

INTRODUCTION TO BIG DATA ANALYTICS TOOLS

CKME 134 – BIG DATA ANALYTICS TOOLS
RYERSON UNIVERSITY
SPRING 2015

Instructor: Shaohua Zhang

General Course Information

□ Instructor

- Shaohua Zhang
- Ryerson shaohua.zhang@ryerson.ca
- Personal shaohua.zhang@live.com

□ GA

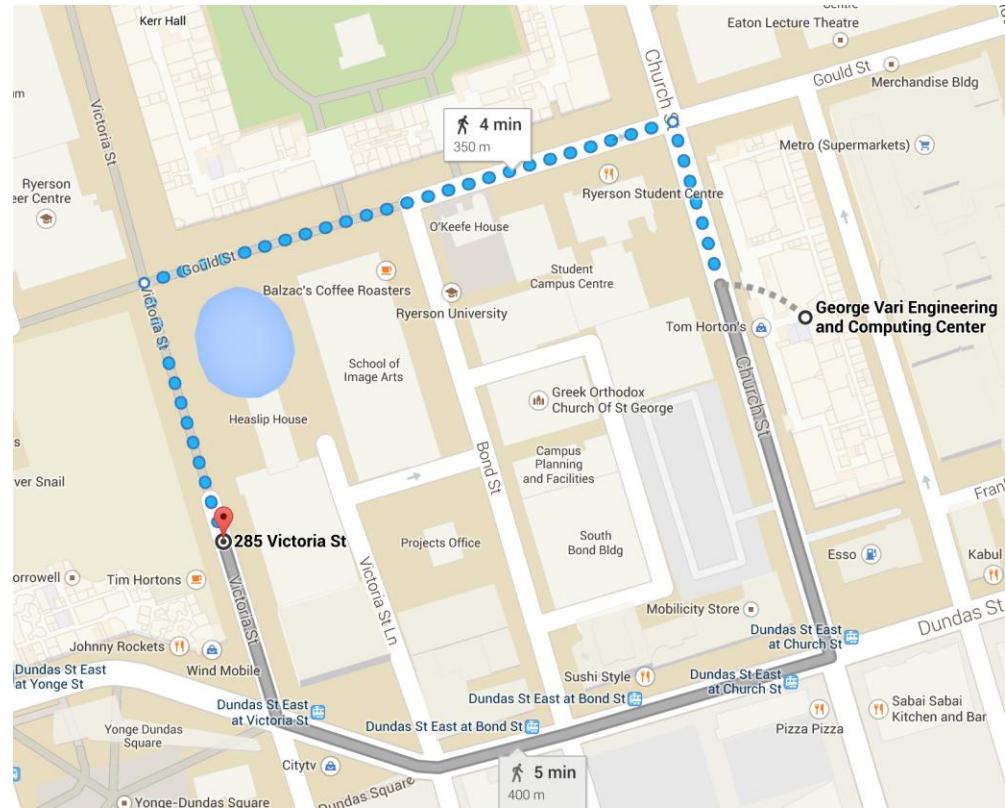
- Behjat Soltanifar
- behjat.soltanifar@ryerson.ca

□ Lectures

- 6:30~8:30
- ENGLG06

□ Lab

- 8:30~9:30
- 285 Victoria St (403/404)
 - Take the elevator to 4FL



About This Course

3

- **This course focuses on practicality**
 - It follows industry trends and job market trends
 - You'll be hands-on with several popular tools
 - You'll learn practical use cases and how to choose the right tools
 - You'll learn enough to be able to extend on your own after this course
- **This course teaches big data tools related to analytics**
 - Will focus less on Infrastructure, ETL and BI
- **It mainly focuses on batch processing tools**
 - Will introduce streaming processing
- **Lab work will be done via virtual machines**

Course Outline (*subject to change*)

1. **Intro to Big Data**
2. Distributed Computing and MapReduce
3. Hadoop Ecosystem
4. Intro to Hive
5. Pig
6. Advanced Pig
7. Hadoop Performance Optimization
8. Big Data Use Cases: Location Intelligence and Marketing Analytics
9. Big Data Use Cases: Recommendation Engine and Computational Advertising
10. Hadoop In Action: Building Data Pipelines
11. Beyond Hadoop: Spark
12. Beyond Hadoop: Real-Time Analytics

Course Project/Assignments

- Analysis of Geo-tagged Tweets
 - Twitter data analysis
 - Simple NLP
 - Location clustering
 - User segmentation
 - Census demographic

Exams

- No mid-term (probably)
- Final Exam
 - Hands-on lab test
 - Apply everything you've learned to complete a small real-life task
 - Labs and assignments are important

Survey

- **Self-introduction**
 - How many of you are already data analysts?
 - What background do you come from?
 - What kind of programming languages you know?
 - Size of data you worked with?
- **Why do you take this course?**
 - Curious about big data; want to gain knowledge
 - Starting career
 - Switching career

Lecture 1 - Outline

1. Big Data Introduction
2. Big Data Use Cases
3. Data Analytics Tooling
4. Big Data Job Market
5. Big Data Challenges



1. Intro to Big Data

Big Data Is Hot!



Data Scientist: *The Sexiest Job of the 21st Century*

McKinsey&Company

McKinsey Global Institute

June 2011

Big data: The next frontier for innovation, competition, and productivity

Big Data – Why Now?

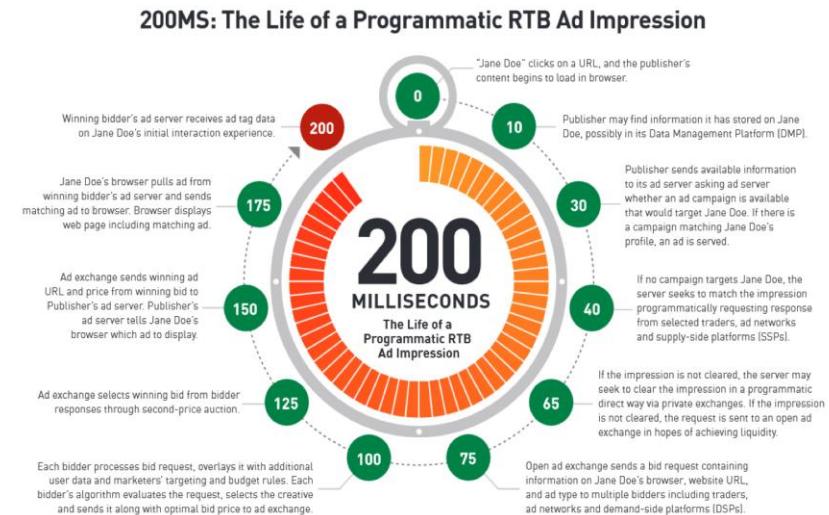
- Data at scale (Volume)
 - Since when was 1TB not big data any more :-)
 - Speed (Velocity)
 - Near Realtime response is key to the modern web/mobile experience
 - Data in many forms (Variety)
 - Structured
 - Unstructured
 - Location
 - Text
 - Image
 - Video
 - Semi-structured
 - Graph
- | | |
|--|---|
| <ul style="list-style-type: none">• Internet | <ul style="list-style-type: none">◦ 2.5 exabytes (2.5×10^{18}) per day – 2012◦ 2.3 zettabytes (2.3×10^{21}) per day - 2014 |
| <ul style="list-style-type: none">• Facebook | <ul style="list-style-type: none">◦ 500+ terabytes per day◦ 100+ petabytes in a single Hadoop cluster |
- “More data cross the internet every second than were stored in the entire internet just 20 years ago” - Big Data: The Management Review (HBR)
- 

Four V's of Big Data: <http://www.ibmbigdatahub.com/infographic/four-vs-big-data>

Big Data (wiki): http://en.wikipedia.org/wiki/Big_data

Big Data – Why Now?

