

# CIS 490: Introduction to Programming Techniques for Data Science and Analytics

## Fall 2018

**Hours:** 3 hours; extended course project option (CIS 597/598, 690) available

**Prerequisite:** CIS 200 (Fundamentals of Programming) or CIS 209 (C Programming for Engineers) or instructor permission

**Textbook:** Drabas, T., & Lee, D. (2017). *Learning PySpark*. Birmingham, UK: Packt Publishing. URL: <http://bit.ly/learning-pyspark-mapt>

**Venue:** MWF 12:30 – 13:20 U.S. Central Time, 1116 Engineering Hall (CIS 490: Reference #15130)

**Instructor:** William H. Hsu, Department of Computer Science

Office: 2164 Engineering Hall (DUE) Google Voice (office/home/cell): +1 785 236 8247

TA: CS graduate TA (TBD), 1119 DUE Instructional alias: [bigdataclassta@listserv.ksu.edu](mailto:bigdataclassta@listserv.ksu.edu)

URL: <http://bit.ly/hsu-calendar-color> E-mail: [bhsu@ksu.edu](mailto:bhsu@ksu.edu)

**Office hours:** 10:30 – 11:30, 15:30 – 16:30 Mon, Wed; 08:00 – 09:00 Fri; 09:00 – 10:00 Tue; 18:30 Thu for distance students; by appointment

**Web page:** <http://bit.ly/kstate-datascience-class> (public)

**MediaSite lectures:** Linked from both K-State Canvas (official) and mirror (public)

### Course Description

This is an introductory course on programming techniques for working with big data, including data sets that are terascale and larger as well as complex, heterogeneous data from various domains. It is intended for students who have had at least one introductory programming course and are prepared to learn a new programming language and one or more libraries. No additional background is assumed. The course will survey programming concepts that underlie the MapReduce architecture, their implementation using platforms such as Apache Hadoop, and specific tools for data integration and data transformation such as Apache Hive and Scalding. Basic NoSQL databases and query processing will be presented in the context of real-world problems and students will be given full-scale data sets and problems to work with, along with the opportunity to bring data and problems to work on from other disciplines.

### Course Requirements

Component	Components	Grade Value	Total Value
Exams and quizzes	2 online hour exams	20% (10% each)	45%
	1 online final exam	20%	
	5 of 6 online quizzes	5% (1% each)	
Homework	5 of 6 programs	10% (2% each)	20%
	5 of 6 exercise sets (written)	10% (2% each)	
Term project	Planning/design, interview	8%	25%
	Intermediate milestone	8%	
	Implementation, report	8%	
	Peer review	1%	
Class participation	Forum participation	6%	10%
	Answering questions/discussion	2%	
	10 of 15 labs	2%	

### Recommended text (on reserve in K-State CIS Library)

Scott, J. A. (2015). *Getting Started with Apache Spark: From Inception to Production*. San Jose, CA, USA: MapR Technologies, Inc. URL: <https://mapr.com/ebooks/spark/>

Miner, D., & Shook, A. (2012). *MapReduce Design Patterns: Building Effective Algorithms and Analytics for Hadoop and Other Systems*. Sebastopol, CA, USA: O'Reilly Media. ISBN-10: 1449327176. ISBN-13: 978-1449327170.

Lin, J., Dyer, C., & Hirst, G. (2010). *Data-Intensive Text Processing with MapReduce (Synthesis Lectures on Human Language Technologies)*. San Rafael, CA, USA: Morgan & Claypool Publishers. ISBN-10: 1608453421. ISBN-13: 978-1608453429. URL: <http://bit.ly/linbook>

