

Introduction to Computational Thinking and Data Science

USC Viterbi School
of Engineering

INF 549

Term: Fall 2018

Syllabus

Term: Fall 2018

Units: 4

Time: Tues-Thur 10am-11:50pm

Location: [Waite Phillips Hall \(WPH\)](#) 207

Instructor: Dr. Yolanda Gil

Office: GER 207

Office Hours: Tuesdays 12pm-1pm

Contact: gil@isi.edu

Instructor: Dr. Gale Lucas

Office: GER 207

Office Hours: Thursdays 12pm-1pm

Contact Info: lucas@ict.usc.edu

Grader: Mengyue Huan

Contact Info: mhuan@usc.edu

Catalogue Course Description

Introduction to data analysis techniques and associated computing concepts for non-programmers. Topics include foundations for data analysis, visualization, parallel processing, metadata, provenance, and data stewardship.

Expanded Course Description

This course will teach non-programmers to think in computing terms about modern topics, and to approach real-world phenomena through data science. The course will enable students to:

- Acquire computational thinking skills that will enable students to represent and reason about complex problems in the digital arena
- Understand different kinds of data in terms of their possibilities and limitations to approach complex problems cast in terms of the emerging field of data science
- Become data science scholars through best practices in data documentation and dissemination

The course is intended for students in disciplines outside of computer science, so no prior experience with computer science is assumed. The course topics will be particularly relevant to students interested in physical sciences and social sciences.

This class will include eight homework assignments and a final exam.

Learning Objectives

This course teaches non-programmers to think in computing terms about modern topics, and to approach real-world phenomena through data science. The course introduces different kinds of data and corresponding approaches to data analysis, including geospatial data, time series, networks, and multimedia data. Students learn to run multi-step analysis through a graphical workflow interface, and will experience first hand complex concepts in data science such as parallel computing, provenance, and visualization. Students also learn to use ontologies and logic representations to capture metadata and other knowledge about complex data. The course includes practical lessons to use workflow and ontology development toolkits, as well as best practices for data stewardship and dissemination.

Prerequisite(s): none

Co-Requisite (s): none

Recommended Preparation: Mathematics and logic undergraduate courses.

Required Readings and Supplementary Materials

There is no textbook. Handouts of all required readings will be made freely available to students electronically. All required software is freely available for students to install on their personal computers or to access through a web interface.

Representative course readings that will be used include:

- “Computational Thinking.” J. M. Wing. Communications of the ACM, viewpoint, vol. 49, no.3, March 2006.
- “Data Science in the News: Advances and Challenges for the Era of Big Data.” Kate Musen, Alyssa Deng, Taylor Alarcon, Yolanda Gil. Technical Report ISI-TR-702, Information Sciences Institute, University of Southern California. August 24, 2015.
- “Ten Simple Rules for the Care and Feeding of Scientific Data.” Goodman, A.; Pepe, A.; Blocker, A. W.; Borgman, C. L.; Cranmer, K.; Crosas, M.; Stefano, R. D.; Gil, Y.; Groth, P.; Hedstrom, M.; Hogg, D. W.; Kashyap, V.; Mahabal, A.; Siemiginowska, A.; and Slavkovic, A. PLOS Computational Biology, 10, 2014.
- “Intelligent Workflow Systems and Provenance-Aware Software.” Y. Gil. Proceedings of the Seventh International Congress on Environmental Modeling and Software, San Diego, CA, 2014.
- “Data Science for Business”, Foster Provost and Tom Fawcett. O’Reilly Media publishers, 2013.
- “A Primer for the PROV Provenance Model.” Gil, Y.; Miles, S.; Belhajjame, K.; Deus, H.; Garijo, D.; Klyne, G.; Missier, P.; Soiland-Reyes, S.; and Zednik, S. World Wide Web Consortium (W3C) Technical Report, 2013.
- “The Ethics of Data Sharing and Reuse in Biology.” Duke, C. S., & Porter, J. H. BioScience, 63(6), 483–489, 2013. doi:10.1525/bio.2013.63.6.10

Description and Assessment of Homework Assignments

There will be a homework assignment every 3 or 4 lectures. The homeworks include a class project that will be developed by the students independently in 3 separate stages, getting feedback from the instructors at each stage. The assignments must be submitted individually and students will receive individual scores. Students may work in groups to complete the tasks. The homework assignments are expected to take 6-8 hours. Each assignment is graded on a scale of 0-100 and the grading criteria will be specified in each assignment. The homework topics are listed in the Course Schedule.

