

# MIDS, UC Berkeley

## DATASCI W261

### W261 "Machine Learning at Scale"

---

**Course Architect: Dr. James G. Shanahan, 1/1/2015**

Revised on April 3, 2017 by Jimi Shanahan

#### **Course Description:**

Learn the underlying principles required to develop scalable machine-learning pipelines for structured and unstructured data at the petabyte scale, and gain hands-on experience in Apache Hadoop and Apache Spark and related MapReduce frameworks.

#### **Course Summary:**

Until recently, “big data” was very much the purview of database management and summary statistics systems such as Hadoop (HDFS and MapReduce) and was largely underleveraged by machine learning. These systems, though useful, suffered from their limited utility. **This course builds on and goes beyond this collect-and-analyze phase of big data by focusing on how machine learning algorithms can be rewritten and in some cases extended to scale to work on petabytes of data, both structured and unstructured, to generate sophisticated models that can be used for real-time predictions.** Predictive modeling at this scale can lead to huge boosts in performance (typically in the order of 10–20%) over small-scale models running on stand-alone computers that require one to significantly down-sample and, necessarily, to simplify big data. Concretely, this course focuses on how the MapReduce design pattern from parallel computing can be extended and more faithfully leveraged to tackle the somewhat “embarrassingly parallel” task of machine learning (a lot of machine learning algorithms fit this mold). In this course, this is accomplished via the Apache Spark project and its many related subprojects. Apache Spark is an open-source cluster-computing framework. It has emerged as the next-generation big data processing engine, overtaking Hadoop MapReduce, which helped ignite the big data revolution. **Spark maintains MapReduce’s linear scalability and fault tolerance but extends it in a few important ways: it is much faster (100 times faster for certain applications); much easier to program in, due to its rich APIs in Python, Java, and Scala (and R) and its core data abstraction, the distributed data frame; and it goes far beyond batch applications to support a variety of compute-intensive tasks, including interactive queries, streaming, machine learning, and graph processing.**

This course will provide an accessible introduction to MapReduce frameworks and to Spark and its potential to revolutionize academic and commercial data science practices through scale. Conceptually,

the course is divided into two parts. The first part will cover fundamental concepts of MapReduce parallel computing via Hadoop, MRJob, and Spark, while diving deep into Spark core, data frames, the Spark shell, Spark streaming, Spark SQL, MLlib, and more. The second part will focus on hands-on algorithmic design and development in parallel computing environments such as Spark; developing algorithms from scratch, such as decision-tree learning; graph-processing algorithms such as PageRank and shortest path; gradient descent algorithms such as support vectors machines; and matrix factorization. Industrial applications and deployments of MapReduce parallel compute frameworks from various fields, including advertising, finance, healthcare, and search engines, will also be presented. Examples and exercises will be made available in Python notebooks (Hadoop streaming, MRJob, and PySpark).

**Prerequisites:** Introductory machine learning course or equivalent. Intermediate programming capabilities in an object-oriented language (such as Python).

## Course Assignments and Exams:

There will be homework associated with each week of class in this course. These assignments will be completed individually. The midterm exam will cover all the major topics in the course to that point. Finally, there is a group project that you will complete with a small team, using one of the course data sets.

% of Final Grade	Task
40%	Homework+Projects (iterate/agile): Projects will have phases (all structured)
20%	Midterm Exam (Week 8 of the semester)
20%	Class participation (fact-to-face and in the online forums; please answer each other's questions; collaborate; communicate)
20%	End term exam (Week 14 of the Semester)

## Homework Schedule and policy

- Homework submissions are due on Tuesday morning by 8AM, West Coast Time (this is a strict deadline as we have a tough schedule to keep up with each week during semester.
- Each day you are late you will incur a **33.33% penalty**.
- Grade each HW and give written feedback to students by Sunday night of each week. This is a great point of calibration. TAs should make instructors aware of student performances and areas where students did not do well.
- Publish HW results on Monday, along with master solution
- E.g., (with all times in US West Coast Time)
  - HW1 is due on the Tuesday of Week 2 at 8AM (Sharp)
  - HW1 will be graded by Sunday night of Week 2
  - Grades and master solution for HW1 will be published (on ISVC) by Sunday midnight of Week 2.

## Homework Location

- Github Repo:
- <https://classroom.github.com/assignment-invitations/562761d1f7fde3700357401ff5ad25fa>

## Office Hours

Office hours are an on an on-demand basis. This need will be met largely by extra time after live sessions. If required we can do additional office on **Sundays, at 2PM West Coast Time.**

**NOTE:** You can only attend if you submit a question(s) 24 hours before the start of office hours to the online forum (SLACK). Use the [STAR](#) methodology to present your Questions/Answers for both Knowledge based questions and for problem solving questions

## Slack Channel

Instructors Channel: <https://ucbyschool.slack.com/messages/w261-instructors>

Students Channel: <https://ucbyschool.slack.com/messages/w261-summer-2017>

- Join Slack Channel for questions/discussion
- JimiUCB@ucbyschool.slack.com, Mike.Tamir, KyleHamilton
  - #2017-0109\_w261
  - [https://ucbyschool.slack.com/messages/2017-0109\\_w261/](https://ucbyschool.slack.com/messages/2017-0109_w261/)

