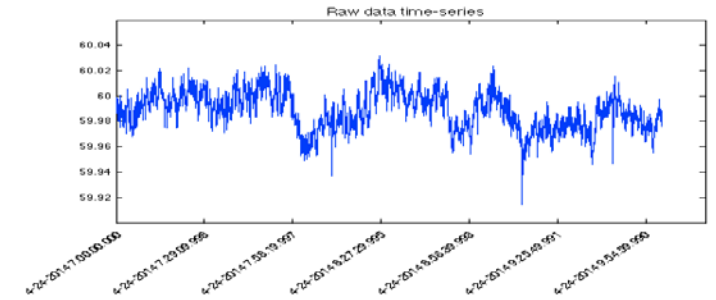  - Power outage prevention

## Gas Pipelines



- Challenge:
  - Increased compliance regulations
  - Rising demand
  - Aging infrastructure

- Goals
  - Improved safety
  - Asset optimization
  - Situational awareness
  - Proactive decisions
  - Resource allocation
  - Data driven risk prediction

# PMU Data Analysis


Raw data time-series

- Frequency, Magnitude, and Angle for two PMU's
- Collected 30 times a second
- Very sensitive and noisy measurements
- Goal: detect and predict event


Normal Operations


Event


PCA - 2D

Outliers detection


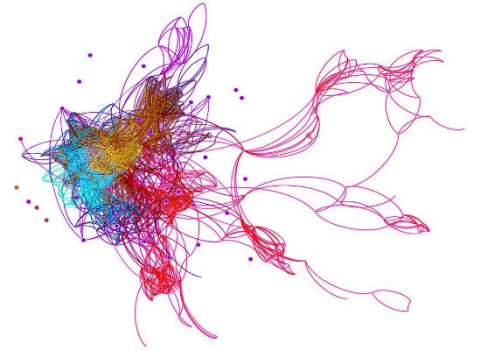Input FFT data

# Big Data for Cyber Security

- Identify risks and anticipate problems before they occur

- Traditional security mechanisms leverage rule, pattern, signature and algorithm-based approaches to detect threats

- Analyze changes in behavior and predict risks and breaches before they happen
    - Malicious Code Detection
    - Network intrusion detection
    - Anomaly detection
    - Data Stream Mining
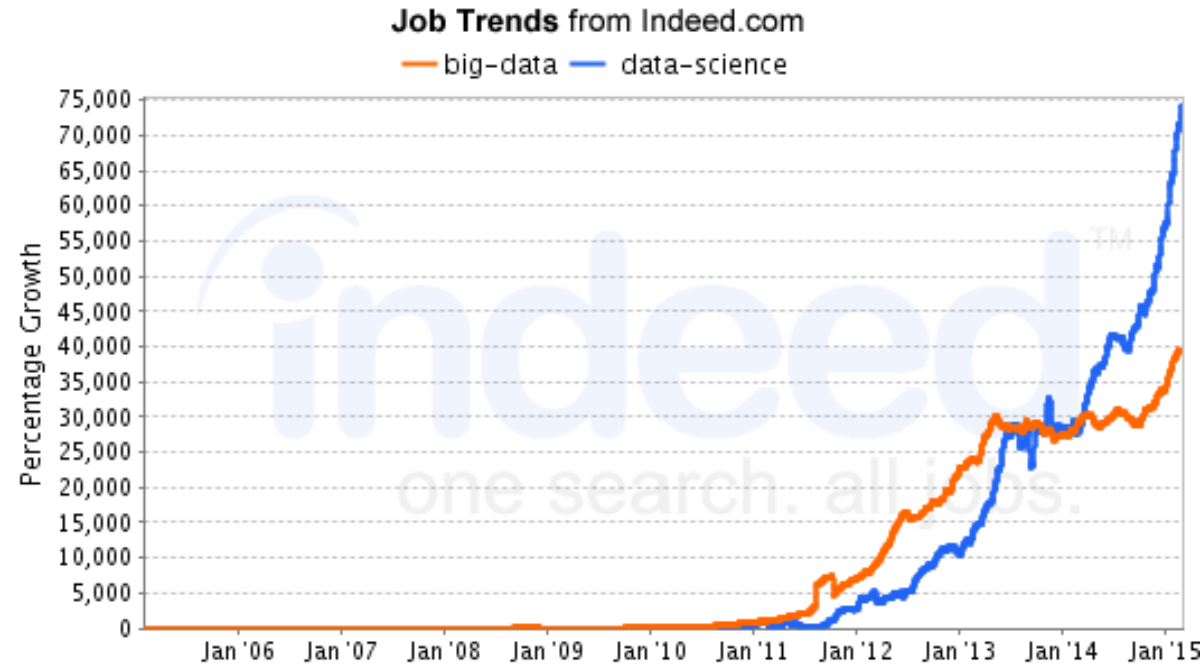
# Big Data – Big Data Science

