

A decorative graphic on the left side of the slide, consisting of a dark grey vertical band with a pattern of thin, light brown lines and small circles, resembling a circuit board or data flow diagram.

Lecture 1

Prof. Shannon Quinn

CSCI 8360: Data Science Practicum

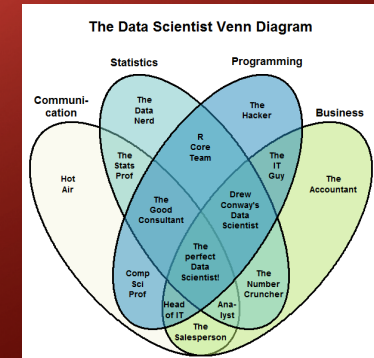
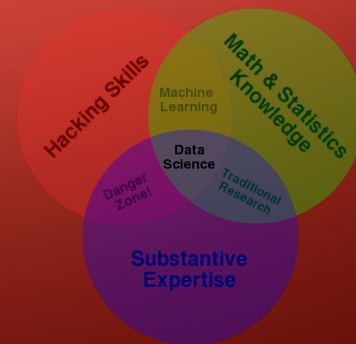
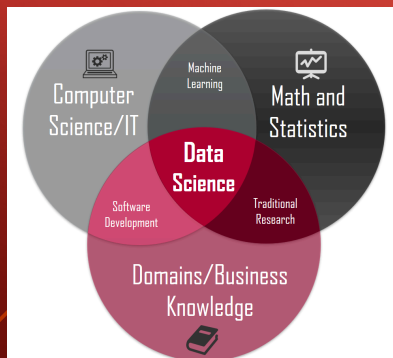


Part I: Lightning Overview

CSCI 8360 DATA SCIENCE PRACTICUM

Data Science

- What is it?
- Why is it important?
- How does one learn it?



CSCI 8360: What Is It?

- What this class **IS**:
 - Hands-on data science
 - Team-based problem solving
 - “Kaggle in the Classroom”



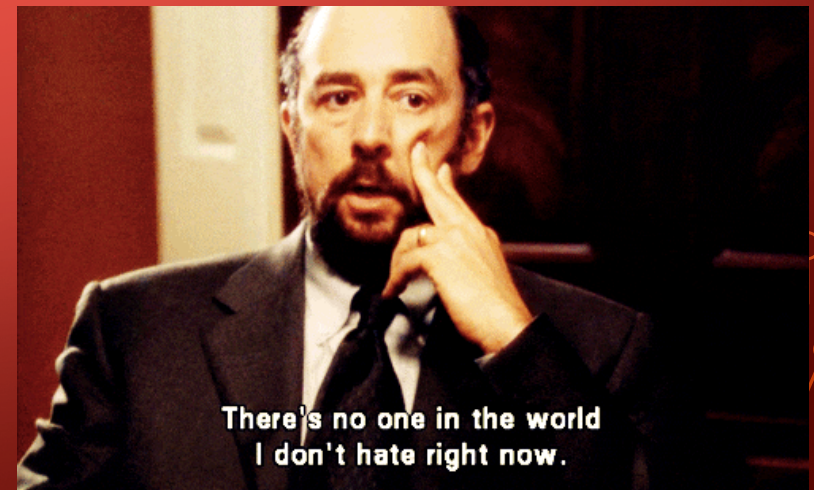
CSCI 8360: What Is It?

- What this class **IS NOT**:
 - Introduction to Machine Learning
 - Introduction to Distributed Systems
 - Introduction to Software Engineering



CSCI 8360 Requirements

- Thorough understanding of machine learning and statistics
 - (or teammates who can bring you up to speed very quickly)
- Good software engineering skills
 - (working on teams)
- An ability to learn fast
 - (definition of “graduate student”)



