

Schedule (2016-09-20)

All required readings should be completed by the following week.

Date	Topic / Readings	Deadlines
2016-08-30	<p>Introductions; computing setup: Jupyter notebook and command line shell basics; Git and GitHub basics.</p> <p><u>Readings for next week:</u> Required: Software Carpentry Lesson: The Unix Shell, http://swcarpentry.github.io/shell-novice/</p> <p>Required: Roger Peng on Reproducible Research (three videos): http://tinyurl.com/jhu-reproducible-research</p> <p>Optional: Software Carpentry Lesson: Version Control with Git, http://swcarpentry.github.io/git-novice/</p>	Exercise #1, Friday, 9/2, 12pm
2016-09-06	<p>The command line shell: input, output, and pipelines; csvkit; data types.</p> <p><u>Readings</u> Required: Wickham, "Tidy Data." http://vita.had.co.nz/papers/tidy-data.pdf</p> <p>Optional: Data Science at the Command Line, chapters 1-5</p>	Exercise #2, Friday, 9/9, 12pm
2016-09-13	<p>Command line filters in the shell and Python; parallel processing in the shell.</p> <p><u>Readings</u> Required: Software Carpentry Lesson: Using Databases and SQL, Topics 1-5, http://software-carpentry.org/lessons.html</p> <p>Optional: Data Science at the Command Line, chapters 6-8</p>	Project #1, Friday, 9/23, 12pm
2016-09-20	<p>RDBMS: schema, keys, basic SQL operations, aggregate functions.</p> <p><u>Readings</u> Required: Software Carpentry Lesson: Using Databases and SQL, Topics 6-10, http://software-carpentry.org/lessons.html</p> <p>Optional: Learning SQL, chapters 1-4; Database System Concepts, chapters 1-3</p>	Review #1, Tuesday, 9/27, 7pm
2016-09-27	<p>RDBMS: subqueries, joins, integrity, transactions, functions, triggers, schema design and E-R models, normal forms.</p> <p><u>Readings</u> Optional: Learning SQL, chapters 5, 6, 7, 9, 10</p>	Exercise #3, Friday 9/30, 12pm

	<p>Optional: A Gentle Introduction to Algorithm Complexity Analysis (online at http://discrete.gr/complexity/)</p> <p>Optional: Visualizing Algorithms (online at http://bost.ocks.org/mike/algorithms/)</p>	
2016-10-04	<p>RDBMS: advanced SQL, indexes, query processing, analysis, and optimization.</p> <p>Note: no office hours on Tuesday, October 4.</p> <p><u>Readings</u> Required: Star Schema, chapters 1-5</p> <p>Optional: Learning SQL, chapters 12, 13, 14</p>	<p>Exercise #4, Friday 10/7, 12pm</p> <p>Project #2, Friday 10/15, 12pm</p>
2016-10-11	No class	Review #2, Tuesday, 10/18, 7pm
2016-10-18	<p>Warehouses: facts and dimensions, architectures, schemas</p> <p><u>Readings</u> Required: Star Schema, chapters 4-7</p>	Exercise #5, Friday, 10/21, 12pm
2016-10-25	No class (fall break)	
2016-11-01	<p>Warehouses: dimension design</p> <p><u>Readings</u> Required: Star Schema, chapter 11</p> <p>Required: AWS Redshift. https://aws.amazon.com/redshift/</p>	Exercise #6, Friday, 11/4, 12pm
2016-11-08	<p>Warehouses: fact table design</p> <p><u>Readings</u> Required: Dean and Ghemawat, "MapReduce: Simplified Data Processing on Large Clusters." http://research.google.com/archive/mapreduce.html</p> <p>Required: Drake, "Command-line tools can be 235x faster than your Hadoop cluster." http://aadrake.com/command-line-tools-can-be-235x-faster-than-your-hadoop-cluster.html</p> <p>Optional: Chang et al. "Bigtable: A Distributed Storage System for Structured Data." http://research.google.com/archive/bigtable.html</p>	Project #3, Friday, 11/18, 12pm

	Optional: DeCandia et al. "Dynamo: Amazon's Highly Available Key-value Store", http://www.read.seas.harvard.edu/~kohler/class/cs239-w08/de-candia07dynamo.pdf	
2016-11-15	Contemporary data management tools: Hadoop, map/reduce, Dynamo, Trifacta <u>Readings</u> Required: Apache Spark. https://spark.apache.org/ Required: Lambda Architecture. http://lambda-architecture.net/	Exercise #7, Friday, 11/18, 12pm Review #3, Tuesday, 11/22, 7pm
2016-11-22	Contemporary data management tools: Spark introduction <u>Readings</u> Required: CAP theorem. https://en.wikipedia.org/wiki/CAP_theorem Required: Kudu. http://getkudu.io/ Required: AWS Kinesis. https://aws.amazon.com/kinesis/	Exercise #8, Tuesday 11/29, 7pm
2016-11-29	Contemporary data management tools: Spark SQL, DataFrames, MLib, Streaming	Final Project, Friday 12/9, 12pm
2016-12-06	Final Project presentations, course wrap-up	