- "Data Scientist"
  - **The "Hot new gig in town"**
    - O'Reilly report
  - **Data Scientist: The Sexiest Job of the 21st Century**
    - Harvard Business Review, October 2012
    - The next sexy job in next 10 years will be statistician" – Hal Varian, Google Chief Economist
    - Geek Chic – Wall Street Journal – new cool kids on campus
  - The future belongs to the companies and people that turn data into products
- *"The human expertise to capture and analyze big data is both the most expensive and the most constraining factor for most organizations pursuing big data initiatives" –* Thomas Davenport
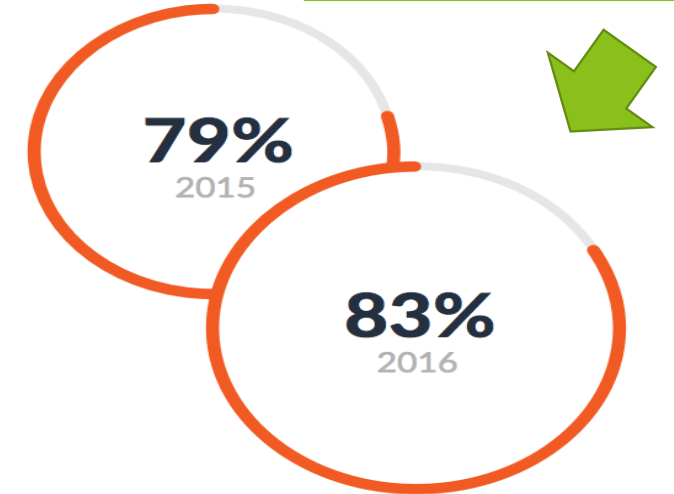
# Data scientist:  The hot new gig in tech

- Article in *Fortune* - *"nerdy-cool job that companies are scrambling to fill: data scientist"*

- Gartner in 2012 said there would be a shortage of 100,000 data scientists in the United States by 2020.

- *McKinsey Global Institute* "Big data Report" in 2011

  - By 2018, the United States alone could face a shortage of 140,000 to 190,000 people with deep analytical skills as well as 1.5 million managers and analysts with the know-how to use the analysis of big data to make effective decisions - demand that's 60 percent greater than supply

- In 2014, the consulting firm Accenture found that more than 90 percent of its clients planned to hire people with data science expertise

- Gartner says the current demand for data scientist exceeds the current supply by factor of three