CSCI 8360 Links

- Course website
 - <https://dsp-uga.github.io/sp21>
 - Lectures and assignments will be posted here
- Discord
 - This is where **all course communication** will happen
 - (yes, we used Slack in previous years; testing this out this semester)
- GitHub
 - <https://github.com/dsp-uga/>
 - All team development will happen here
- AutoLab
 - <https://autolab.cs.uga.edu>
 - Project submissions for grading and evaluation
- Google Compute Platform (GCP)
 - Everyone will get credits



Course Outline

- 3.5 Projects (+ a pseudo-project), each 3(+) weeks
- Lecture every Wednesday (except the first two weeks)
- Office hours Tuesday/Thursday (except the first two weeks)
- No exams!
- No final project!
- **No grades!**
- Attendance



Part II: Administrative Details

CSCI 8360 DATA SCIENCE PRACTICUM

Lectures, Revisited

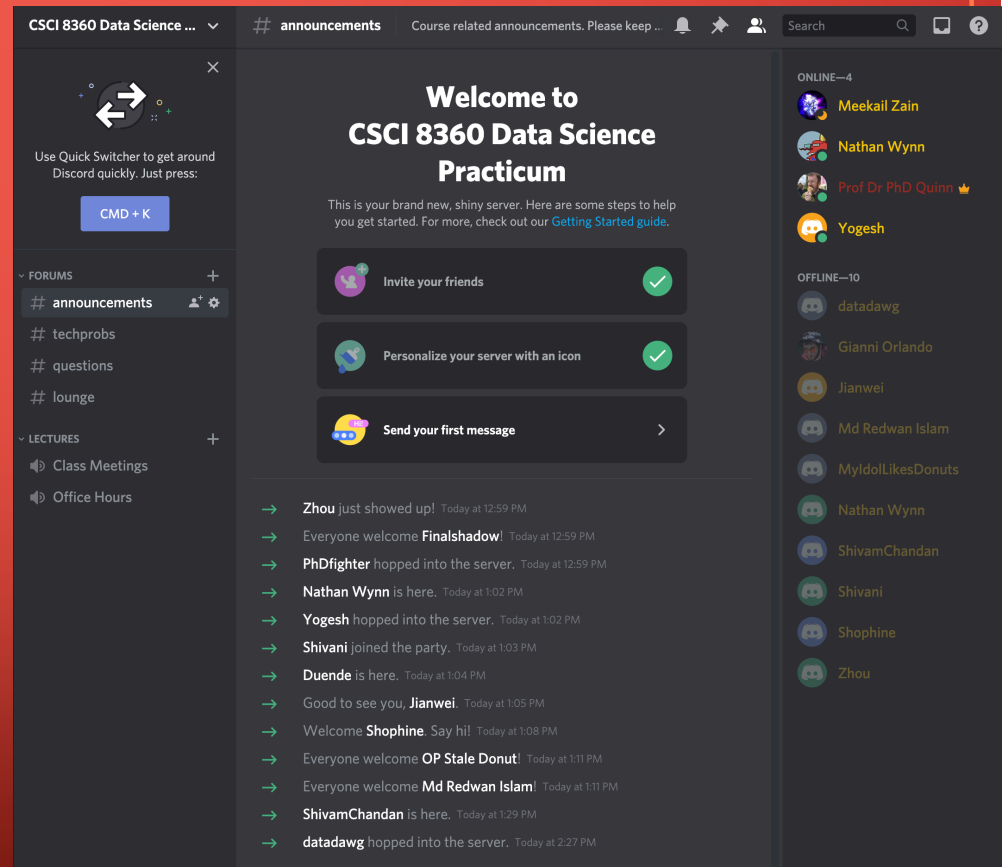
- Locations:
 - In-person: **Boyd 208 (Wednesdays); Food Science 202 (Tues/Thurs)**
 - Virtually: **Discord Server**
- Time
 - Today, 11:30am – 12:20pm
 - Tomorrow, 11:10am – 12:25pm
 - Next Tuesday (Jan 19), Wednesday (Jan 20), and Thursday (Jan 21)
 - **NO OTHER LECTURE TIMES** (unless announced in Discord)
- Policy
 - Attend, don't attend, up to you (except for guest lecturers—please come for those)
 - Lectures will be recorded and posted on the website

