



SWE404/DMT413

BIG DATA ANALYTICS

Lecture 1: Introduction

Lecturer: Dr. Yang Lu

Email: luyang@xmu.edu.my

Office: AI-432

Office hour: TBD

Lecturer Information

Yang Lu (Jason), Ph.D.

■ Experience:

- 2020~Present: Assistant Professor, School of Electrical and Computer Engineering, Xiamen University Malaysia.
- 2019~Present: Assistant Professor, Department of Computer Science, School of Informatics, Xiamen University.
- 2018~2019: Research Intern, Tencent.
- 2014~2019: Ph.D. student, Department of Computer Science, Hong Kong Baptist University.
- 2012~2014: Master student, Department of Software Engineering, University of Macau.
- 2008~2012: Bachelor student, Department of Software Engineering, University of Macau.

■ Research interests:

- Artificial intelligence, machine learning.

■ Personal website:

- <https://jasonyanglu.github.io/>

Consultation Hours

- Office location: AI-432
- Consultation Hours:
 - TBD
- For other time slot:
 - Send email to confirm my presence first.
 - Or take your luck to knock the door (I will be in my office for most of the time in this semester).

Class Rules

- You can do anything except:
 - Make noises (chatting, singing...)
 - Catch others' eyes (dancing, push-ups...)
 - Eat food with strong smell (instant noodles, durian...)
- Feel free to interrupt me if you have questions (no need to raise hand).
- According to the university policy, taking attendance is needed.
 - **Important: you are required to have an 80% attendance to be able to seat for the final exam.**

Course Assessment

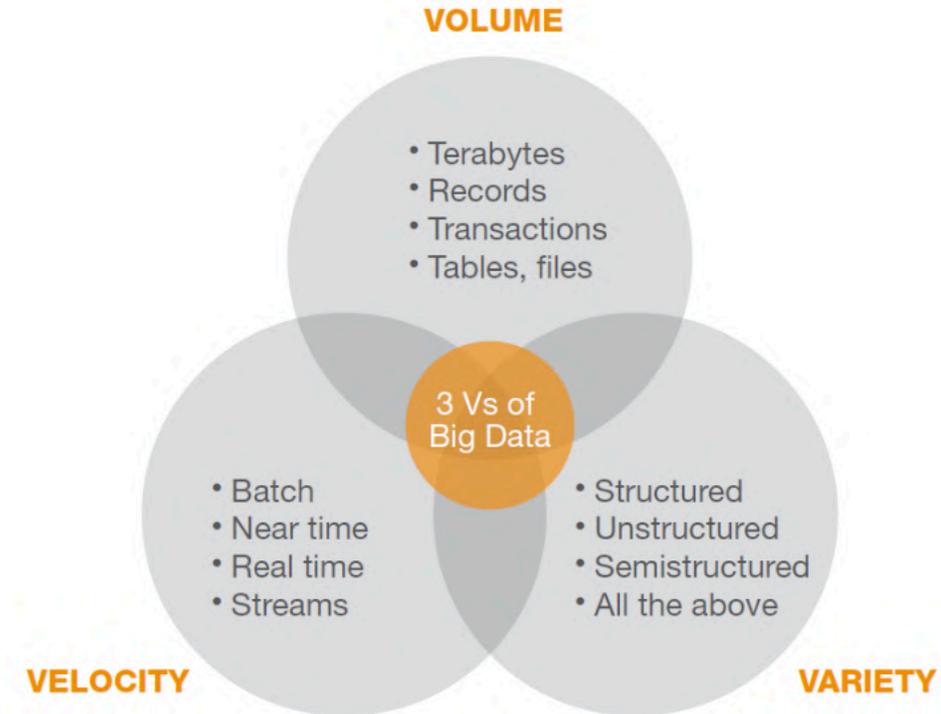
- Temporary according to the situation:
 - Final exam: 50%
 - Assignment: 20%, individually
 - Project: 30%, 2-3 members per group, report and presentation are required.
- **Important: cheating and plagiarism will get no marks.**

Textbook

- No textbook used on this course.
- Links of necessary materials will be shown in the slides.

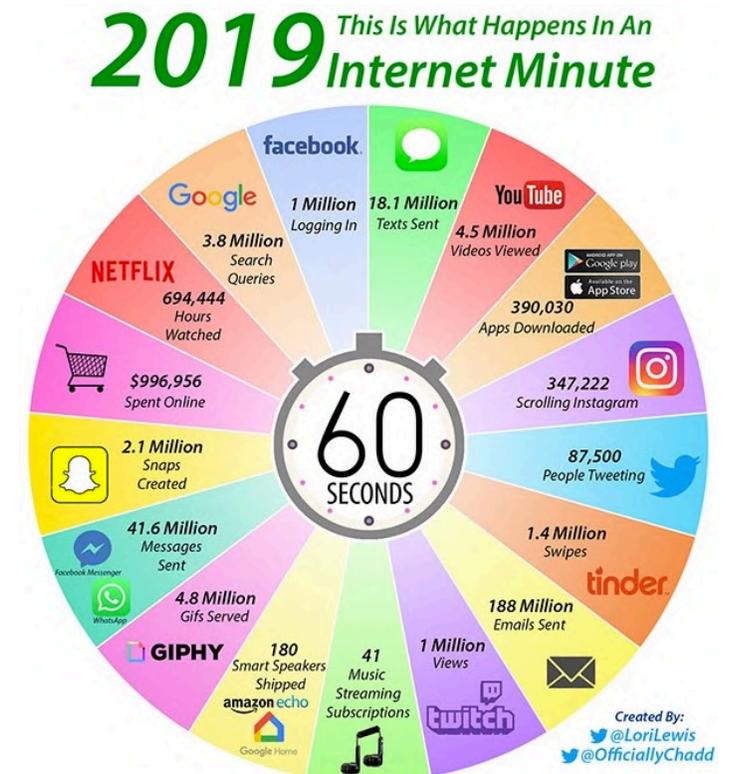
What Is Big Data?

- The term “big data” refers to data that is so **large**, **fast** or **complex** that it’s difficult or impossible to process using traditional methods.



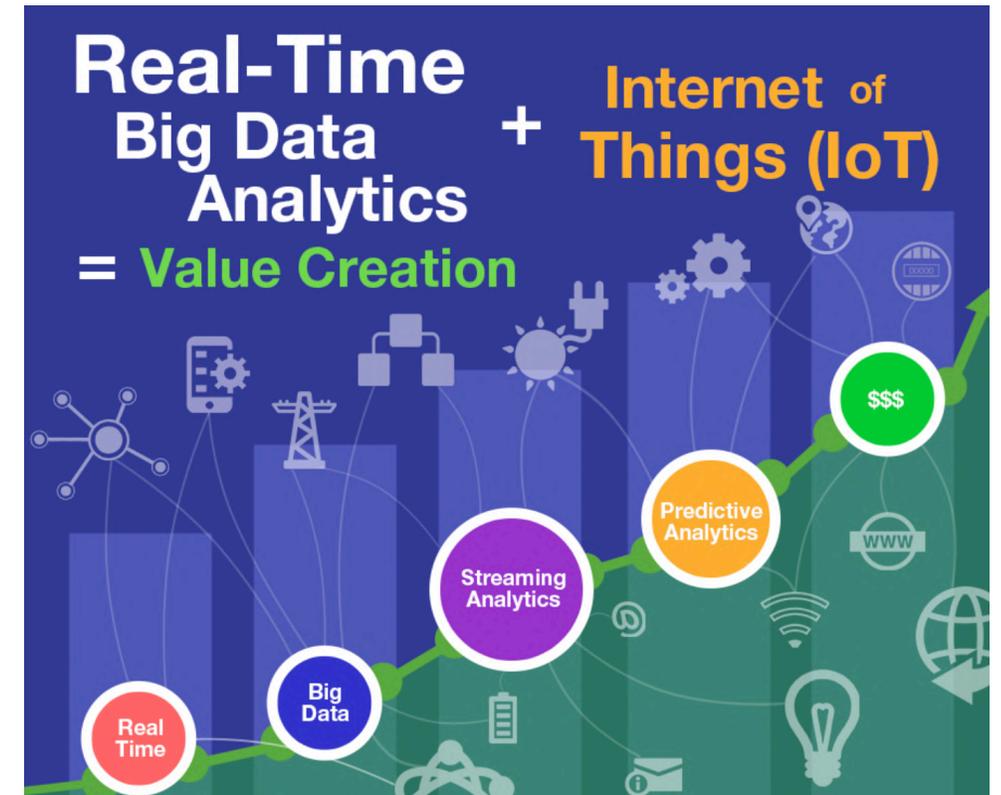
Volume of Big Data

- Organizations collect data from a variety of sources, including
 - business transactions
 - smart (IoT) devices
 - industrial equipment
 - videos
 - social media
 - ...
- In the past, storing it would have been a problem – but cheaper storage on platforms like data lakes and Hadoop have eased the burden.