- Data at scale (Volume)
 - Since when was 1TB not big data any more :-)
- Speed (Velocity)
 - Near Realtime response is key to the modern web/mobile experience
- Data in many forms (Variety)
 - Structured
 - Unstructured
 - Location
 - Text
 - Image
 - Video
 - Semi-structured
 - Graph



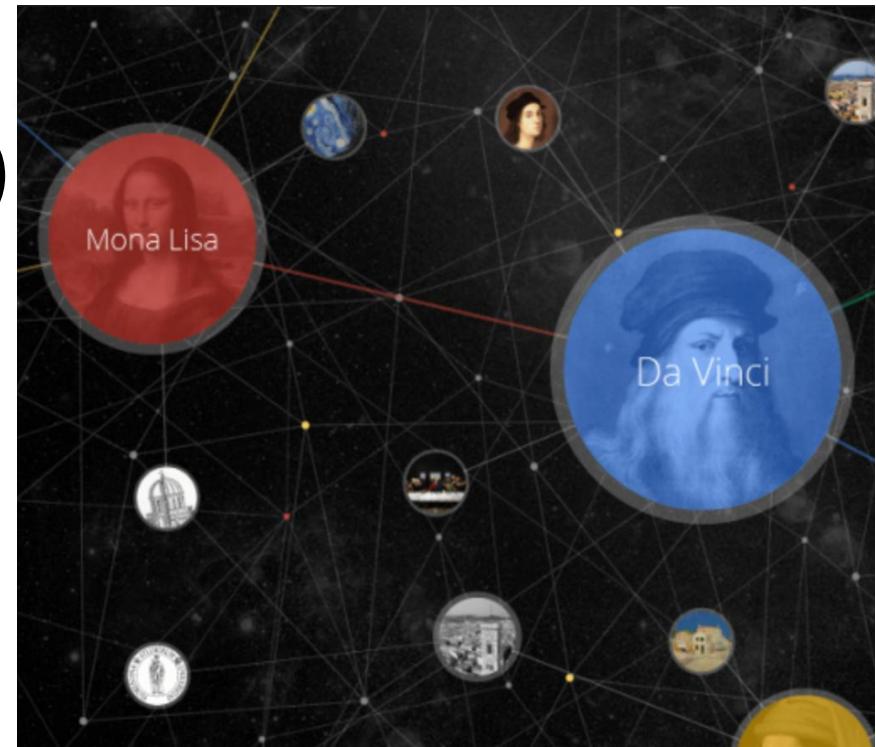
[video](#)

Four V's of Big Data: <http://www.ibmbigdatahub.com/infographic/four-vs-big-data>

200MS: The Life of a Programmatic RTB Ad Impression: <http://bit.ly/1iPqAlt>

Big Data – Why Now?

- Data at scale (Volume)
 - Since when was 1TB not big data any more :-(
- Speed (Velocity)
 - Near Realtime response is key to the modern web/mobile experience
- Data in many forms (Variety)
 - Structured
 - Unstructured
 - Location
 - Text
 - Image
 - Video
 - Semi-structured
 - Graph



Applications Driving the Need for Big Data

- Data-driven Applications
 - Location-based services
 - Social media apps
 - Image/voice recognition
 - Personal Assistant
 - Advertising
 - Internet of Things
- Internet Economies
 - Monetization needs - data is the new oil



Communities Driving the Need for Big Data

□ Hadoop vendors

- Cloudera
- Hortonworks (public)
- MapR

□ Traditional Vendors

- Oracle
- SAS
- IBM
- Revolution Analytics

□ Open Source Communities

- RHadoop
- RapidMiner
- NoSQL

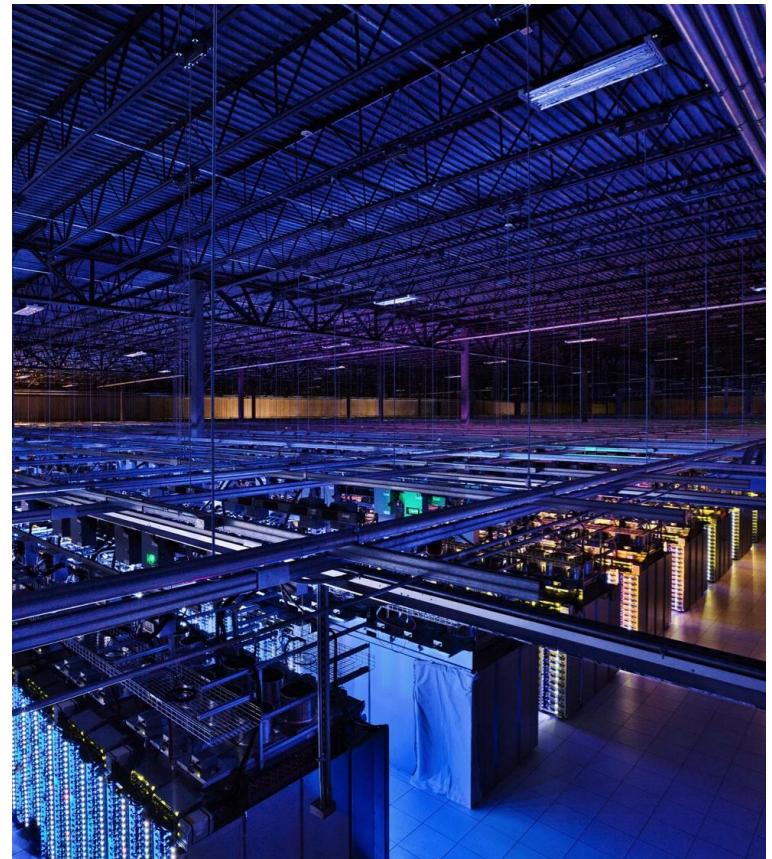
Big Data Made Possible

□ Hardware

- Big cluster of commodity machines at lower cost
 - Faster processor
 - Cheaper memory
 - Bigger hard drive space
 - Faster network bandwidth

