

## Schedule (updated 2016-11-22)

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, <a href="http://swcarpentry.github.io/shell-novice/">http://swcarpentry.github.io/shell-novice/</a></p> <p>Required: Roger Peng on Reproducible Research (three videos): <a href="http://tinyurl.com/jhu-reproducible-research">http://tinyurl.com/jhu-reproducible-research</a></p> <p>Optional: Software Carpentry Lesson: Version Control with Git, <a href="http://swcarpentry.github.io/git-novice/">http://swcarpentry.github.io/git-novice/</a></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." <a href="http://vita.had.co.nz/papers/tidy-data.pdf">http://vita.had.co.nz/papers/tidy-data.pdf</a></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, <a href="http://swcarpentry.github.io/sql-novice-survey/">http://swcarpentry.github.io/sql-novice-survey/</a></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, <a href="http://swcarpentry.github.io/sql-novice-survey/">http://swcarpentry.github.io/sql-novice-survey/</a></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>	Exercise #3, Friday 9/30, 12pm

	<u>Readings</u> Optional: Learning SQL, chapters 5, 6, 7, 9, 10  Optional: A Gentle Introduction to Algorithm Complexity Analysis (online at <a href="http://discrete.gr/complexity/">http://discrete.gr/complexity/</a> )  Optional: Visualizing Algorithms (online at <a href="http://bost.ocks.org/mike/algorithms/">http://bost.ocks.org/mike/algorithms/</a> )	
2016-10-04	RDBMS: advanced SQL, ETL, indexes, query processing, analysis, and optimization, SQL from Python.  <b>Note:</b> no office hours on Tuesday, October 4.  <u>Readings</u> Required: Star Schema, chapters 1-5  Optional: Learning SQL, chapters 12, 13, 14	
2016-10-11	<b>No class</b>  <b>Note:</b> no office hours on Tuesday, October 11.	
2016-10-18	Warehouses: facts and dimensions, architectures, schemas  <u>Readings</u> Required: Star Schema, chapters 4-7	Exercise #4, Monday, 10/24, 12pm
2016-10-25	<b>No class (fall break)</b>	
2016-11-01	Warehouses: dimension design  <u>Readings</u> Required: Star Schema, chapter 11  Required: AWS Redshift. <a href="https://aws.amazon.com/redshift/">https://aws.amazon.com/redshift/</a>	Project #2, Friday, 11/11, 12pm
2016-11-08	Warehouses: fact table design  <u>Readings</u> Required: Dean and Ghemawat, "MapReduce: Simplified Data Processing on Large Clusters." <a href="http://research.google.com/archive/mapreduce.html">http://research.google.com/archive/mapreduce.html</a>  Required: Drake, "Command-line tools can be 235x faster than your Hadoop cluster." <a href="http://aadrake.com/command-line-tools-can-be-235x-faster-than-your-hadoop-cluster.html">http://aadrake.com/command-line-tools-can-be-235x-faster-than-your-hadoop-cluster.html</a>	Review #2, Tuesday, 11/15, 7pm

	<p>Optional: Chang et al. "Bigtable: A Distributed Storage System for Structured Data."  <a href="http://research.google.com/archive/bigtable.html">http://research.google.com/archive/bigtable.html</a></p> <p>Optional: DeCandia et al. "Dynamo: Amazon's Highly Available Key-value Store",  <a href="http://www.read.seas.harvard.edu/~kohler/class/cs239-w08/de-candia07dynamo.pdf">http://www.read.seas.harvard.edu/~kohler/class/cs239-w08/de-candia07dynamo.pdf</a></p>	
2016-11-15	<p>Contemporary data management tools: Map/Reduce, Hadoop, Trifacta</p> <p><u>Readings</u>  Required: Apache Spark. <a href="https://spark.apache.org/">https://spark.apache.org/</a>  Required: Lambda Architecture.  <a href="http://lambda-architecture.net/">http://lambda-architecture.net/</a></p>	<p>Exercise #5,  Friday, 11/18,  12pm</p> <p>Final Project,  Friday 12/9,  12pm</p>
2016-11-22	<p>Contemporary data management tools: Spark introduction</p> <p><u>Readings</u>  Required: Spark SQL and DataFrames Programming Guide.  <a href="http://spark.apache.org/docs/latest/sql-programming-guide.html">spark.apache.org/docs/latest/sql-programming-guide.html</a>  Optional: CAP theorem.  <a href="https://en.wikipedia.org/wiki/CAP_theorem">https://en.wikipedia.org/wiki/CAP_theorem</a>  Optional: Kudu. <a href="http://getkudu.io/">http://getkudu.io/</a>  Optional: AWS Kinesis. <a href="https://aws.amazon.com/kinesis/">https://aws.amazon.com/kinesis/</a></p>	
2016-11-29	<p>Contemporary data management tools: Spark SQL, DataFrames, MLib, Streaming, Column-oriented storage</p>	<p>Exercise #6,  Friday 12/2,  12pm</p> <p>Final Project presentations  Tuesday, 12/6</p>
2016-12-06	<p>Final Project presentations, course wrap-up</p>	<p>Review #3,  Tuesday, 12/12  12pm</p>