Velocity of Big Data

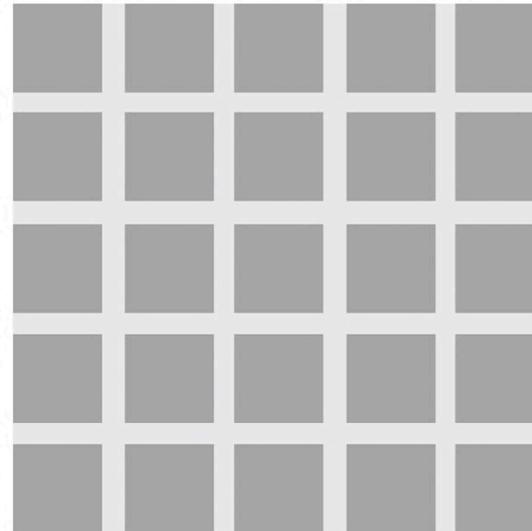
- With the growth in the Internet of Things, data streams in to businesses at an unprecedented speed and must be handled in a timely manner.
- Big data techniques are driving the need to deal with these torrents of data in near-real time.



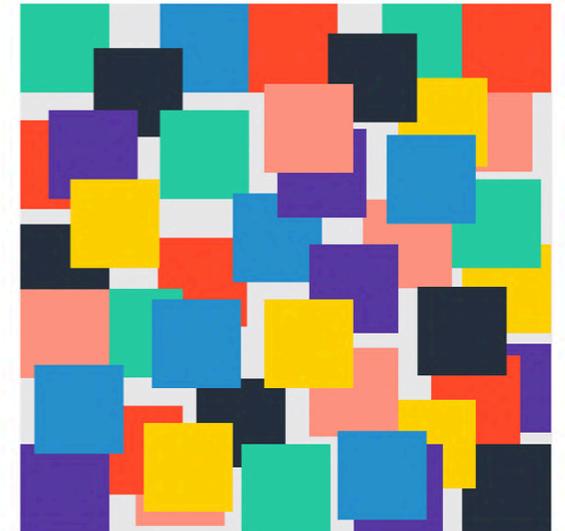
Variety of Big Data

- Data comes in all types of formats
 - Structured: quantitative data, fits neatly within fixed fields and columns in relational databases and spreadsheets
 - Numeric data
 - Database
 - ...
 - Unstructured: qualitative data, cannot be processed and analyzed using conventional tools and methods.
 - Text documents
 - Videos
 - Financial transactions
 - ...

