# Introduction to Data Science: CptS 483–04
## Syllabus

## Course Information

Credit Hours: 3
Semester: Fall 2017
Meeting times and location: MWF, 12:10–13:00, Sloan 38
Course website: https://scads.eecs.wsu.edu/index.php/data-science/

The course website will be used to post relevant course material—including this syllabus—and course related resources. Additionally, the online portal OSBLE+ (https://plus.osble.org) will be used for posting lecture material, assignments, announcements, and messages; and for handling student submissions and instructor feedbacks.

## Instructor Information

Assefaw Gebremedhin
Office: EME B43
Email: assefaw AT eecs DOT wsu DOT edu
Homepage: www.eecs.wsu.edu/~assefaw

**Office hours**: Wednesdays 1:00–2:30pm, or by appointment.

## Course Description

Data Science is the study of the generalizable extraction of knowledge from data. Being a data scientist requires an integrated skill set spanning computer science, mathematics, statistics, and domain expertise along with a good understanding of the art of problem formulation to engineer effective solutions. The purpose of this course is to introduce students to this rapidly growing field and equip them with some of its basic principles and tools as well as its general mindset, using primarily the statistical computing language R.

**Topics to be covered include**: the data science process, exploratory data analysis, linear regression, classification, clustering, principal components analysis, data wrangling, data visualization, time-series data mining, recommender systems, social network mining, deep learning, and data and ethics. The focus in the treatment of these topics is on breadth, rather than depth, and emphasis is placed on integration and synthesis of concepts and their application to solving problems. Necessary theoretical abstractions (mathematical and algorithmic) are introduced as and when needed.

## Audience

The course is suitable for upper-level undergraduate or graduate students in computer science, engineering, applied mathematics, the sciences, business, and related analytic fields.

## Prerequisites

Students are expected to have basic knowledge of algorithms and reasonable programming experience (equivalent to completing a data structures course such as CptS 223), and some familiarity with basic linear algebra (e.g. solution of linear systems and eigenvalue/vector computation) and

basic probability and statistics. *If you are interested in taking the course, but are not sure if you have the right background, talk to the instructor. You may still be allowed to take the course if you are willing to put in the extra effort to fill in any gaps.*

## Course Work

The course consists of several elements: lectures (three times a week, 50 min each) that incorporate labs as needed; in-class exercises; a set of homework assignments; a substantial semester project; and an exam. Below is how the coursework and assessment is broken down.

- **Assignments (30%)**. There will be about four assignments spread through the semester. Each assignment will have one major topic of emphasis. Assignments are to be completed and submitted individually. Each assignment will carry equal weight. Together all assignments account for 30% of final grade.
- **Semester Project (30%)**. Students, working in teams of two or three, will complete a semester project. A project could take one of several forms: analyzing an interesting dataset using existing methods and software tools; building your own data product; or creating a visualization of a complex dataset. Students will be given an opportunity to choose from a list of projects the instructor provides or propose their own project. Guidelines for what constitutes a project will be provided by the instructor. A project will culminate in a written report and a presentation in class. General guidelines for how to prepare a report will be provided by the instructor.
- **Exam (30%)**. There will be one exam designed to a) cover most of the material in class and b) complement the assignments and semester project.
- **Class participation (10%)**. Attendance and active class participation (in discussions, in-class exercises and thought experiments) is required. It will count to 10% of the final grade.

## Expectations for Student Effort

For each hour of lecture equivalent, students should expect to have a minimum of two hours of work outside class.

## Grading

Letter grades will be given according to the following ranges:

A (93–100%), A- (90–92.99%), B+ (87–89.99%), B (83–86.99%), B- (80–82.99%), C+ (77–79.99%), C (70–76.99%), C- (67–69.99%), D (60–66.99%), F (less than 60%).

## Learning Outcomes and Assessment

| Student Learning Outcomes. *By the end of the course, students should be able to:* | Course Topics/Dates. *The following topics/dates will address this outcome:* | Evaluation. *This outcome will be evaluated primarily by:* |
| --- | --- | --- |
| Describe what Data Science is and the skill sets needed | What is Data Science? (week 1); Statistical Learning (weeks 2, 3) | Assignments; Exam |
| Describe the Data Science Process | EDA and the Data Science Process (weeks 3, 4 ) | Assignments; Exam; Project |
| Use R to carry out basic statistical modeling and analysis | Intro to R (week 2); Most subsequent topics throughout the semester | Assignments; Project |
| Carry out exploratory data analysis | EDA (weeks 3, 4) | Assignments; Project |
| Apply basic machine learning algorithms for predictive modeling | Linear Regression (week 4); Classification I, II (weeks 5, 6); | Assignments; Project; Exam |
| Apply learning methods to discover patterns, trends and anomalies in data | Unsupervised learning (week 7); Time Series Data Mining (week 10) | Assignments; Project; Exam |
| Use effective data wrangling approaches to manipulate data | Data Wrangling (week 8) | Assignments; Project |
| Identify and explain fundamental mathematical and algorithmic ingredients that constitute a Recommender System | Recommender Systems (week 11) | Exam; Project |
| Create effective visualization of data (to communicate or persuade) | Data Visualization (week 9) | Assignments; Project |
| Reason around ethical and privacy issues in data science conduct and apply ethical practices | Data and Ethics (week 13) | In-class exercise |
| Work effectively in teams on data science projects | | Project |
| Apply knowledge gained in the course to carry out a project and write a technical report | | Project |