□ Software

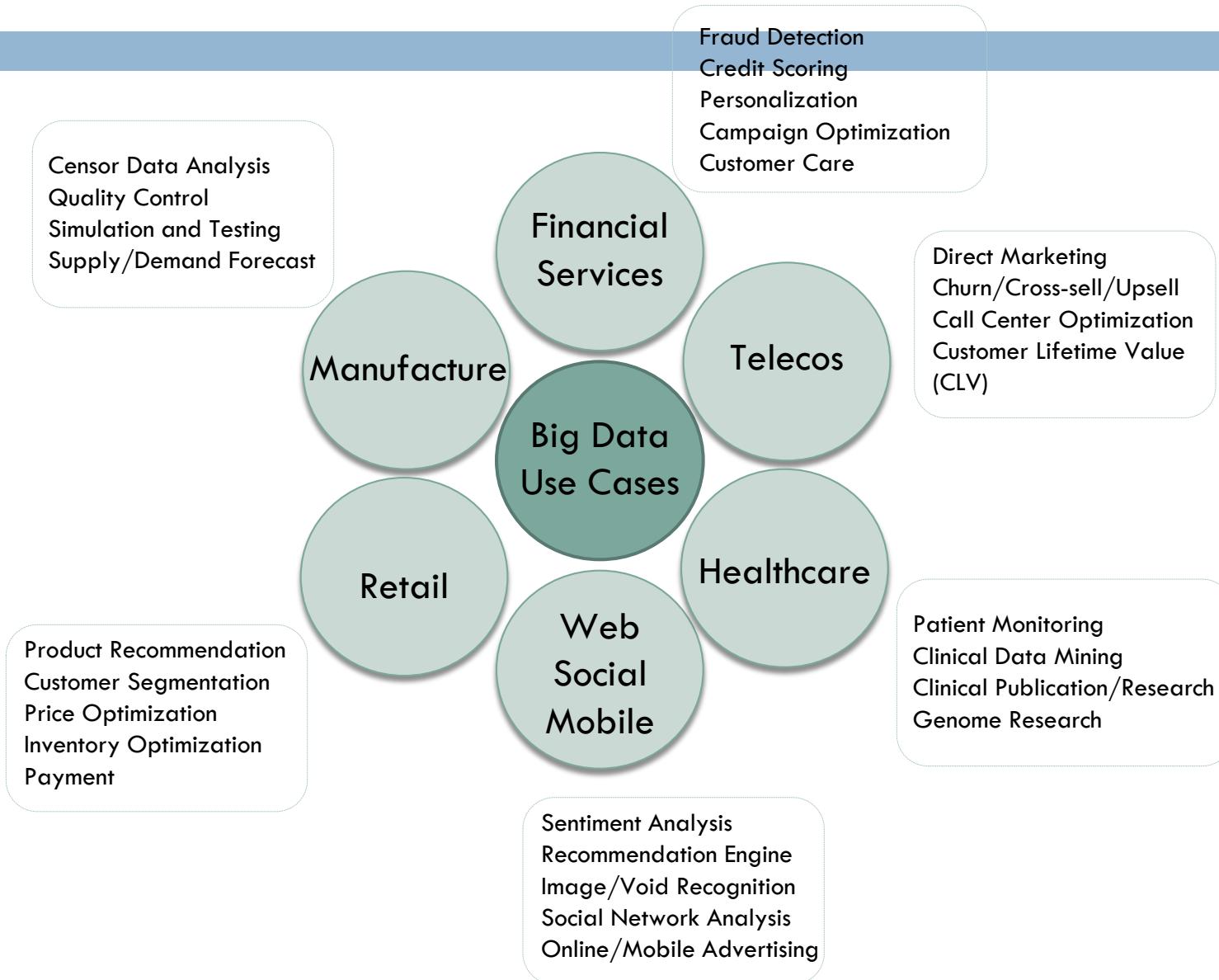
- Algorithms to allow parallel computing (map-reduce)





2. Big Data Use Cases

Big Data Use Cases



Big Data Use Case – Search & Media

□ Google

- Original MapReduce paper
2004
- Search & Advertising
- Image Recognition
- Google Voice
- Etc.

□ Yahoo!

- Hadoop 2005 (Doug Cutting)
- Page personalization
- Flickr Image Recognition
- Advertising



Recommended links

+79% clicks
vs. randomly selected

News Interests

+160% clicks
vs. one size fits all

Top Searches

+43% clicks
vs. editor selected

Google MapReduce Paper 2004: <http://bit.ly/1ADYtjC>

Google, Stanford build hybrid neural networks that can explain photos: <http://bit.ly/11owuZ2>

Big Data Use Case - Banking

□ Hadoop is the new backbone of American Express

- Recommender systems
- Graph algorithms
- Machine learning for Fraud and Marketing
- Data products
- Experiments

Website Personalization



Big Data Use Case – Social Media

□ Facebook (designed Hive, Giraph)

- New Feeds
- Friend recommendation
- Ads
- Graph Search

□ Twitter

- Follower recommendation
- Tweet search
- Timeline

□ LinkedIn

- PYMK
- Job recommendation

The screenshot shows a LinkedIn interface titled 'People You May Know'. At the top, there's a navigation bar with links for Home, Profile, Contacts, Groups, Jobs, Inbox, Companies, News, and More. Below the navigation is a search bar. The main content area is titled 'People You May Know' with a sub-instruction 'See people from different parts of your professional life'. It displays a grid of profile cards for various users, each with a thumbnail, name, title, location, and a 'Connect' button. The profiles include Brett Petersel (Director of Business Development at The Next Web), Wytze de Haan (Events at The Next Web), Becca Colbaugh (Vice President of Production & Operations at Voice123, LLC), Morgan Levy (Senior Artist Services Manager at Artist Arena), David Lecinski (Casting Advisory Manager at Voice123, LLC), Dan Rockwell (Program Manager of the Software Prototyping Center at The Ohio State University), Allison Grant (Press Coordinator & Researcher at Americans Elect), and Mitch Neff (Digital Marketing Manager at Cisco Systems).

Big Data Use Case – LBS

□ Foursquare

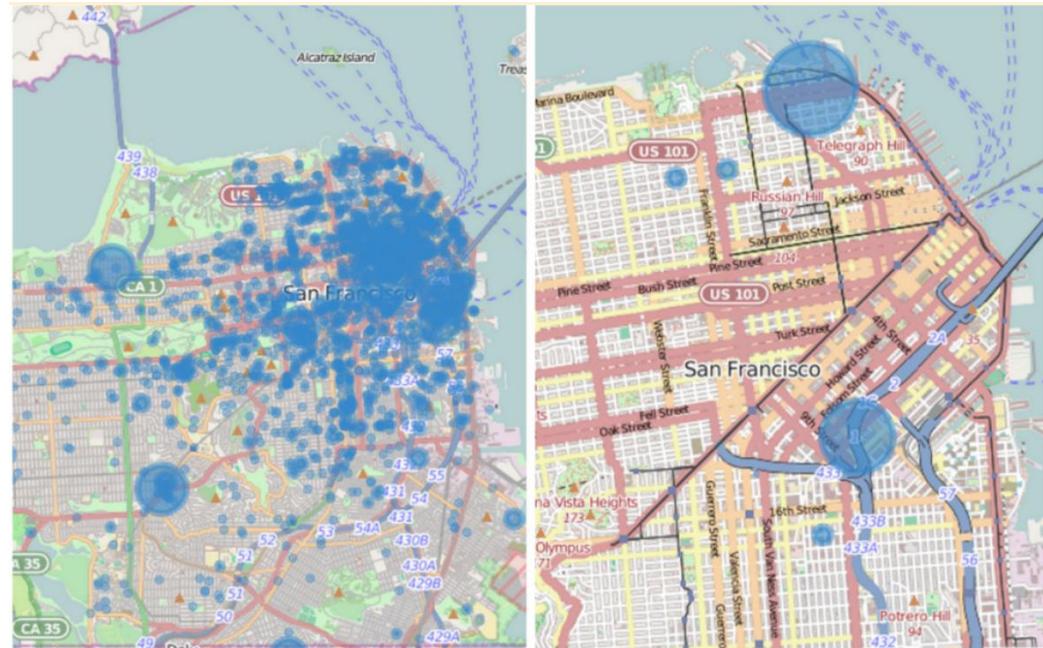
- Location Recommendation
- Local Search
- Location-based Social Network
- Location-based Interest

□ Yelp

- Sentiment Analysis
- POI Recommendation
- Text Classification
- Personalized Star Rating

□ Uber

- Trip Prediction
- Location-based User Segmentation
- Location-based Demographic Prediction



Big Data Use Case – E-Commerce

□ Amazon

- Product Recommendation
 - People who bought this also ...
- Fire Phone Image Recognition

□ Ebay

- Product Tagging
- User Taste/Interest Graph
- Fraud Detection
- Personalization