# Data Science Job Growth



**Job Trends** from Indeed.com
— big-data — data-science



**79%** 2015

**83%** 2016

Not enough
Data Scientist

Crowdflower survey report:
A full 83% of respondents said there weren't
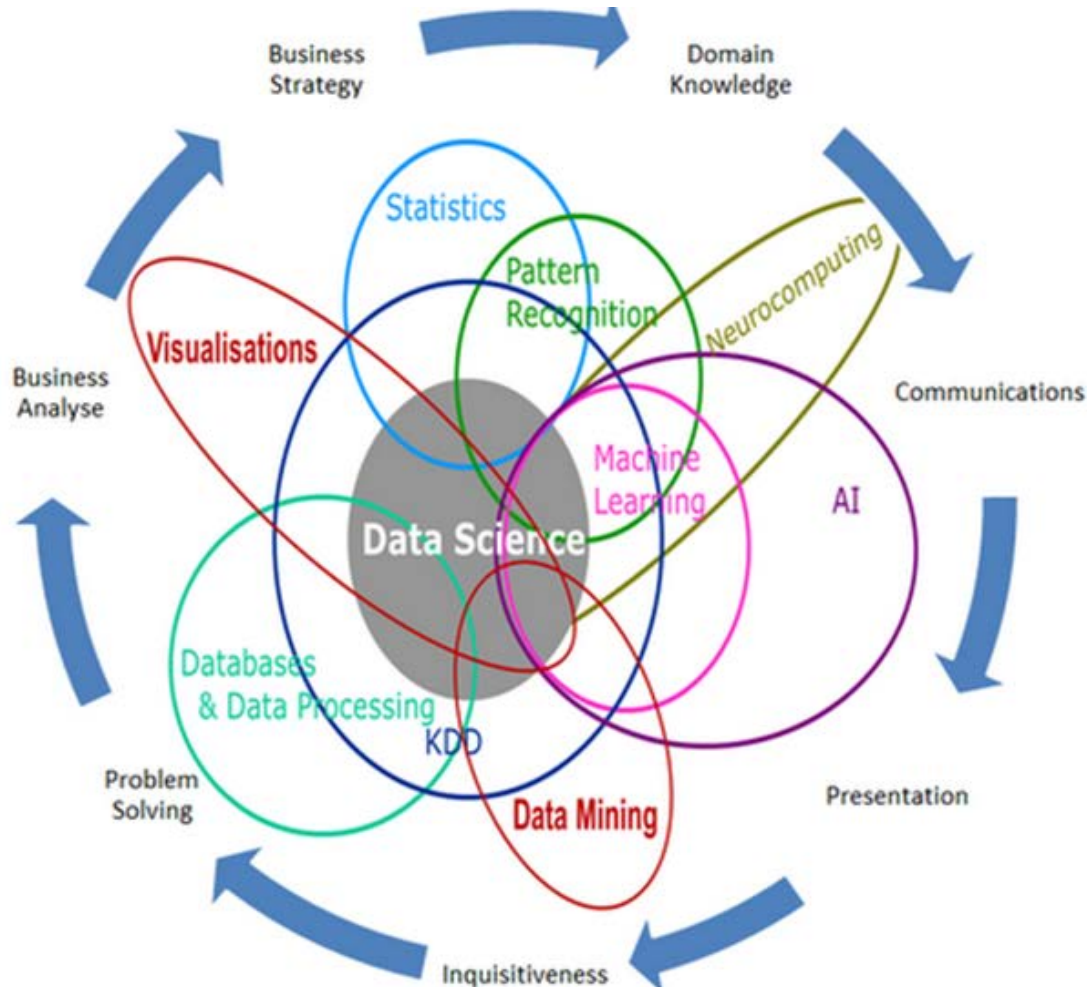enough data scientists to go around

**By 2018 shortage of 140-190,000 predictive analysts and 1.5M managers / analysts in the US**

# Data Miners: Past and Present

- Traditional approaches have been for DM experts: "White-coat PhD statisticians"

  - DM tools also fairly expensive

- Today: approach is designed for those with *some* Database/Analytics  skills

  - DM built into DB, easy to use GUI, Workflows

  - Many jobs available from Statistical analyst to Data Scientist!

- Data Science:  The Art of mathematically sophisticated data engineers delivering insights from data into business decisions and systems