## Detailed Topics and Course Outline

1. Introduction: What is Data Science?
   - Big Data and Data Science hype – and getting past the hype
   - Current landscape of perspectives
   - Skill sets needed

2. Statistical Learning and Intro to R
   - Statistical learning overview
   - Assessing model accuracy
   - Intro to R (combined with lab sessions):
     Basic commands, Graphics, Indexing data, Loading data,
     Graphical and numerical summaries

3. Exploratory Data Analysis and the Data Science Process
   - Basic tools (plots, graphs and summary statistics) of EDA
   - Philosophy of EDA
   - The Data Science Process

4. Linear Regression
   - Simple linear regression
   - Multiple linear regression
   - Extensions of the linear model
   - Lab session on linear regression

5. Classification
   - Overview of classification
   - K-Nearest Neighbors (KNN)
   - Logistic regression
   - Naive Bayes classifier
   - Decision Trees
   - Lab session on classification methods

6. Unsupervised Learning
   - K-means clustering
   - Hierarchical clustering
   - Principal Components Analysis (PCA)
   - Lab session on clustering and PCA

7. Data Wrangling
   - Data cleaning, data reshaping, data integration
   - dplyr, tidyr

8. Data Visualization
   - Telling story with data
   - Choosing tools to visualize data
   - Visualizing patterns over time
   - Visualizing proportions
   - Visualizing relationships
   - Visualizing text information

9. Time Series Data Mining Basics
   - Examples of areas where time series data arise
   - Distance measures
   - Transformations (dimensionality reduction)
   - Algorithms (motif discovery, anomaly detection, segmentation, classification, clustering)
   - Tools: Matrix Profile and SAX

10. Recommender Systems and Social Network Mining
    - Collaborative filtering models
    - Content-based recommender systems
    - Knowledge-based recommender systems
    - Demographic recommender systems
    - Dimension Reduction
    - Social Network Mining

11. Intro to Deep Learning
    - What is deep learning?
    - Overview of deep networks
    - TensorFlow

12. Data Science and Ethical Issues
    - Discussions on privacy, security, ethics
    - A look back at Data Science

## Books

There is no required "textbook" for this course. Select chapters from the followings books will be used as a starter to guide the discussions, but they will be supplemented with instructor-developed lecture notes and reading assignments from other sources. The lecture notes and reading material will be made available on the OSBLE+ page of the course as the course proceeds.

- Gareth James, Daniela Witten, Trevor Hastie and Robert Tibshirani. *An Introduction to Statistical Learning with Applications in R.* Springer, 2013. ISBN 978-1461471370.
  The book is freely available online at: http://www-bcf.usc.edu/~gareth/ISL/).
- Cathy O'Neil and Rachel Schutt. *Doing Data Science, Straight Talk From The Frontline.* O'Reilly. 2014. ISBN 978-1-449-35865-5.
- Jure Leskovek, Anand Rajaraman and Jeffrey Ullman. *Mining of Massive Datasets.* v2.1, Cambridge University Press. 2014.
  The book is freely available online at: http://www.mmds.org/#ver21).
- Jiawei Han, Micheline Kamber and Jian Pei. *Data Mining: Concepts and Techniques.* Third Edition. Morgan Kaufmann Publishers. 2012. ISBN 978-0-12-381479-1.
- Ethem Alpaydin. *Introduction to Machine Learning.* Third Edition. MIT Press, 2014. ISBN 978-0-262-02818-9.
- Nathan Yau. *Visualize This: The FlowingData Guide to Design, Visualization, and Statistrics.* Wiley Publications, 2011. ISBN-13: 978-0470944882.
- Ian Goodfellow, Yoshua Bengio and Aaron Courville. *Deep Learning.* MIT Press, 2016. ISBN 9780262035613. The book is freely available online at: http://www.deeplearningbook.org

## Weekly Schedule

See Table 1 for a weekly schedule of topics and assignments.

