



Computational Structures in Data Science

UC Berkeley EECS
Lecturer
Michael Ball



Lecture #1: Welcome to CS88!



<http://cs88.org>

January 24, 2020

1

CS88 Team - Michael

- Michael Ball**
 - ball@Berkeley.edu – You're best off by using Piazza! ☺
 - 625 Soda Hall
 - <http://michaelball.co> – I don't update this much...
 - » It was great procrastination when I was a CS student.
 - Office hours: TBD @ 625 Soda
 - A few minutes after class
- Things I do:**
 - Intro CS Research
 - » Tools, curriculum
 - Training TAs
 - Building Educational Software (Gradescope)
 - Tools for web accessibility



1/24/20 UCB CS88 Sp20 L1 2

2

CS88 Team - Gerald

- Dr. Gerald Friedland**
 - fractor@Berkeley.edu – You're best off by using Piazza! ☺
- Projects you might want to check out:**
 - <http://mmcommons.org>
 - Work with 100M images, 1M videos in your own Amazon instance.
 - <http://www.teachingprivacy.org>
 - Creating teaching materials informing about data over sharing.





1/24/20 UCB CS88 Sp20 L1 3

3

CS88 — Our Awesome Team!

Head Teaching Assistants

		
Email: alexkasel@berkeley.edu	Email: bri@berkeley.edu	Email: julyayu@berkeley.edu

Teaching Assistants

			
Email: aleckan@berkeley.edu	Email: cmaloy@berkeley.edu	Email: shreyakannan@berkeley.edu	Email: sophia.qn@berkeley.edu
			
Email: srgt@berkeley.edu	Email: vandanag@berkeley.edu		

1/24/20 UCB CS88 Sp20 L1 4

4

Welcome

- We are all here to learn:
Knowledge (end) – Knowledge (start)

9/9/19 UCB CS88 Sp20 L1 5

5

Goals today

- Introduce you to**
 - the field
 - the course
 - the team
- Answer your questions**
- Big Ideas:**
 - Abstraction
 - Data Types





1/24/20 UCB CS88 Sp20 L1 6

6

Logistics

- **Labs are Mon & Tues**
 - Please see the signup form on Piazza
 - Go to any lab section next week
 - Attendance is tracked, but not required. ("Bonus" points.)
- **There are 2 sections on CalCentral, it does not matter which you are enrolled in.**
- **1-2pm lecture is webcast, 3pm is not.**

UCB CS88 Sp20 L1



7

Data Science

Nearly every field of discovery is transitioning from "data poor" to "data rich"



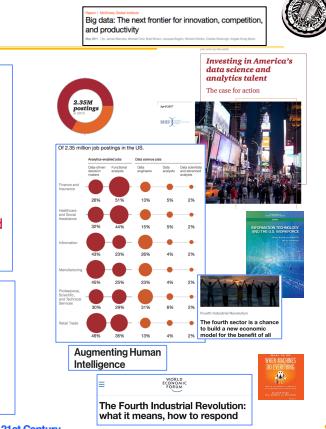
8

A National Challenge

Increasingly US jobs require data science and analytics skills. Can we meet the demand? The current shortage of skills in the national job pool demonstrates that business-as-usual strategies won't satisfy the growing need. If we are to unlock the promise and potential of data and all the technologies that depend on it, employers and educators will have to transform.

By 2021, **69%** of employers expect candidates with DSA skills to get preference for jobs in their organizations.

Only **23%** of college and university leaders say their graduates will have those skills.

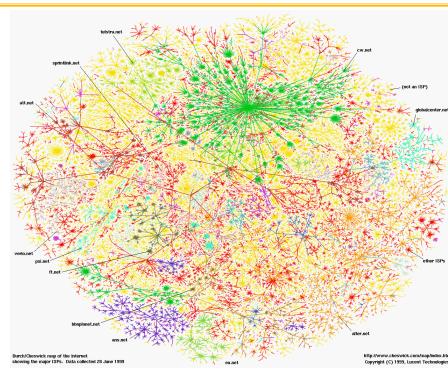


5/24/18

21st Century

8

Greatest Artifact of Human Civilization ...



UCB CS88 Sp20 L1



10

WORLD INTERNET USAGE AND POPULATION STATISTICS DEC 31, 2017 - Update

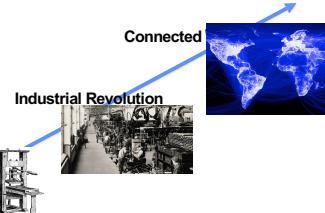
World Regions	Population (2018 Est.)	Population % of World	Internet Users 31 Dec 2017	Penetration Rate (% Pop.)	Growth 2000-2018
Africa	1,287,914,329	16.9 %	453,329,534	35.2 %	9,941 %
Asia	4,207,588,157	55.1 %	2,023,630,194	48.1 %	1,670 %
Europe	827,650,849	10.8 %	704,833,752	85.2 %	570 %
Latin America / Caribbean	652,047,996	8.5 %	437,001,277	67.0 %	2,318 %
Middle East	254,438,981	3.3 %	164,037,259	64.5 %	4,893 %
North America	363,844,662	4.8 %	345,660,847	95.0 %	219 %
Oceania / Australia	41,273,454	0.6 %	28,439,277	68.9 %	273 %
WORLD TOTAL	7,634,758,428	100.0 %	4,156,932,140	54.4 %	1,052 %