Syllabus and Class Schedule

	Topic	Material Covered	Homework assigned
Section I: Introduction to Computational Thinking and Data Science			
1	Computational thinking and data science	<ul style="list-style-type: none"> • What is computational thinking • Computational thinking for reasoning and analysis • What is data science • Data scientists • The context of data science 	
2	Data	<ul style="list-style-type: none"> • What is data • What is not (yet) data • Time series data • Networked data • Geospatial data • Text data • Labeled and annotated data • Big data 	HW1: Project part 1 – Finding data
3	Data analysis software	<ul style="list-style-type: none"> • Programs for data analysis • Inputs and Outputs • Program Parameters • Programming Languages • Programs as Black Boxes • Algorithms versus software 	
4	Multi-step data analysis as workflows	<ul style="list-style-type: none"> • Building workflows by composing software • Pre-processing and post-processing data • Workflows for data analysis • Workflow inputs and parameters • Executing workflows • Exploring data through workflows • Workflows in practice 	
5	Workflow practicum	<ul style="list-style-type: none"> • The WINGS workflow system • Workflows in practice 	Homework HW2: Exploring data analysis workflows
Section II: Data Analysis			
6	Basic statistics	<ul style="list-style-type: none"> • Descriptive statistics • Inferential statistics • Consuming statistical results 	
7	Data analysis tasks (I)	<ul style="list-style-type: none"> • Data analysis tasks in data mining, statistics, and machine learning • Supervised learning 	

		<ul style="list-style-type: none"> ○ Classification tasks ○ Classification algorithms ○ Evaluation of classifiers 	
8	Data analysis tasks (II)	<ul style="list-style-type: none"> ● Unsupervised learning <ul style="list-style-type: none"> ○ Clustering ○ Pattern detection ○ Anomaly detection ● Simulation and prediction 	
9	Data analysis tasks (III)	<ul style="list-style-type: none"> ● Causality <ul style="list-style-type: none"> ○ Probabilistic graphical models ○ Bayesian networks ○ Causal models 	
Section III: Data Analysis in Practice			
10	Analyzing different kinds of data (I)	<ul style="list-style-type: none"> ● Analyzing text data <ul style="list-style-type: none"> ○ Pre-processing text ○ Document classification ○ Document clustering ○ Topic detection ○ Sentiment analysis 	
11	Analyzing different kinds of data (II)	<ul style="list-style-type: none"> ● Analyzing time series data <ul style="list-style-type: none"> ○ Collecting time series data ○ Pre-processing time series data ○ Event detection ○ Granger causality 	Homework HW3: Analyzing data with workflows
12	Analyzing different kinds of data (III)	<ul style="list-style-type: none"> ● Analyzing network data <ul style="list-style-type: none"> ○ Network structure ○ Dynamic networks ○ Scale-free networks ○ Network analysis 	
13	Analyzing different kinds of data (IV)	<ul style="list-style-type: none"> ● Analyzing multimedia data <ul style="list-style-type: none"> ○ Pre-processing images ○ Segmentation ○ Edge detection ○ Object detection ○ Video analysis ● Analyzing geospatial data <ul style="list-style-type: none"> ○ Coordinate systems ○ GIS systems 	
Section IV: User interfaces and user studies			
14	Data visualization	<ul style="list-style-type: none"> ● Quality of visualizations ● Major types of visualizations 	Homework HW3: Data visualization

		<ul style="list-style-type: none"> • Time series visualizations • Geospatial visualizations • Multi-dimensional spaces • Network visualizations 	
15	User experience and user interfaces	<ul style="list-style-type: none"> • UX/UI Design Principles • AB testing • Basics of user study design 	Homework HW5: Project part 2 – Design of data analysis approach
16	User studies	<ul style="list-style-type: none"> • User study design • Null hypothesis significance testing • Advanced analysis for experiments 	
17	Causal claims from user studies	<ul style="list-style-type: none"> • Correlational research • Comparing correlational research to experiments • Ensuring internal validity 	
Section V: Data analysis at large scale			
18	Parallel and distributed computing for big data (I)	<ul style="list-style-type: none"> • Cost of computation • Divide and conquer • Speedup with parallel processing • Limits of speedup: Critical path • Amdahl's law • When problems are not parallelizable 	
19	Parallel and distributed computing for big data (II)	<ul style="list-style-type: none"> • Multi-core computing • Distributed computing • Cluster computing • Cloud computing • Grid computing • Virtual machines • Web services • Practical concerns in distributed computing • Parallel programming languages <ul style="list-style-type: none"> ○ MapReduce/Hadoop 	Homework HW6: Data analysis with parallel processing
Section VI: Metadata			
20	Semantic metadata	<ul style="list-style-type: none"> • What is metadata • Basic metadata versus semantic metadata • Metadata about data collection • Metadata about data processing • Metadata for search and retrieval • Metadata standards 	