## Other important Links

- Calendar and Google Drive:
  - [https://docs.google.com/spreadsheets/d/1RRmGU5aymKWbFjxvK\\_ElQJnfTaB-G2U3ByNv2b0wWO0/edit#gid=1565460380](https://docs.google.com/spreadsheets/d/1RRmGU5aymKWbFjxvK_ElQJnfTaB-G2U3ByNv2b0wWO0/edit#gid=1565460380)
- Slides and other folders
  - <https://drive.google.com/open?id=0B0b2FfTVAWAMTDZOWnN3VDhUOWc>

## Course Software:

Python and iPython Notebooks will be the common link between all the MapReduce frameworks we will be using to teach with and to complete projects. Having said that, the core focus of this class is large-scale machine learning. As such, we will be using frameworks that allow the handling of petabytes

of data (both structured and unstructured) from the following perspectives: storage, network transfer, exploratory data analysis, modeling, and experimentation. We will primarily focus on the following MapReduce frameworks:

- Hadoop streaming [Week 1]
- MRJob [Week 3]
- Spark and some of its specialty libraries (SparkSQL, MLlib, GraphX) [Week 10]
- AltaScale [Week 4 introduction]

## Textbooks and Reader:

This course will use a combination of textbooks and some online readings. You can buy the paper or electronic copies of the textbooks. All students are expected to do the readings before watching a asynchronous lecture and before all synchronous (also known as live) sessions. In many cases, the lecture materials and the weekly discussion sessions will not make sense if you do not have knowledge of the examples and issues from the readings.

### Recommended Textbooks:

Karau, Holden, Konwinski, Andy, Wendell, Patrick, & Zaharia, Matei. (2015). *Learning Spark: Lightning-fast big data analysis*. Sebastopol, CA: O'Reilly Publishers.

Lin, Jimmy, & Dyer, Chris. (2010). *Data-intensive text processing with MapReduce*. San Rafael, CA: Morgan & Claypool Publishers. (Free online)

Hastie, Trevor, Tibshirani, Robert, & Friedman, Jerome. (2009). *The elements of statistical learning: Data mining, inference, and prediction* (2nd ed.). Stanford, CA: Springer Science+Business Media. (Free online)

Ryza, Sandy, Laserson, Uri, Owen, Sean, & Wills, Josh. (2015). *Advanced analytics with Spark: Patterns for learning from data at scale*. Sebastopol, CA: O'Reilly Publishers.

Leskovec Jure, Rajaraman Anand, Ullman jeff, (2014). *Mining of Massive Datasets*, Cambridge University Press. Book available online at <http://www.mmids.org/>

Online references: Other references such as research pages, webpages and examples are included in the lectures directly.

### Course Logistics:

Every week you will cover lecture 1 by yourself asynchronously. The goal of lecture 2 (also known as the live session) is for us to meet and discuss questions that you may have, review homework and also to reinforce the topics in lecture 1 with additional examples or background material. If time permits we

may occasionally build on the lecture 1 topics and introduce new and related topics.

# Syllabus

Unit 1	Introduction and Motivation for Machine Learning at Scale
Lecture 1	Introduction and Motivation for Machine Learning at Scale
References	<p>Data Science</p> <ul style="list-style-type: none"> <li>• Selections from, Doing data science, Schutt and O'Neil; (Chapters 1 and 2)</li> <li>• Youtube video: The human face of big data: Rick Smolan at TEDxMidwest (<a href="https://www.youtube.com/watch?v=h8FLkRK-wF4">https://www.youtube.com/watch?v=h8FLkRK-wF4</a>)</li> </ul> <p>Excellent articles with code in matlab for bias-variance analysis of squared error loss</p> <ul style="list-style-type: none"> <li>• <a href="https://theclevermachine.wordpress.com/2013/04/21/model-selection-underfitting-overfitting-and-the-bias-variance-tradeoff/">https://theclevermachine.wordpress.com/2013/04/21/model-selection-underfitting-overfitting-and-the-bias-variance-tradeoff/</a></li> <li>• <a href="http://insidebigdata.com/2014/10/22/ask-data-scientist-bias-vs-variance-tradeoff/">http://insidebigdata.com/2014/10/22/ask-data-scientist-bias-vs-variance-tradeoff/</a></li> </ul> <p>Hastie, Trevor, Tibshirani, Robert, &amp; Friedman, Jerome. (2009). <i>The elements of statistical learning: Data mining, inference, and prediction</i> (2nd ed.). Stanford, CA: Springer Science+Business Media. (Chapter 7)</p> <p>Also other referenced materials (online links are embedded in slides); this is the case for all lectures</p>
Topics	<p>Big data</p> <p>Sources of big data</p> <p>Data Scientist</p> <p>Large scale machine learning</p> <p>Data modeling pipeline</p> <p>Bias-variance tradeoff</p> <p>Back of envelope calculations</p> <p>A MapReduce framework using Unix Command line</p>
Lecture 2	A command line map-reduce environment (Lecture 2 corresponds to the weekly synchronous time period; we may cover some or all of the topics listed here depending on the available time and needs of the students)
References	<p>Introduction to information retrieval; Manning, Raghavan, &amp; Schutz, Chapter 13 (Naïve Bayes)</p> <p><a href="https://en.wikipedia.org/wiki/Reservoir_sampling">https://en.wikipedia.org/wiki/Reservoir_sampling</a></p>
Topics	<p>Self-Introductions (Each student: please prepare 100-200 word bio of yourself to paste into the chat pod in Adobe Connect)</p> <p>Bernoulli Naive Bayes</p> <p>Multinomial Naive Bayes</p> <p>Coding up a distributed Naive Bayes using Python (and Unix command-line utilities)</p>

