# [320] Welcome + First Lecture [reproducibility]

Tyler Caraza-Harter

### Welcome to the first ever offering of Data Programming II!

Builds on CS 301 220. <u>https://stat.wisc.edu/undergraduate-data-science-studies/</u>

CS 220	CS 320	
getting results	getting reproducible results	
writing correct code	writing efficient code	
using objects	designing new types of objects	
functions: f(obj)	<pre>methods: obj.f()</pre>	
lists+dicts	graphs+trees	
analyzing datasets	collecting+analyzing datasets	
plots	animated visualizations	
tabular analysis	simple machine learning	



# Who am I?

Tyler Caraza-Harter

- Long time Badger
- Email: <u>tharter@wisc.edu</u>
- Just call me "Tyler"

Industry experience

- Worked at Microsoft on SQL Server and Cloud
- Other internships/collaborations: Qualcomm, Google, Facebook, Tintri







# Who are You?

Year in school?

• Ist year? 2nd? Junior/senior? Grad student?

Area of study

• Natural science, social science, engineering, business, statistics, data science, other?

What CS courses have people taken before?

• CS 220/301 (the import one here)? CS 200? CS 300? CS 354?

Please fill this form: <a href="https://forms.gle/SdFHQq37txmqkkgq7">https://forms.gle/SdFHQq37txmqkkgq7</a>. Why?

- Help me get to know you
- Get participation credit

# Course Logistics

# Class organization

### Teams

- you'll be randomly assigned to a team of 4-7 students
- teams will last the whole semester
- some types of collaboration with team members are allowed (not required) on graded work, such as projects+quizzes
- most collaboration with non-team members in not allowed

### Staff

- I. Instructor
- 2. Teaching Assistants
- 3. Mentors

we all provide office hours, and you can attend any that you prefer!

# Class organization

### Teams

- you'll be randomly assigned to a team of 4-7 students
- teams will last the whole semester lacksquare
- some types of collaboration with team members are allowed (not required) on graded work, such as projects+quizzes
- most collaboration with non-team members in not allowed lacksquare

### Staff

- Instructor
   Teaching Assistants
   head TA: in charge of projects
   team TA: primary contact for team, same whole semester
   grader TA: reviews projects (rotates weekly)

we all provide office hours, and you can attend any that you prefer!

# Course Website

It's here: https://tyler.caraza-harter.com/cs320/f20/schedule.html

Data Programming II Schedule Sy	llabus Projects Resources 🕶 Tools 🕶		
Course Schedu	le		
Part 1: Performance Week 1			read syllabus carefully and checkout other content
[Mon] No Class (Aug 31)	[Wed] Reproducibility 1 (Sep 2)	[Fri] Reproducibility 2 (Sep 4)	
	<ul> <li>Course Overview</li> <li>Hardware, OS, Interpreters</li> <li>Read: Syllabus</li> <li>WEEKLY LAB: Cloud Setup</li> </ul>	<ul> <li>versioning</li> <li>git</li> <li>Read: Git Tutorial</li> </ul>	
Week 2			

I'll also use Canvas for four things:

- general announcements
- quizzes
- help you keep track of your progress through lectures, labs, etc
- simple grade summaries (not feedback or exam answers)

# Other Communication

#### Piazza

- find link on site
- don't post >5 lines of project-related code (considered cheating)
- pinned post will list office hours (me,TAs, mentors)

### Forms

- <u>https://tyler.caraza-harter.com/cs320/f20/surveys.html</u>
- Who are you? Feedback Form. Thank you!

### Email

- me: <u>tharter@wisc.edu</u>
- TAs: <u>https://tyler.caraza-harter.com/cs320/f20/contact.html</u>

# Course Etiquette

### Meetings

- I. office hours are drop-in (no need to reserve)
- 2. email me to schedule individual meetings

### Email

- 3. let us know your NetID (if not from <a href="mailto:netid@wisc.edu">netid@wisc.edu</a>)
- 4. don't start new email thread if topic is the same
- 5. unless urgent, please give me 48 hours to respond before following up (I'll try to be faster usually)
- 6. use your judgement about whether to email me or TA first
- 7. if general question, consider using piazza instead

# Graded Work

## 7 Projects - 8% each

- **format**: notebook, module, or program
- part I: you can collaborate with team
- part 2: must be individualy (only help from 320 staff)
- still a test.py, but more depends on TA evaluation (more plots)
- ask for specific feedback (giving constructive criticism is a priority in CS 320)

## 14 Quizzes - 2% each

- after each week, anytime before deadline
- on Canvas, open book/notes
- can take together AT SAMETIME with team members (no other human help)

# Graded Work

# I Final - 10%

- short, open-ended project on topic of your choosing
- due at originally scheduled exam time
- more details/constraints when it gets closer...

## Participation - 6%

- class surveys
- interacting with posted discussions
- active in weekly team meetings
- doing labs before projects
- etc.

# Academic Misconduct

Read syllabus to make sure you know what is and isn't OK.

It's not obvious! Especially this semester...

### In Fall 2019, I made the following misconduct reports:

- 23 students for cheating on projects
- 2 past students for sharing solutions from past semesters
- 7 students for cheating on exams

### How we'll keep the class fair

- run MOSS on submissions
- randomize exam question order

Please talk to me if you're feeling overwhelmed with 320 or your semester in general!

# Reading: same as 301 and some others...



I'll post links to other online articles and my own notes

Lectures don't assume any reading prior to class

# Any questions?

Drop an email, or better, post to piazza

# Today's Lecture: **Reproducibility**

# Reproducibility (Fall 19 Grading for CS 301)



why was project 9 so problematic?