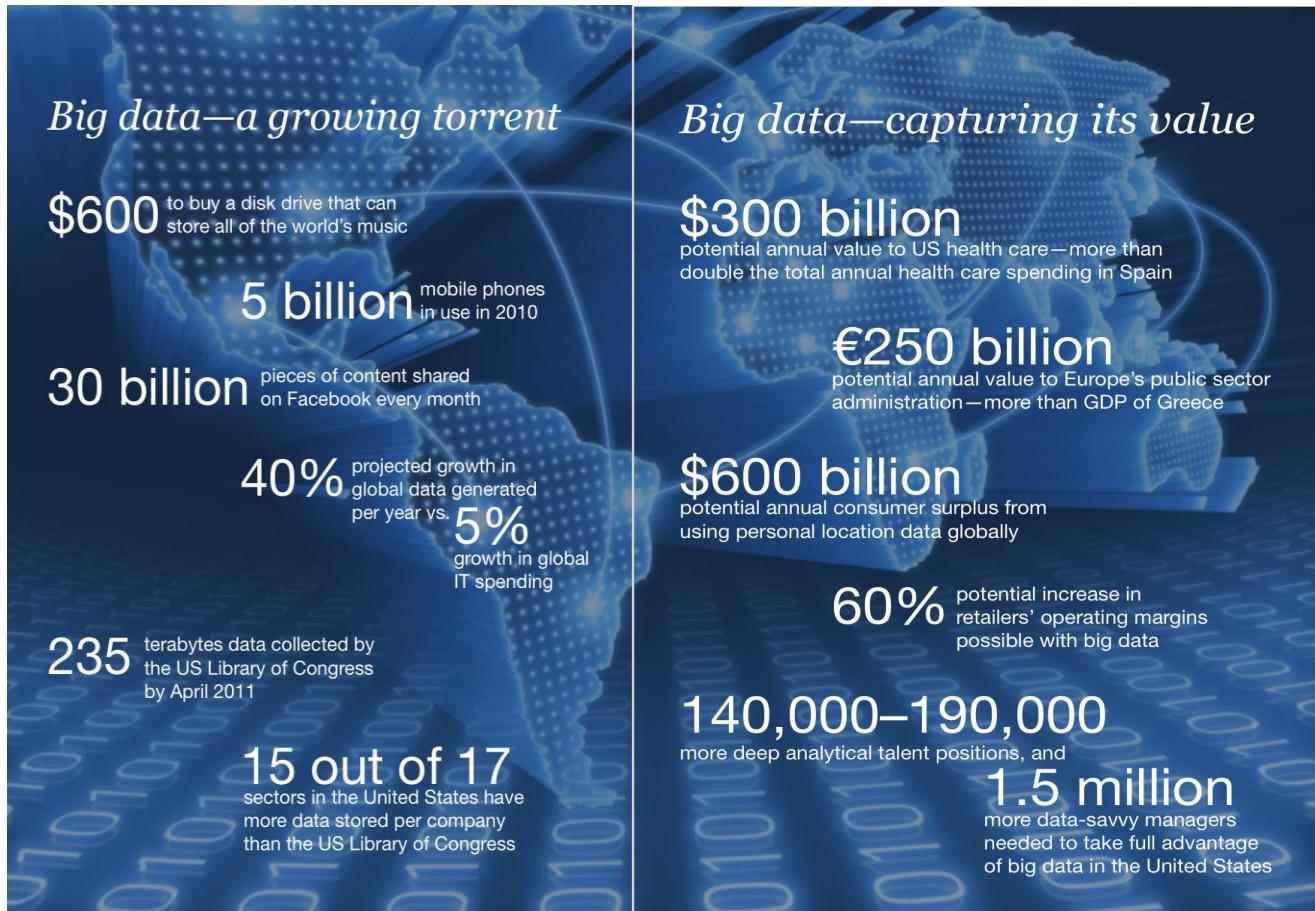


3. Big Data & Analytics Tooling

Big Data – Exciting Future

25

Great! Cool! Promising! Exciting!