Structured data

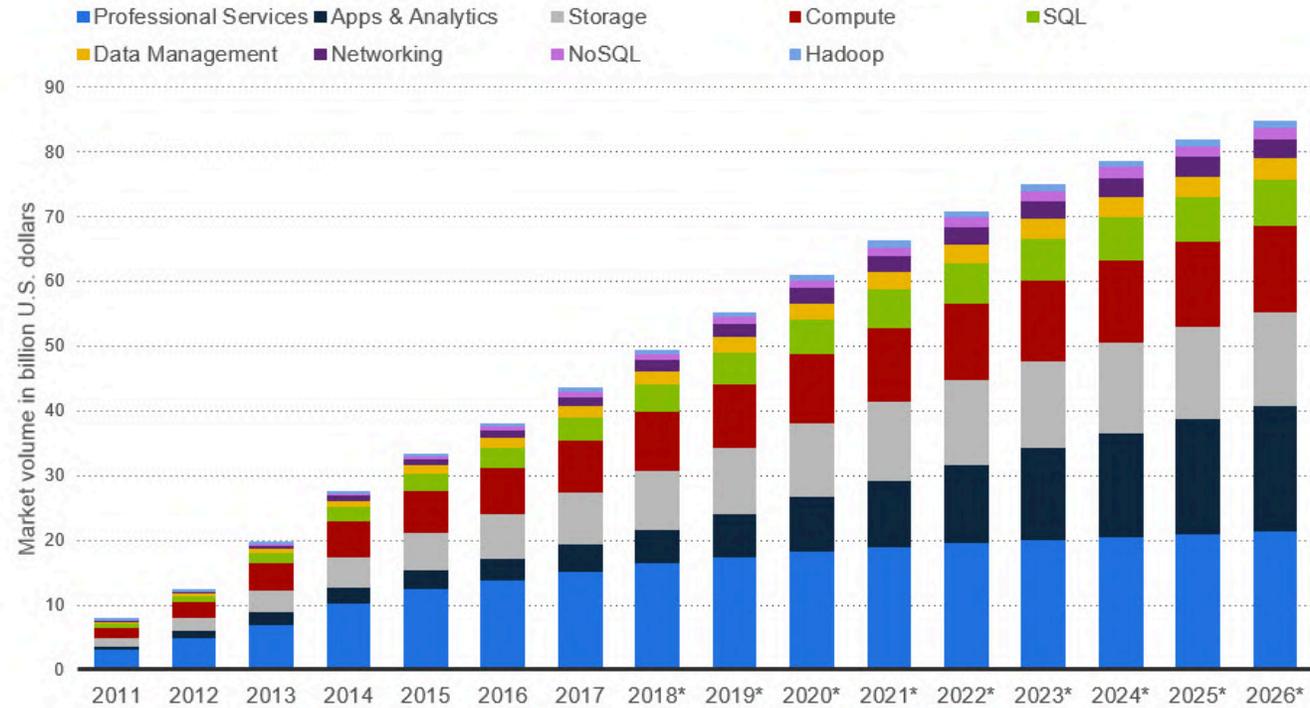


Unstructured data

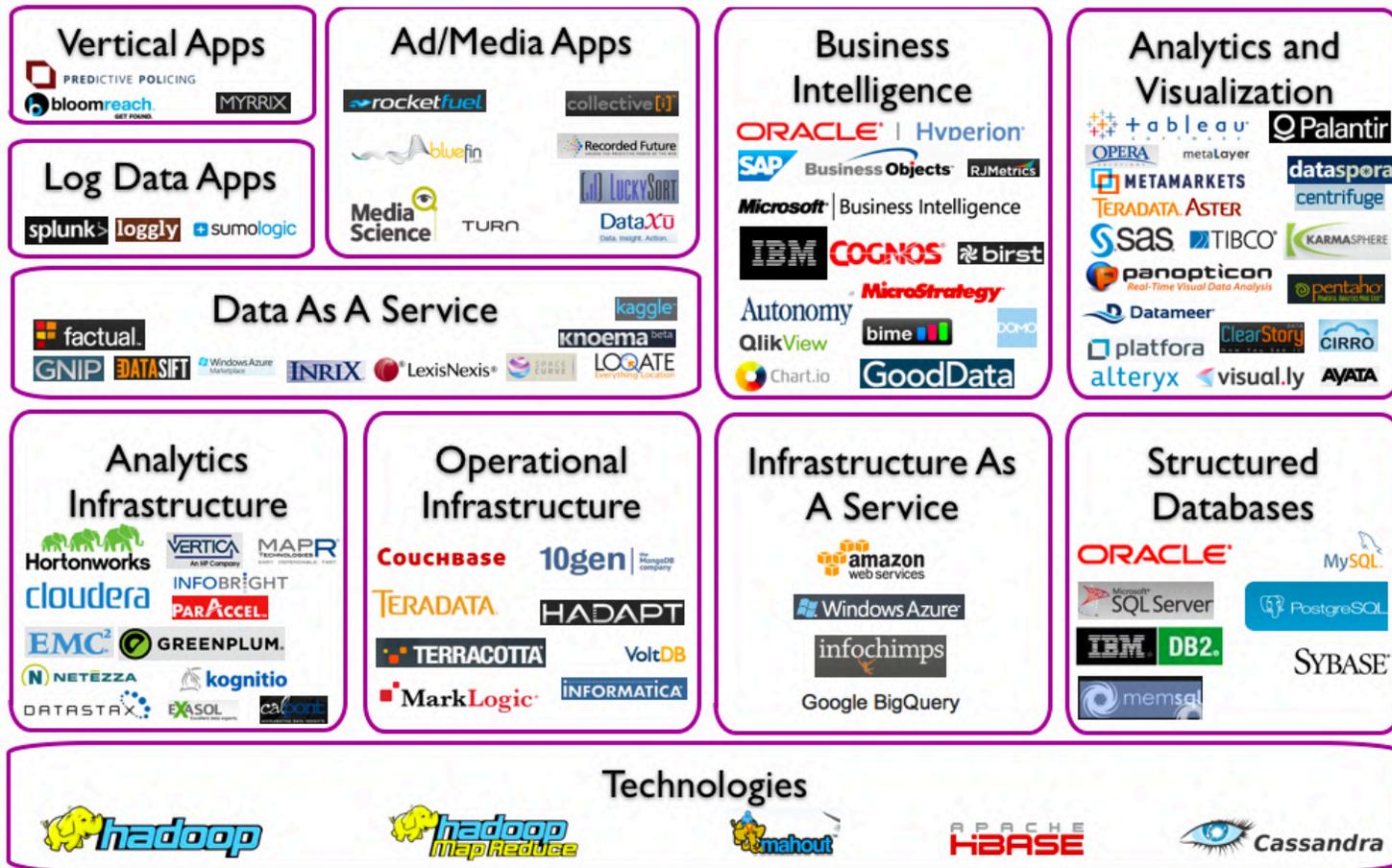


Big Data Market Worldwide Segment Revenue Forecast 2011-2026

Big Data Market Forecast Worldwide from 2011 to 2026, by segment (in billion U.S. dollars)



Big Data Landscape 2012

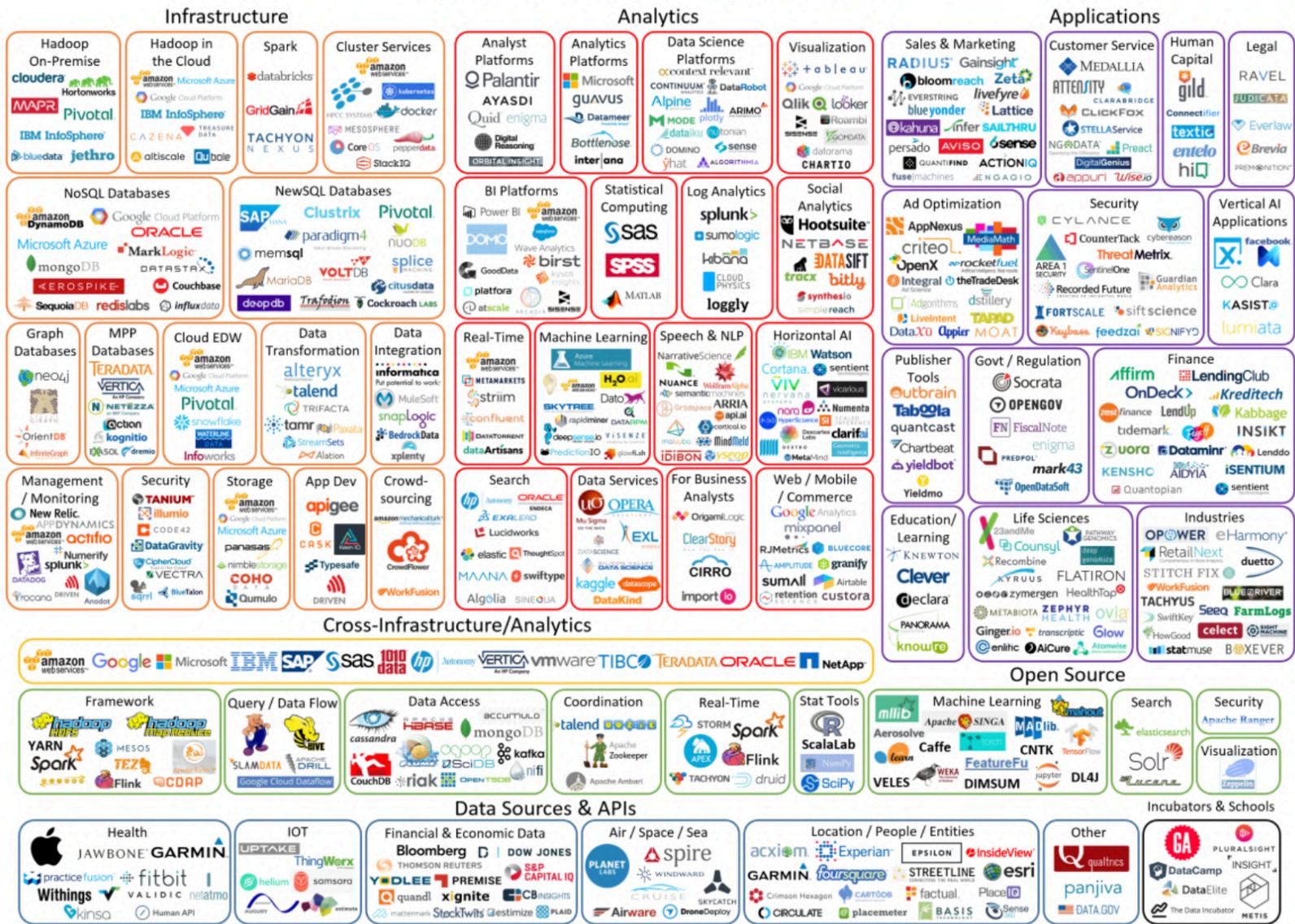


Copyright © 2012 Dave Feinleib

dave@vcdave.com

blogs.forbes.com/davefeinleib

Big Data Landscape 2016 (Version 3.0)



Last Updated 3/23/2016

© Matt Turck (@mattturck), Jim Hao (@jimhao), & FirstMark Capital (@firstmarkcap)

FIRSTMARK

Image source: <https://mattturck.com/big-data-landscape-2016-v18-final/>

What Is Analytics

- The importance of big data doesn't revolve around how much data you have, but **what you do with it.**
- Analytics is the scientific process of transforming data into insight for making better decisions, offering new opportunities for a competitive advantage.

Types of Analytics

Predictive Analytics

- Predicting the future based on historical patterns.
- What could happen?

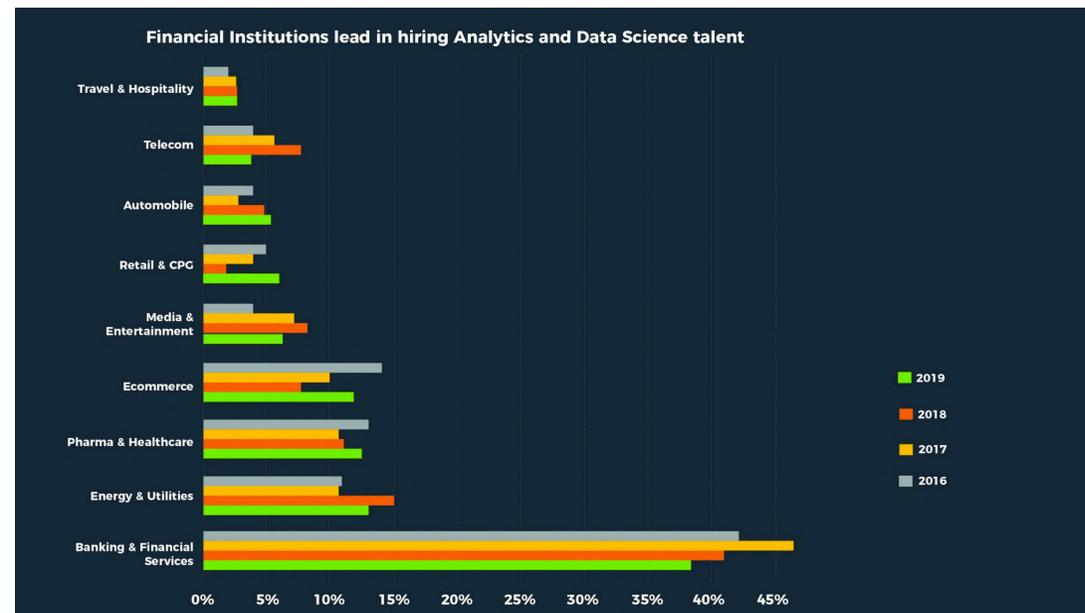
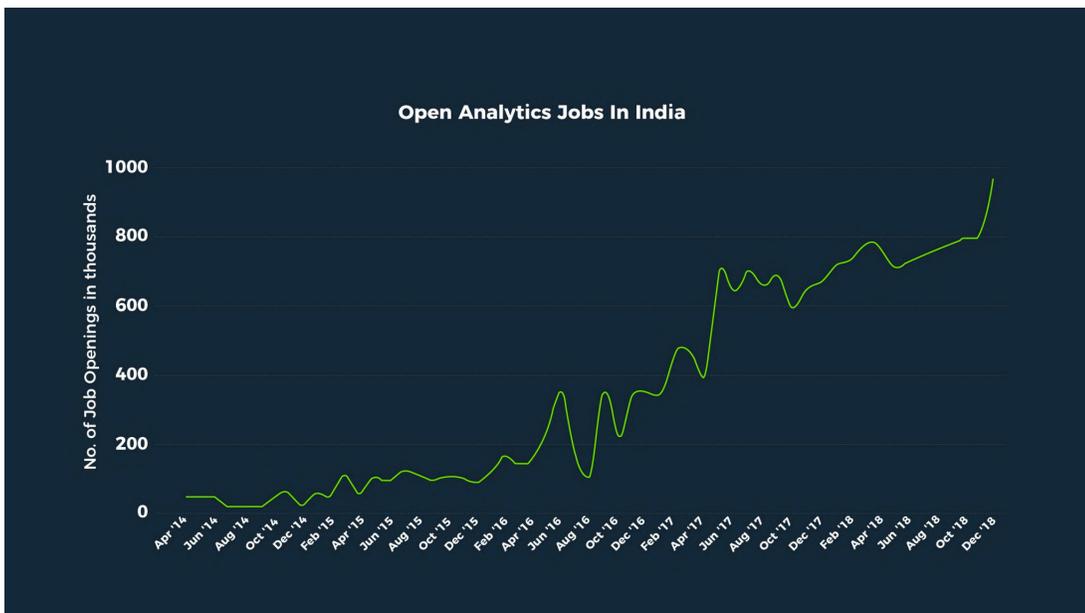
Descriptive Analytics