# Data Scientist Skill and Characteristics



- Intellectual curiosity, Intuition
  - Find needle in a haystack
  - Ask the right questions – value to the business
- Communication and engagements
- Presentation skills
  - Let the data speak but tell a story
  - Story teller – drive business value not just data insights
- Creativity
  - Guide further investigation
- Business Savvy
  - Discovering patterns that identify risks and opportunities
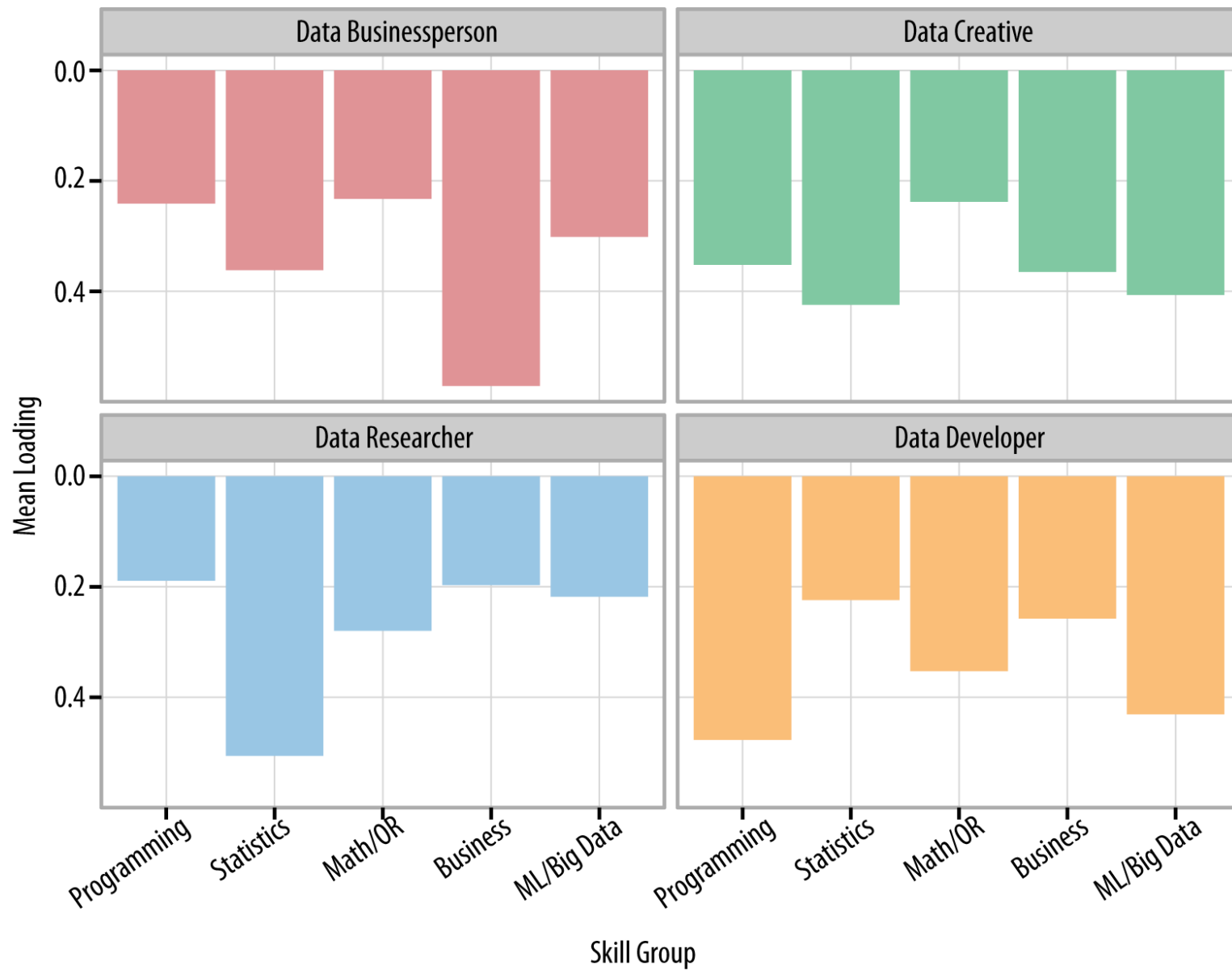  - Measure

# Data Scientist Self-ID

| | | | |
|---|---|---|---|
| **Data Developer** | Developer | Engineer | |
| **Data Researcher** | Researcher | Scientist | Statistician |
| **Data Creative** | Jack of All Trades | Artist | Hacker |
| **Data Businessperson** | Leader | Businessperson | Entrepeneur |