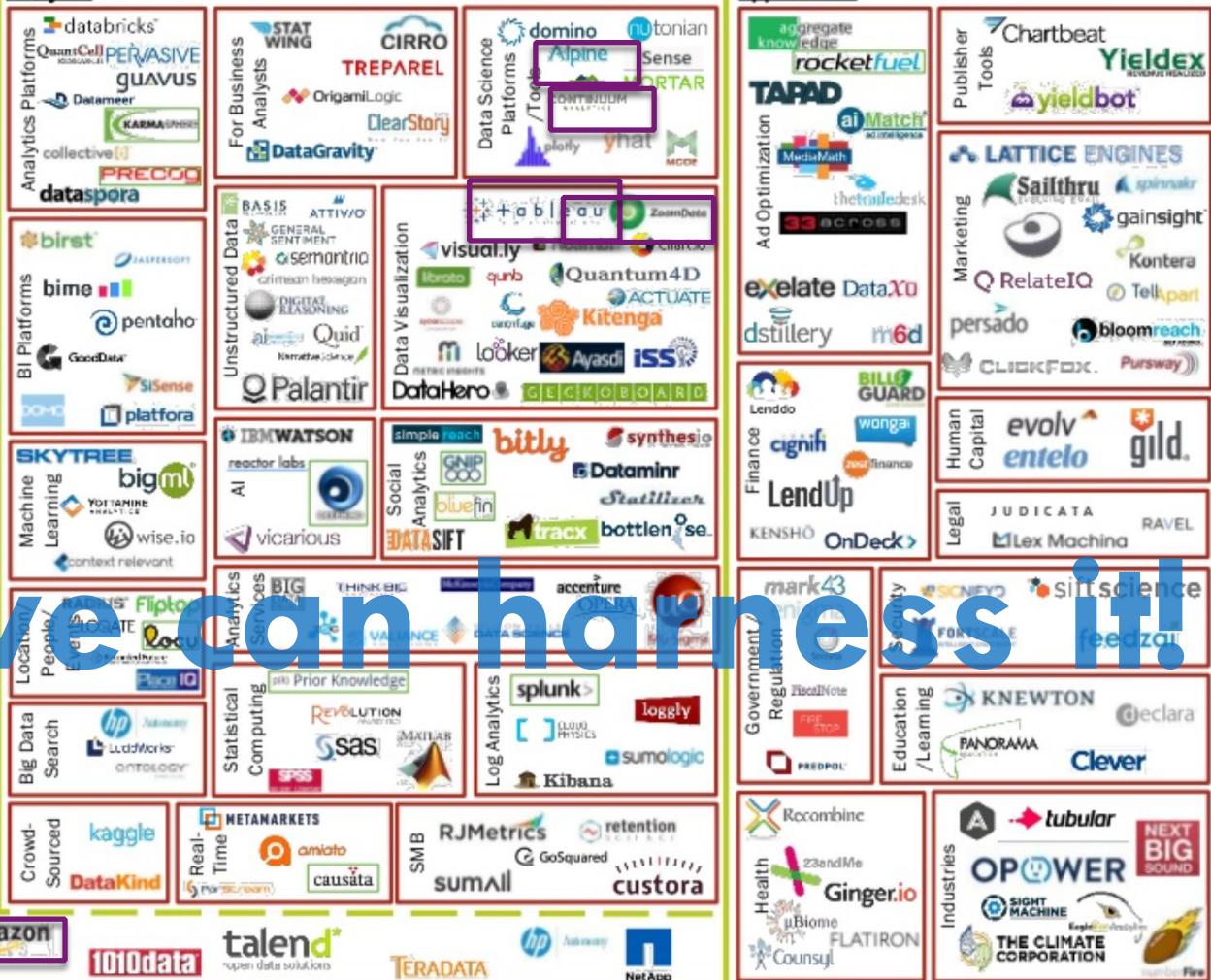
BIG DATA LANDSCAPE, VERSION 3.0

Exited: Acquisition or IPO

Infrastructure



Analytics



Applications

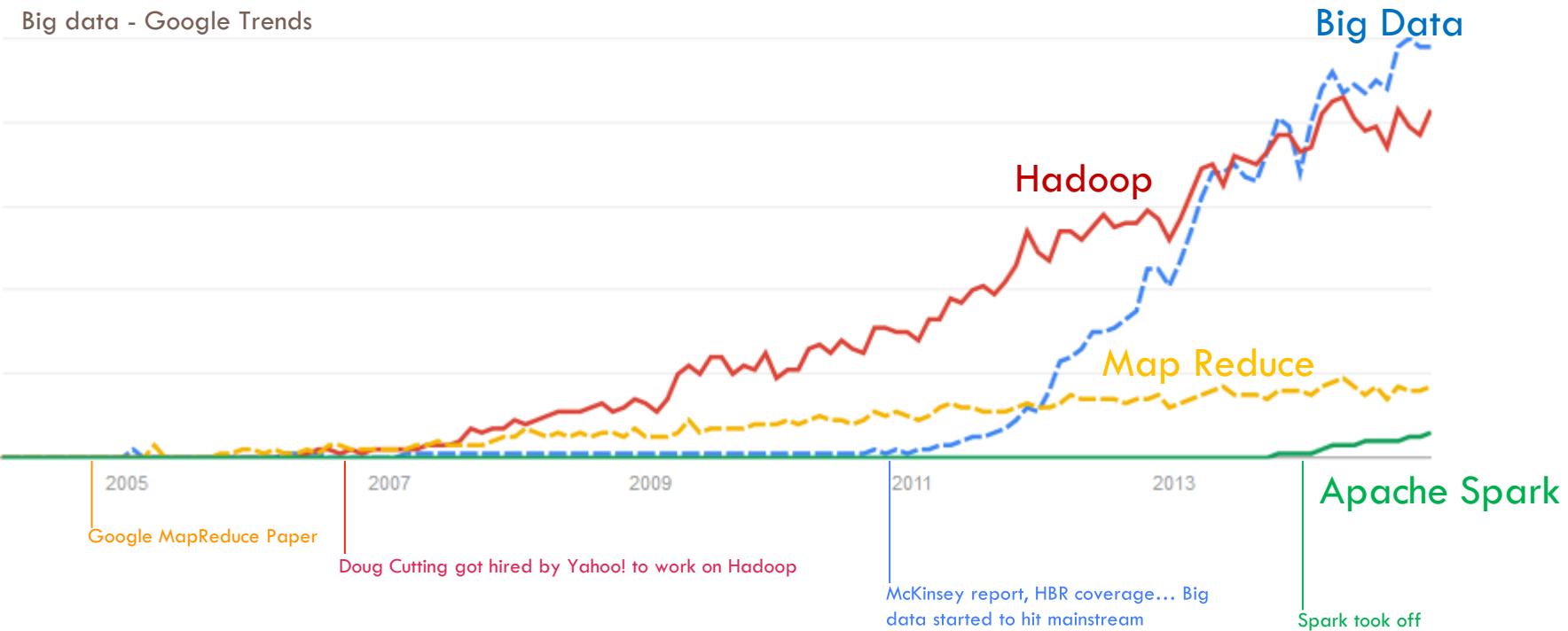


Data Sources



If only we can harness it!

Big Data History



Google MapReduce Paper: <http://research.google.com/archive/mapreduce.html>

Big data: The next frontier for innovation, competition, and productivity: <http://bit.ly/1pCOqom>

Big Data Landscape - Simplified

28

	Open Source	Commercial	Comments
Big Data Platform	Hadoop (MR, Pig, Hive etc.)	Cloudera, Hortonworks, MapR	<i>Hadoop is going mainstream</i>
	Spark	DataStax	<i>Spark is HOT! considered as next-generation big data platform</i>
		AWS	<i>Elastic MapReduce (EMR) and EC2 from AWS is most popular among startups.</i>
Machine Learning & Statistical Learning	Mahout		<i>Mahout was one of the earliest ML libraries for MapReduce. It is being revamped to take advantage of Spark currently</i>
	MLlib (Spark)		<i>MLlib is Spark's machine learning library. It's written in Scala and also provides Python and Java API</i>
	H2O		<i>H2O is the latest buzzing big data machine learning tool, backed by Oxdata. It works with a Hadoop cluster but also works on Standalone cluster. It has an amazing lineup of algorithms and even supports Deep Learning. The GUI-based predictive analytics suites works like a charm</i>
		SAS	<i>SAS integration with Hadoop will be very powerful. Imaging writing your data steps that runs procedures on hadoop</i>
		Revolution Analytics	<i>Commercial version of open source R. Enterprise-class big data analytics capability</i>
		Alpine	<i>World's first code-free in-cluster web analytics platform to analyze big data and hadoop</i>

Big Data Landscape - Simplified

29

	<i>Open Source</i>	<i>Commercial</i>	<i>Comments</i>
Graph Processing	Giraph		Graph processing framework on top of Hadoop. Used extensively at Facebook for large-scale graph algorithms
	GraphLab		Developed at CMU by Dr. Carlos and his team. Superior graph processing performance. Building the tools to make data scientists' lives easier. Great as a standalone graph processing and machine learning tool but won't fit well into the existing hadoop cluster
	GraphX (Spark)		Graph processing on Spark platform
Search	Solr		Open source search server based on Lucene Java library
	Elastic Search		Open source search and analytics engine
Stream Processing	Storm		Real-time stream processing framework developed at Twitter. Most popular streaming processing tool
	Spark Streaming		Streaming processing on Spark. Less mature than Storm at the moment but growing rapidly
Visualization	d3.js		Fantastic javascript library for visualization
		Tableau, Qlikview, Zoomdata	Popular visualization tools widely adopted
	Kibana		Log and time series data visualization tool from Elasticsearch

Big Data Analytics Tooling

□ Choosing the right tools – considerations

- Legacy systems
- Data scientists/engineers preference
- Scalability
- Availability
- Data manipulation capability
- Algorithms/libraries supported
- Operations (use in production)
- Cost
- Industry/vertical standards
 - Security
 - Support and service
- University programs

Big Data Analytics Tooling

Case Study

Data size

- 100TB

Formats

- Structured
- Unstructured

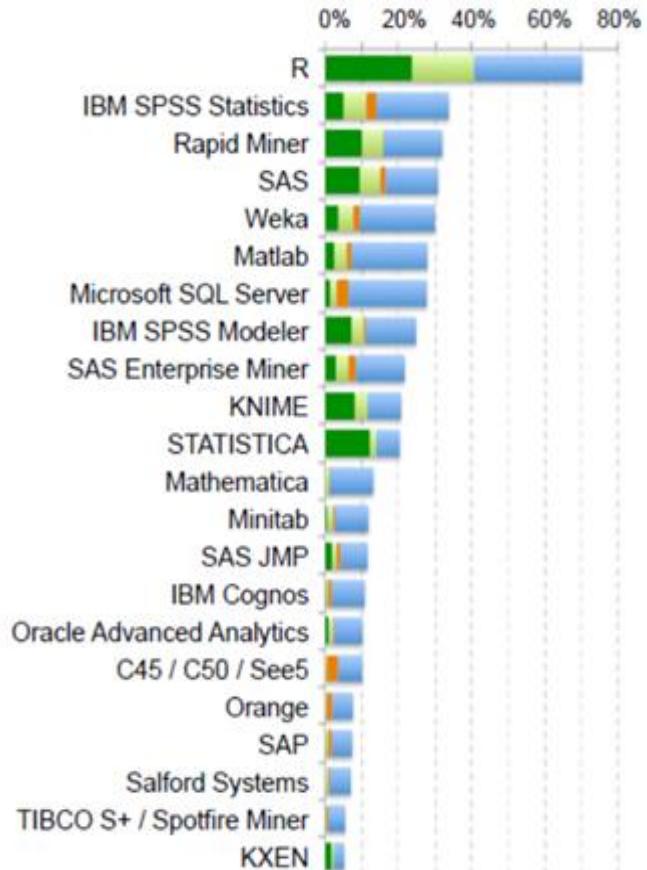
Tasks

- ETL
- Data analysis
- Machine Learning

Considerations	SAS	R	Python	Java/Hadoop	Pig/Hive	Spark	Oracle
Legacy System	●	●	●	●	●	●	●
Scalability	●	●	●	●	●	●	●
Availability	●	●	●	●	●	●	●
Ease of data manipulation	●	●	●	●	●	●	●
Algorithms/Libraries	●	●	●	●	●	●	●
Operations/Production Readiness	●	●	●	●	●	●	●
Cost	●	●	●	●	●	●	●
Support & Service	●	●	●	●	●	●	●
Business/Data analyst	●	●	●	●	●	●	●
Statistician	●	●	●	●	●	●	●
Data engineer	●	●	●	●	●	●	●
Data scientist	●	●	●	●	●	●	●

