

SE256:Jan16 (2:1) L1:Introduction Scalable Systems for Data Science

Yogesh Simmhan Partha Talukdar



©Yogesh Simmhan & Partha Talukdar, 2016 This work is licensed under a <u>Creative Commons Attribution 4.0 International License</u> Copyright for external content used with attribution is retained by their original author



About This Course

www.serc.iisc.in/courses/se256/



- Instructors: Yogesh Simmhan, Partha Talukdar
- Teaching Assistant: Ravikant Dindokar



What is this course about?

- Fundamental systems aspects of big data platforms,
- How these platforms can be used to build large-scale data intensive applications.
- Why Big Data platforms are necessary?
- How they are designed? What are the programming abstractions (e.g. MapReduce) that are used to compose data science applications?
- How the programming models are translated to scalable runtime execution on clusters and Clouds (e.g. Hadoop)?
- How do you design algorithms for analyzing large datasets?
- How do you map them to Big Data platforms?

What is this course about?

- Platforms used for developing applications over large columnar/tuple based data,
 - Map Reduce/Apache Hadoop,
 - Streaming data like Apache Storm and Spark Streaming,
 - Graph/linked data like Apache Giraph
- Declarative tools
 - NoSQL databases and TensorFlow.
- Scale up of data analysis techniques
 - clustering, collaborative filtering, frequent itemset mining, classification, graph analytics,
 - Apply them over large datasets
- Why do platforms scale, and how can you use them to scale data science applications?



Pre-requisites

 Data Structures, Programming and Algorithm concepts. Programming experience required, preferably in Java.

Other Data Science Courses

- Foundations of Data Science
- Data Analytics
- Data Analysis and Visualization
- Introduction to Cloud Computing, Parallel Programming
- Topics in Web-scale Knowledge Harvesting, Data Mining, Deep Learning, Machine Learning, Artificial Intelligence, Probability & Statistics
- Bioinformatics, Chemoinformatics, Video Analytics

Lecture Schedule

- Wed 2-4PM, SE202
- Big Data & Platform Design Goals (YS)
 - Big Data & other computing platforms
 - Runtime Platforms: Hadoop, HDFS.
- Data-intensive algorithms and applications (PPT)
 - Analytics: clustering, collaborative filtering, frequent itemset mining, Classification, graph analytics, etc.
- Programming models data streams and graphs (YS)
 - Dataflow, BSP, Pregel. Runtime Platforms: Apache Storm, Apache Giraph, Apache Spark
- Big Tools: NoSQL databases, SparkQL, TensorFlow, etc (YS/PPT)

Tutorials

Time TBD with Ravikant, Students

- Check room availability
- 330-430PM Mondays
- Programming aspects
- Assignment guides
- Installation issues
- Bring laptops if possible

Text Book

 Select chapters from <u>Mining of Massive</u> <u>Datasets</u>, Jure Leskovec, Anand Rajaraman and Jeff Ullman, 2nd Edition (v2.1), 2014.



Assignments & Exams (tentative)

- Assignment 0: Jan 12
- Assignment 1: Jan 26
- Assignment 2: Feb 23
- Midterm: Mar 2
- Assignment 3: Mar 15
- Project Topic Finalized: Mar 16
- Project Midterm Report: Apr 5
- Project Final Report: Apr 26
- Final Exam: TBD

Interaction

- http://www.serc.iisc.in/courses/se256/
- se256.jan16@mailman.serc.iisc.in
- Class/lab discussions

- Students must uphold <u>IISc's Academic Integrity</u> <u>guidelines</u>. While these are common sense, it is helpful to review them since failure to follow them will lead to sanctions and penalties.
- This includes a reduced or failing grade in the course. Severe cases of academic violations will be reported to the Institute and may lead to an expulsion.

http://www.iisc.ernet.in/students-corner/existingstudents-academicintegrity.php

- Learning takes place both within and outside the class. Hence, discussions between students and reference to online material is <u>encouraged</u> as part of the course to achieve the intended learning objectives.
- However, while you may learn from any valid source, you must form your own ideas and complete problems and assignments by yourself.
- All works submitted by the student as part of their academic assessment must be their own!