- Mining historical data to provide business insights.
- What has happened?

Prescriptive Analytics

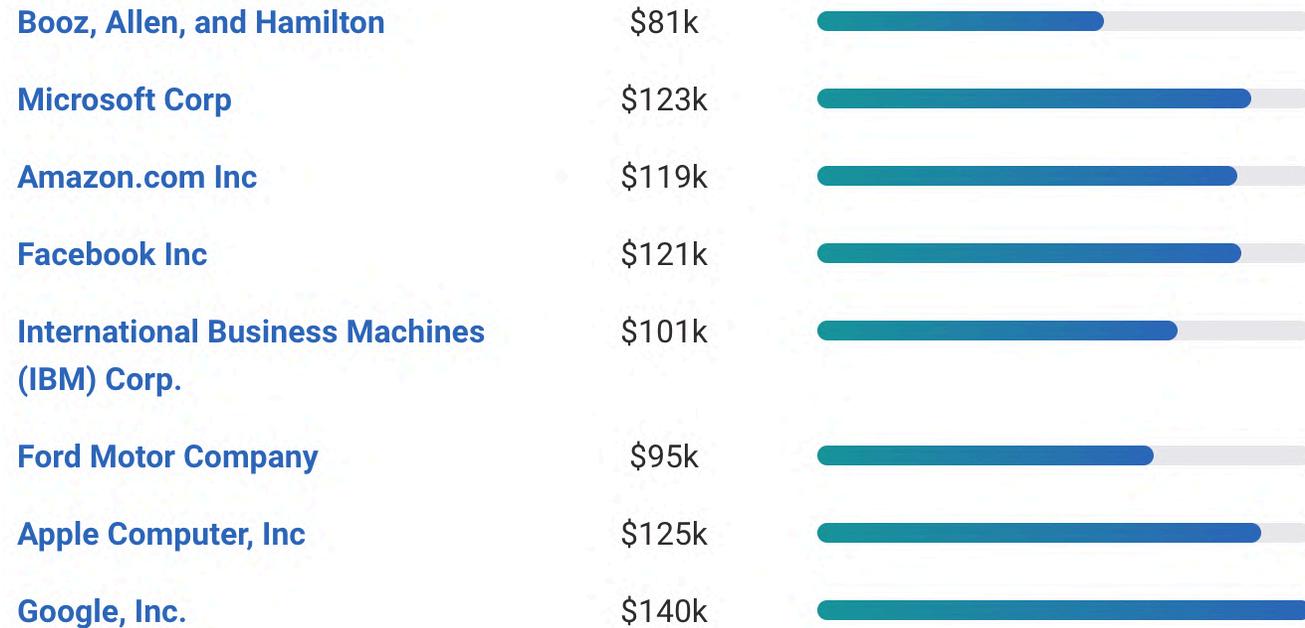
- Enabling smart decisions based on data.
- What should we do?

Job Market for Analytics



Salary

Popular Employer Salaries for Data Scientist

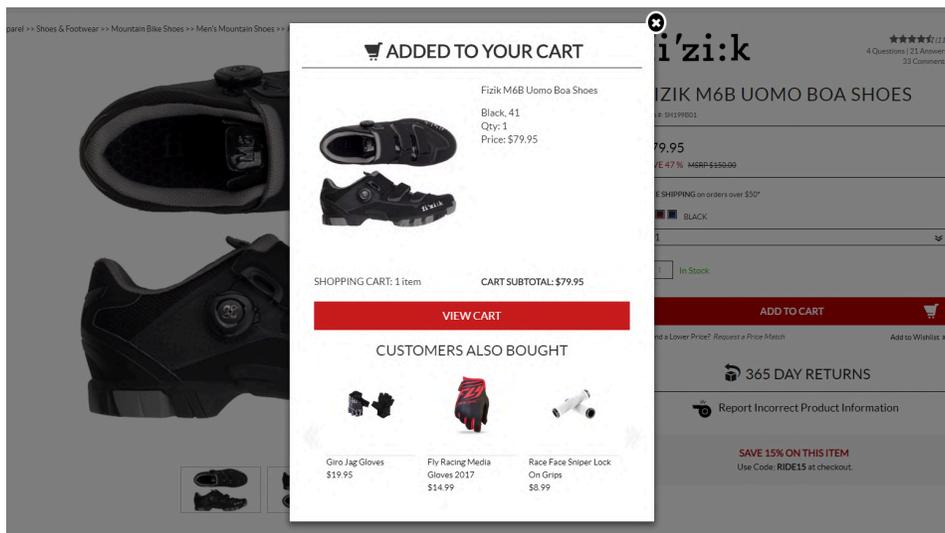


Top 10 Applications of Big Data Across Industries

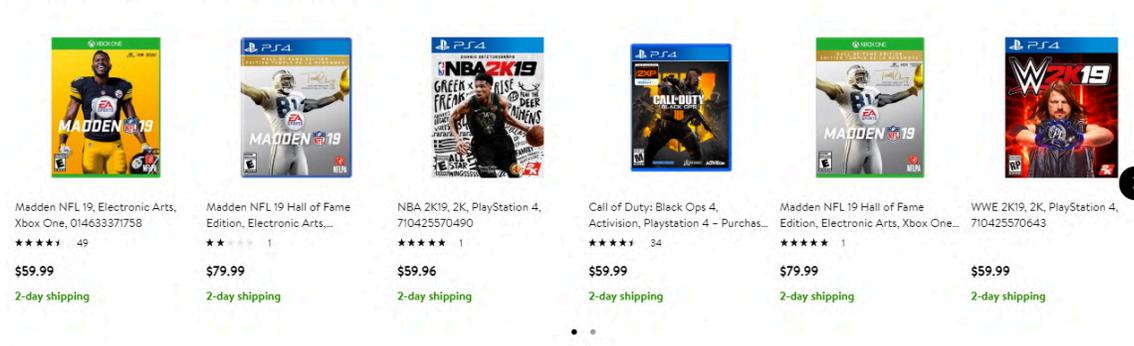
- **Banking and Securities**
 - Trade analytics used in high-frequency trading, pre-trade decision-support analytics, sentiment measurement.
- **Communications, Media and Entertainment**
 - Create content for different target audiences, recommend content on demand, measure content performance.
- **Healthcare Providers**
 - Evidence-based medicine test, faster identification, tracking the spread of chronic disease.
- **Education**
 - Track study activity, measure teacher's effectiveness, correct astray students.
- **Manufacturing and Natural Resources**
 - Geospatial data analysis, seismic interpretation.
- **Government**
 - Energy exploration, financial market analysis, fraud detection, health-related research, and environmental protection.
- **Insurance**
 - provide customer insights for transparent and simpler products, fraud detection.
- **Retail and Wholesale Trade**
 - Optimized staffing, reduced fraud, timely analysis of inventory.
- **Transportation**
 - Traffic control, route planning, intelligent transport systems, congestion management.
- **Energy and Utilities**
 - Faster meter readers, better asset and workforce management.

