

15-388/688 - Practical Data Science: Introduction

J. Zico Kolter
Carnegie Mellon University
Spring 2021

Outline

What is data science?

What is data science not?

(A few) data science examples

Course objectives and topics

Course logistics

Outline

What is data science?

What is data science not?

(A few) data science examples

Course objectives and topics

Course logistics

Some possible definitions

Data science is the application of computational and statistical techniques to address or gain insight into some problem in the real world

Some possible definitions

Data science is the application of **computational** and **statistical** techniques to address or gain insight into some problem in the **real world**

Some possible definitions

Data science = statistics +
data processing +
machine learning +
scientific inquiry +
visualization +
business analytics +
big data + ...

Data science is the best job in America

The screenshot shows the Glassdoor website interface. At the top, there is a green navigation bar with the Glassdoor logo and links for Jobs, Companies, Salaries, and Interviews. A search bar is located in the center of the navigation bar. To the right, there are links for Sign In, a notification bell, and a plus sign. Below the navigation bar, there is a large banner image with the text "25 Best Jobs in America".

On the left side, there is a sidebar with the following categories:

- Employees' Choice Awards
- Other Lists
- Oddball Interview Questions
- Best Jobs
- Best Cities for Jobs
- Trends
- Additional Resources

The main content area displays the "25 Best Jobs in America" list for 2016. The first job listed is "Data Scientist". The list includes the following information for each job:

- Job Openings
- Median Base Salary
- Career Opportunity
- Job Score

The "Data Scientist" job has the following statistics:

- Job Openings: 1,736
- Median Base Salary: \$116,840
- Career Opportunity: 4.1
- Job Score: 4.7

The page also includes a "United States" dropdown menu and a "2016" dropdown menu. There are social media sharing icons for Facebook, Twitter, LinkedIn, and Email, along with a "2.5k Shares" indicator.

Outline

What is data science?

What is data science not?

(A few) data science examples

Course objectives and topics

Course logistics

Data science is not machine learning

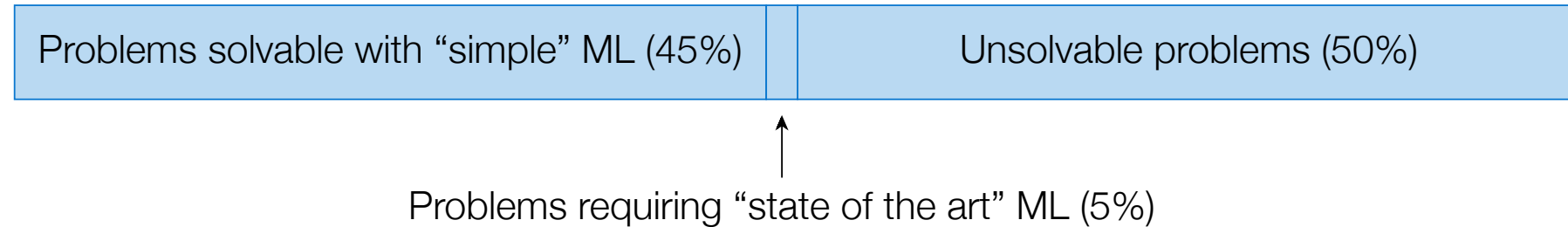
Machine learning involves computation and statistics, but has not (traditionally) been very concerned about answering *scientific questions*

Machine learning has a heavy focus on fancy algorithms...

