

# Data, Environment and Society (ER131)

## Lecture 1: Introduction

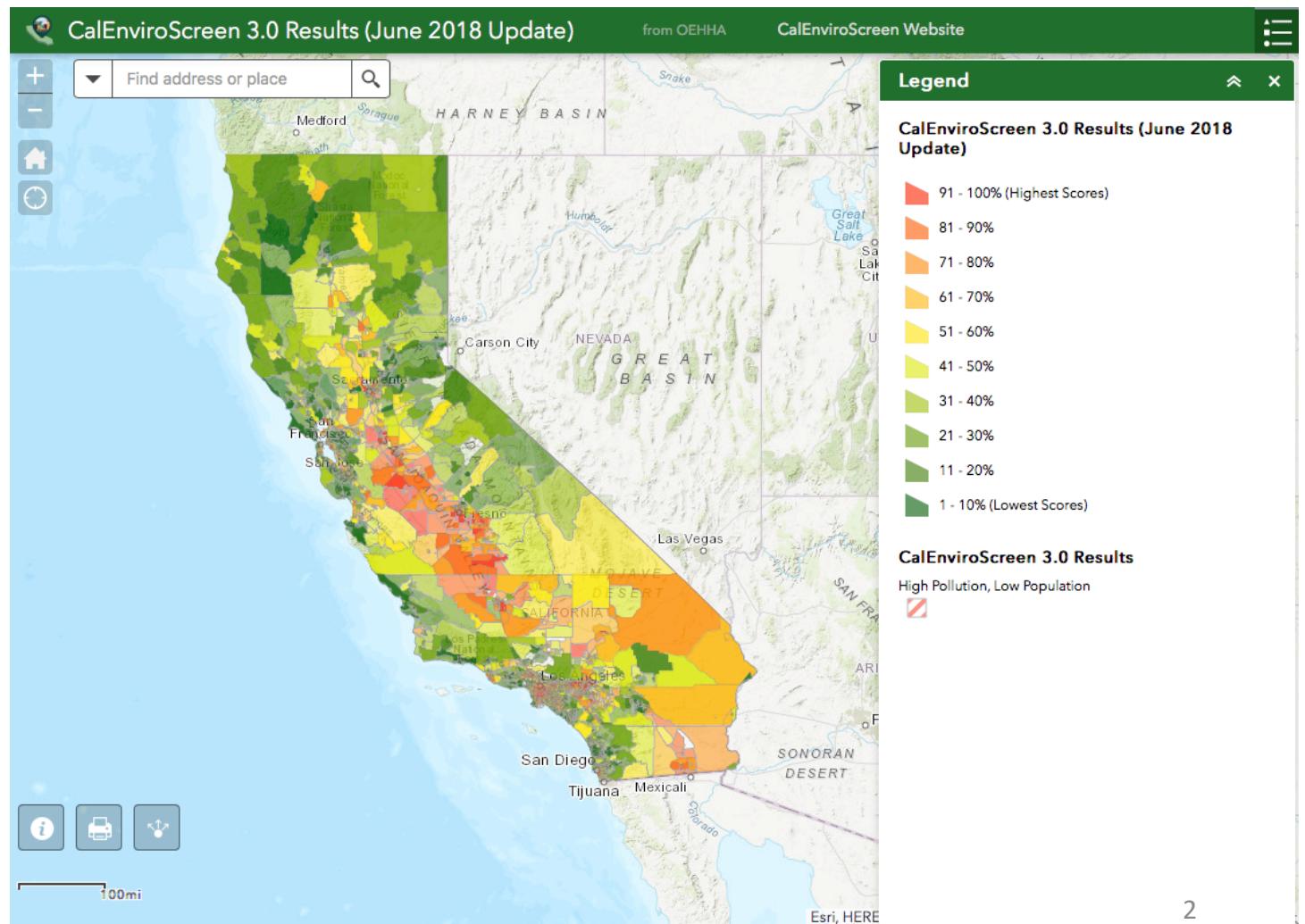
August 29, 2019

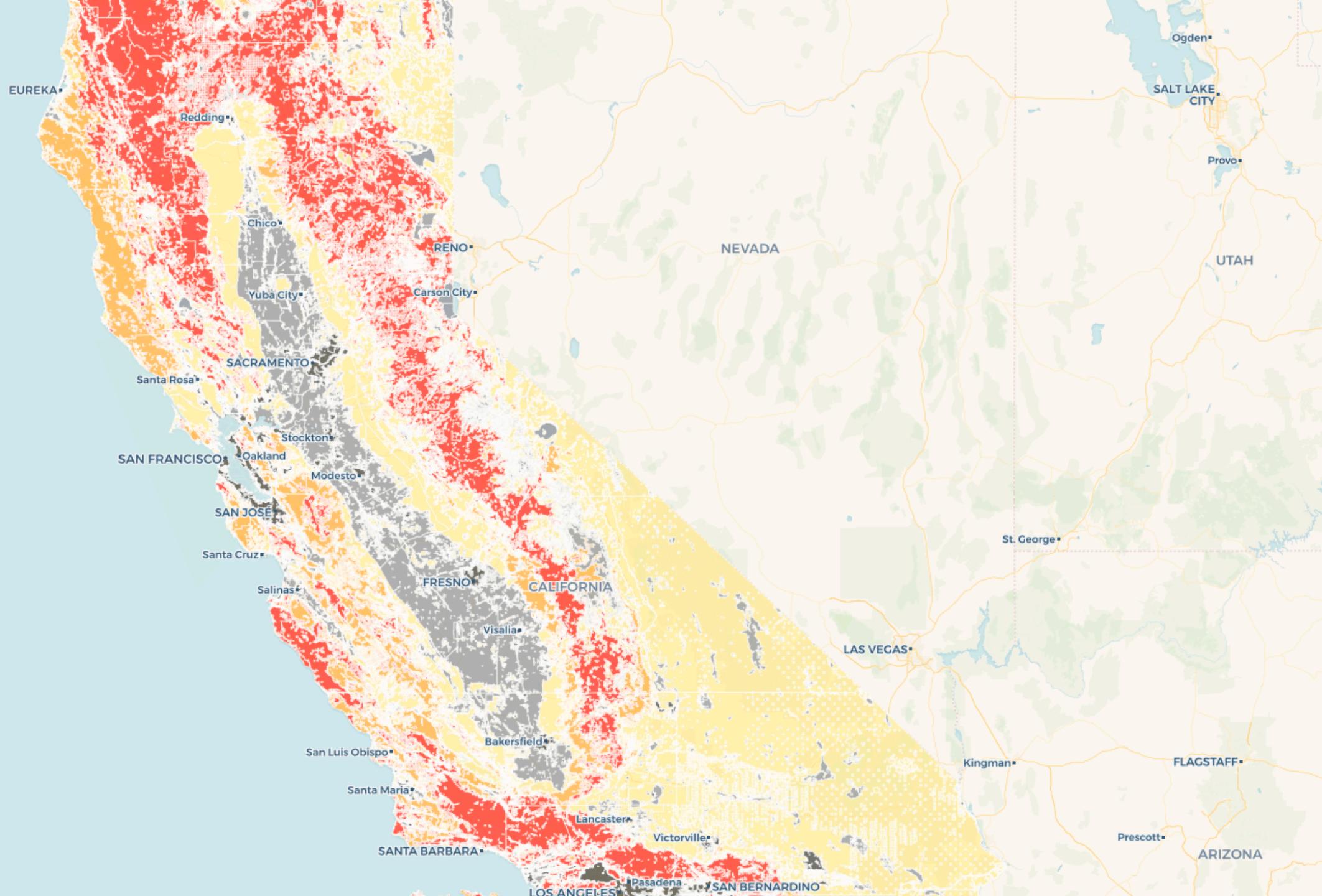
Instructor: Duncan Callaway

GSI: Salma Elmallah

# Why data, environment and society?

1. Energy consumption and other industrial activities generate pollution. Its distribution and impact are heterogeneously distributed.





