

CS 6350 Big Data Analytics and Management (Graduate Level) Spring 2015

People:

Instructor: Dr. Latifur Khan

Office: ECSS (ES) 3.228

Phone: (972) 883 4137

E-mail: lkhan@utdallas.edu

Office Hours: Tuesday: 11.45 a.m. to 1.00 p.m. & Friday: 6.45 p.m. to 7.30 p.m.

URL: <http://www.utdallas.edu/~lkhan/Spring2015>

Class Time

15S	CS 6350.001 26923	Big Data Management and Analytics (3 Credits)	Tues & Thurs : 1:00pm-2:15pm ECSS 2.415
15S	CS 6350.002 26999	Big Data Management and Analytics (3 Credits)	Fri : 4:00pm-6:45pm ECSS 2.410

Teaching Assistants (TA):

TBD;

Office Hour: TBD

Office Location: TBD

Course Summary

Popular relational database systems like [IBM DB2](#), [Microsoft SQLServer](#), [Oracle](#), and [Sybase](#) are struggling to handle massive scale of data introduced by the Web, Social network and cyber physical systems. Now-a-days, companies have to deal with extremely large datasets. For example, on one hand, Facebook handles 15 TeraBytes of data each day into their [2.5 PetaByte Hadoop-powered data warehouse](#) and on the other hand, eBay maintains a 6.5 PetaByte data warehouse. To handle emerging data at massive scale, "big data analytics" and "big data management" areas are emerging. Many traditional assumptions are not working, instead, [new query and programming interfaces are required](#), and new computing models are emerging.

The course will focus on data mining and machine learning algorithms for analyzing very large amounts of data or Big data. Map Reduce and No SQL system will be used as tools/standards for creating parallel algorithms that can process very large amounts of data.

The course material will be drawn from textbooks as well as recent research literature. The following topics will be covered this year: Hadoop, Mapreduce, NoSQL systems (Cassandra, Pig, Hive, BigTable, HBASE, SPARK), Storm, Large scale supervised machine learning, Data streams, Clustering, and Applications including recommendation systems, Web and security.

Requirements

Two exam (in March (after Spring break), May), a couple of assignments (preferably 4) and a project. Your course grade will be based on 45% of the exam, 42.5% of assignments, and 12.5% on the project. The project includes 1/2 page proposal, implementation and demonstration on April (at the end).

Class Learning Outcomes	Number of Students			Material U
	Below Expectations	Meets Criteria	Exceeds Criteria	
Ability to understand of conceptual, logical and physical organization of big data				HW 1, 2, 3, 4
Ability to understand of large data processing using Map-Reduce				HW 1, 2, 3, 4
Ability to understand of NoSQL models, theory and practices				HW 1.2 3, 4,
Ability to understand of data modeling, indexing, query processing for big data				HW 1, 2, 3, 4,
Ability to understand of recommendation methods for big data				HW 4, EXA
Ability to understand of unsupervised learning for big data				HW 4, EXA
Ability to Understand of supervised learning for big data				HW 4, EXA
Ability to communicate and work on team software project				Project

Perquisite:

Database Management Systems, JAVA (intermediate/advanced), Linux OS, Machine Learning/AI (co-requisite)

Course Materials

The following textbook can be used this semester to augment the material presented in lectures:

- B1: Jimmy Lin and Chris Dyer, Data-Intensive Text Processing with MapReduce, Morgan & Claypool Publishers, 2010. <http://lntool.github.com/MapReduceAlgorithms/> [Mandatory]
- B2: Pang-Ning Tan, Michael Steinbach, and Vipin Kumar, Introduction to Data Mining, Addison-Wesley April 2005. [Mandatory]
- B3: Anand Rajaraman and Jeff Ullman, Mining of Massive Datasets, Cambridge Press, <http://infolab.stanford.edu/~ullman/mmds/book.pdf> [Mandatory]
- B4: Jiawei Han and Micheline Kamber, Data Mining: Concepts and Techniques, The Morgan Kaufmann Series in Data Management Systems, Jim Gray, Series Editor Morgan Kaufmann Publishers, August 2000. 550 pages. ISBN 1-55860-489-8. [Optional]
- B5: Chuck Lam, Hadoop in Action, December, 2010 | 336 pages ISBN: 9781935182191, <http://netlab.ulusofoona.pt/cp/HadoopinAction.pdf> [Optional]

Papers Related to Big Data Analytics and Management [May Expand Further]