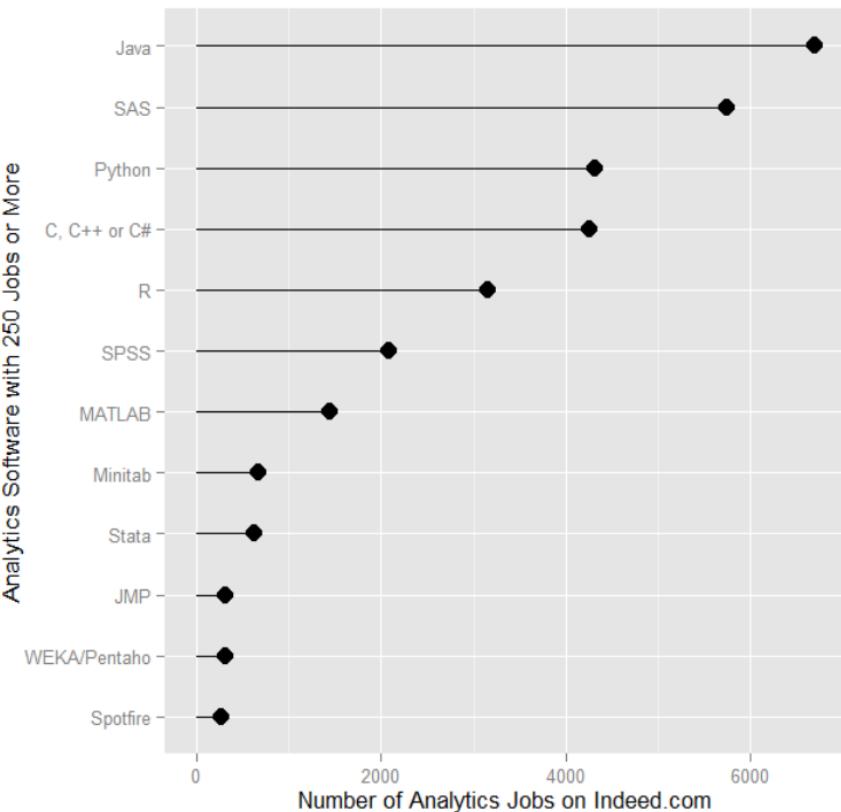
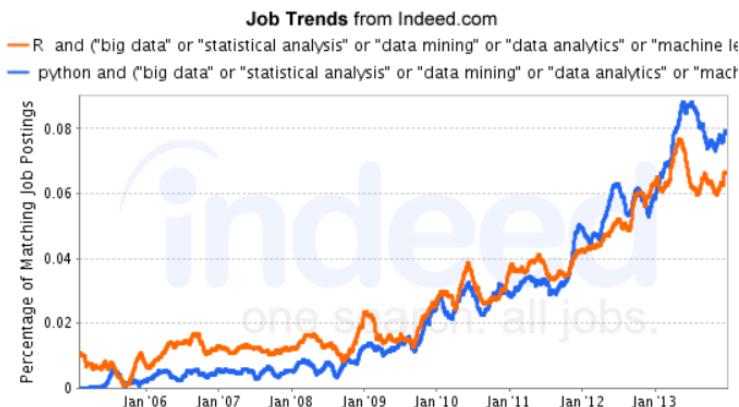
Analytics Tools Comparison

- Ranked by usage
- Overall trend is insightful
- But don't get too serious about the results
 - Should probably break down by industry and business functions
- Statistically significant?
 - Depends on how you interpret it!
- R gets all the love
- SAS is expensive, but still popular in big companies



Analytics Tools – Popular Tools

□ SAS vs. R vs. Python



The Popularity of Data Analysis Software 2012 <http://r4stats.com/articles/popularity/>

Indeed.com Job Trends: SAS vs R <http://www.indeed.com/jobtrends/R+statistics,+SAS+statistics.html>

Choosing The Right Tools - Previously

	SAS	R	Python
Prototyping	SAS Base, SAS EG	R (requires sampling)	Python (requires sampling)
Data Manipulation	SAS, Oracle SQL	R, Oracle	Python
Modeling	Enterprise Miner, SAS Base, SAS EG	R	Scikit-learn
Scoring	Enterprise Miner, SAS Base, PMML	R, PMML	Python

Choosing The Right “Big Data” Tools - Today

	Java	SAS	R	Python	Spark
Prototyping	Weka, Java	SAS Base, SAS EG	R	Python	Spark/R
Data Manipulation	Hadoop, Pig/Hive	SAS Connector for Hadoop	RHadoop	Hadoop Streaming Pig/Hive	Spark
Modeling	Weka, Mahout	Enterprise Miner, SAS Hadoop	RHadoop	Hadoop Streaming	Mllib, GraphX
Scoring	Hadoop, Mahout	Enterprise Miner, SAS Base, Hadoop PMML	RHadoop	Hadoop Streaming, Pig	Spark



4. Big Data Job Market

Job Market

□ Where are big data jobs?

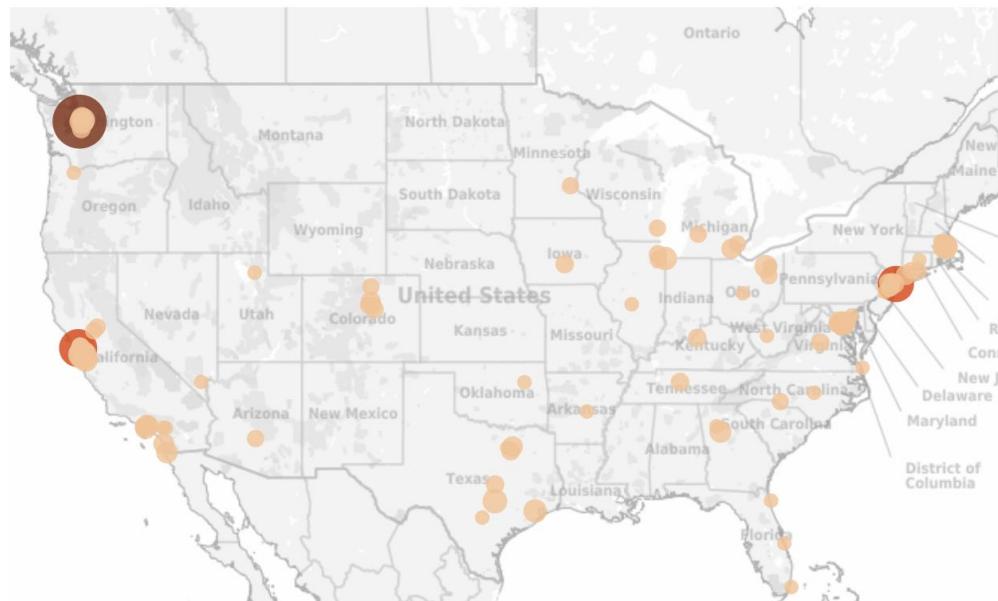
□ North America

- Silicon Valley, Seattle, NYC, Toronto

□ India/China



Total Jobs
822
Approximate Avg. Salary
\$93,205



Big data jobs around the nation: <http://www.tableausoftware.com/public/gallery/big-data-jobs>