Case I: Personalized Recommendation

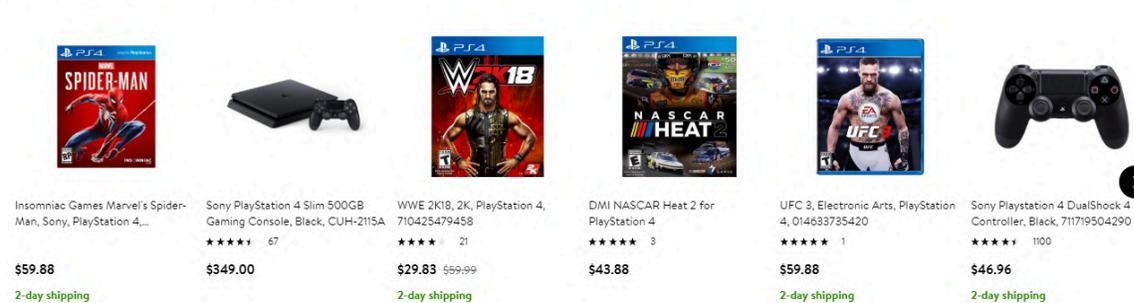
- Recommend based on your profile:
 - Demographic characteristics, historical order, recent clicks, ...



Customers also viewed these products

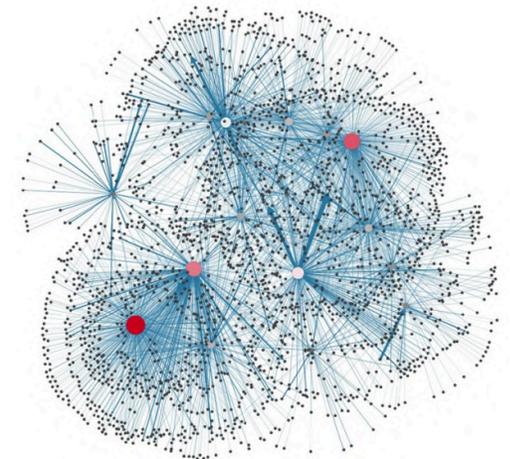


Customers also bought these products

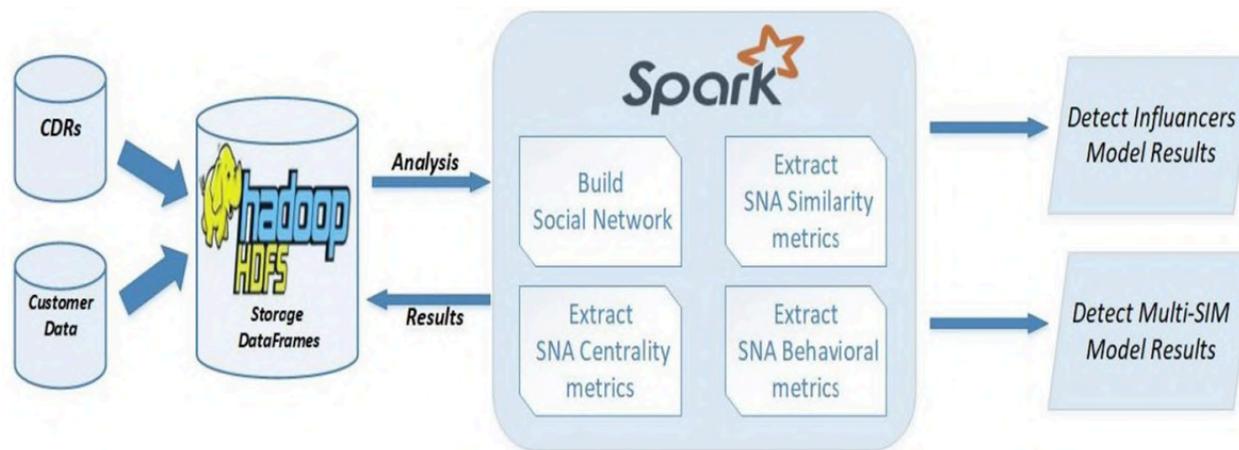


Case 2: Social Network Analysis by Call Detail Records (CDR)

- Construct a weighed graph to model social network, representing how close two subscribers are to each other.
 - Predict the relationship between subscribers.
- Usage: market campaign, churn prediction, prevent phone fraud,...



Visualization for a sample of the Telecom social network



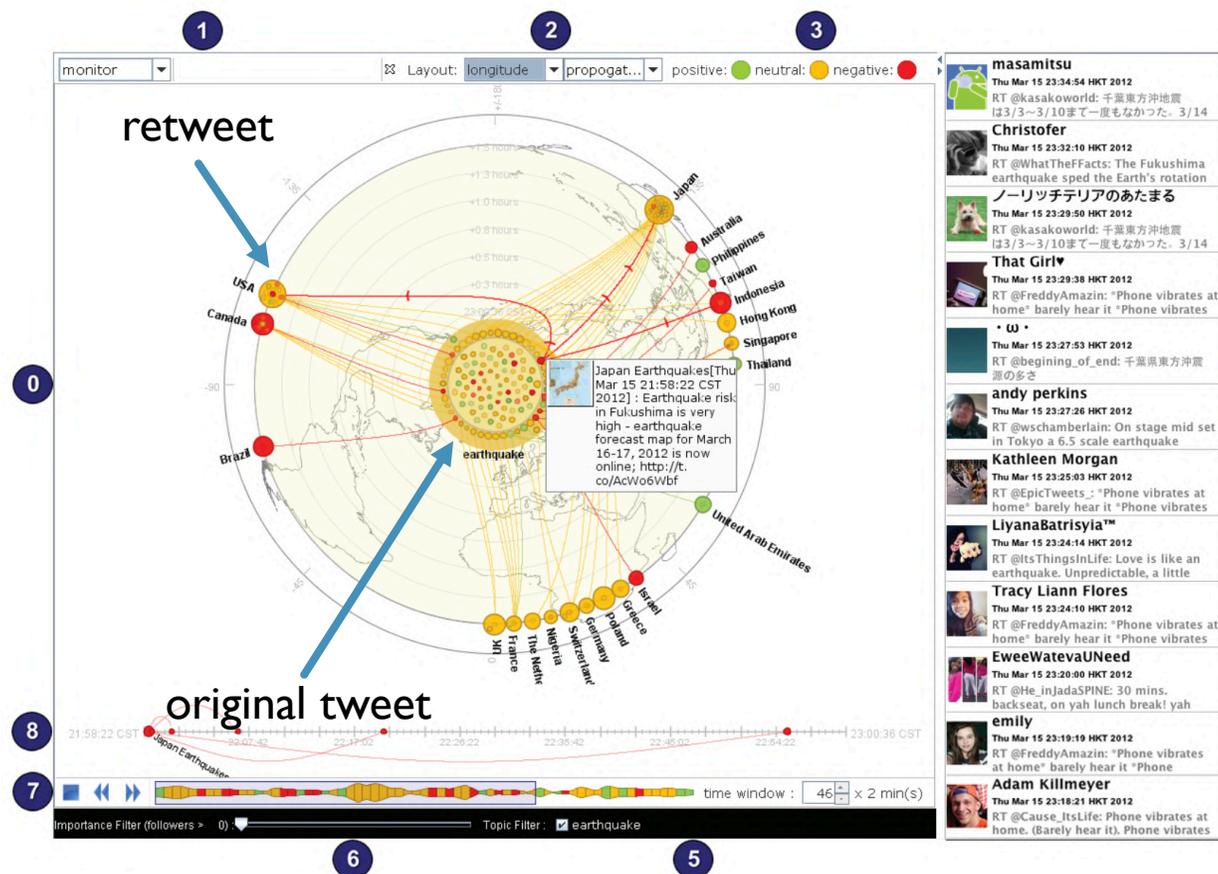
Solution architecture cycle

Source	Destination	Weight
963-9*****14	963-9*****22	0.0425
963-9*****31	963-9*****62	0.0496
963-9*****94	963-9*****11	0.0272
963-9*****34	963-9*****78	0.0335