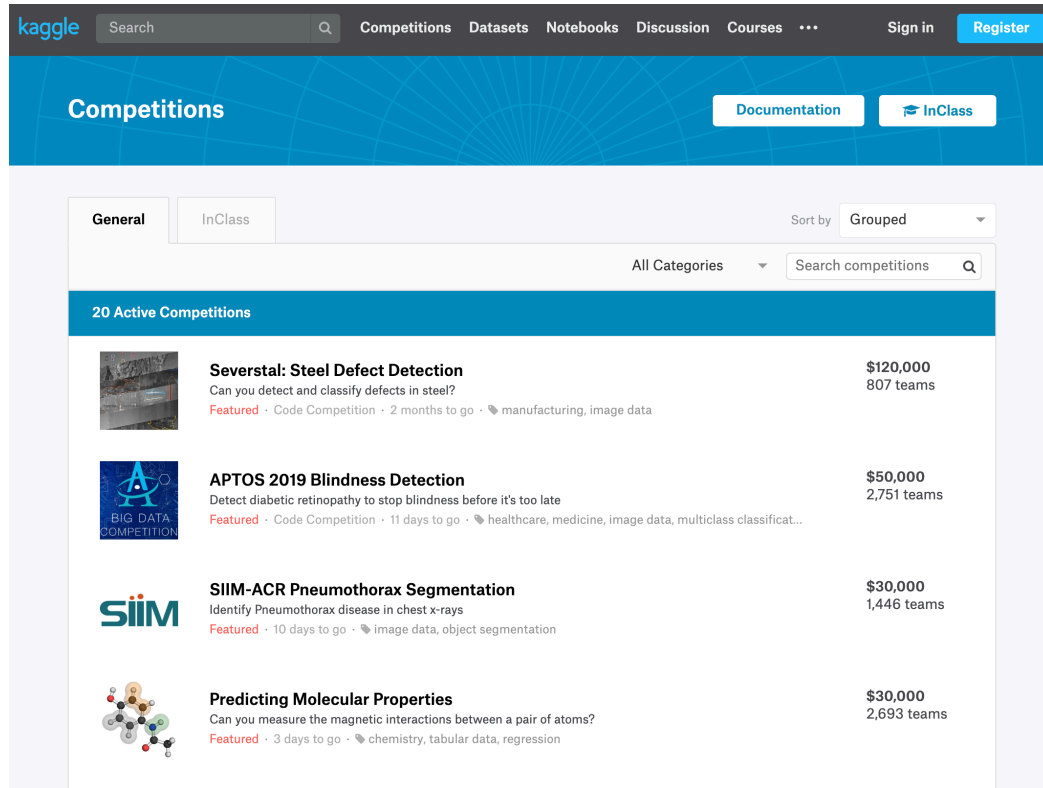
... but sometimes the best way to solve a problem is just by visualizing the data, for instance

Data science is not machine learning

Universe of machine learning problems



Data science is not machine learning competitions



The screenshot shows the Kaggle website's 'Competitions' page. The header includes the Kaggle logo, a search bar, and navigation links for Competitions, Datasets, Notebooks, Discussion, Courses, Sign in, and Register. Below the header, there are buttons for 'Documentation' and 'InClass'. The main content area is titled '20 Active Competitions' and lists several competitions with their respective details:

Competition Name	Prize	Teams
Severstal: Steel Defect Detection	\$120,000	807 teams
APTOS 2019 Blindness Detection	\$50,000	2,751 teams
SIIM-ACR Pneumothorax Segmentation	\$30,000	1,446 teams
Predicting Molecular Properties	\$30,000	2,693 teams

Data science competitions like Kaggle ask you to optimize a metric on a fixed data set

This may or may not ultimately solve the desired business/scientific problem

Data science is the iterative cycle of designing a concrete problem, building an algorithm to solve it (or determining that this is not possible), and evaluating what insights this provides for the real underlying question

Data science is not statistics

“Analyzing data computationally, to understand some phenomenon in the real world, you say? ... that sounds an awful lot like statistics”

Statistics (at least the academic type) has evolved a lot more along the mathematical/theoretical frontier

Not many statistics courses have a lecture on e.g. web scraping, or a lot of data processing more generally

Plus, statisticians use R, while data scientists use Python ... clearly these are completely different fields

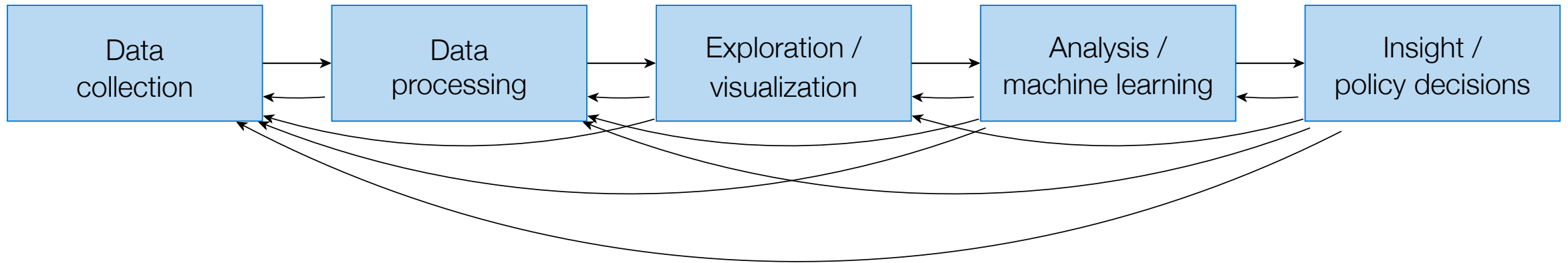
Data science is not big data

Sometimes, in order to truly understand and answer your question, you need massive amounts of data...