| Week | Topics | Assignments |
|------|--------|-------------|
| 01 (Aug 21) | What is Data Science | Survey out |
| 02 (Aug 28) | Statistical Learning, R | Survey due, Assignment 1 out |
| 03 (Sep 04) | Exploratory Data Analysis, R | Assignment 1 due |
| 04 (Sep 11) | Linear Regression | Assignment 2 out |
| 05 (Sep 18) | Classification I | Assignment 2 due |
| 06 (Sep 25) | Classification II | Assignment 3 out |
| 07 (Oct 02) | Unsupervised Learning | Assignment 3 due |
| 08 (Oct 09) | Data Wrangling, Project setup | Project proposal out |
| 09 (Oct 16) | Data Visualization | |
| 10 (Oct 23) | Time Series Data Mining | Project proposal due |
| 11 (Oct 30) | Recommender Systems, Social Networks | Mid-term Exam |
| 12 (Nov 06) | Intro to Deep Learning | Assignment 4 out |
| 13 (Nov 13) | Ethics, Wrap-up | Assignment 4 due |
| 14 (Nov 20) | Thanksgiving break | |
| 15 (Nov 27) | Project presentations | |
| 16 (Dec 04) | Project presentations | Final project report due |

Table 1: Tentative week-by-week schedule of topics and assignments.

## Policies

### Conduct

Students are expected to maintain a professional and respectful classroom environment. In particular, this includes:

- silencing personal electronics
- arriving on time and remaining throughout the class

### Correspondence

All class related correspondence with the instructor will be made via OSBLE+. I will check check messages sent to my Inbox or posted to the Dashboard on a regular basis, and will do my best to respond promptly. Students are encouraged to choose their OSBLE+ settings so that they get emails notifications when messages are sent or posted.

### Attendance

Regular attendance is required. While students may miss class for urgent reasons, absences that are not cleared with the instructor will factor into the Class Participation portion of the semester grade.

**Missing or late work**

Submissions will be handled via the OSBLE page of the course. Students are expected to submit assignments by the specified due date and time. Assignments turned in up to 48 hours late will be accepted with a 10% grade penalty per 24 hours late. Except by prior arrangement, missing or work late by more than 48 hours will be counted as a zero.

**Academic Integrity**

Academic integrity is the cornerstone of higher education. As such, all members of the university community share responsibility for maintaining and promoting the principles of integrity in all activities, including academic integrity and honest scholarship. Academic integrity will be strongly enforced in this course. Any student who violates the University's standard of conduct relating to academic integrity will receive an F as a final grade in this course, will not have the option to withdraw from the course and will be reported to the Office of Student Standards and Accountability. Cheating is defined in the Standards for Student Conduct WAC 504-26-010 (3). You can learn more about Academic Integrity on the WSU campus at http://conduct.wsu.edu. Please also read this link carefully: EECS Academic Integrity Policy (http://www.eecs.wsu.edu/~schneidj/Misc/academic-integrity.html). Use these resources to ensure that you do not inadvertently violate WSU's standard of conduct.

**Safety on Campus**

Washington State University is committed to enhancing the safety of the students, faculty, staff, and visitors. It is highly recommended that you review the Campus Safety Plan (http://safetyplan.wsu.edu/) and visit the Office of Emergency Management web site (http://oem.wsu.edu/) for a comprehensive listing of university policies, procedures, statistics, and information related to campus safety, emergency management, and the health and welfare of the campus community.

**WSU Classroom Safety**

Classroom and campus safety are of paramount importance at Washington State University, and are the shared responsibility of the entire campus population. WSU urges students to follow the "Alert, Assess, Act" protocol for all types of emergencies and "Run, Hide, Fight" response for an active shooter incident. Remain ALERT (through direct observation or emergency notification), ASSESS your specific situation, and act in most appropriate way to assure your own safety (and the safety of others if you are able).

Please sign up for emergency alerts on your account at MyWSU. For more information on this subject, campus safety and related topics, please view the FBI's Run, Hide, Fight video (https://www.fbi.gov/about-us/cirg/active-shooter-and-mass-casualty-incidents/run-hide-fight-video) and visit the WSU safety portal (https://faculty.wsu.edu/classroom-safety).

**Students with Disabilities**

Reasonable accommodations are available for students with a documented disability. If you have a disability and need accommodations to fully participate in this class, please either visit or call the

Access Center (Washington Building 217; 509-335-3417) to schedule an appointment with an Access Advisor. All accommodations must be approved through the Access Center. For more information, consult the webpage http://accesscenter.wsu.edu or email at Access.Center@wsu.edu.

## Important Dates and Deadlines

Students are encouraged to refer to the academic calendar often to be aware of critical deadlines throughout the semester. The academic calendar can be found at http://registrar.wsu.edu/academic-calendar.

## Weather Policy

For emergency weather closure policy, consult: http://alert.wsu.edu.

## Changes

This syllabus is subject to change. Updates will be posted on the course website.