<b>Unit 2</b>	<b>Parallel Computing, MapReduce, and Hadoop (Data Storage and Algorithms)</b>
<b>Lecture 1</b>	<b>Parallel Computing, MapReduce, and Hadoop (Data Storage and Algorithms)</b>
References	<ul style="list-style-type: none"> <li>• Lin, Jimmy, &amp; Dyer, Chris. (2010). Data-intensive text processing with MapReduce. San Rafael, CA: Morgan &amp; Claypool Publishers. (Chapter 2)</li> <li>• <a href="http://en.wikipedia.org/wiki/Barrier_(computer_science)">http://en.wikipedia.org/wiki/Barrier_(computer_science)</a></li> <li>• <a href="http://en.wikipedia.org/wiki/Embarrassingly_parallel">http://en.wikipedia.org/wiki/Embarrassingly_parallel</a></li> <li>• <a href="http://gerardnico.com/wiki/data_storage/shared_disk_nothing">http://gerardnico.com/wiki/data_storage/shared_disk_nothing</a></li> <li>• <a href="http://en.wikipedia.org/wiki/Shared_nothing_architecture">http://en.wikipedia.org/wiki/Shared_nothing_architecture</a></li> <li>• <a href="http://en.wikipedia.org/wiki/Message_Passing_Interface">http://en.wikipedia.org/wiki/Message_Passing_Interface</a></li> <li>• <a href="http://en.wikipedia.org/wiki/Map_(parallel_pattern)">http://en.wikipedia.org/wiki/Map_(parallel_pattern)</a></li> <li>• <a href="http://www.michael-noll.com/tutorials/writing-an-hadoop-mapreduce-program-in-python/">http://www.michael-noll.com/tutorials/writing-an-hadoop-mapreduce-program-in-python/</a></li> <li>• <a href="http://hadoop.apache.org/core/docs/current/">http://hadoop.apache.org/core/docs/current/</a></li> <li>• See also other referenced material</li> </ul>
Topics	<ul style="list-style-type: none"> <li>• Motivation for parallel computing</li> <li>• Parallel computing</li> <li>• Parallel Computing (PC) Definition and Communication Synchronization Types of PC Tasks</li> <li>• Architectures for Parallel Computation</li> <li>• Developer Frameworks for Parallel Computation</li> <li>• Hadoop Background and History</li> <li>• Hadoop Distributed File System (HDFS)</li> <li>• MapReduce: Functional Programming</li> <li>• Hadoop: MapReduce</li> <li>• Animated Examples</li> </ul>
<b>Lecture 2</b>	<b>Installing and Programming Hadoop</b>
References	Hadoop online documentation
Topics	<ul style="list-style-type: none"> <li>• MapReduce: Runtime environment</li> <li>• Hadoop 2.0</li> <li>• Install Hadoop; run locally; run in the cloud?</li> <li>• Run wordcount notebook and explain each step</li> <li>• 10 most popular words</li> <li>• Examples of problem solving in Hadoop</li> <li>• Naive Bayes in Hadoop</li> </ul>

<b>Unit 3</b>	
<b>MapReduce Algorithm Design and Design Patterns</b>	
<b>Lecture 1</b>	<b>MapReduce Algorithm Design and Design Patterns</b>
References	<ul style="list-style-type: none"> <li>• Lin, Jimmy, &amp; Dyer, Chris. (2010). Data-intensive text processing with MapReduce. San Rafael, CA: Morgan &amp; Claypool Publishers. (Chapter 3).</li> <li>• <a href="https://en.wikipedia.org/wiki/Priority_queue">https://en.wikipedia.org/wiki/Priority_queue</a></li> <li>• <a href="http://en.wikipedia.org/wiki/Merge_algorithm">http://en.wikipedia.org/wiki/Merge_algorithm</a></li> </ul> <p>See also other referenced material</p>
Topics	<ul style="list-style-type: none"> <li>• Background</li> <li>• RAM vs. disk vs. bandwidth</li> <li>• Priority queues and merge sort</li> <li>• The internals of the Hadoop shuffle</li> <li>• Local aggregation</li> <li>• Combiners and in-mapper combining</li> <li>• Algorithmic correctness with local aggregation</li> <li>• How many Mappers and Reducers?</li> <li>• Pairs and stripes</li> <li>• Computing relative frequencies</li> </ul>
<b>Lecture 2</b>	<b>Hadoop at Scale and Installing MRJob</b>
References	<ul style="list-style-type: none"> <li>• Introduction to information retrieval; Manning, Raghavan, &amp; Schutz, Chapter 13</li> <li>• Lin, Jimmy, &amp; Dyer, Chris. (2010). Data-intensive text processing with MapReduce. San Rafael, CA: Morgan &amp; Claypool Publishers. (Section 3.4).</li> <li>• <a href="https://en.wikipedia.org/wiki/Apriori_algorithm">https://en.wikipedia.org/wiki/Apriori_algorithm</a></li> <li>• MMDS Book on Data mining, <ul style="list-style-type: none"> <li>o Mining of Massive Datasets by Jure Leskovec, Anand Rajaraman, Jeff Ullman, Book available online at <a href="http://www.mmms.org/">http://www.mmms.org/</a></li> <li>o <a href="https://www.dropbox.com/s/yoapjbpoiwmfssy/MMDS-Book-Chapter6.pdf?dl=0">https://www.dropbox.com/s/yoapjbpoiwmfssy/MMDS-Book-Chapter6.pdf?dl=0</a></li> </ul> </li> </ul>
Topics	<ul style="list-style-type: none"> <li>• Secondary sorting</li> <li>• Coding and debugging Hadoop</li> <li>• A-Priori algorithm for mining frequent itemsets (single brute force map-reduce job for itemsets 1... n, and for association rule extraction as opposed to a multistep pipeline with that has iterative part)</li> <li>• Installing MRJob</li> </ul>