Sample of Telecom social network data

Case 3: Tracing Information Diffusion with Data Visualization

- Social media, like Twitter, has been increasingly used for exchanging information, opinions and emotions about events that are happening across the world.
- Three major characteristics of diffusion processes in social media:
 - temporal trend
 - social-spatial extent
 - community response of a topic of interest



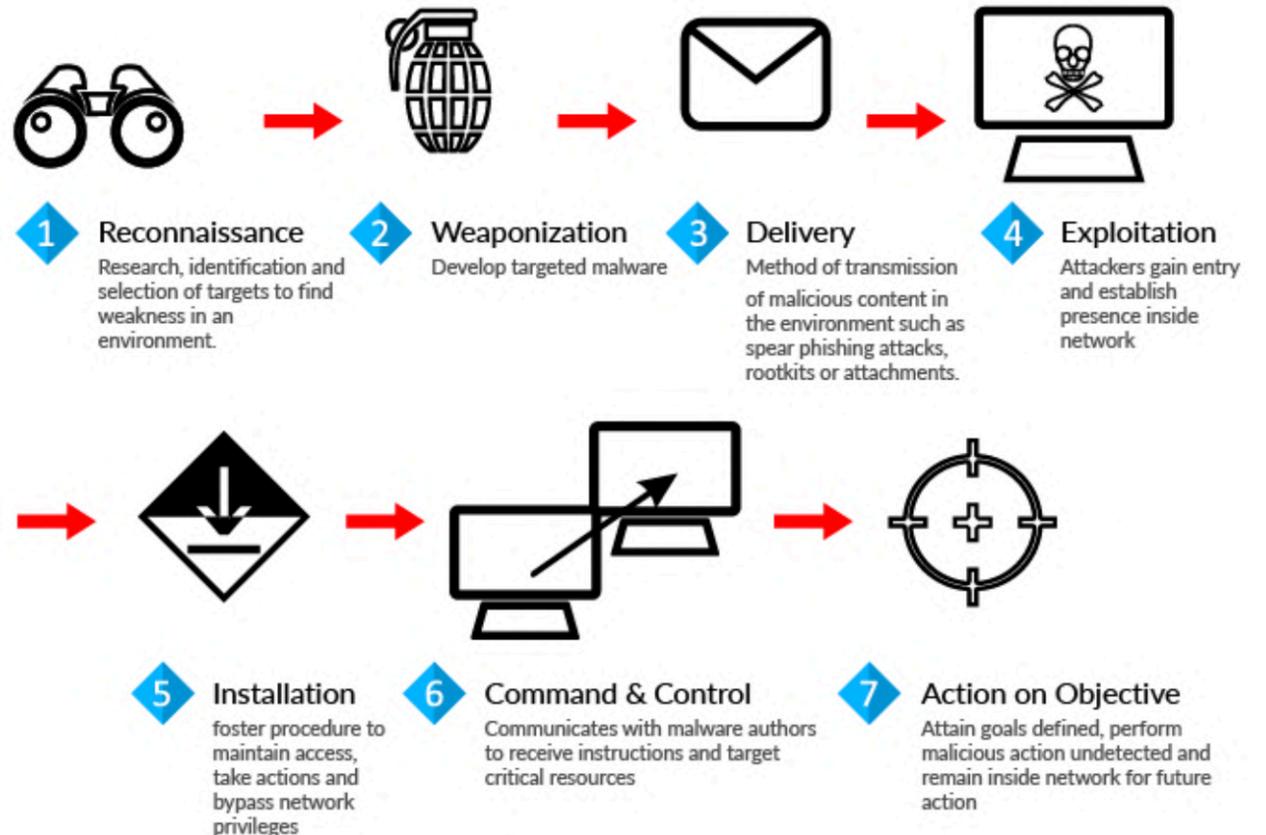
Case 4: Financial Analytics

- Expected outcome:
 - Rejected transactions
 - Real time alerts
 - Real time dashboard
 - Automated learning and improvement
 - Audit trails and analytics



Case 5: Detecting Cyber Attack

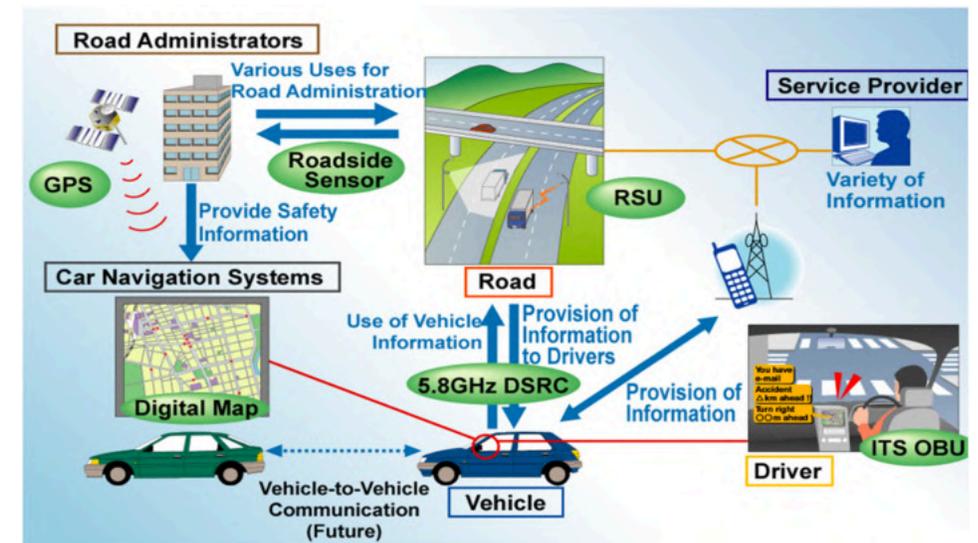
- Using big data technology to cope with such substantial amounts of data allows us to quickly perform the following analysis to assess alerts and determine if they are suspicious or benign.
 - Risk scoring
 - Graph analysis
 - Time series analysis
 - Determine causal relationships between entities



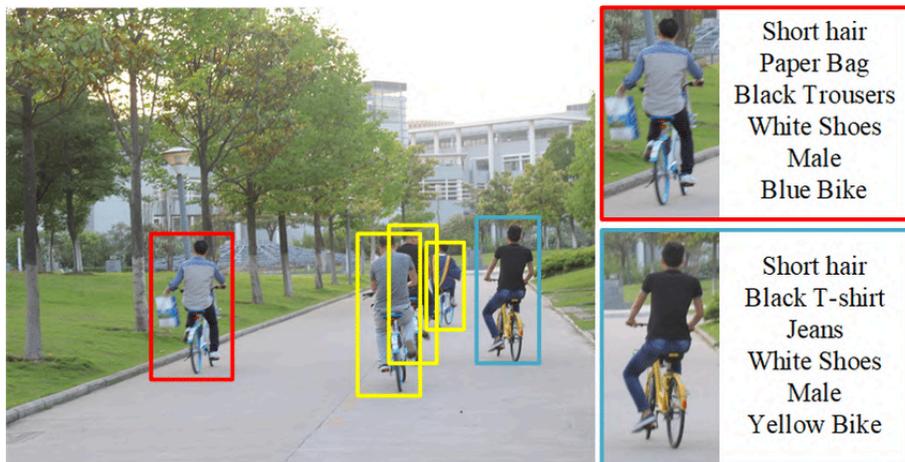
Cyber kill chain

Case 6: Smart Traffic Management

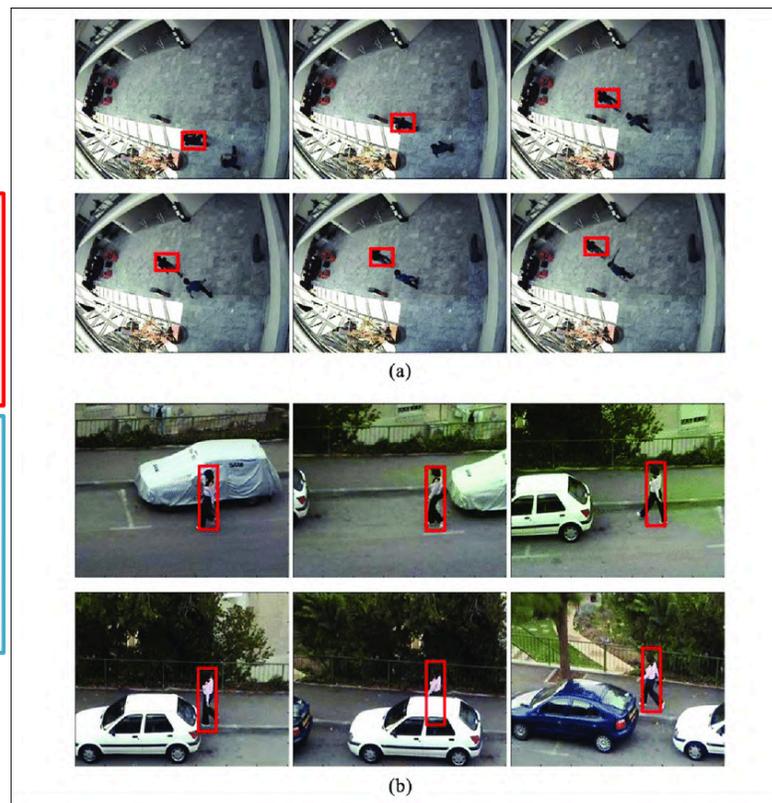
- Use sensors and traffic signals to monitor, control and respond to traffic conditions.
- The aim of smart traffic management systems:
 - Reduce day-to-day congestion by improving traffic flow
 - Prioritize traffic according to real-time changes in traffic conditions
 - Reduce pollution by limiting traffic jams
 - Prioritize buses entering intersections
 - Improve traffic incident response time