...But sometimes you don't

Don't create more work for yourself than you need to

Back to what data science is



Outline

What is data science?

What is data science not?

(A few) data science examples

Course objectives and topics

Course logistics

Gendered language in professor reviews

Gendered Language in Teacher Reviews

This interactive chart lets you explore the words used to describe male and female teachers in about 14 million reviews from RateMyProfessor.com.

You can enter any other word (or two-word phrase) into the box below to see how it is split across gender and discipline: the x-axis gives how many times your term is used per million words of text (normalized against gender and field). You can also limit to just negative or positive reviews (based on the numeric ratings on the site). For some more background, see [here](#).

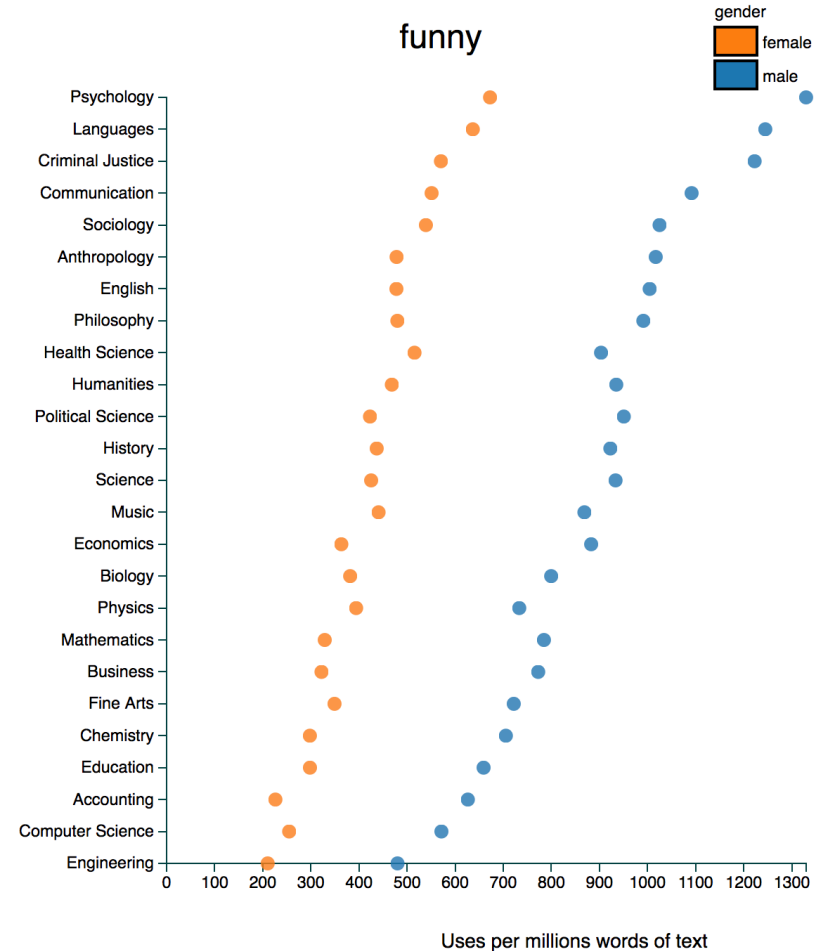
Not all words have gender splits, but a surprising number do. Even things like pronouns are used quite differently by gender.

Search term(s) (case-insensitive):
use commas to aggregate multiple terms

All ratings

Only positive

Only negative



<http://benschmidt.org/profGender/>

Obligatory quote

The greatest value of a picture is when it forces us to notice what we never expected to see.

-John Tukey

FiveThirtyEight

Search for a race or candidate

How do you like your House forecast?