<b>Unit 4</b>	<b>Unsupervised Learning at Scale in MapReduce: Clustering, K-Means, Expectation Maximization</b>
<b>Lecture 1</b>	<b>MRJob, Unsupervised Learning at Scale: Clustering and K-Means</b>
References	<ul style="list-style-type: none"> <li>• Introduction to information retrieval; Manning, Raghavan, &amp; Schutz, Chapters 16 and 17</li> <li>• MRJob online documentation (<a href="https://pythonhosted.org/mrjob/">https://pythonhosted.org/mrjob/</a> )</li> <li>• MRJob examples <ul style="list-style-type: none"> <li>◦ <a href="http://pythonhosted.org/mrjob/guides/writing-mrjobs.html%23writing-basics">http://pythonhosted.org/mrjob/guides/writing-mrjobs.html%23writing-basics</a></li> <li>◦ <a href="https://github.com/uchicago-cs/cmsc12300/tree/master/examples/data_analysis/bin">https://github.com/uchicago-cs/cmsc12300/tree/master/examples/data_analysis/bin</a></li> </ul> </li> <li>• Microsoft customer visitor log file data <ul style="list-style-type: none"> <li>◦ <a href="http://archive.ics.uci.edu/ml/machine-learning-databases/anonymous/">http://archive.ics.uci.edu/ml/machine-learning-databases/anonymous/</a></li> </ul> </li> <li>• Google N-Gram data <ul style="list-style-type: none"> <li>◦ <a href="http://storage.googleapis.com/books/ngrams/books/datasetsv2.html">http://storage.googleapis.com/books/ngrams/books/datasetsv2.html</a></li> </ul> </li> <li>• <b>Yelp review challenge:</b> <ul style="list-style-type: none"> <li>◦ <a href="http://www.yelp.com/dataset_challenge/">http://www.yelp.com/dataset_challenge/</a></li> <li>◦ <a href="https://github.com/jblomo/pycon-mrjob/blob/master/code/unique_review.py">https://github.com/jblomo/pycon-mrjob/blob/master/code/unique_review.py</a></li> </ul> </li> <li>• See also other referenced material</li> </ul>
	<ul style="list-style-type: none"> <li>• MRJob Installation</li> <li>• MRJob fundamentals and concepts</li> <li>• Writing MRJob code</li> <li>• Log-file processing</li> <li>• Serializable, JSON, and other MRJob information</li> <li>• Clustering overview</li> <li>• K-means algorithm</li> </ul>
<b>Lecture 2</b>	<b>Model-Based Clustering: Expectation Maximization</b>
References	<ul style="list-style-type: none"> <li>• Introduction to information retrieval; Manning, Raghavan, &amp; Schutz, Chapter 16</li> <li>• MMDS Book on Data mining, <ul style="list-style-type: none"> <li>◦ Mining of Massive Datasets by Jure Leskovec, Anand Rajaraman, Jeff Ullman, Book available online at <a href="http://www.mmids.org/">http://www.mmids.org/</a></li> <li>◦ <a href="https://www.dropbox.com/s/yoapjbpoiwmfssy/MMDS-Book-Chapter6.pdf?dl=0">https://www.dropbox.com/s/yoapjbpoiwmfssy/MMDS-Book-Chapter6.pdf?dl=0</a></li> </ul> </li> </ul>
Topics	<ul style="list-style-type: none"> <li>• Distributed k-means in MRJob</li> <li>• Distributed k-means in MRJob (code)</li> <li>• Data and distributions</li> <li>• Model-based clustering</li> <li>• EM algorithm for a mixture of Bernoulli distributions</li> <li>• Distributed EM clustering algorithm in MRJob</li> <li>• <i>Initialization (canopy clustering) if time permits</i></li> </ul>

<b>Unit 5</b>	<b>Big Data Pipelines</b>
---------------	---------------------------

Lecture 1	Big Data Pipelines
References	<ul style="list-style-type: none"> <li>• Lin, Jimmy, &amp; Dyer, Chris. (2010). Data-intensive text processing with MapReduce. San Rafael, CA: Morgan &amp; Claypool Publishers. (Chapter 3);</li> <li>• Shanahan, James G., &amp; Kurra, Goutham. (2011). Digital advertising: An information scientist's perspective. In Massimo Melucci &amp; Ricardo Baeza-Yates (Eds.), Advanced topics in information retrieval, pp. 209–238. Berlin and Heidelberg: Springer-Verlag</li> <li>• <a href="http://en.wikipedia.org/wiki/Data_warehouse">http://en.wikipedia.org/wiki/Data_warehouse</a></li> <li>• <a href="https://en.wikipedia.org/wiki/Join_(SQL)">https://en.wikipedia.org/wiki/Join_(SQL)</a></li> </ul> <p>See also other referenced material</p>
Topics	<ul style="list-style-type: none"> <li>• Mobile advertising example</li> <li>• Data warehousing, OLTP/OLAP</li> <li>• Database relational joins</li> <li>• Databases Operations in MapReduce</li> <li>• Relational Joins</li> <li>• Reduce Side Join</li> <li>• Reduce Side Join Example</li> <li>• Memory Backed Map Side Join</li> <li>• Map Side Join Merge Join</li> </ul>
Lecture 2	Exploratory Data Analysis and Evaluation and Distributed Information Retrieval
References	<ul style="list-style-type: none"> <li>• Lin, Jimmy, &amp; Dyer, Chris. (2010). Data-intensive text processing with MapReduce. San Rafael, CA: Morgan &amp; Claypool Publishers. (Chapter 4);</li> <li>• <a href="https://en.wikipedia.org/wiki/Reservoir_sampling">https://en.wikipedia.org/wiki/Reservoir_sampling</a></li> <li>• <a href="https://en.wikipedia.org/wiki/Column-oriented_DBMS">https://en.wikipedia.org/wiki/Column-oriented_DBMS</a></li> </ul>
Topics	<ul style="list-style-type: none"> <li>• Counters in MrJob/Hadoop</li> <li>• Online advertising <ul style="list-style-type: none"> <li>◦ Demand-side platform challenges</li> <li>◦ Case study (DSP competition)</li> </ul> </li> <li>• Column oriented databases</li> <li>• Distributed information retrieval? Why a new database structure?</li> <li>• Summary statistics</li> <li>• Sampling</li> </ul>

