

**San José State University**  
**Department of Computer Science**  
**CS185C, Solving Big Data Problems, Section 02, Fall 2017**

**Course and Contact Information**

<b>Instructor:</b>	James Casaletto
<b>Office Location:</b>	DH 282
<b>Telephone:</b>	(408) 394-5748
<b>Email:</b>	james.casaletto@sjsu.edu
<b>Office Hours:</b>	Tuesday Thursday 07:00-7:30
<b>Class Days/Time:</b>	Tuesday Thursday 7:30-8:45
<b>Classroom:</b>	SCI 311
<b>Prerequisites:</b>	CS146 or equivalent Java programming experience

**Faculty Web Page and MYSJSU Messaging**

Course materials such as syllabus, handouts, notes, assignment instructions, etc. can be found on the Canvas learning management system course website. You are responsible for regularly checking with the messaging system through MySJSU (or other communication system as indicated by the instructor) to learn of any updates.

**Course Description**

This course is a comprehensive overview of solving big data problems using Apache Hadoop and is comprised of three main parts. The first part of this course explores the core of Apache Hadoop. The second part of the course explores the Apache Hadoop ecosystem. The third part of the course explores topics in machine learning using Apache Spark. All programming assignments and coding examples are in Java.

**Course Goals**

The goal of this course is to develop a working knowledge of solving big data problems from end to end using Apache Hadoop and its ecosystem.

**Course Learning Outcomes (CLO)**

Upon successful completion of this course, students will be able to:

1. Use HDFS to store and retrieve data at scale
2. Write MapReduce programs in Java to transform, filter, and enrich data at scale.
3. Use Hadoop ecosystem components to ingest, transform, store, and analyze big data.
4. Use machine learning algorithms to derive insights from big data

## Required Texts/Readings

### Textbook

No textbook is required for this class.

### Other Readings

A list of other readings will be provided on the CANVAS page associated with this class.

### Other equipment / material

Students are required to have a 64-bit laptop running either Windows, MacOS, or Linux with at least 8GB memory installed, 2 CPU cores, and approximately 30GB disk space free. All labs are to be developed using a virtual machine installed on the laptop with these and other requirements.

## Course Requirements and Assignments

SJSU classes are designed such that in order to be successful, it is expected that students will spend a minimum of forty-five hours for each unit of credit (normally three hours per unit per week), including preparing for class, participating in course activities, completing assignments, and so on. More details about student workload can be found in [University Policy S12-3](http://www.sjsu.edu/senate/docs/S12-3.pdf) at <http://www.sjsu.edu/senate/docs/S12-3.pdf>.

- 6 x individual labs
- 1 x midterm exam
- 1 x team project
- 1 x final exam

### Final Examination or Evaluation

The final exam will cover all the material discussed in class after the midterm exam. The format for the final exam will be all multiple-choice.

## Grading Information

### Determination of Grades

The final grade will be calculated as follows:

$$0.2 * (\text{midterm exam score}) + 0.4 * (\text{average score of labs}) + 0.2 * (\text{final exam score}) + 0.2 * (\text{project score})$$

98-100 => A+, 93-97 => A, 90-92 => A-  
87-89 => B+, 83-86 => B, 80-82 => B-  
77-79 => C+, 73-76 => C, 70-72 => C-  
67-69 => D+, 63-66 => D, 60-62 => D-  
below 60 => F

- There will not be any extra credit provided
- 10% will be deducted for each day that any lab is turned in after it's due.

### Grading Information

“Passage of the Writing Skills Test (WST) or ENGL/LLD 100A with a C or better (C- not accepted), and completion of Core General Education are prerequisite to all SJSU Studies courses. Completion of,

or co-registration in, 100W is strongly recommended. A minimum aggregate GPA of 2.0 in GE Areas R, S, & V shall be required of all students.”

### Classroom Protocol

Class begins promptly at 7:30 and ends abruptly at 8:45. Please silence all cell phones during class. Your active participation in the lecture discussions is greatly encouraged.

### University Policies

Per University Policy S16-9, university-wide policy information relevant to all courses, such as academic integrity, accommodations, etc. will be available on Office of Graduate and Undergraduate Programs’ [Syllabus Information web page](http://www.sjsu.edu/gup/syllabusinfo/) at <http://www.sjsu.edu/gup/syllabusinfo/>”

## CS185C / Solving Big Data Problems, Sp17, Course Schedule

*The following schedule is subject to change. All changes will be verbally communicated in class as well as messaged using the CANVAS system.*

### Course Schedule

Week	Date	Topics, Readings, Assignments, Deadlines
1	24.aug	Introduction to this course
2	29.aug	Introduction to big data (lab 1 assigned)
2	31.aug	Building and managing a 1-node virtual machine Hadoop cluster
3	5.sep	Developing big data applications using Eclipse and Maven (lab 2 assigned)
3	7.sep	Introduction to distributed file systems I: HDFS architecture
4	12.sep	Introduction to distributed file systems II: HDFS CLI and API (lab 3 assigned)
4	14.sep	MapReduce programming I: conceptual overview
5	19.sep	MapReduce programming II: writing a MapReduce program (lab 4 assigned)
5	21.sep	Introduction to the Hadoop ecosystem I: overview
6	26.sep	Introduction to the Hadoop ecosystem II: solution design
6	28.sep	Apache Spark I: RDD, data sets, data frames
7	3.oct	Apache Spark II: transformations and actions
7	5.oct	Apache Spark III: developing Spark applications (lab 5 assigned)
8	10.oct	Apache Kafka I: consumers, producers, topics, keys, values
8	12.oct	Apache Kafka II: partitioning, cursor management, offsets
9	17.oct	Putting it all together: using Spark with Kafka on Hadoop (lab 6 assigned)
9	19.oct	Midterm exam
10	24.oct	Introduction to data science I: brief history of machine learning

<b>Week</b>	<b>Date</b>	<b>Topics, Readings, Assignments, Deadlines</b>
10	26.oct	Introduction to data science II: overview of ML algorithms
11	31.oct	Linear regression I: overview, use cases, demo
11	2.nov	Linear regression II: algorithm
12	7.nov	Decision trees and random forests: overview, use cases, demo (group project proposal due)
12	9.nov	Decision trees and random forests: algorithm
13	14.nov	K-means clustering I: overview, use cases, demo
13	16.nov	K-means clustering II: algorithm
14	21.nov	Putting it all together: using Spark MLlib pipelines
14	23.nov	Thanksgiving holiday – no class
15	28.nov	Group project presentations I
15	30.nov	Group project presentations II
16	5.dec	Group project presentations III
16	7.dec	Group project presentations IV
17	15.dec	Final exam 07:15 – 09:30