# Why data, environment and society?

2. Wildfire is caused by a range of factors and has a range of impacts.

The temporality, scope and complexity of data available for analyzing historical fires and *predicting* future fires is enormous



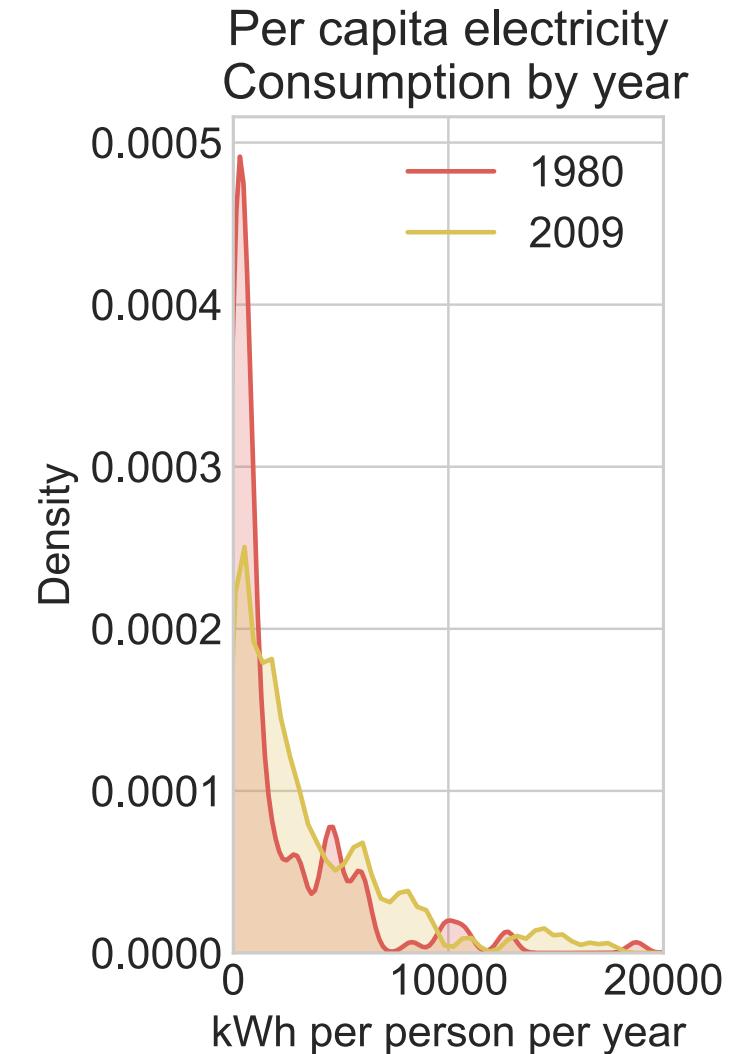
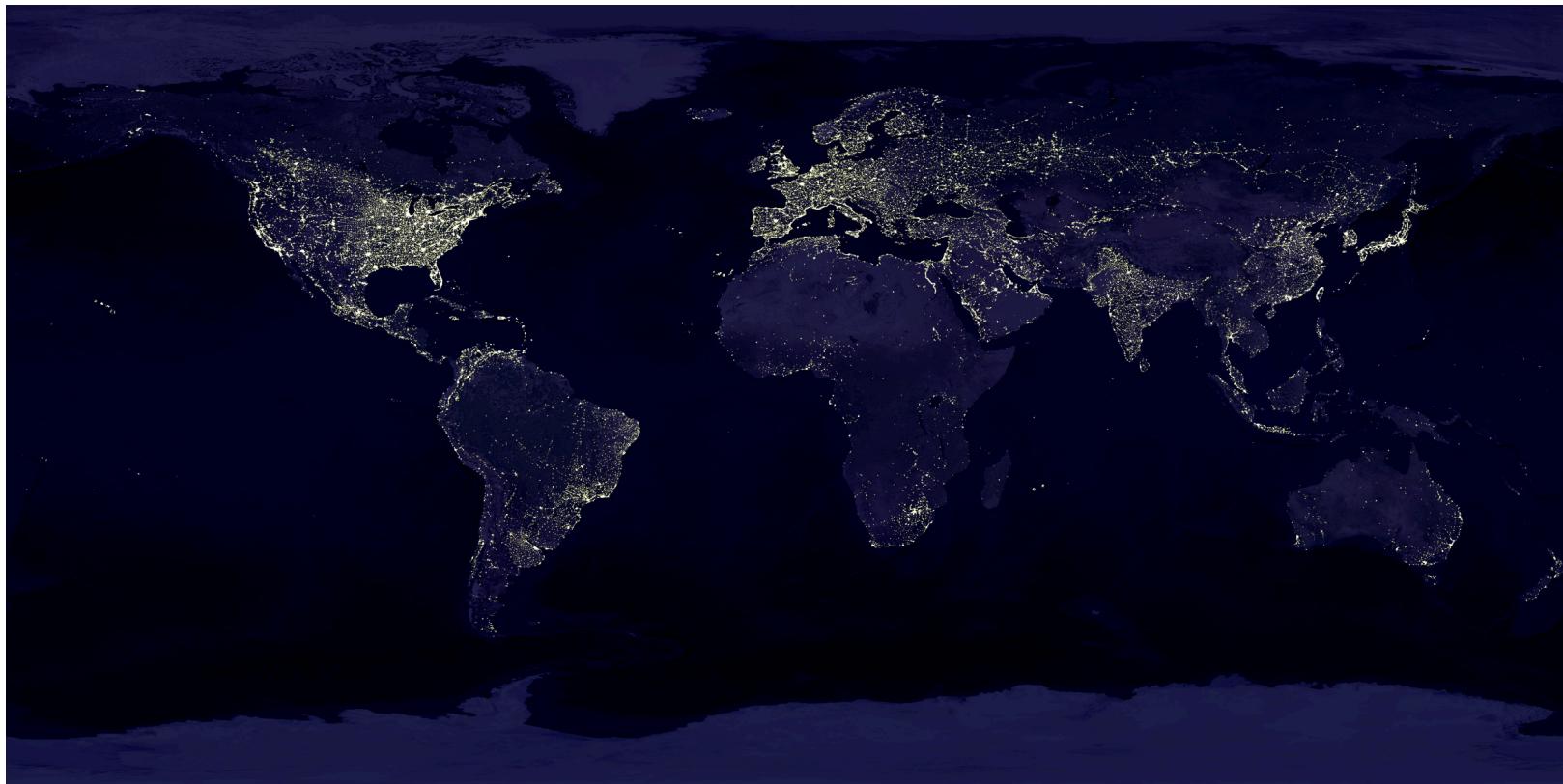
(sfgate)