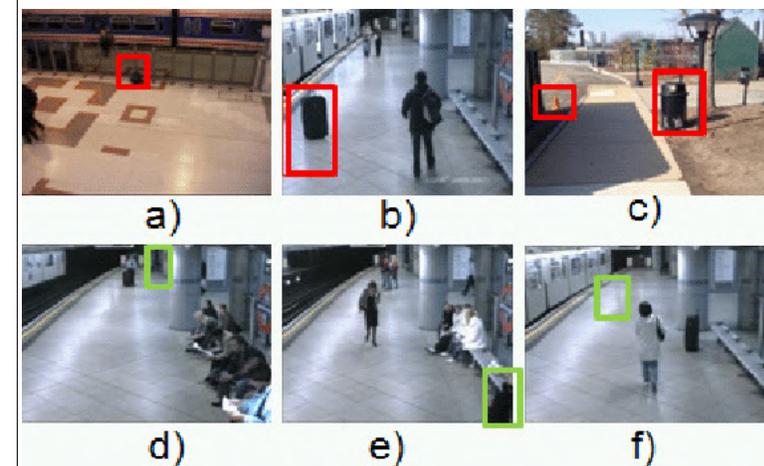
Case 7: Video Surveillance



Pedestrian attributes recognition

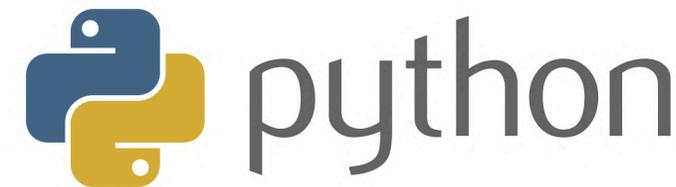


Pedestrian tracking



Abandoned object alert

Big Data Tools



SAS

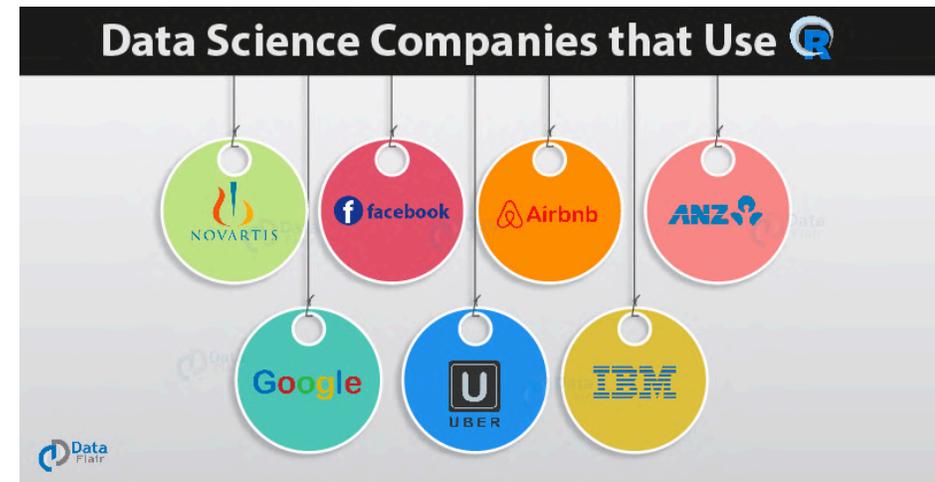
- No. 1 market leader in analytics.
 - The largest independent vendor in the business intelligence market.
 - The industry standard for Clinical Data Analysis.
- Integrated platform for end to end solutions.
 - SAS provides an integrated set of software products and services and integrated technologies for information management, advanced analytics and reporting.
- Business solutions across domains and industries.
 - Unmatched domain specific industry focused analytics solutions.



Used in
60,000+
companies in
over 135
countries

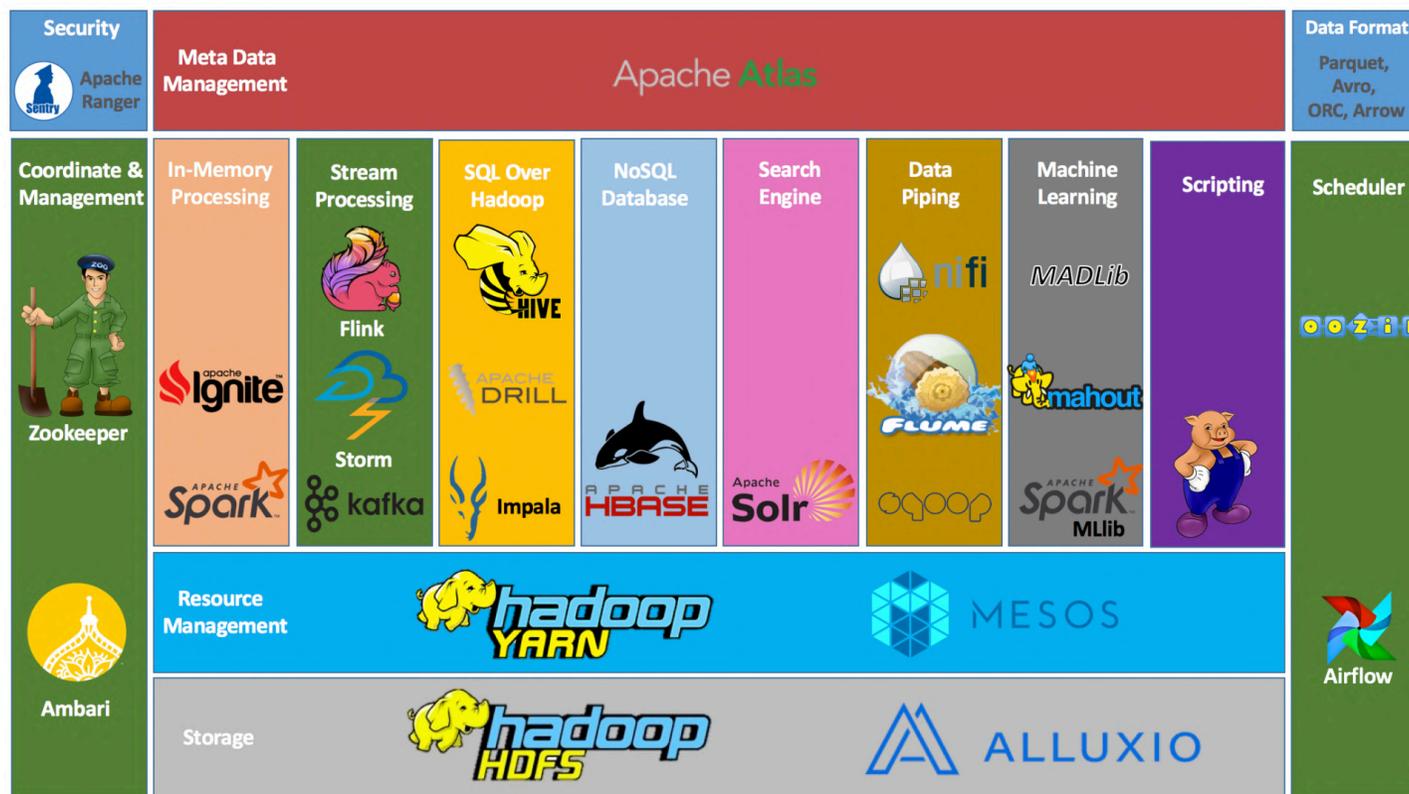
R

- R is a language and environment for statistical computing and graphics.
- R provides a wide variety of statistical and graphical techniques, and is highly extensible.