		<ul style="list-style-type: none"> • Domain metadata and ontologies 	
21	Ontologies (I)	<ul style="list-style-type: none"> • What is an ontology • Taxonomies and class inheritance • Properties • Logical constraints 	
22	Ontologies (II)	<ul style="list-style-type: none"> • Logical reasoning and inference • Expressivity and computation • The Semantic Web 	
23	Ontologies (III)	<ul style="list-style-type: none"> • Practicum: the PROTÉGÉ ontology editor 	Homework HW7: Developing ontologies
Section VII: Data dissemination			
24	Provenance	<ul style="list-style-type: none"> • What is provenance • Provenance concerning objects • Provenance concerning people and institutions • Provenance concerning processes • Provenance models • Provenance standards 	
25	Data standards and data stewardship	<ul style="list-style-type: none"> • Data formats and standards • Data repositories and services • Data sharing • Data identifiers • Licenses for data • Data citation and attribution • Software and other work products 	
26	Data lifecycle	<ul style="list-style-type: none"> • Data collection and storage • Data cleaning • Data extraction and querying • Data preparation • Quality control • Data integration 	Homework HW8: Project part 3 – Management plan and final report
Section VIII: Advanced Topics			
27	Advanced topics (I)	<ul style="list-style-type: none"> • Multidisciplinary collaborations 	
28	Advanced topics (II)	<ul style="list-style-type: none"> • Privacy and ethics 	
29	Review	<ul style="list-style-type: none"> • Selected topics 	

Final Exam

The final exam will be on Tuesday December 11 at 8am-10am. The last lecture will be a final review of the material.

Assignment Submission Policy

Homework assignments are due at 11:59pm on the due date and should be submitted in Blackboard. Homework will be accepted up to one week late as long as the student requested a late submission ahead of the deadline, and in that case the assignment will be graded at 20% less than the possible points for the assignment. After one week, the assignment will not be graded.

Grading Breakdown

Quizzes: There will be weekly quizzes based on the material from the week before. There is no mid-term for this class.

Homework: There will be eight homework assignments throughout the course.

Final Exam: There is a final exam at the end of the semester covering all of the material covered in the class.

Grading Schema:

Quizzes	20%
Homework assignments	50%
Class participation	10%
Final:	20%
<hr/>	
Total	100%

Grades will range from A through F. The following is the breakdown for grading:

94 - 100 = A	74 - 76 = C
90 - 93 = A-	70 - 73 = C-
87 - 89 = B+	67 - 69 = D+
84 - 86 = B	64 - 66 = D
80 - 83 = B-	60 - 63 = D-
77 - 79 = C+	Below 60 is an F

Academic Conduct and Support Systems

Honor Code

In response to recommendations made by the Academic Integrity Task Force to the Dean, the USC Viterbi School of Engineering now has an Honor Code. The Code was developed by Viterbi students, and its text is as follows:

Engineering enables and empowers our ambitions and is integral to our identities. In the Viterbi community, accountability is reflected in all our endeavors.

Engineering+ Integrity.

Engineering+ Responsibility.

Engineering+ Community.

Think good. Do better. Be great.

These are the pillars we stand upon as we address the challenges of society and enrich lives.

During your time here at Viterbi, please know that academic and personal resources are available to help you:

- The student-driven and student-written Honor Code is here: <http://viterbi.usc.edu/academics/integrity/>.

- An introductory video is posted at <https://myviterbi.usc.edu/> under the link "Academic Integrity Introduction" and serves as a reminder of the school's emphasis in maintaining a high level of academic integrity.
- Master's and PhD students can contact the GAPP office in OHE 106 (<https://gapp.usc.edu/>) for other helpful resources.
- The Viterbi Academic and Resource Center (VARC) (<http://viterbi.usc.edu/students/undergrad/varc>) has a variety of services available.