(nasa)

# Why data, environment and society?

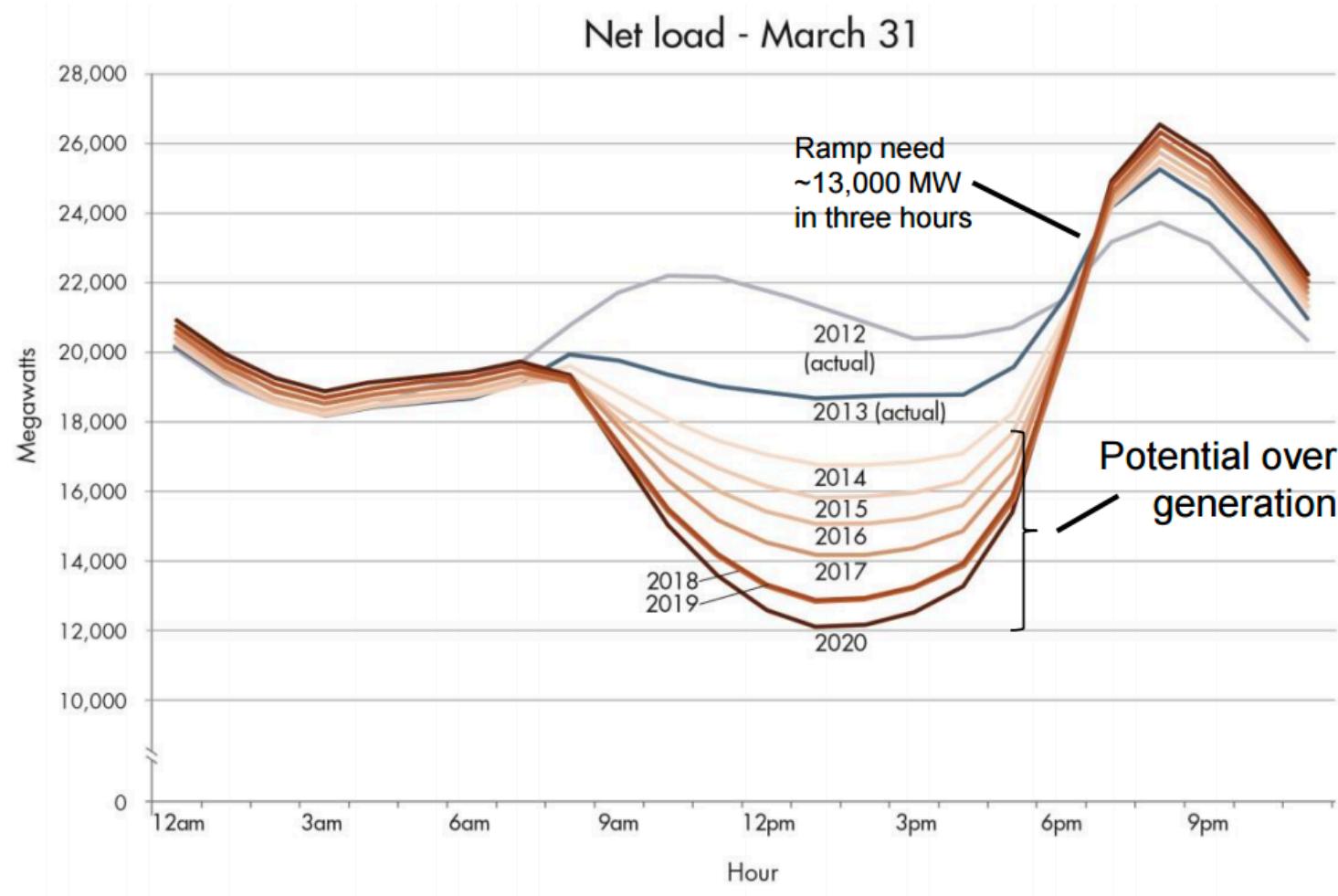
## 3. Energy conversion and use is at the core of nearly all societies