Big Data Salary: <http://goo.gl/4et998>

Oreilly Media Data Science Salary Survey: <http://www.oreilly.com/data/free/files/stratasurvey.pdf>

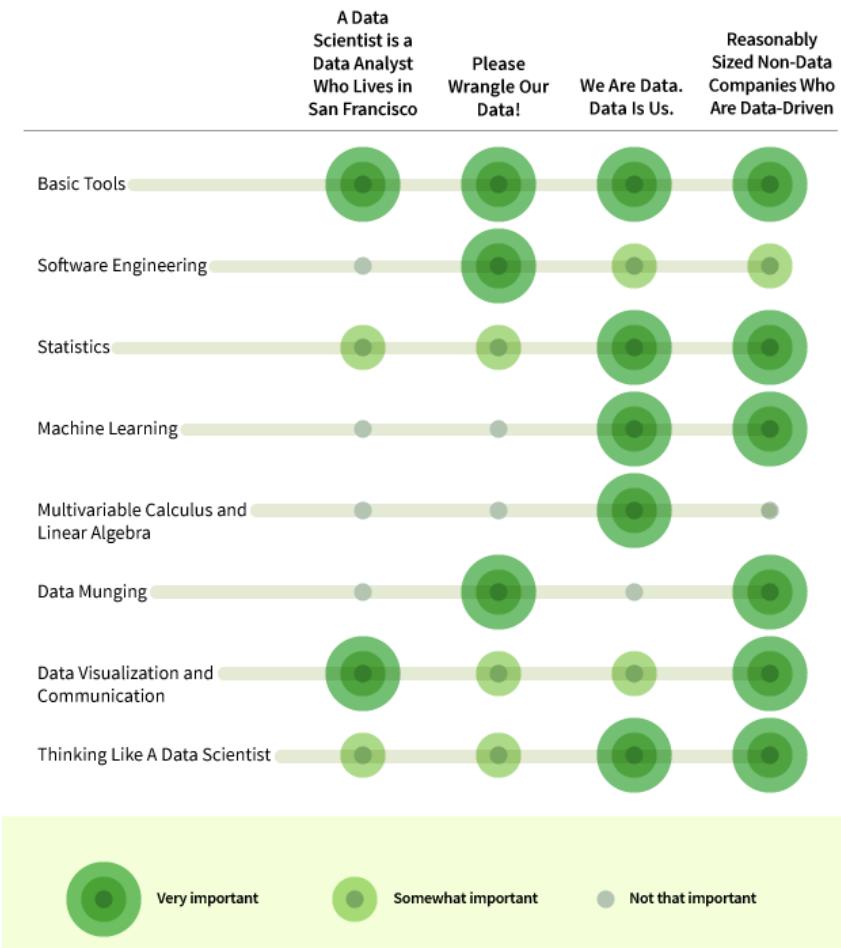
KDNuggets 2014 Analytics/Data Science Salary Poll: <http://goo.gl/VhO9IW>

Big Data/Analytics Jobs (Toronto)

- Banks
 - RBC, TD, CIBC, Scotiabank, AMEX, CapitalOne, ING Direct
- Telcommunications
 - Rogers, Telus, Bell, etc.
- Technology
 - BlackBerry, Huawei, CGI
- Manufacture/Services
 - GM, Canada Post, Workopolis
- Insurance
 - SunLife, Manulife
- Web/Mobile/Startup
 - Google, Mozilla
- Digital Media/Agencies
 - Globe and Mail, Kobo
- Consulting
 - Accenture, IBM, Deloitte, SAS
- Retail/e-commerce
 - Amazon, HR, Hudson Bay, Sears, Shoppers, Canadian Tire, Sobeys
- Pharmaceutical/Healthcare
 - Hospitals, Clinical Research Companies etc.

Data Organizations

- Data Analyst Who Lives in San Francisco
- Please Wrangle Our Data!
- We Are Data. Data Is Us
- Reasonably Sized Non-Data Companies Who Are Data-Driven



Reference:

Big Data Skill Sets (Analytics)

- **Data munging**
 - SAS, Oracle, Hadoop Pig/Hive
- **Database**
 - SQL
 - Hive
- **Programming**
 - Java, Scala, Python
- **Machine learning, predictive modeling at scale**
 - SAS Enterprise Miner
 - Mahout, Spark Mllib, H₂O
 - R, Python
- **Algorithms**
 - Regression, Trees, Clustering, SVM, KNN, etc.
- **Data Engineering**
 - linux, shell, cron, automation
- **Reporting**
 - Pivot tables
 - Dashboards
- **Visualization/presentation**
 - d3.js, Tableau, Qlikview, Excel

Job Functions/Titles

- Teams
 - Platform/Infrastructure
 - Data Science
 - Business Intelligence
 - Analytics (Marketing, Sales, Product etc.)
 - Startup - Jack of all trades
- ✖ Titles
 - Data Engineer (Database Developer, Data System Developer, etc.)
 - Data Scientist/Statistician
 - Big Data Researcher
 - Machine Learning Researcher/Scientist
 - Platform Engineer
 - Platform/Infrastructure Architect
 - Infrastructure Specialist
 - Data Analyst
 - Software Engineer
 - Database Administrator
 - ETL Developer/Specialist
 - Web developer
 - Data Ninja
 - Data Evangelist

Software
Engineer

Data
Engineer

Data
Scientist

Applied
Scientist

Research
Scientist



Salary

- Great compensation!
- But range is wide...

Big data jobs around the nation: <http://www.tableausoftware.com/public/gallery/big-data-jobs>

Big Data Salary: <http://goo.gl/4et998>

Oreilly Media Data Science Salary Survey: <http://www.oreilly.com/data/free/files/stratasurvey.pdf>

KDNuggets 2014 Analytics/Data Science Salary Poll: <http://goo.gl/VhO9IW>

Big Data Salaries: An Inside Look Data Scientist Salaries, Data Analyst Salaries, DBA Salaries, etc: <https://datajobs.com/big-data-salary>



4. Big Data Challenges

Getting Over The Big Data Hype

- “Big Data” is NOT about “big”
 - we’ve done it for many many years (costly)
 - isn’t it expected anyway with the growth of the Internet
 - it is a mentality
- You don’t need big data sometimes
- Having big data and Hadoop cluster doesn’t solve your problems... it may create new problems if you can’t harness it
 - you need the right tools, right talent, right management support and team structure
- Just a different tool or platform
 - How you do analytics haven’t fundamentally changed
- Bigger doesn’t mean better
 - Big data vs small data

Big Data – The Challenges

- Reality is that Hadoop is still hard to use (usability for business analysts)
 - Requires low-level Map Reduce programming to achieve sophisticated task
 - Mostly command line, GUI is not user friendly (improving)
- SQL-on-hadoop not delivering the promise yet
 - The SQL vs. NoSQL war
 - NewSQL (Google's F1 paper)
- Rapid growth causes confusions
 - Emerging stack such as Spark
 - Uncertainty
 - Vendors and confusions

Big Data – The Challenges

46

- Most companies are still in very early stage, leveraging hadoop for data storage and ETL, not really taking the full advantage of the stacks
- Building data pipelines is hard
 - pipelines are the glues
 - different platforms/tools cause frictions
 - Hadoop stack works
 - Spark is the challenger
- Talent gap
 - High quality data scientists/engineers hard to find
 - Unicorns are rare

Lecture 1 - Summary

- Data science/analytics is a competitive market, you need to master a set of new tools to stay competitive
- Data size is growing exponentially. You need to choose the right tools for your analytics needs
- Tools you'll learn in this course
 - Hadoop – basic MapReduce concept
 - Pig/Hive large-scale data processing
 - Building automated data pipelines
 - Apache Spark Introduction
- Use cases you'll learn in this course
 - Location analytics
 - Marketing analytics
 - Recommendation engine
 - Computational advertising
 - Real-time analytics

Big Data Analytics Resources

48

- Blogs/Talks
 - Datasciencecentral
 - DataTau
 - DataScience Weekly
 - Meetups (HakkaLab) – youtube/slideshare
 - SF Machine Learning
 - Engineering blogs - FB, Yahoo, Twitter, 4SQ etc.