- Plagiarism: Verbatim reproduction of material from external sources (web pages, books, papers, etc.) is not acceptable.
- If you are *paraphrasing* external content (or even your own prior work) or were otherwise *influenced* by them while completing your assignments, projects or exams, **you must clearly acknowledge them**.
- When in doubt, add a citation!

- Cheating: While you may discuss lecture topics and broad outlines of homework problems and projects with others, you cannot collaborate in completing the assignments, copy someone else's solution or falsify results.
- You <u>cannot use</u> notes or unauthorized resources during exams, or copy from others.
- The narrow exception to collaboration is between team-mates when competing the project, and even there, the contribution of each team member for each project assignment should be <u>clearly documented</u>.

- Classroom Behaviour: Ensure that the course atmosphere, both in the class, outside and on the online forum, is conducive for learning.
- Participate in discussions but do not dominate or be abusive. There are no "stupid" questions
 Be considerate of your fellow students and avoid disruptive behaviour.



Big Data & Why is it Important



What is Big Data?

Data whose characteristics exceeds the capabilities of conventional algorithms, systems and techniques to derive useful value. **CDS.IISc.ac.in** | **Department of Computational and Data Sciences**





06-Jan-16

Wordle of "Thought Leaders'" definition of Big Data, © Jennifer Dutcher, 2014 https://datascience.berkeley.edu/what-is-big-data/



So, where does Big Data come from?

World Wide Web & Search engines



World Wide Web & Search engines





Social Networks & Micro-blogs



Facebook monthly users

06-Jan-16

http://www.theguardian.com/news/datablog/2014/feb/04/facebook-innumbers-statistics

Social Networks & Micro-blogs



Youtube





The Netflix Prize (2006-2009)

- "An open competition for the best collaborative filtering algorithm to predict user ratings for films, based on previous ratings"
- Training set: 99,072,112 ratings
 - <user, movie, date of grade, grade>
- Qualifying set: 2,817,131 ratings
 - <user, movie, date of grade>
- Siri, Cortana, Now, Echo

CDS.IISc.ac.in | **Department of Computational and Data Sciences**

Enterprises & Government

Online retail & eCommerce



Source: CRISIL Research



http://blogs.ft.com/beyond-brics/2014/02/28/online-retail-in-india-learning-to-evolve/

http://www.peridotcapital.com/2014/04/amazon-sales-growth-projections-for-next-two-years-appear-overly-optimistic.html

06-Jan-16

CDS.IISc.ac.in | **Department of Computational and Data Sciences**

Enterprises & Government

Online retail & eCommerce



Source: CRISIL Research



http://blogs.ft.com/beyond-brics/2014/02/28/online-retail-in-india-learning-to-evolve/

http://www.peridotcapital.com/2014/04/amazon-sales-growth-projections-for-next-two-years-appear-overly-optimistic.html

06-Jan-16

Enterprises & Government

Logistics

- FedEx handles ~2.2 billion transactions a day.
- Indian Railways
 - 150,000 Tickets Per Hour

Aadhaar Data Streams



- Aadhaar packet sizes and rates
 - Hourly rates for a 24 hour period, Sep, 2015
- Mean packet size is ~3MB
- Bi-modal rate distribution
 - After morning & evening sessions
 - ~600K/day now. Peak was 1.3M/day in 2013
- Soon >20 Millions of authentications per day

06-Jan-16

Science: Pan-STARRS Sky Survey Sc. org

- Discover & characterize Earthapproaching objects that might pose a danger to our planet.
- One of the largest telescopes
 - 1.4 Gigapix camera world's largest!
- Scan ²/3^{rds} of sky, 3 times/month
 - 1 PB of images, 30 TB of processed data/year
 - 150 M detections / night
 - 5.5 Billion objects, 350 Billion detections







SDSS, PS, LSST, SKA, ...