<b>Unit 6</b>	<b>Supervised Machine Learning in MapReduce (Part 1: Linear Regression)</b>
<b>Lecture 1</b>	<b>Distributed Supervised Machine Learning (Distributed Linear Regression via Two Approaches)</b>
References	<ul style="list-style-type: none"> <li>• Hastie, Trevor, Tibshirani, Robert, &amp; Friedman, Jerome. (2009). <i>The elements of statistical learning: Data mining, inference, and prediction</i> (2<sup>nd</sup> ed.). Stanford, CA: Springer Science+Business Media. (Chapter 3).</li> <li>• <a href="https://en.wikipedia.org/wiki/Machine_learning">https://en.wikipedia.org/wiki/Machine_learning</a></li> <li>• Lin, Jimmy, &amp; Dyer, Chris. (2010). Data-intensive text processing with MapReduce. San Rafael, CA: Morgan &amp; Claypool Publishers. (Section 3.4).</li> <li>• <a href="https://en.wikipedia.org/wiki/Mathematical_optimization">https://en.wikipedia.org/wiki/Mathematical_optimization</a></li> <li>• <a href="https://en.wikipedia.org/wiki/Convex_optimization">https://en.wikipedia.org/wiki/Convex_optimization</a></li> <li>• <a href="http://en.wikipedia.org/wiki/Gradient">http://en.wikipedia.org/wiki/Gradient</a></li> <li>• MapReduce Study (<a href="http://www.cs.stanford.edu/people/ang/papers/nips06-mapreducemulticore.pdf">http://www.cs.stanford.edu/people/ang/papers/nips06-mapreducemulticore.pdf</a>)</li> </ul> <p>See also other referenced material</p>
Topics	<ul style="list-style-type: none"> <li>• Supervised machine learning (parametric vs. nonparametric)</li> <li>• Applications of matrices at scale</li> <li>• Distributed matrix and vector multiplication</li> <li>• Optimization theory</li> <li>• Gradient descent</li> <li>• Convex optimization</li> <li>• Linear regression</li> <li>• Distributed closed-form linear regression</li> <li>• Complexity, scaling out, and case study</li> <li>• Distributed linear regression via gradient descent</li> </ul>
<b>Lecture 2</b>	<b>Regression Diagnostics, Experiments, and Extensions of Linear Regression</b>
References	Hastie, Trevor, Tibshirani, Robert, & Friedman, Jerome. (2009). <i>The elements of statistical learning: Data mining, inference, and prediction</i> (2 <sup>nd</sup> ed.). Stanford, CA: Springer Science+Business Media. (Chapter 3).
Topics	<ul style="list-style-type: none"> <li>• Ridge and Lasso regression (part 1)</li> <li>• Nonparametric approaches to machine learning at scale</li> <li>• Evaluation of models (internal diagnostics and metrics)</li> <li>• A-B testing</li> </ul>

<b>Unit 7</b>	<b>Introduction to Graph Algorithms at Scale: Single Shortest Path Algorithm</b>
<b>Lecture 1</b>	<b>Graph-Based Algorithms</b>
References	<ul style="list-style-type: none"> <li>• Lin, Jimmy, &amp; Dyer, Chris. (2010). <i>Data-intensive text processing with MapReduce</i>. San Rafael, CA: Morgan &amp; Claypool Publishers. (Chapter 5).</li> <li>• <b>Dijkstra's_algorithm</b> <ul style="list-style-type: none"> <li>o <a href="http://www.cs.cornell.edu/courses/cs312/2002sp/lectures/lec20/lec20.htm">http://www.cs.cornell.edu/courses/cs312/2002sp/lectures/lec20/lec20.htm</a></li> <li>o <a href="http://en.wikipedia.org/wiki/Dijkstra's_algorithm">http://en.wikipedia.org/wiki/Dijkstra's_algorithm</a></li> </ul> </li> <li>•</li> </ul> <p>See also other referenced material</p>
Topics	<ul style="list-style-type: none"> <li>• Networks Introduction and Motivation</li> <li>• Applications of Graphs</li> <li>• Graph Definitions</li> <li>• Shortest Path Introduction</li> <li>• Single Source Shortest Path (SSSP) Unweighted Graphs</li> <li>• BFS for Unweighted Graphs Directed and Animations</li> <li>• SSSP for Weighted Graphs BFS</li> <li>• Dijkstras Algorithm</li> </ul>
<b>Lecture 2</b>	<b>Implementing graph algorithms</b>
References	<a href="https://en.wikipedia.org/wiki/Connected_component_(graph_theory)">https://en.wikipedia.org/wiki/Connected_component_(graph_theory)</a>
Topics	<ul style="list-style-type: none"> <li>• SSSP</li> <li>• Connected-component</li> <li>• Implementing Connected components</li> <li>• Q&amp;A and planning for mid term</li> </ul>