Big Data Analytics Resources

49

□ Conference/Workshop

- Oreilly Strata/Hadoop World
- Hadoop Summit
- Cassandra Summit
- H2O World
- Solr Revolution
- GraphLab Conference
- Qcon
- MLconf
- PAW – Predictive Analytics World
- SAS Conference

Lab & Assignments

□ Lab Computer

- Username: datastudent
- Password: datastudent

□ Lab 1

- Virtualbox
- Hadoop setup

□ Assignment 1

- Pick a big data job title that fits your interests
- Find 3 companies (U.S. or Canada) that offer the job position you identified
- Read the job postings and summarize the skillsets required by the position
 - Send to Instructor's email: shaohua.zhang@ryerson.ca by putting "[Assignment1]" in the title



Lab 1 – Environment Setup

Lab 1 – Install VirtualBox

□ Installation of VirtualBox

- <https://www.virtualbox.org/wiki/Downloads>

Our Lab Director Shannon has already downloaded and installed VirtualBox for us 😊

Lab 1 – Install Putty

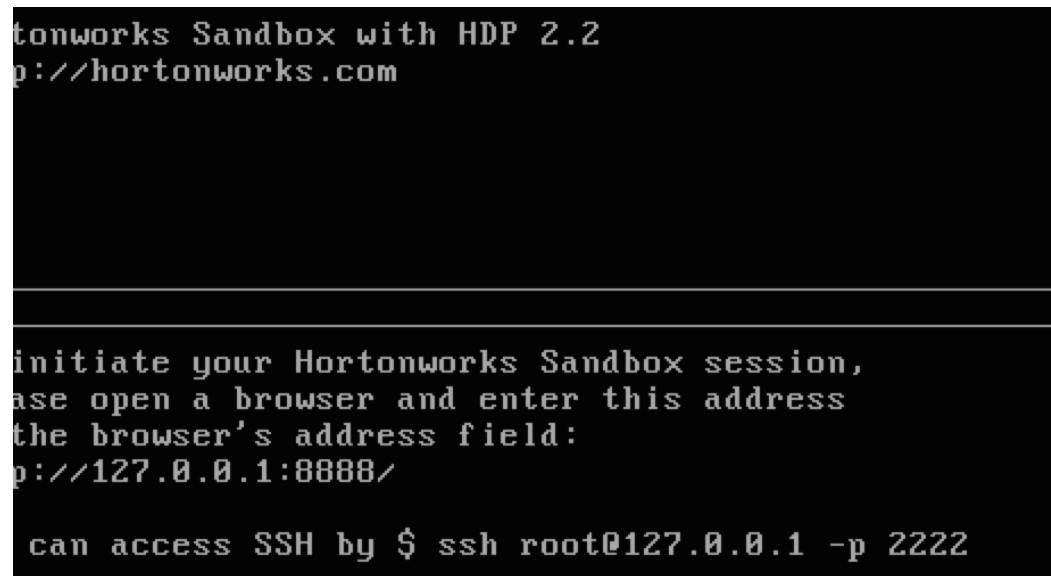
- Windows users need to install Putty in order to use SSH to access hadoop
 - <http://www.chiark.greenend.org.uk/~sgtatham/putty/>

Our Lab Director Shannon has already downloaded and installed Putty for us ☺

Lab 1 – Install Hortonworks HDP Sandbox

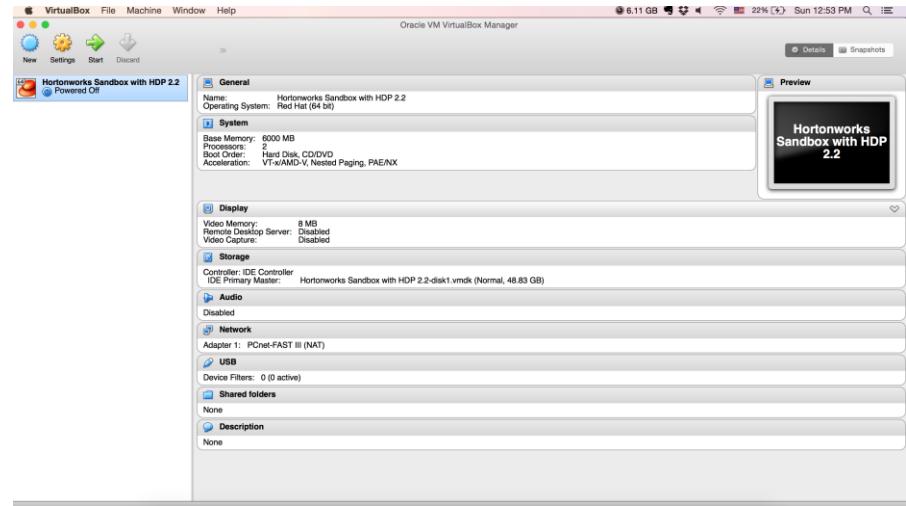
- Download HDP hadoop vm
image: <http://hortonworks.com/products/hortonworks-sandbox/>
- Tutorial for Windows: <http://hortonworks.com/wp-content/uploads/unversioned/pdfs/InstallingHortonworksSandbox2onWindowsusingVB.pdf>

- The HDP VM Sandbox has been downloaded to C drive on windows.
- After installation and importing the instance, you start the vm and will see the screen on the right



Lab 1 – Log On and Test Hive

1. Open VirtualBox
2. Start the HDP 2.2 VM
3. Open Putty and connect to 127.0.0.1 port 2222
 - Username: root
 - Password: hadoop
4. Test Hive
 - \$ hive



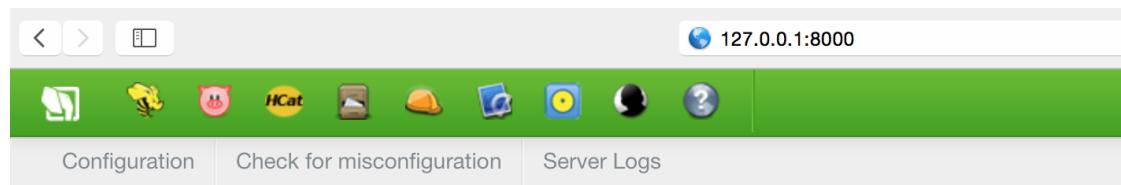
```
Hortonworks Sandbox with HDP 2.2
http://hortonworks.com

initiate your Hortonworks Sandbox session,
please open a browser and enter this address
in the browser's address field:
http://127.0.0.1:8888/
You can access SSH by $ ssh root@127.0.0.1 -p 2222
```

Lab 1 – Access Hadoop via Browser Using Hue

□ Access Hue via Browser

- <http://127.0.0.1:8000>



Hortonworks Sandbox with HDP 2.2

The logo features three stylized green elephants standing in a row, with the word 'Hortonworks' in a bold, sans-serif font below them. At the bottom of the box, there is a green link labeled 'Leave Feedback'.

Component	Version
Hue	2.6.1-2041
HDP	2.2.0
Hadoop	2.6.0
Pig	0.14.0
Hive-Hcatalog	0.14.0
Oozie	4.1.0
Ambari	1.7-169

At the bottom right of the table, there is a button labeled 'Enable'.