Hadoop

- Hadoop is the most popular big data ecosystem.
- Hadoop is highly scalable, that is designed to accommodate computation ranging from a single server to a cluster of thousands of machines.



Hadoop ecosystem

Python

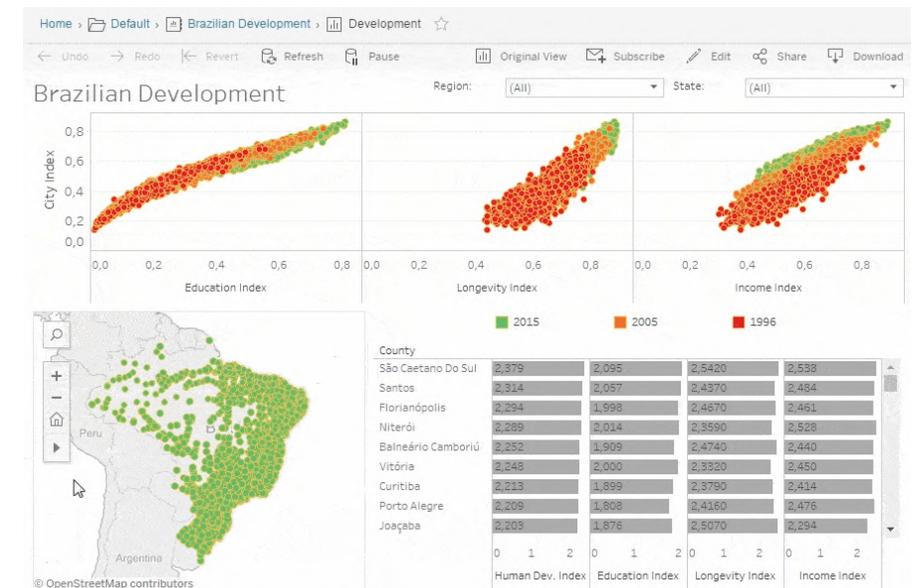
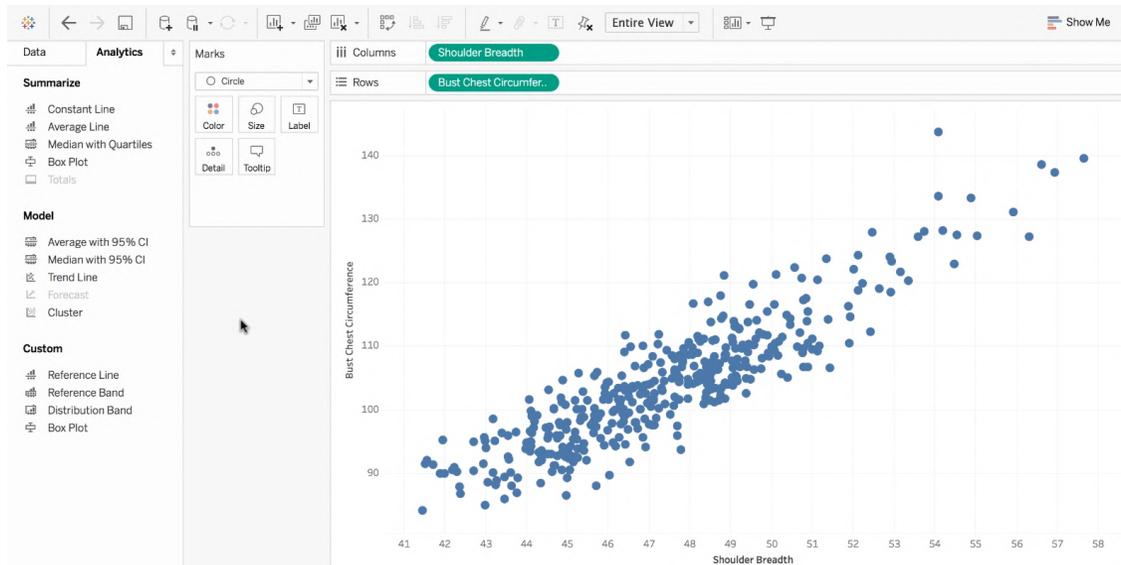
- Python is an interpreted, high-level, general-purpose programming language.
- One of the most popular programming language in recent years.
- Ten areas that uses Python most frequently:
 - Web Development
 - Game Development
 - Machine Learning and Artificial Intelligence
 - Data Science and Data Visualization
 - Desktop GUI



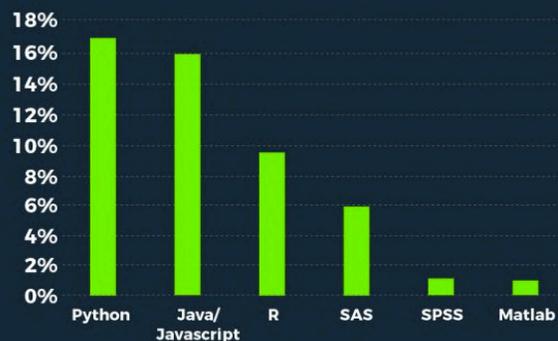
- Web Scraping Applications
- Business Applications
- Audio and Video Applications
- CAD Applications
- Embedded Applications

Tableau

- Tableau is a data visualization tool that is widely used for business intelligence.
- Create interactive graphs and charts in the form of dashboards and worksheets to gain business insights.



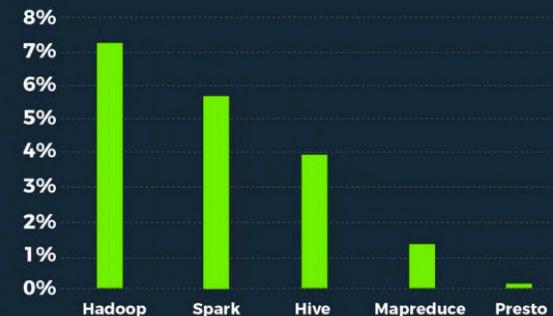
Skills in demand in 2019



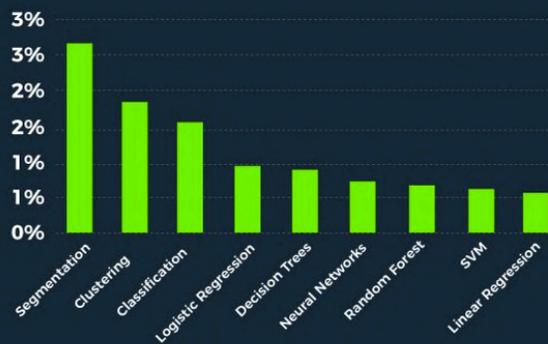
Programming language



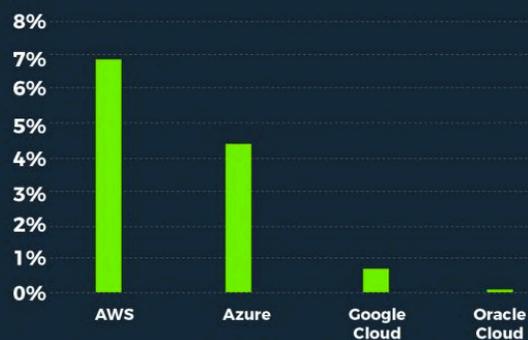
Data visualization



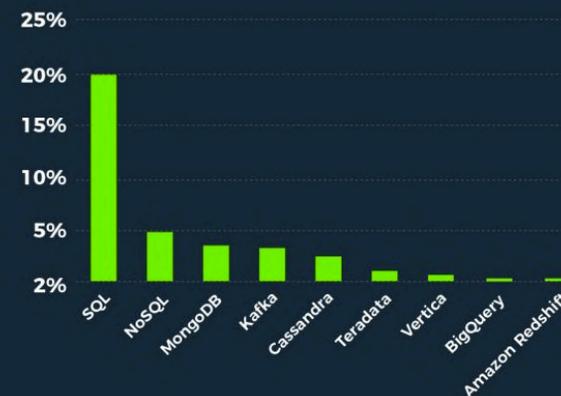
Distributed framework



ML algorithm



Cloud computing



DB skill

Thank you!

- Any question?
- Don't hesitate to send email to me for asking questions and discussion. 😊

Acknowledgement: Thankfully acknowledge slide contents shared by Prof. Yongxuan Lai