<b>Unit 8</b>	<b>Supervised Machine Learning in MapReduce (Part 1: Linear Regression)</b>
<b>Lecture 1</b>	<b>Mid-Term Exam</b>
References	Weeks 1–7 reading material
Topics	A three-hour open-book exam will be administered based on the course material to date
<b>Lecture 2</b>	<b>Mid-Term Exam</b>
References	----
Topics	----

<b>Unit 9</b>	<b>Large Scale Graph Processing: Random Walks, PageRank, Personalized PageRank</b>
<b>Lecture 1</b>	<b>Large Scale Graph Processing: Random Walks, PageRank, Personalized PageRank</b>
References	<p>Lin, Jimmy, &amp; Dyer, Chris. (2010). <i>Data-intensive text processing with MapReduce</i>. San Rafael, CA: Morgan &amp; Claypool Publishers. (Chapter 5).</p> <p>Overview article on PageRank</p> <ul style="list-style-type: none"> <li>• Page, Lawrence; Brin, Sergey; Motwani, Rajeev and Winograd, Terry (1999). The PageRank citation ranking: Bringing order to the Web</li> <li>• <a href="http://www.dam.brown.edu/people/mchb/la/GooglePageRank.pdf">http://www.dam.brown.edu/people/mchb/la/GooglePageRank.pdf</a></li> </ul> <p>PageRank in Python/MrJob</p> <ul style="list-style-type: none"> <li>• <a href="https://github.com/Yelp/mrjob/blob/master/mrjob/examples/mr_page_rank.py">https://github.com/Yelp/mrjob/blob/master/mrjob/examples/mr_page_rank.py</a></li> <li>• <a href="http://www.cs205.org/resources/05-MapReduceDesignPatterns.pdf">http://www.cs205.org/resources/05-MapReduceDesignPatterns.pdf</a></li> <li>• The IR Book Chapter 21, <a href="http://nlp.stanford.edu/IR-book/pdf/21link.pdf">http://nlp.stanford.edu/IR-book/pdf/21link.pdf</a></li> </ul>
Topics	<ul style="list-style-type: none"> <li>• Web Search Introduction</li> <li>• WebGraph History</li> <li>• Random Walk</li> <li>• Random Walk as a Markov Chain</li> <li>• PageRank</li> <li>• PCA and Power Iteration</li> <li>• Distributed PageRank</li> <li>• Distributed PageRank Example</li> <li>• PageRank and Dangling Nodes</li> <li>• Topic-Specific PageRankcent</li> </ul>
<b>Lecture 2</b>	<b>Implementation and applications of Distributed Pagerank and install Spark</b>
References	Karau, Holden, Konwinski, Andy, Wendell, Patrick, & Zaharia, Matei. (2015). <i>Learning Spark: Lightning-fast big data analysis</i> . Sebastopol, CA: O'Reilly Publishers. (Chapters 1—2)
Topics	<ul style="list-style-type: none"> <li>• Implementations of Distributed Pagerank and its derivatives</li> <li>• HITS Algorithm</li> <li>• TextRank</li> <li>• Applications of pagerank (e.g., Graph analysis of movies)</li> <li>• Project planning</li> <li>• Install Spark</li> <li>• Final project planning</li> </ul>

<b>Unit 10</b>	<b>Apache Spark: From Basics to Advanced</b>
<b>Lecture 1</b>	<b>Introduction to Apache Spark</b>
References	<p>Karau, Holden, Konwinski, Andy, Wendell, Patrick, &amp; Zaharia, Matei. (2015). Learning Spark: Lightning-fast big data analysis. Sebastopol, CA: O'Reilly Publishers. (Chapters 1–3)</p> <p>See also other referenced material</p>
Topics	<ul style="list-style-type: none"> <li>● Functional Programming Review and lazy evaluation</li> <li>● Background on Spark</li> <li>● Spark basics</li> <li>● Programming with Base RDDs</li> <li>● Transformation and actions</li> <li>● Animated Example Data Flow</li> <li>● Pair RDDs</li> <li>● Basic WordCount in Spark</li> </ul>
<b>Lecture 2</b>	<b>Advanced Programming in Spark</b>
References	<p>Karau, Holden, Konwinski, Andy, Wendell, Patrick, &amp; Zaharia, Matei. (2015). Learning Spark: Lightning-fast big data analysis. Sebastopol, CA: O'Reilly Publishers. (Chapter 4)</p>
Topics	<ul style="list-style-type: none"> <li>● Accumulators and broad cast variables</li> <li>● Joins</li> <li>● Linear regression in Spark</li> <li>● Spark at scale (EC2 cluster)</li> <li>● Scala, tools for EDA, SparkSQL</li> <li>● Join, statistics in Spark, streaming, Naive Bayes</li> <li>● Project planning</li> </ul>