Office Hours, Revisited

- Location: **Food Science 202** (and Discord server)
- Time: Tuesdays / Thursdays, 11:10am – 12:25pm
- (yep, when we'd otherwise have lecture, so I know you can't possibly have conflicts)
- Happy to set up appointments if you need them

Discord, Revisited

- Discord: free team messaging platform
- Web-based and mobile apps
- Teams can set up private channels to coordinate
- Can also send individual DMs
- **All course announcements will be made here**



GitHub, Revisited

- Most popular code repository in the world
- Uses the *git* concurrent versioning system (itself an open source project)
- Lots of useful team-based tools (issue tracker, wiki, GUI)
- **All projects will be sourced in the DSP-UGA GitHub organization**



The screenshot shows a GitHub repository page for the organization "dsp-uga" and repository "sp18". The page includes navigation tabs for Code, Issues (0), Pull requests (0), Projects (0), Wiki, Insights, and Settings. The repository name is "Spring 2018 rendition of CSCI 8360." and it has 4 commits, 1 branch, 0 releases, and 1 contributor. The repository is licensed under MIT. The file list shows:

File	Commit Message	Time
docs	Added date, time, and location.	2 months ago
LICENSE	Initial commit	2 months ago
README.md	Updated README.	2 months ago

The README.md file content is visible below the file list:

```
Spring 2018: CSCI 8360

Spring 2018 rendition of CSCI 8360 Data Science Practicum.
```

AutoLab, Revisited

- Assignment submission and autograder
- Also has leaderboards!
- **All project outputs will be submitted to AutoLab for ranking**



Google Cloud Platform

- (comparable to Amazon Web Services, or AWS)
- Spin up elastic compute resources on-demand
- Every student gets \$50 in credits (usable across ALL services)
- “Cloud Dataproc” contains APIs for specifically spinning up Spark and Hadoop clusters
- **Details to come**



Google Cloud Platform



Part III: Projects

CSCI 8360 DATA SCIENCE PRACTICUM

Project Overview

- Solving real-world machine learning problem
 - Classify a large corpus of documents
 - Convex optimization over a huge dataset
 - Dimensionality reduction over a high-dimensional matrix
 - etc.
- Each project varies in length from 3 to 4 weeks
 - “Introductory” Project 0 out **tomorrow**, will be 1.5 weeks long
 - Project 1 (P1) out the following Tuesday (Jan 26), will be 3 weeks long

Project Requirements: Teams

- Teams (3-4 people per team)
 - Assigned *completely randomly* (by me)
 - Will change for each project
- Each team member should have a **clear, well-defined role**
 - Not everyone has to be a coder!
 - But 1 person should not be carrying the whole project
 - **The only way you can fail this course is by being a bad teammate who either ghosts their team or does everything themselves**



Project Requirements: Code

- Use good coding practices
 - Documentation (in code, in GitHub wiki, in README, in commit comments)
 - Well-organized structure (should be easy for me to understand)
- Use organizational GitHub account
 - <https://github.com/dsp-uga>
- Recommended additional practices
 - License your code with a permissive open license (<https://choosealicense.com/>)
 - Add a continuous integration module
 - Implement unit testing for your code
 - Create a website for your project (see GitHub documentation; makes this easy)



Project Requirements: AutoLab

- Submit to AutoLab before the deadline
 - One submitter per team (can submit as many times as you like)
 - Unless otherwise specified, submission will always be a text file with your code's predictions on a test dataset
 - If your submission is correctly formatted, your performance should show up on the leaderboard in short order
- AutoLab submission **shuts down after the deadline**



Project Requirements: Lightning Talks

- Wednesdays *after* a project deadline, all teams will give a lightning talk (~5 minutes long)
- Talks will outline the problem, the team's approach, their results, and any other discussion points
- Creativity welcome—code examples, live demos, interactive slides, etc
- One person from each team will speak