**Discuss:** how might we define "reproducibility" for a data scientist?

# 15 new terms to learn today...

reproducibility: others can run our analysis code and get same results process: byte: process memory: address: encoding: CPU: how many terms do you know already? instruction set: operating system: resource: allocation: abstraction: virtual machine: cloud: ssh:

**Big question:** will my program run on someone else's computer? (not necessarily written in Python)



## Hardware: Mental Model of Process Memory

Imagine...

- one huge list, **per each** running program process
- every entry in the list is an integer between 0 and 255 (aka a "byte")





- multiple lists
- variables and other references > data
- strings
- code



Is this really all we have for state?

- multiple lists
- variables and other references
- strings
- code



the [11,22,33] list starts at address 12 in the giant list

- multiple lists
- variables and other references
- strings
- code



# fast
L2.append(44)

implications for performance...

- multiple lists
- variables and other references
- strings
- code



# fast
L2.append(44)

implications for performance...

- multiple lists
- variables and other references
- strings
- code



implications for performance...

# fast
L2.append(44)

# slow
L2.pop(0)

- multiple lists
- variables and other references
- strings
- code



implications for performance...

# fast
L2.append(44)

# slow
L2.pop(0)

- multiple lists
- variables and other references
- strings
- code



We'll think more rigorously about performance in CS 320 (big-O notation)

# fast
L2.append(44)

# slow
L2.pop(0)

- multiple lists
- variables and other references
- strings
- code





PythonTutor's visualization

- multiple lists
- variables and other references
- strings
- code



- multiple lists
- variables and other references
- strings
- code



- multiple lists
- variables and other references
- strings
- code





	code	operation
	5	ADD
Instruction Set	8	SUB
	33	JUMP
	•••	•••

CPUs interact with memory:

- keep track of what instruction we're on
- understand instruction codes
- much more





	code	operation
	5	ADD
Instruction Set	8	SUB
	33	JUMP
	•••	•••

CPU

- line that just executed
- next line to execute

CPUs interact with memory:

- keep track of what instruction we're on
- understand instruction codes •



	code	operation
	5	ADD
Instruction Set	8	SUB
	33	JUMP
	•••	•••

CPUs interact with memory:

• keep track of what instruction we're on



	code	operation
	5	ADD
Instruction Set	8	SUB
	33	JUMP
	•••	•••

0

15

CPUs interact with memory:

- keep track of what instruction we're on
- understand instruction codes
- much more



	code	operation
	5	ADD
Instruction Set	8	SUB
	33	JUMP
	•••	•••

CPU



for CPUY

•••

8

33

•••

undefined

•••

for	C	ΡU	X

code	operatior
5	ADD
8	SUB
33	JUMP
	•••

# A Program and CPU need to "fit"





# A Program and CPU need to "fit"



## why haven't we noticed this yet for our Python programs?

## Interpreters



Interpreters (such as python.exe) make it easier to run the same code on different machines

A compiler is another tool for running the same code on different CPUs

## Interpreters



Interpreters (such as python.exe) make it easier to run the same code on different machines

**Discuss:** if all CPUs had the instruction set, would we still need a Python interpreter?

**Big question:** will my program run on someone else's computer? (not necessarily written in Python)



**Big question:** will my program run on someone else's computer? (not necessarily written in Python)



# OS jobs: Allocate and Abstract Resources

[like CPU, hard drive, etc]



## Parallelism -- more later this semester...





most modern CPUs actually contain multiples CPUs (called "cores") on a single chip

Later: how can we write programs that run in parallel, going faster by using multiple cores?

# OS jobs: Allocate and Abstract Resources

[like CPU, hard drive, etc]



## Harder to reproduce on different OS...



f = open("/data/file.txt")

• • •

The Python interpreter mostly lets you [Python Programmer] ignore the CPU you run on.

But you still need to work a bit to "fit" the code to the OS.

# Harder to reproduce on different OS...



f = open("c:\data\file.txt")

• • •

The Python interpreter mostly lets you [Python Programmer] ignore the CPU you run on.

But you still need to work a bit to "fit" the code to the OS.

# Harder to reproduce on different OS...



# solution 1:
f = open(os.path.join("data", "file.txt"))
...

# solution 2: tell anybody reproducing your results to use the same OS!

tradeoffs?

The Python interpreter mostly lets you [Python Programmer] ignore the CPU you run on.

But you still need to work a bit to "fit" the code to the OS.

# VMs (Virtual Machines)

#### popular virtual machine software



With the right virtual machines created and operating systems installed, you could run programs for Mac, Linux, and Windows -- at the same time without rebooting!

# The Cloud



# Lecture Recap: Reproducibility

**Big question:** will my program run on someone else's computer?

Things to match:



# Recap of 15 new terms

reproducibility: others can run our analysis code and get same results process: a running program byte: integer between 0 and 255 process memory: a big "list" of bytes, per process, for all state address: index in the big list encoding: pairing of letters characters with numeric codes CPU: chip that executes instructions, tracks position in code instruction set: pairing of CPU instructions/ops with numeric codes operating system: software that allocates+abstracts resources resource: time on CPU, space in memory, space on SSD, etc. allocation: the giving of a resource to a process abstraction: hiding inconvenient details with something easier to use virtual machine: "fake" machine running on real physical machine allows us to running additional operating systems cloud: place where you can rent virtual machines and other services ssh: secure shell -- tool that lets you remotely access another machine