<b>Unit 11</b>	
<b>Distributed Supervised Machine Learning Part 2 (in Spark)</b>	
<b>Lecture 1</b>	<b>Distributed Supervised Machine Learning Part 2 (in Spark)</b>
References	<ul style="list-style-type: none"> <li>• Hastie, Trevor, Tibshirani, Robert, &amp; Friedman, Jerome. (2009). The elements of statistical learning: Data mining, inference, and prediction (2nd ed.). Stanford, CA: Springer Science+Business Media. (Section 4.4).</li> <li>• Izenman, Alan J. (2013). Modern multivariate statistical techniques: Regression, classification, and manifold learning. Springer Texts in Statistics. Stanford, CA: Springer Science+Business Media. (Chapter 11)</li> <li>• Loss functions <ul style="list-style-type: none"> <li>o <a href="http://www.ics.uci.edu/~dramanan/teaching/ics273a_winter08/lectures/lecture14.pdf">http://www.ics.uci.edu/~dramanan/teaching/ics273a_winter08/lectures/lecture14.pdf</a></li> </ul> </li> <li>• <b>Distributed Perceptron</b> <ul style="list-style-type: none"> <li>o <a href="http://static.googleusercontent.com/external_content/untrusted_dlcp/research.google.com/en//pubs/archive/36266.pdf">http://static.googleusercontent.com/external_content/untrusted_dlcp/research.google.com/en//pubs/archive/36266.pdf</a></li> <li>o <a href="http://cui.unige.ch/~gesmundo/papers/gesmundo-eacl2012.pdf">http://cui.unige.ch/~gesmundo/papers/gesmundo-eacl2012.pdf</a></li> </ul> </li> <li>•</li> </ul>
Topics	<ul style="list-style-type: none"> <li>• Loss Functions and General Framework for Gradient Descent</li> <li>• Logistic Regression at Scale</li> <li>• Perceptrons</li> <li>• Distributed Perceptron and case study</li> <li>• Support vector machines (SVM)</li> <li>• SGD Based SVMs at Scale</li> </ul>
<b>Lecture 2</b>	<b>Regression Diagnostics, Experiments, and Extensions of Linear Regression</b>
References	<ul style="list-style-type: none"> <li>• Hastie, Trevor, Tibshirani, Robert, &amp; Friedman, Jerome. (2009). The elements of statistical learning: Data mining, inference, and prediction (2nd ed.). Stanford, CA: Springer Science+Business Media. (Section 4.4).</li> <li>• Izenman, Alan J. (2013). Modern multivariate statistical techniques: Regression, classification, and manifold learning. Springer Texts in Statistics. Stanford, CA: Springer Science+Business Media. (Chapter 11)</li> </ul>
Topics	<ul style="list-style-type: none"> <li>• Derive dual form of SVMs</li> <li>• SMO algorithm</li> <li>• Distributed SMO</li> <li>• Kernel-AdaTron</li> <li>• QP-based algorithms</li> <li>• Nonparametric approaches to machine learning at scale</li> <li>• Ridge and Lasso logistic regression</li> <li>• Project work</li> </ul>



<b>Unit 12</b>	<b>Non-Gradient-Based Machine Learning Algorithms in Spark and Ensembles (Special Focus on Decision Trees)</b>
<b>Lecture 1</b>	<b>Non-Gradient-Based Machine Learning Algorithms in Spark (Special Focus on Decision Trees)</b>
References	<ul style="list-style-type: none"> <li>● Hastie, Trevor, Tibshirani, Robert, &amp; Friedman, Jerome. (2009). The elements of statistical learning: Data mining, inference, and prediction (2nd ed.). Stanford, CA: Springer Science+Business Media. (section 9.2)</li> <li>● <a href="https://en.wikipedia.org/wiki/Decision_tree_learning">https://en.wikipedia.org/wiki/Decision_tree_learning</a></li> <li>● Planet: Massively Parallel Learning of Tree Ensembles <ul style="list-style-type: none"> <li>○ <a href="http://static.googleusercontent.com/media/research.google.com/en//pubs/archive/36296.pdf">http://static.googleusercontent.com/media/research.google.com/en//pubs/archive/36296.pdf</a></li> </ul> </li> <li>● Ryza, Sandy, Laserson, Uri, Owen, Sean, &amp; Wills, Josh. (2015). Advanced analytics with Spark: Patterns for learning from data at scale. Sebastopol, CA: O'Reilly Publishers. (Chapter 4)</li> </ul>
Topics	<ul style="list-style-type: none"> <li>● Introduction to Decision Trees</li> <li>● Decision Trees Modes Medians Means and Machine Learning</li> <li>● Decision Tree Prediction</li> <li>● Decision Trees Learning a Regression Tree Over Real Valued Variables</li> <li>● Decision Trees Learning a Classification Tree</li> <li>● Decision Trees Understanding Input Variables</li> <li>● Decision Trees Representation</li> <li>● Learning a Decision Tree</li> <li>● Distributed Decision Tree Learning</li> <li>● Planet Overview</li> </ul>
<b>Lecture 2</b>	<b>Ensembles and Machine learning Case Studies in Spark</b>
References	<ul style="list-style-type: none"> <li>● Planet: Massively Parallel Learning of Tree Ensembles <ul style="list-style-type: none"> <li>○ <a href="http://static.googleusercontent.com/media/research.google.com/en//pubs/archive/36296.pdf">http://static.googleusercontent.com/media/research.google.com/en//pubs/archive/36296.pdf</a></li> </ul> </li> <li>● Ryza, Sandy, Laserson, Uri, Owen, Sean, &amp; Wills, Josh. (2015). Advanced analytics with Spark: Patterns for learning from data at scale. Sebastopol, CA: O'Reilly Publishers. (chapter 4)</li> </ul>
Topics	<ul style="list-style-type: none"> <li>● Decisions trees at scale</li> <li>● Ensembles of decision trees</li> <li>● MLlib introduction</li> <li>● Data frames in Spark</li> <li>● Case study: Airline delay prediction</li> <li>● Other case studies (advanced Spark book)</li> <li>● Project work</li> </ul>