5/24/18

21st Century

11

Era of Transformation



5/24/18
21st Century



12

A Connected World of Data

- The world's knowledge at our finger tips
- Digitalization* of life, industry and society
- Intimately connected to billions of us, globally
- Explosion of observational instruments
 - Genomics, Microscopy, Astronomical, ...
- Vast Computational power to do analytics
- Synthetic design exploration thru simulation
- Machine reading of everything
- Statistical machine learning algorithms to "discover" structure

5/24/18 21st Century 13

13

What if I could ... ?

- See the world's digital footprints?
- Read everything that's ever been written?
- Take it all in and dive down anywhere as far as the science can take me?
- Learn the physical/chemical/biological /sociological/neurological... models from the data?
- Explore billions of designs and pick the one I want?
- ... ?

5/24/18 21st Century 14

14

A Connected World

Internet Users in the world: 3,293,151,639 (3.0 B 11/15)

Google searches today: 2,652,887,737

Videos viewed today: 5,835,884,253

ARPANet (1969) - RFC 675 TCP/IP (1974)

Internet (1983) - 2.0 B 1/26/11

United States (2010) - 3.0 B 11/15

1/24/20 15

15

Data 8 – Foundations of Data Science

- Computational Thinking + Inferential Thinking in the context of working with real world data
- Introduce you to several computational concepts in a simple data-centered setting
 - Authoring computational documents
 - Tables
 - Within Python3 and "SciPy"

UCB CS88 Sp20 L1

1/24/20 16

16

CS88 – Computational Structures in Data Science

- Deeper understanding of the computing concepts introduced in C8
 - Hands-on experience => Foundational Concept
 - How would you create what you use in c8 ?
- Extend your understanding of the structure of computation
 - What is involved in interpreting the code you write ?
 - Deeper CS Concepts: Recursion, Objects, Classes, Higher-order Functions, Declarative programming, ...
 - Managing complexity in creating larger software systems through composition
- Create complete (and fun) applications
- In a data-centric approach

1/24/20 UCB CS88 Sp20 L1 17

17

How does CS88 relate to CS61A ?

Interpretation CS Concepts and Techniques Intro Programming & Tools	Thinking w/ Data Statistics Concepts in a Computational Approach Intro Programming	Working w/ Data CS Concepts and Techniques & Tools
CS61A	DATA8	CS88

1/24/20 UCB CS88 Sp20 L1 18

18

Opportunities for students

The diagram illustrates the progression of opportunities for students:

- CS minor:** c8, c8 CS88
- CS major:** c8 CS88, c8 CS61B, ***
- CS61a:** c8 cs61a, cs61a

UCB CS88 Sp16 L1 19

19

The Data Science Major

The Data Science Major curriculum is structured as follows:

- Foundational Lower Division (16 units):**
 - Mathematics
 - Computing
 - Data 8: Foundations of Data Science
- Individualized Upper Division (30 units):**
 - Computational & Inferential Depth
 - Modeling, Learning & Decision Making Probability
 - Domain Emphasis
 - Human Contexts & Ethics
 - Electives
- College Breadth & Electives:** Domain Emphasis

UCB CS88 Sp16 L1 19

20

Course Structure

- 2 hours lecture, 2 hours lab, ~5 hours HW+study
- Lecture introduces concepts (quickly!), answers why questions.
- Lab provides concrete detail hands-on
- Homework (12) cements your understanding
 - Out Tuesdays, Due Next Thurs (-9 days)
- Projects (2) put your understanding to work in building complete applications
 - Maps
 - Ants vs Some Bees
 - Maybe: open-ended final?
- Readings: <http://composingprograms.com>
 - Same as CS61a

UCB CS88 Sp20 L1 21

21

Course Culture

- Learning
- Community
- Respect
- Collaboration
- Peer Instruction

UCB CS88 Sp20 L1 22

Piazza for {ask,answer}ing questions

The Piazza interface shows a question from a TA about professor office hours, with a response from a student named Luke Segars. The interface includes features like "Good Question", "Good Answer", and "Ask a Followup".

UCB CS88 Sp20 L1 23

23

Where will we work?