Lite

Keep it simple, please — give me the best forecast you can based on what local and national polls say

Classic

I'll take the polls, plus all the "fundamentals": fundraising, past voting in the district, historical trends and more

Deluxe

Gimme the works — the Classic forecasts plus experts' ratings

Forecasting the race for the House



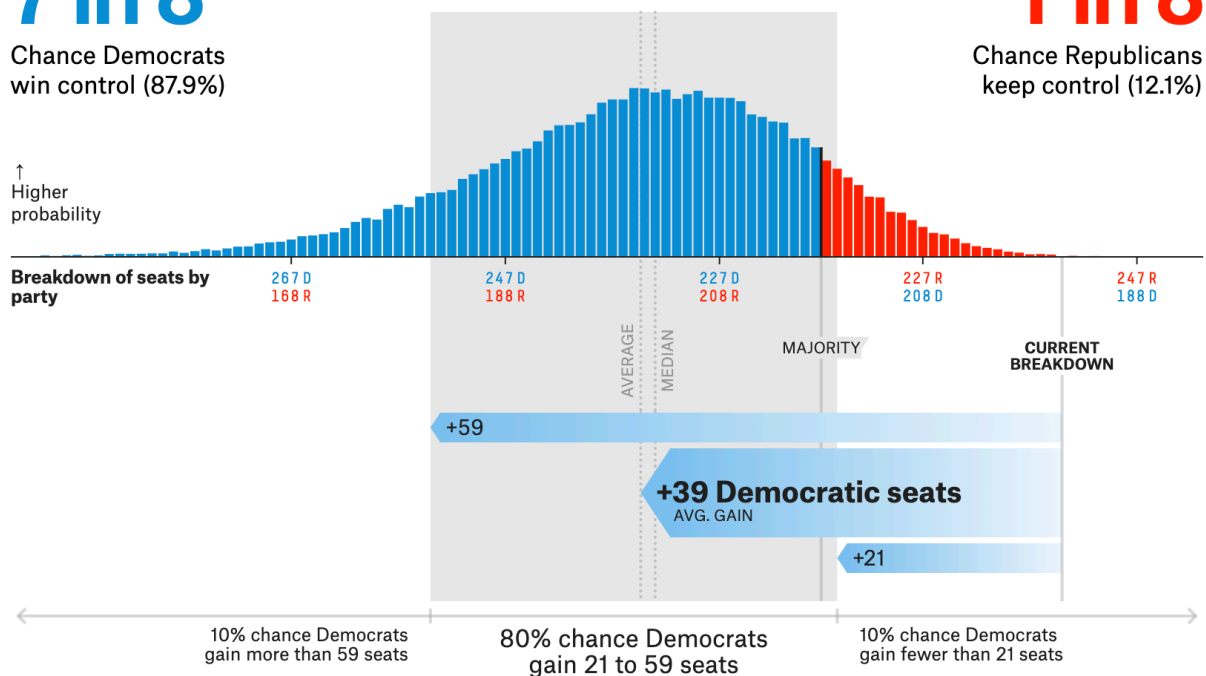
Updated Nov. 6, 2018, at 11:06 AM

7 in 8

Chance Democrats win control (87.9%)

1 in 8

Chance Republicans keep control (12.1%)



Poverty Mapping



Figure 2: Example of metal roof in center of satellite image.

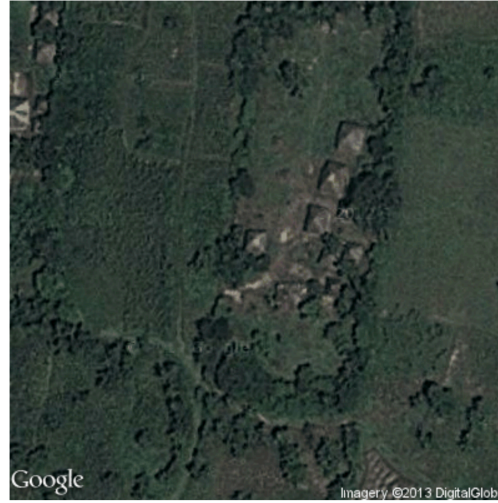


Figure 3: Example of thatched roof in center of satellite image.