- P1: Jeffrey Dean and Sanjay Ghemawat, MapReduce: Simplified Data Processing on Large Clusters, *OSDI'04: Sixth Symposium on Operating System Design and Implementation*, San Francisco, CA, December, 2004. http://static.googleusercontent.com/external_content/untrusted_dlcp/research.google.com/en/us/archive/mapreduce-osdi04.pdf (shortened version: <http://dl.acm.org/citation.cfm?doid=1327452.1327492>)
- P2: Michael Stonebraker, Daniel Abadi, David J. DeWitt, Sam Madden, Erik Paulson, Andrew Pavlo, and Alexander Rasin. (2010) [MapReduce and Parallel DBMSs: Friends or Foes?](#) *Communications of the ACM*, 53(1):64-71.
- P3: Jeffrey Dean and Sanjay Ghemawat. (2010) [MapReduce: A Flexible Data Processing Tool](#). *Communications of the ACM*, 53(1):72-77.
- P4: Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach, Michael Burrows, Tushar Chandra, Andrew Fikes, and Robert Gruber. (2006) [Bigtable: A Distributed Storage System for Structured Data](#). *Proceedings of the 7th Symposium on Operating System Design and Implementation (OSDI 2006)*, pages 205-218.
- P5: Christopher Olston, Benjamin Reed, Utkarsh Srivastava, Ravi Kumar, and Andrew Tomkins. (2008) [Pig Latin: A Not-So-Foreign Language for Data Processing](#). *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, pages 1099-1110.
- P6: Mohammad M. Masud, Jing Gao, Latifur Khan, Jiawei Han, and Bhavani Thuraisingham, Classification and Novel Class Detection in Concept-Drifting Data Streams under Time Constraints, *IEEE Transactions on Knowledge & Data Engineering (TKDE)*, 2011, IEEE Computer Society, June 2011, Vol. 23, No. 6, Page 859-874.
- P7: [Mohammad M. Masud](#), [Clay Woolam](#), [Jing Gao](#), Latifur Khan, [Jiawei Han](#), [Kevin W. Hamlen](#), [Nikunj C. Oza](#): Facing the reality of data stream classification: coping with scarcity of labeled data. *Knowl. Inf. Syst.* 33(1): 213-244 (2011)

- P8: Mohammad Masud, Tahseen Al-Khateeb, Latifur Khan , Charu Aggarwal, and Jiawei Han. Recurring and Novel Class Detection using Class-Based Ensemble, In *Proc. of IEEE International Conference on Data Mining(ICDM)*, 2012, Belgium, Dec 2012.
- P9: Haoyuan Li, Yi Wang, Dong Zhang, Ming Zhang, and Edward Y. Chang. PFP: Parallel FP-Growth for Query Recommendation. In Proceedings of the 2008 ACM conference on Recommender systems.
- P10: Avinash Lakshman, Prashant Malik, Cassandra: a decentralized structured storage system, ACM SIGOPS Operating Systems Review archive, Volume 44 Issue 2, April 2010, Pages 35-40, ACM New York, NY, USA
- P11: Alan F. Gates, Olga Natkovich, Shubham Chopra, Pradeep Kamath, Shravan M. Narayanamurthy, Christopher Olston, Benjamin Reed, Santhosh Srinivasan, Utkarsh Srivastava, Building a High-Level Dataflow System on top of Map-Reduce: The Pig Experience, VLDB 2009.

Hive:

- P12: Thusoo, A.; Sarma, J.S.; Jain, N.; Zheng Shao; Chakka, P.; Ning Zhang; Antony, S.; Hao Liu; Murthy, R., "Hive - a petabyte scale data warehouse using Hadoop," *Data Engineering (ICDE), 2010 IEEE 26th International Conference on* , vol., no., pp.996,1005, 1-6 March 2010
doi: 10.1109/ICDE.2010.5447738

Software

Mahout: <http://mahout.apache.org/>

Hive: <https://cwiki.apache.org/confluence/display/Hive/Home>

Piglatin: <http://pig.apache.org/docs/r0.7.0/tutorial.html>

Hadoop: <http://hadoop.apache.org/>

Cassandra: <http://cassandra.apache.org/>

Kafka: <https://cwiki.apache.org/confluence/display/KAFKA/Kafka+papers+and+presentations>

Storm: <http://storm-project.net/>

Spark <https://spark.apache.org/>

Lectures

Topic	Chapters/Papesr	Homework/Lecture Notes
Hadoop+ Mapreduce	Chapter 1, 2, 3 [B1], [B3]. Paper: P1, P2, P3, P4	Hadoop with Mapreduce
Big Data Management: NoSQL, Pig Latin, Hive, Cassandra, BigTable, HBASE	Paper: P5, P10, P11	
Stream Data Management: Storm Spark & Kafka		
Clustering Analysis	Chapter 8 [B2] Chapter 9 [B2]	Lecture Note in Clustering
Classification, Prediction, Stream Mining	Chapter 5 [B2] Chapter 4 [B3] Paper: P6, P7, P8	
Recommendation System	Chapter 9 [B3], P9	
Graph Processing Applications	Papers TBD	
Exam I, II	March, May	
Project Demonstration	TBD	

Assignment (Tentative)

Assignments will be based on:

- 1. Map Reduce Programming—Basic & Hands on HDFS setup**
- 2. Map Reduce Programming – advanced**
- 3. Big Data Management Using Hive, Pig, Cassandra, HBASE (including Hands on setup)**
- 4. Stream Processing and Analytics using Spark, and Mahout; Problem Solving Questions/Exercise Problems from Books; and Large Scale Machine Learning Using Mahout**

Projects (Sample—Expanded further)

- Tweeter Data Management**
- Google Glass**
- Anomaly Detection**
- Stream Mining for Tweets**
- Text Mining with LDA**
- Sentiment Analysis**