- Your laptop
 - Using an editor and a terminal
- cs88.org
- datahub.berkeley.edu
 - “Jupyter Notebooks”. An industry standard tool.
 - Not as often, but an option

UCB CS88 Sp20 L1 24

iClicker Check In

- Are you enrolled in Data 8?
- A. I took it Fall 2018 or earlier
- B. I took it Spring 2019
- C. I'm taking it right now
- D. I am trying to enroll in Data 8
- E. I am not taking Data 8

UCB CS88 Sp20 L1



26

26

Pro-student Grading Policies

- EPA
 - Rewards good behavior
 - Effort
 - » E.g., Office hours, doing every single lab, hw, reading Piazza pages
 - Participation
 - » E.g., Raising hand in lec or discussion, asking questions on Piazza
 - Altruism
 - » E.g., helping other students in lab, answering questions on Piazza
- You have 3 “Slip Days”
 - Homework and Projects
 - You use them to extend due date, 1 slip day for 1 day extension
 - You can use them one at a time or all at once or in any combination
 - They follow you around when you pair up (you are counted individually)
 - » E.g., A has 2, B has 0. Project is late by 1 day. A uses 1, B is 1 day late

27

27

Pro-student Grading Policies

- Turning in work will never harm you! Even if it's late or you can't get it completely correct.
- Labs are based primarily on effort.
- If you're struggling, come talk to us!
 - Especially your TAs. They have been there too.

01/28/19 UCB CS88 Sp19 L1



28

28

Abstraction

- Detail removal
 - “The act of leaving out of consideration one or more properties of a complex object so as to attend to others.”
- Generalization
 - “The process of formulating general concepts by abstracting common properties of instances”
- Technical terms:
Compression, Quantization, Clustering, Unsupervised Learning



01/28/19 UCB CS88 Sp19 L1

29

29

Experiment



01/28/19

UCB CS88 Sp19 L1

30

30

Where are you from?

Possible Answers:

- Planet Earth
- Europe
- California
- The Bay Area
- San Mateo
- 1947 Center Street, Berkeley, CA
- $37.8693^\circ \text{ N}, 122.2696^\circ \text{ W}$



All correct but different levels of abstraction!

01/19/18

UCB CS88 Sp18 L1

31

31

Abstraction gone wrong!

I Can Stalk U
Raising awareness about inadvertent information sharing

Home How Why About Us Contact Us

What are people *really* saying in their tweets?

- denisluque:** I am currently nearby http://maps.google.com
1 minute ago · Map Location · View Tweet · View Picture · Reply to denisluque
- mikosofficial:** I am currently nearby http://maps.google.com
5 minutes ago · Map Location · View Tweet · View Picture · Reply to mikosofficial
- dilmanaredet:** I am currently nearby http://maps.google.com
7 minutes ago · Map Location · View Tweet · View Picture · Reply to dilmanaredet
- downtowndan:** I am currently nearby http://maps.google.com
10 minutes ago · Map Location · View Tweet · View Picture · Reply to downtowndan
- MommaGooseBC:** I am currently nearby 15745 Weaver Lake Rd Maple Grove MN

How did you find me?

Did you know that a lot of smart phones encode the location of where pictures are taken? Anyone who has a copy can access this information.

01/28/19 UCB CS88 Sp19 L1 32

32

Detail Removal (in Data Science)

You'll want to look at only the interesting data, leave out the details, zoom in/out...

Abstraction is the idea that you focus on the essence, the cleanest way to map the messy real world to one you can build

Experts are often brought in to know what to remove and what to keep!

The London Underground 1928 Map & the 1933 map by Harry Beck.

01/28/19 UCB CS88 Sp19 L1 33

33

The Power of Abstraction, Everywhere!

- Examples:**
 - Functions (e.g., $\sin x$)
 - Hiring contractors
 - Application Programming Interfaces (APIs)
 - Technology (e.g., cars)
- Amazing things are built when these layer**
 - And the abstraction layers are getting deeper by the day!

We only need to worry about the interface, or specification, or contract *NOT how (or by whom) it's built*

Above the abstraction line

Abstraction Barrier (Interface) (the interface, or specification, or contract)

Below the abstraction line

This is where / how / when / by whom it is actually built, which is done according to the interface, specification, or contract.

1/24/20 UCB CS88 Sp20 L1 34

34

Abstraction: Pitfalls

- Abstraction is not universal without loss of information (mathematically provable). This means, in the end, the complexity can only be “moved around”
- Abstraction makes us forget how things actually work and can therefore hide bias. Example: AI and hiring decisions.
- Abstraction makes things special and that creates dependencies. Dependencies grow longer and longer over time and can become unmanageable.

1/24/20 UCB CS88 Sp20 L1 35

35

Abstraction in CS: Data Type

- What's this?

Real (or ideal) world

42

Computer representation

1/24/20 UCB CS88 Sp20 L1 36

36

Data Types and Operations

- Set of elements**
 - with some internal representation
 - E.g. Integers, Floats, Booleans, Strings, ...
- Set of operations on elements of the type**
 - e.g. $+$, $*$, $-$, $/$, $==$, $<$, $>$, $<=$, $>=$
- Properties**
 - Commutative, Associative, ...
- Expressions are valid well-defined sets of operations on elements that produce a value of a type**

1/24/20 UCB CS88 Sp20 L1 37

37

Lab and HW next week



- Lab will get you setup with Python locally
 - You'll use your own computer
 - Learn about text files and the terminal

1/24/20

UCB CS688 Sp20 L1

38

38

Thoughts for the Wandering Mind



A binary digit (bit) is a symbol from {0,1}.

- How many things can you represent with N bits?
- How many things can you represent with 1 digit (0-9)?
- 2 digits? 6 digits?

1/24/20

UCB CS688 Sp20 L1

39

39