- Sloan Digital Sky Survey...first all digital survey...precursor to PS
- LSST: 10M events/night. 37B objects, 7T obs, 30T measurements per year*
- Square Km Array
 - Radio Telescope, 300 dishes
 - 68Tb/s to be sampled
 - 10PB images/day. 1PB catalog.
 Graphs
- Next Gen Sequencers
- TB of data per run

Large Hadron Collider

Petabytes of data per run







Internet of Everything

- Personal Devices: Smart Phones, Fitbit
- Smart Appliances
- Smart Cities: Power, Water, Transportation, Environment
- Smart Retail









2015-09-26





06-Jan-16





06-Jan-16





06-Jan-16





06-Jan-16







Why is Big Data Difficult?



CDS.IISc.ac.in | **Department of Computational and Data Sciences**

Evolving Nature of Science Data

Large Hadron Collider



Illumina NGS @ IISc



Bluetooth Mote @ IISc



Smart Meter @ LADWP



Few Instruments, Large Data Volume

Many Devices, Volume & Velocity 10² Sources TB's Data Days to Proc.

10⁵ Sources GB's Data Hours to Proc.

Numerous Sensors, High data **Velocity** 10⁸ Sources MB's Data <Mins to Proc.

2015-04-15









1 TB OF TRADE







Modern cars have close to 100 SENSORS

that monitor items such as fuel level and tire pressure But what exactly is big of massive amounts of data

As a leader in the sect break big data into for Velocity, Variety and Vera

Depending on the indust data encompasses info internal and external sour social media, enterprise mobile devices. Comparadapt their products and customer needs, optiinfrastructure, and find (

By 2015 4.4 MILLION IT JOI

will be created globally with 1.9 million in the



The New York Stock Exchange captures

1 TB OF TRADE INFORMATION

during each trading session



By 2016, it is projected there will be

18.9 BILLION NETWORK CONNECTIONS

- almost 2.5 connections per person on earth

ANALYSIS OF STREAMING DATA

Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, GAS http://www.ibmbigdatahub.com/infographic/four-vs-big-data



CDS.IISc.ac.in | **Department of Computational and Data Sciences**

1441

R V's ig 1

nd music downloads to web records, data is recorded, to enable the technology world relies on every day, g data, and how can these ata be used?

ctor, IBM data scientists four dimensions: Volume, tracity

ustry and organization, big iformation from multiple burces such as transactions, rise content, sensors and As of 2011, the global size of data in healthcare was estimated to be

150 EXABYTES 161 BILLION GIGABYTES 1



30 BILLION PIECES OF CONTENT

are shared on Facebook every month



Variety

DIFFERENT

FORMS OF DATA

By 2014, it's anticipated there will be

420 MILLION WEARABLE, WIRELESS HEALTH MONITORS

4 BILLION+ HOURS OF VIDEO

are watched on YouTube each month



400 MILLION TWEETS

are sent per day by about 200 million monthly active users



Poor data quality costs the US economy around



lata, and how can these i be used?

or, IBM data scientists ir dimensions: Volume, city

try and organization, big rmation from multiple cessuch as transactions, a content, sensors and ies can leverage data to services to better meet mize operations and new sources of revenue.

3S

to support big data, United States





400 MILLION INCLIS

are sent per day by about 200 million monthly active users

1 IN 3 BUSINESS LEADERS

don't trust the information they use to make decisions



Poor data quality costs the US economy around

\$3.1 TRILLION A YEAR



27% OF RESPONDENTS

in one survey were unsure of how much of their data was inaccurate

Veracity UNCERTAINTY OF DATA

IBM.

http://www.ibmbigdatahub.com/infographic/four-vs-big-data

Data Platforms

- Acquire, manage, process Big Data
- At large scales
- To meet application needs



The Platform Stack: Cloudera

Cloudera's Distribution for Hadoop

UI Framework Hue				SDK Hue SDK			
Workflow Oozie		Scheduling <i>oozie</i>			Metada	ata Hive	
Data Integration <i>Flume, Sqoop</i>		Languages	, Compilers	Pig/ Hive	Fast read/w access	rite HBase	
		Coord	ination		Z	Zookeeper	



The Platform Stack: BDAS



Reading & Assignment

- Textbook chapter 2.1, 2.2
- Assignment 0
 - Install Hadoop 2.6.x+ in pseudo distributed mode on your laptop