

5460. Big Data Scaling

Maizie (Xin) Zhou, PhD

Syllabus

Class Information

Class Hours: Tuesday and Thursday 4:00 - 5:15pm

Room: Sony Building 2071

Office Hours: By appointment (maizie.zhou@vanderbilt.edu)

TA/Grader: Xiaoxing Qiu (xiaoxing.qiu@vanderbilt.edu)

Office hours on Campuswire: Wednesday 3 - 4 pm

Siqi (Christine) Zhou (siqi.zhou@vanderbilt.edu)

Office hours on Campuswire: Monday 1 - 2 pm

Grades: Grading will be based on the following elements:

Homework Assignments, In-class live coding participation, Midterm

Final project (writeup and presentation).

Assessment:

- **Homework Assignments 30%:** There will be homework assignments nearly for every topic. They are intended primarily to help you prepare for the exam and project. Because of the frequency of the assignments, I will not accept late homeworks for any reason. However, I will drop the lowest homework score when calculating your overall grade in the course. You are allowed to work in groups on the homework, but you must write up your own solutions in your own words. **ASSIGNMENTS ARE DUE AT 11:59 midnight OF THE DUE DATE THROUGH [BRIGHTSPACE](#).**
- **In class live coding participation 10%:** There will be living coding participation nearly for every topic on Thursday.
- **Midterm 30%:** Midterm will be a take-home exam.
- **Final Project 30%:** Projects are required to be related in a substantive way to at least one of the central topics of the course. Final projects can be done in groups of 1 - 3 people. We encourage you to form a group of 3 members, since groups of 3 usually lead to the best outcomes. We will talk about more details in class.

Class Announcements: All students are held responsible for all announcements made in the class and [Campuswire](#).

Campuswire: We will use Campuswire for all homework/midterm/project announcements, questions and public communications, holding office hours. You should have already received the invite by email!

Course Materials: The course relies on Jupyter notebooks, Powerpoint presentations, and online resources (<https://spark.apache.org/docs/latest/>). The lecture slides and Jupyter notebooks for each week will be found in the [box](#) as the course progresses. The textbook (Mining of Massive Dataset: <http://www.mmds.org/>) is useful and recommended, but not required. You can download it for free or purchase the hardcopy from Cambridge University Press.

All academic work at Vanderbilt is done under the Honor System.

The course will discuss data mining and machine learning algorithms, and the emphasis will be on MapReduce and Spark as tools for creating parallel algorithms that can process very large amounts of data.

Topics include: Cloud Computing, Distributed File Systems, MapReduce, Apache Spark, SparkSQL, Regression, Tree Methods, Locality Sensitive Hashing (LSH), Clustering, Recommendation Systems, Link Analysis, Natural Language Processing (NLP).

Schedule:

Week	Topic / Contents
1	Cloud Computing Intro to Clouds, Cloud Computing, Popular Commercial Cloud Architectures Parallel processing in python https://www.vanderbilt.edu/accre/jupyter/ https://jupyter.accre.vanderbilt.edu/
2	Big Data Processing Distributed File Systems, HDFS, Hadoop, MapReduce Review course for Linux, ACCRE
3	Apache Spark Spark Architecture, Resilient Distributed DataSets (RDDs), Transformations and actions (slides) SparkContext and RDD Basics.ipynb (notebook)
4	Spark DataFrames <ul style="list-style-type: none"> ● Spark DataFrames Section Introduction (slides) ● Spark DataFrame Basics (notebook) ● Spark DataFrame Basic Operations (notebook) ● Groupby and Aggregate Functions (notebook) ● Missing Data (notebook) ● Dates and Timestamps (notebook)
5	Spark SQL (slides) Tutorial: PySpark, SparkSQL in Google Colab
6	Machine Learning Overview (slides) Midterm
7	Regression <ul style="list-style-type: none"> ● Regression (slides) ● Introduce Linear Regression in PySpark (notebook)

	<ul style="list-style-type: none"> • Data Transformations (Data_Transformations.ipynb, notebook) • Linear Regression Example (Ecommerce Customer Data) (notebook) • Linear Regression Consulting Project • Introduce Logistic Regression in PySpark (notebook) • Data Transformations and Pipeline (MLlib_features_and_pipeline.ipynb, notebook) • Logistic Regression Example (Titanic Data, notebook) • Logistic Regression Consulting Project (notebook)
8	<p>Tree Methods</p> <ul style="list-style-type: none"> • Large Scale Machine Learning: Decision Tree (slides) • Introduce Tree Methods in PySpark (notebook) • Three Tree Models Comparison Example (College data, notebook) • Random Forest Classification Consulting Project (Dog Food, notebook) • Tree Methods Consulting Project (notebook)
9	<p>Locality-Sensitivity Hashing, Clustering</p> <ul style="list-style-type: none"> • Locality-Sensitivity Hashing, Clustering (slides) • Introduce Clustering in PySpark (notebook) • Clustering Example (Iris Dataset, notebook) • Clustering Consulting Project - Customer Segmentation (Fake Data, notebook)
10	<p>Recommendation System</p> <ul style="list-style-type: none"> • Introduction to Recommender Systems and Collaborative Filtering (slides) • Recommendation System Example (MovieLens Dataset, notebook) • Recommendation System Consulting Project (notebook)
11	<p>PageRank, Link Analysis</p> <ul style="list-style-type: none"> • PageRank (slides) • Introduce NetworkX for PageRank (notebook)
12	<p>Natural Language Processing</p> <ul style="list-style-type: none"> • Tools for NLP (notebook) • NLP Example (Naive Bayes Model, Spam Filter, notebook) <p>Panel: Big Data in Industry (Freenome (Healthcare), Oracle, Facebook, Apple, Apollo Education Group (NLP))</p>
13	Final Project Presentations
14	Project Writeup Due