**O'Reilly Strata Survey suggested Self-ID Group, along with the self-ID categories most strongly associated with each Group**
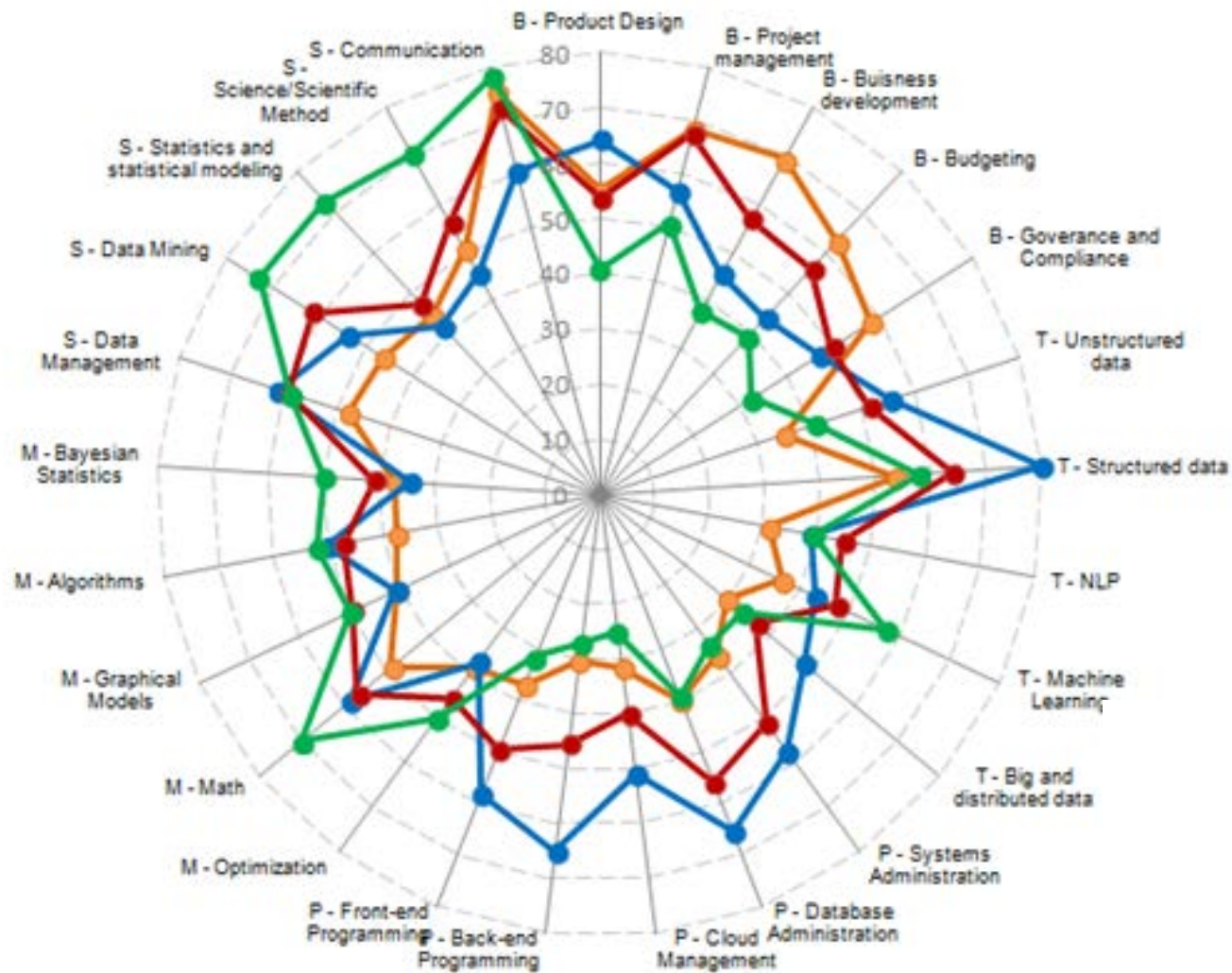
# Strata Survey Skills

| Business | ML / Big Data | Math / OR | Programming | Statistics |
|----------|---------------|-----------|-------------|------------|
| Product Developement | Unstructured Data | Optimization | Systems Administration | Visualization |
| Business | Structured Data | Math | Back End Programming | Temporal Statistics |
| | Machine Learning | Graphical Models | Front End Programming | Surveys and Marketing |
| | Big and Distributed Data | Bayesian / Monte Carlo Statistics | | Spatial Statistics |
| | | Algorithms | | Science |
| | | Simulation | | Data Manipulation |
| | | | | Classical Statistics |