Abelson, Varshney, and Sun. "Targeting Direct Cash Transfers to the Extremely Poor," 2012

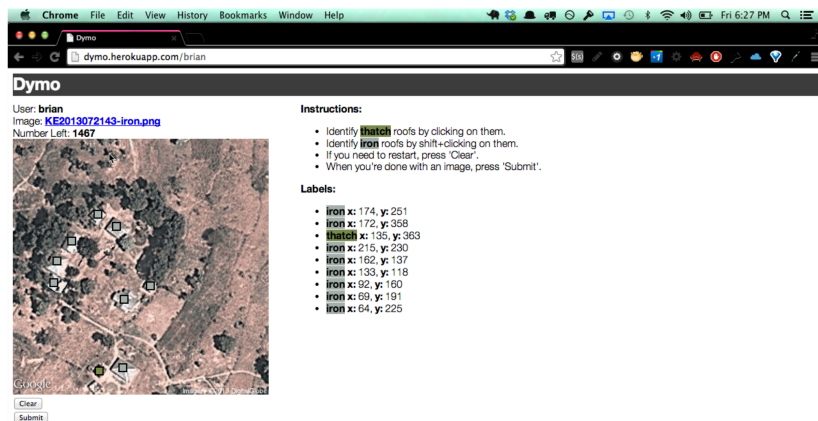


Figure 6: Screen shot of application deployed for crowdsourced labeling of roofs in satellite images.

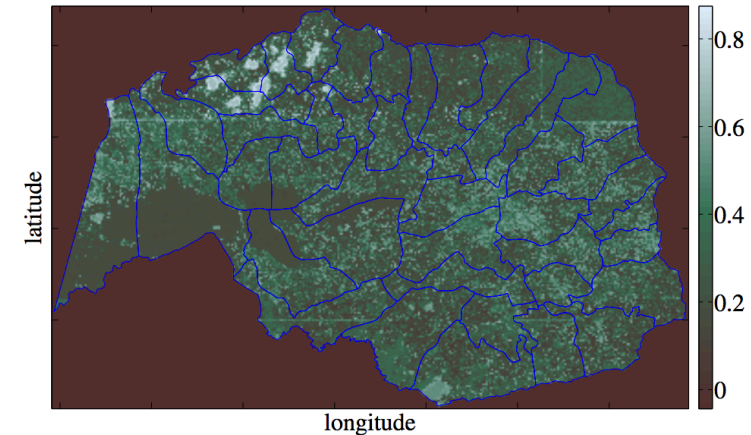


Figure 11: Heat map of proportion of roofs that are metal in the region of interest.

Outline

What is data science?

What is data science not?

(A few) data science examples

Course objectives and topics

Course logistics

Learning objectives of this course

After taking this course, you should...

... understand the full data science pipeline, and be familiar with programming tools to accomplish the different portions

... be able to collect data from unstructured sources and store it using appropriate structure such as relational databases, graphs, matrices, etc

... know to explore and visualize your data

... be able to analyze your data rigorously using a variety of statistical and machine learning approaches

Topics covered (subject to change)

Data collection and management: relational data, matrices and vectors, graphs and networks, free text processing, geographical data

Statistical modeling and machine learning: linear and nonlinear classification and regression, regularization, data cleaning, hypothesis testing, kernel methods and SVMs, boosting, clustering, dimensionality reduction, recommender systems, deep learning, probabilistic models, scalable ML

Visualization: basic visualization and data exploration, data presentation and interactivity

Philosophy: tools and deeper understand

Most of the techniques we will teach in this course have mature tools that you will likely use in practice

But, the philosophy of this course is that you will use these tools most effectively when you understand what is going on under the hood

This course will teach you some of the more common tools, but (especially in 15-688 problem sets), you will also need to implement some of the underlying methods

Example: we'll teach you how to run machine learning algorithms using scikit-learn library, but you'll also need to implement some of the algorithms yourself

Differences between 15-388/688 and XX

There are many courses that cover similar or related material (10-601, 10-701, 11-663, 05-839, 36-402, etc)

In general, this course puts a high emphasis on exploring and analyzing real (unprepared) data, managing the entire data science pipeline

Compared to other machine learning or statistics courses, there is relatively little theory, higher emphasis on implementation and use on practical data sets

Recommended background

The only formal prerequisite for this course is an intro to programming (if you have taken one at another university, this is fine)

We strongly recommend that students have **experience with Python**, ideally some background in **probability and statistics, and linear algebra**

If you don't have background in these areas, you may still sign up, but be aware that you will probably need to learn some of these items as the class goes on (we will be providing pointers to references)

General rule of thumb: If the homework seems hard, but you have ideas about how to proceed, you probably have the right level of background; if the homework seems hard and you have no idea how to proceed, this may be the wrong course

Outline

What is data science?

What is data science not?