Project Grading



Project Grading

- COVID-19 sucks
- Peer evaluation: each team is assigned another team to evaluate their work
- Evaluation split into three categories
 1. Theory (the approach you use as implemented by the code)
 2. Engineering (everything around the implementation)
 3. “Extras”
- Go above and beyond—extra points
- Note shortcomings (approach is flawed or too simple, code poorly documented, one person did almost all of the project, poor performance), as well as strengths (innovative approach—even if it did not work—good division of labor, unit testing, easy to understand code, good use of tickets and milestones)
- Evaluation reports will be issued to each team between projects

Project Grading

- In all seriousness, if you
 - Work with your teammates (don't ghost them)
 - Be a good classmate (answer other questions in Discord, come to occasional office hours / lecture, generally interact)
 - Conduct yourself as a decent human being
- **You get an A this semester**
- I want you all to have a little fun this semester without worrying about grades

Final Project

- Have done this in previous years
- Not doing it this semester



Part IV: The Next Step

CSCI 8360 DATA SCIENCE PRACTICUM

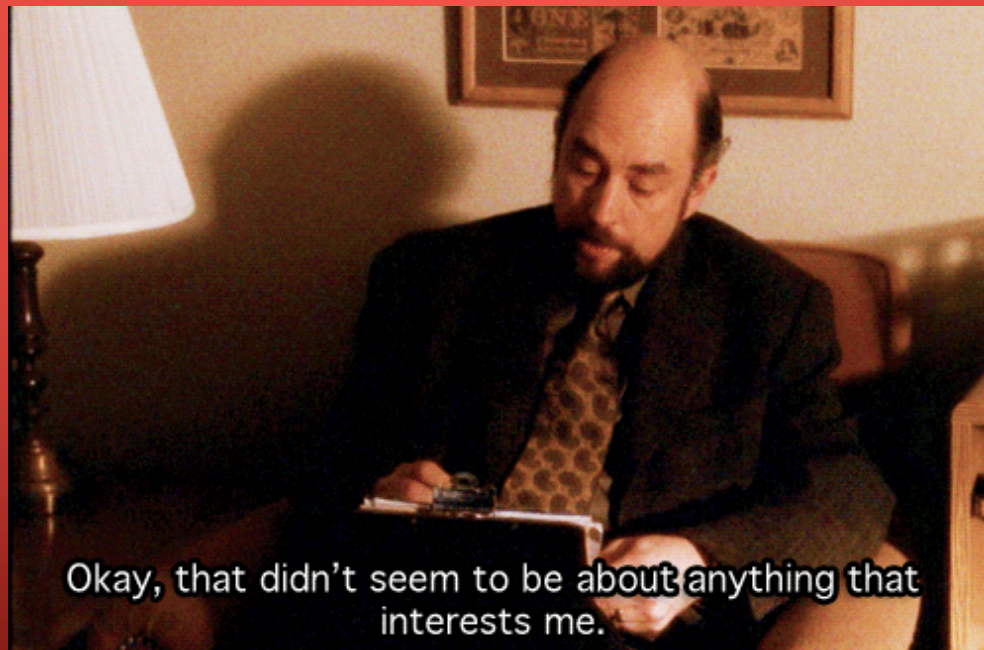
Project α

1. Join the Discord server (as of writing this lecture, 13 of 20 have already joined)
2. Send me your GitHub username over Discord (create an account if you don't have one). Join the GitHub "Data Science Practicum" team.
3. Start looking at Apache Spark and/or dask (for Project 0 tomorrow).

Tomorrow: Project 0

- The only individual project of the semester
- Mainly to familiarize you with Apache Spark or dask (used for Project 1), AutoLab, GitHub, and Discord

QUESTIONS?



Okay, that didn't seem to be about anything that interests me.

Finally...

- What large-scale problems do *you* want to work on?
- Yes, this an opportunity to suggest Projects. If you have an idea, send me:
 1. The problem to be solved (optimization, dimensionality reduction, classification, etc)
 2. How the solutions should be evaluated
 3. Training and validation datasets



Your idea
could be
featured as a
full project!