Academic Integrity

The Viterbi School takes academic integrity violations seriously. Most of the violations that have been reported in the past fall into four categories: unauthorized collaboration, plagiarism, code sharing, and cheating on an exam. Specifically:

- Unauthorized collaboration - Unauthorized collaboration on a project, homework or other assignment. (section 11.14.B) All homework assignments must be individually developed. Students that collaborate on assignments will be referred to the Academic Integrity Coordinator.
- Plagiarism - presenting someone else's ideas as your own, either verbatim or recast in your own words - is a serious academic offense with serious consequences.
- Code sharing - Obtaining for oneself or providing for another person a solution to homework, a project or other assignment, without the knowledge and expressed consent of the instructor. (section 11.14.A)
- Cheating in an exam - this may involve a number of violations, such as looking at class notes during the exam, looking at other student's exam, "texting" with other students during the exam. See the section titled Two Exams for a list of specific violations.

Please note that that these are only the basic violations that we have encountered in the past, and there are many more. Please familiarize yourself with the discussion of plagiarism in SCampus in Section B.11.00, Behavior Violating University Standards and Appropriate Sanctions available at <https://scampus.usc.edu/b/11-00-behavior-violating-university-standards-and-appropriate-sanctions/>.

All academic integrity violations will be referred to the Academic Integrity Coordinator of the Viterbi School of Engineering. The process for adjudicating these cases is available in SCampus, Part B, Section 13.

Other Misconduct

Other forms of academic dishonesty are equally unacceptable. See additional information in SCampus and university policies on scientific misconduct, <http://policy.usc.edu/scientific-misconduct/>.

Discrimination, sexual assault, and harassment are not tolerated by the university. You are encouraged to report any incidents to the Office of Equity and Diversity <http://equity.usc.edu/> or to the Department of Public Safety <http://capsnet.usc.edu/department/department-public-safety/online-forms/contact-us>. This is important for the safety whole USC community. Another member of the university community - such as a friend, classmate, advisor, or faculty member - can help initiate the report, or can initiate the report on behalf of another person. The Center for Women and Men <http://www.usc.edu/student-affairs/cwm/> provides 24/7 confidential support, and the sexual assault resource center webpage <http://sarc.usc.edu> describes reporting options and other resources.

Support Systems

A number of USC's schools provide support for students who need help with scholarly writing. Check with your advisor or program staff to find out more. Students whose primary language is not English should check with the American Language Institute <http://dornsife.usc.edu/ali> which sponsors courses and workshops specifically for international graduate students. The Office of Disability Services and Programs http://sait.usc.edu/academicsupport/centerprograms/dsp/home_index.html provides certification for students with disabilities and helps arrange the relevant accommodations. If an officially declared emergency makes travel to campus infeasible, USC Emergency Information <http://emergency.usc.edu/> will

provide safety and other updates, including ways in which instruction will be continued by means of blackboard, teleconferencing, and other technology.

Diversity

The diversity of the participants in this course is a valuable source of ideas, problem solving strategies, and engineering creativity. The instructors encourage and support the efforts of all of our students to contribute freely and enthusiastically. As members of an academic community, it is our shared responsibility to cultivate a climate where all students and individuals are valued and where both they and their ideas are treated with respect, regardless of their differences, visible or invisible.

Students with Disabilities

Any student requesting academic accommodations based on a disability is required to register with Disability Services and Programs (DSP) each semester. A letter of verification for approved accommodations can be obtained from DSP. Please be sure the letter is delivered to me (or to TA) as early in the semester as possible. DSP is located in STU 301 and is open 8:30 a.m. - 5:00 p.m., Monday through Friday. Website and contact information for DSP: http://sait.usc.edu/academicsupport/centerprograms/dsp/home_index.html, (213) 740-0776 (Phone), (213) 740-6948 (TDD only), (213) 740-8216 (FAX), ability@usc.edu.

Emergency Preparedness/Course Continuity in a Crisis

In case of a declared emergency if travel to campus is not feasible, USC executive leadership will announce an electronic way for instructors to teach students in their residence halls or homes using a combination of Blackboard, teleconferencing, and other technologies.