(A few) data science examples

Course objectives and topics

Course logistics

Course materials and discussion

All course material (slides, notes, lecture videos, assignments) is available on the course webpage:

<http://www.datasciencecourse.org>

Slides posted before class, videos up ~2-3 hours after, notes posted asynchronously, typically well before lecture

The forum for discussions on the course is linked on the class webpage, but available directly at the link:

<http://forum.datasciencecourse.org>

15-388 vs. 15-688

Two versions of the course: 15-388 (undergrad, 9 unit), 15-688 (graduate, 12 unit)

Courses are identical (same lectures, assignments, etc) except that 15-688 problem sets have an additional question per assignment, usually requiring that students implement some advanced technique

Undergraduates **may take 15-688** for 12 units, but please wait until enrollment shakes out (for now, just start doing the 15-688 questions on the homeworks)

Two sections of 688 (A and B) are identical (they are there for historical reasons)

Course videos

All lectures will be recorded, made available on the course website and via Canvas (this is all that Canvas is used for for the course)

Students from any section may opt to view the class recordings instead of attending class (but of course, you won't be able to ask questions then)

Note that even if you ask a question in class, the video likely will not pick up your voice (I need to repeat questions after they are asked)

Auditing?

Auditing the course is permitted

The requirements to pass an audit are to receive at least **50% of the points on 4 out of the 5 assignments** (out of the whole assignment, so both the two 388 questions and the one additional 688 question)

No tutorial or final project are required for audit

We discourage final projects consisting of some full-credit participants and some auditors, unless you have a very good reason

Grading

Grading breakdown is posted on the web site (updated):

40% homework

20% tutorial

30% class project

10% class participation

Final grades are assigned on a curve (separate for 15-388 and 15-688 versions)

Homeworks

One homework assignment every two weeks: released on Tuesdays by midnight, due the Tuesday two weeks later at midnight

We may miss this deadline sometimes (we are sorry in advance, we will of course also extend the due date)

Work will be largely (solely?) about **writing code** to solve problems

Homeworks are in the form of Jupyter notebooks (accessible via Colab, if desired), **solutions autograded via a new system we are developing, more info with first HW release**

Autograding

The meta-goal for this course is to have a *scalable* introduction to data science

We believe that the current best way to achieve scalability is through heavy use of autograding

This presents additional problem for data science, where part of the process is developing scientific conclusions from the data (this is what the class project is for)

Note: tutorial and class project will be graded manually (by myself)

Late days

Assignments are due at 11:59pm (midnight) on Tuesdays

You have **5 late days** to use over the course of the semester

Each assignment can use a maximum of **2 late days** (midnight Saturday)

You cannot use late days for final project submission

Class participation

Class participation grade is determined by your participation in class discussion forums: you will receive full class participation by *answering* at least 5 questions asked by other students in the course

- We may change this criterion (and make an announcement) if there are issues ... the goal is to encourage participation, not gaming the system

Extra class participation credit for the 5 students who are most active in the forums (in terms of question answering, not asking)

Tutorial

The best way to learn a subject is to teach it

In lieu of a midterm, students will design a mini-tutorial, in the form of a Colab notebook, on a subject of their choice (though we will also provide suggestions)

Your tutorial will be read by the instructors, but also by other students, and peer grading will factor in to your final grade on the tutorial

Class project

A major component of the class: goal is to take a real-world domain that you are interested in, and apply data science methodologies to gain insight into the domain

Work to be done in groups of 2-3 students

Final report will be a Colab Notebook working through the analysis of your data, including code and visual results

Also presented in a video presentation (in lieu of final)

Class projects *must* be focused on some real data problem (ideally one that you collect yourself), not an already-curated data set

Collaboration on homeworks

All submitted content (code and prose for homeworks, tutorials, and and final project) should be your own content, written yourself

However, you *may* (in fact are encouraged to) discuss the homework with others in the class and on the discussion *including posting code*

- This creates some room for undue copying, but please obey the reasonable person principle: discuss as you see fit, but don't simply share answers

You may use snippets of code from sources like Stack Overflow, as long as you cite these properly (put a comment above and below whatever portion of code is copied), but again, be reasonable

Student well-being

CMU and courses like this one are stressful environments

In my experience, most academic integrity violations are the product of these environments and decisions made out of desperation

Please don't let it get to this point (or potentially much worse)

Don't sacrifice quality of life for this course: still make time to sleep, eat well, exercise

Up next

Next class: web scraping and data collection

First homework released hopefully released prior to next Tuesday (without autograding yet), use it as a gauge (after a few of the next lectures) to determine if the course is right for you