**Additional bibliography:** excerpted in online course notes and handouts

**Course Calendar, Syllabus, and Readings (textbook chapters unless otherwise indicated)**

Lecture	Date	Topic	Reading (Before Class)
0	Mon 20 Aug 2018	MapReduce/Hadoop 1 of 3: Overview	Missive; Preface, (Ch.) 1
<b>1</b>	<b>Wed 22 Aug 2018</b>	<b>MapReduce/Hadoop 2 of 3: HDFS/MR</b>	<b>2; 1-2 Lin; BigTable Article</b>
2	Fri 24 Aug 2018	MapReduce/Hadoop 3 of 3: Tools	3 Lin; Stats Handout
3	Mon 27 Aug 2018	Python, Amazon Web Services Intro	Python Handout
<b>4</b>	<b>Wed 29 Aug 2018</b>	<b>Python Streaming</b>	<b>Streaming Handout</b>
5	Fri 31 Aug 2018	Introduction to Queries & SQL	SQL Handout; 5
<b>6</b>	<b>Wed 05 Sep 2018</b>	<b>Pig &amp; Hive 1 of 3: Queries</b>	<b>Pig &amp; Hive Handouts; 1</b>
7	Fri 07 Sep 2018	Pig & Hive 2 of 3: Exercises	Pig & Hive Handouts
8	Mon 10 Sep 2018	Pig & Hive 3 of 3; Scala/sbt Preview	FP Basics Handout
<b>9</b>	<b>Wed 12 Sep 2018</b>	<b>PySpark 1 of 3: Overview</b>	<b>PySpark Handout</b>
<b>10</b>	<b>Fri 14 Sep 2018</b>	<b>PySpark 2 of 3; Exam 1 Review</b>	<b>PySpark Handouts</b>
11	Mon 17 Sep 2018	PySpark 3 of 3; AWS Deployment	PySpark Handouts
<b>12</b>	<b>Wed 19 Sep 2018</b>	<b>Project Topics Review</b>	<b>Projects Handout</b>
<b>13</b>	<b>Fri 21 Sep 2018</b>	<b>Exam 1</b>	<b>Avro Handout</b>
14	Mon 24 Sep 2018	MR Design Patterns 1 of 3: Basics	UMAP Papers; 2
<b>15</b>	<b>Wed 26 Sep 2018</b>	<b>MR Design Patterns 2 of 3: Exercise</b>	<b>RecSys Survey, Papers; 3</b>
<b>16</b>	<b>Fri 28 Sep 2018</b>	<b>MR Design Patterns 3 of 3: Apps</b>	<b>Cloudera Handout; 4</b>
17	Mon 01 Oct 2018	NoSQL 1 of 5: Intro, CAP, Sharding	NoSQL Handout; 2-3
<b>18</b>	<b>Wed 03 Oct 2018</b>	<b>NoSQL 2 of 5: Key-Value Stores</b>	<b>Key-Value Handout; 4</b>
19	Fri 05 Oct 2018	NoSQL 3 of 5: Cassandra (Columnar)	Cassandra Handout; 5
20	Mon 08 Oct 2018	NoSQL 4 of 5: MongoDB (Document)	6; Lin 4; Mongo Handout
<b>21</b>	<b>Wed 10 Oct 2018</b>	<b>NoSQL 5 of 5: Neo4j (Graph)</b>	<b>7; Lin 4; Graph Handout</b>
<b>22</b>	<b>Fri 12 Oct 2018</b>	<b>MR Synopsis; Exam 2 Review</b>	<b>8; Lin 5, ETL Handout</b>
23	Mon 15 Oct 2018	Data Mining 1 of 3: Machine Learning	Mahout Handout
<b>24</b>	<b>Wed 17 Oct 2018</b>	<b>Data Mining 2 of 3: sklearn/MMLib, Mahout</b>	<b>Learning Handout</b>
<b>25</b>	<b>Fri 19 Oct 2018</b>	<b>Data Mining 3 of 3: KDD; Exam 2</b>	<b>KDD, CRM/BI Handouts</b>
26	Mon 22 Oct 2018	Search 1 of 5: Indexing	Search Handout; Lin 4
<b>27</b>	<b>Wed 24 Oct 2018</b>	<b>Search 2 of 5: PageRank</b>	<b>3; Lin 6; Search Handout</b>
<b>28</b>	<b>Fri 26 Oct 2018</b>	<b>Search 3 of 5: TFIDF, Text Analytics</b>	<b>Text Analytics Handout</b>
29	Mon 29 Oct 2018	Search 4 of 5: Solr	Solr Handouts
<b>30</b>	<b>Wed 31 Oct 2018</b>	<b>Search 5 of 5: ElasticSearch</b>	<b>ElasticSearch Handout</b>
31	Fri 02 Nov 2018	Mesos, Lambda Architecture	Shard Basics Handout
<b>32</b>	<b>Mon 05 Nov 2018</b>	<b>Synopsis, Final Review</b>	<b>NoSQL Handout</b>
<b>33</b>	<b>Wed 07 Nov 2018</b>	<b>KVP vs. Columnar DBs</b>	<b>NoSQL Handout</b>
34	Fri 09 Nov 2018	Graph Databases 1 of 3 :More Neo4j	Graph & Neo4j Handouts
<b>35</b>	<b>Mon 12 Nov 2018</b>	<b>Graph Databases 2 of 3 :Examples</b>	<b>Neo4j Handout</b>
36	Wed 14 Nov 2017	Graph Databases 3 of 3 :Applications	Graph & Neo4j Handouts
<b>37</b>	<b>Fri 16 Nov 2018</b>	<b>Statistics &amp; Analytics 1 of 3: Descriptive</b>	<b>Lin 5; Stats Handout</b>
38	Mon 26 Nov 2018	Statistics & Analytics 2 of 3: Big Data	Stats Handout
<b>39</b>	<b>Wed 28 Nov 2018</b>	<b>Statistics &amp; Analytics 2 of 3: Info Vis</b>	<b>Visualization Handout</b>
	Fri 30 Nov 2018	Final Projects & Peer Review 1	N/A (Student Handouts)
	<b>Mon 03 Dec 2018</b>	<b>Final Projects &amp; Peer Review 2</b>	<b>N/A (Student Handouts)</b>
<b>40</b>	<b>Wed 05 Dec 2018</b>	<b>Performance &amp; Scalability Issues</b>	<b>Lin 7; HPC Handout</b>
41	Fri 07 Dec 2018	Advanced Big Data Analytics Topics	Advanced Topics Handout
		<b>PROJECT SUBMISSIONS DUE</b>	

Lightly-shaded entries (with bold text) denote the due date of a written problem set.

Heavily-shaded entries denote the due date of a machine problem (programming assignment).

Green-shaded entries denote the date of a lab (distance students may participate online that week).

Green bold-faced text denotes the date of an exam preparation review.

Red bold-faced text denotes the date of an exam (distance students may use online proctoring that week),

Blue bold-faced text denotes the date of a post-exam (model solution) review.

Project proposals are due on the day of Lecture 23. Interim project interviews will be held right after the second hour exam.

The highlighted date is the due date of the draft project report and demo, with interviews and presentations to be held the last two weeks of class.