Strata Survey Skills

Skill Group

Mean Loading

Data Businessperson · Data Creative · Data Researcher · Data Developer

Programming · Statistics · Math/OR · Business · ML/Big Data
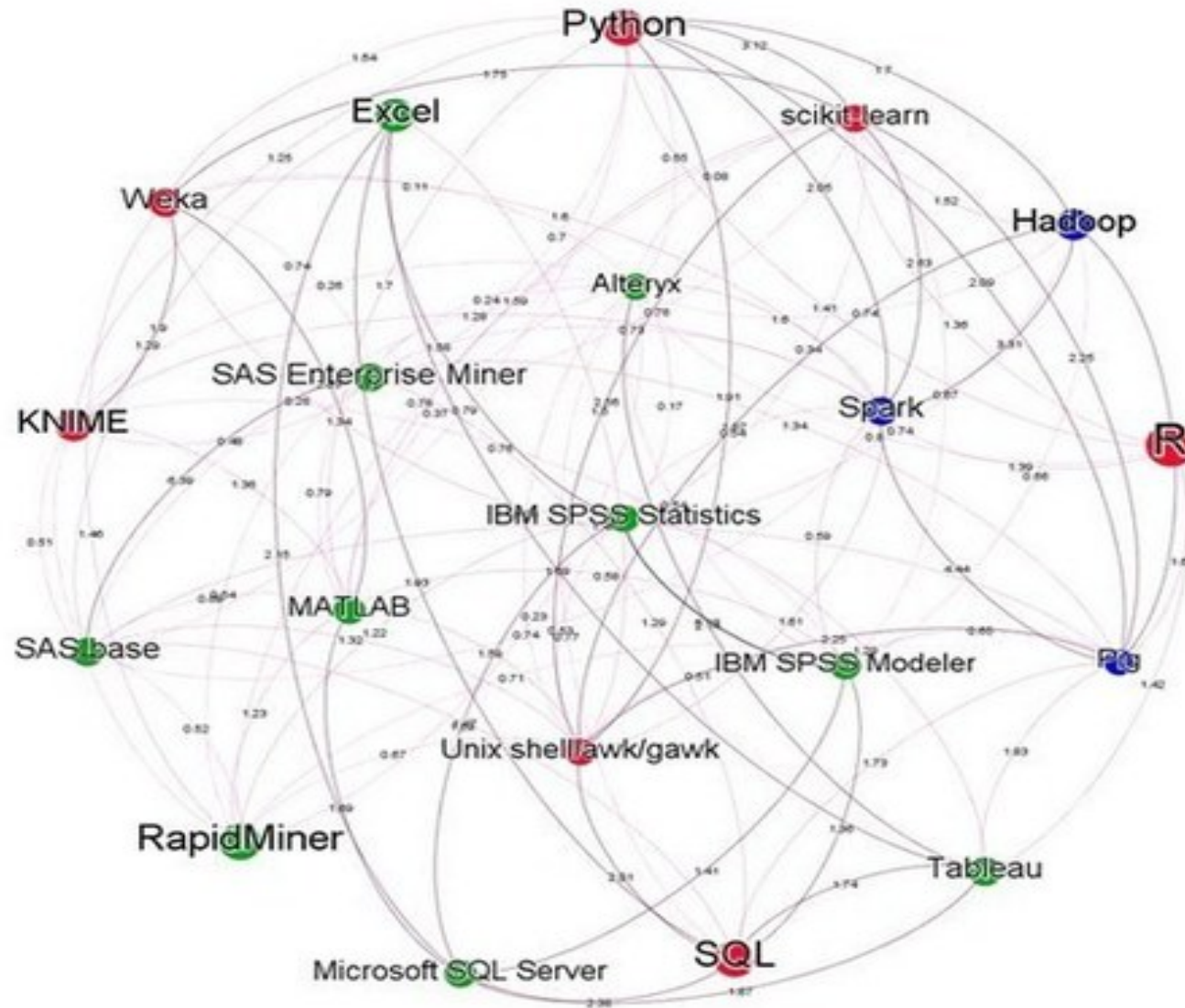
**Strata Survey Skills**

# Learning and Training Opportunities

- Locally: UCSD MAS, UCSD Extension Certificate, PACE Boot Camps, etc.

- Introduction to Data Science Example
  - Part 1: Data Manipulation at scale
    - Databases and the relational algebra
    - Parallel databases, parallel query processing, in-database analytics, MapReduce, Hadoop, relationship to databases, algorithms, extensions, languages
    - Key-value stores and NoSQL; Entity resolution, record linkage
  - Part 2: Analytics, Predictive Analytics, Text mining
  - Part 3: Communicating Results
    - Visualization, data products, visual data analytics
    - Provenance, privacy, ethics, governance

https://www.coursera.org/course/datasci

# World the Data Science Tools

# Successful Data Analytics Project Guidelines

- Don't start a Big Data project without understanding the value – what is my ROI?

- Don't ignore the wider Enterprise story

- A big data project is not just a technology project – it is about business change

- It pays to build the right team – Data Scientist

- Build the support structure from the start

- Build for tomorrow

- Break down the silos - involve the organization

# Thank you!

natashabalac@gmail.com