<b>Unit 13</b>	
<b>Predicting Links and Attributes in Social Networks (OPTIONAL)</b>	
<b>Lecture 1</b>	<b>Predicting Links and Attributes in Social Networks</b>
References	<ul style="list-style-type: none"> <li>• Lars Backstrom and Jure Leskovec. 2011. Supervised random walks: predicting and recommending links in social networks. In Proceedings of the fourth ACM international conference on Web search and data mining (WSDM '11). ACM, New York, NY, USA, 635-644. DOI=10.1145/1935826.1935914 <ul style="list-style-type: none"> <li>o <a href="http://cs.stanford.edu/people/jure/pubs/linkpred-wsdm11.pdf">http://cs.stanford.edu/people/jure/pubs/linkpred-wsdm11.pdf</a></li> </ul> </li> <li>• Label propagation in graphs <ul style="list-style-type: none"> <li>o S. A. Macskassy and F. Provost. A simple relational classifier. In Proceedings of the Multi-Relational Data Mining Workshop (MRDM) at the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2003.</li> <li>o <i>Joint Inference of Multiple Label Types in Large Networks</i>, <a href="#">pdf ppt</a>, by D. Chakrabarti, S. Funiak, J. Chang, and S. A. Macskassy, in ICML 2014.</li> </ul> </li> <li>• TextRank: Bringing Order into Texts, Rada Mihalcea and Paul Tarau, EMNLP Conference 2004 <ul style="list-style-type: none"> <li>o <a href="http://web.eecs.umich.edu/~mihalcea/papers/mihalcea.emnlp04.pdf">http://web.eecs.umich.edu/~mihalcea/papers/mihalcea.emnlp04.pdf</a></li> </ul> </li> </ul>
Topics	<ul style="list-style-type: none"> <li>• A Supervised Random Walk</li> <li>• Supervised Random Walk Case Studies</li> <li>• Label Propagation Introduction</li> <li>• Weighted Vote Relational Neighbor Classifier</li> <li>• Facebook Case Study Inferring Missing Node Attributes</li> <li>• TextRank</li> <li>• Apache Big Data Graph Processing Frameworks</li> <li>• Friends of Friends</li> </ul>
<b>Lecture 2</b>	<b>Regression Diagnostics, Experiments, and Extensions of Linear Regression</b>
References	Ryza, Sandy, Laserson, Uri, Owen, Sean, & Wills, Josh. <i>Advanced analytics with Spark: Patterns for learning from data at scale</i> . Sebastopol, CA: O'Reilly Publishers. (Chapter 7)
Topics	<ul style="list-style-type: none"> <li>• Co-occurrence networks and graphX</li> <li>• Project work</li> </ul>

<b>Unit 14</b>	<b>Recommender Systems, PCA, SVD, Matrix Factorization and MLlib</b>
<b>Lecture 1</b>	<b>Alternating Least Squares and Various Optimization in Spark/Working with Spark MLlib</b>
References	Ryza, Sandy, Laserson, Uri, Owen, Sean, & Wills, Josh. <i>Advanced analytics with Spark: Patterns for learning from data at scale</i> . Sebastopol, CA: O'Reilly Publishers. (Chapter 3)
Topics	<ul style="list-style-type: none"> <li>• ALS Introduction</li> <li>• Single Node ALS Example</li> <li>• ALS Parallelization</li> <li>• Distributed ALS Closed Form</li> <li>• How to Implement ALS in Spark</li> <li>• Distributed ALS Gradient Descent</li> <li>• Pregel GraphX</li> <li>• ALS in Code</li> <li>• MLlib</li> <li>• Course Wrap Up</li> </ul>
<b>Lecture 2</b>	<b>SVD at Scale via ALS and Case Study</b>
References	Ryza, Sandy, Laserson, Uri, Owen, Sean, & Wills, Josh. <i>Advanced analytics with Spark: Patterns for learning from data at scale</i> . Sebastopol, CA: O'Reilly Publishers. (Chapter 6)
Topics	<ul style="list-style-type: none"> <li>• Recommender systems: Overview</li> <li>• Item-to-item algorithm</li> <li>• Project work</li> <li>• Item-to-item algorithm in Spark and MRJob <ul style="list-style-type: none"> <li>o MRJob (<a href="http://aimotion.blogspot.com/2012/08/introduction-to-recommendations-with.html">http://aimotion.blogspot.com/2012/08/introduction-to-recommendations-with.html</a>)</li> </ul> </li> <li>• PCA and Eigenfaces</li> <li>• Matrix factorization</li> <li>• ALS vs. SVD</li> <li>• MovieLens uses all data (not just the people's star ratings but the missing ones also)</li> </ul>

<b>Unit 15</b>	<b>Deep Learning (OPTIONAL)</b>
<b>Lecture 1</b>	<b>Presentations and Discussion</b>
References	---- Deep Learning references are provided within the slides
Topics	<ul style="list-style-type: none"> <li>• Deep Learning lecture</li> </ul>

## **Unit 15: Presentation and Discussion of Final Project**

Lecture 1: Presentations and Discussion

Lecture 2: Presentations and Discussion