# Why data, environment and society?

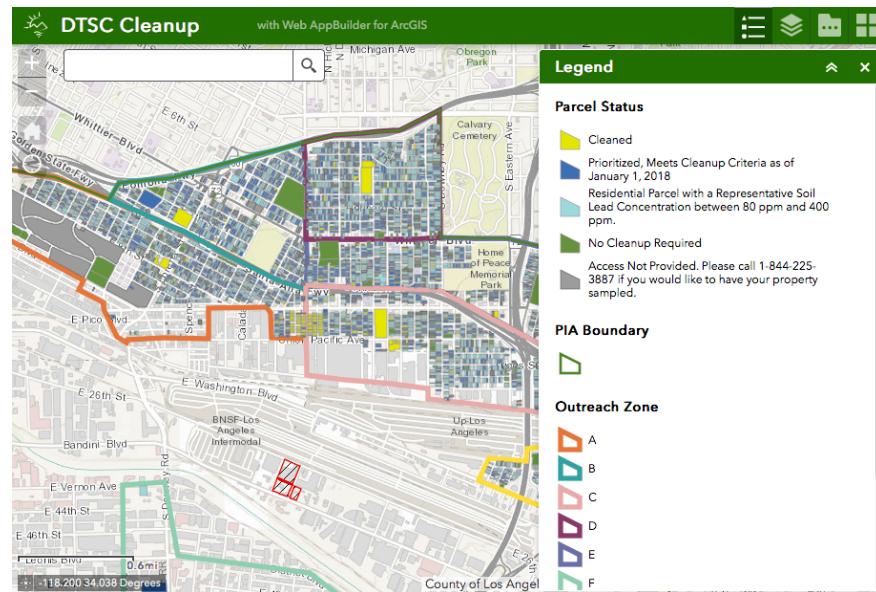
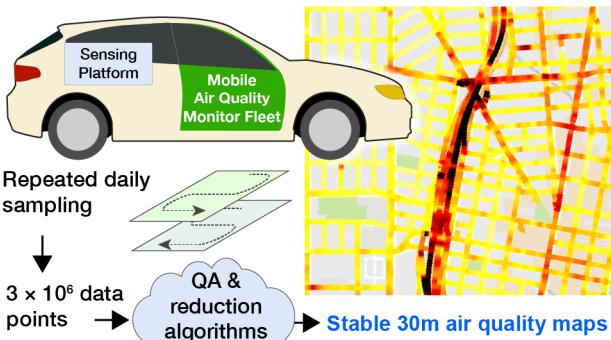
4. As renewable electricity proliferates, electric power systems must transform their operations in response

Analyzing and *predicting* energy production patterns is crucial to economical and reliable future operations



# Why data, environment and society?

## 5. Data (and access to the data) to study these issues is growing



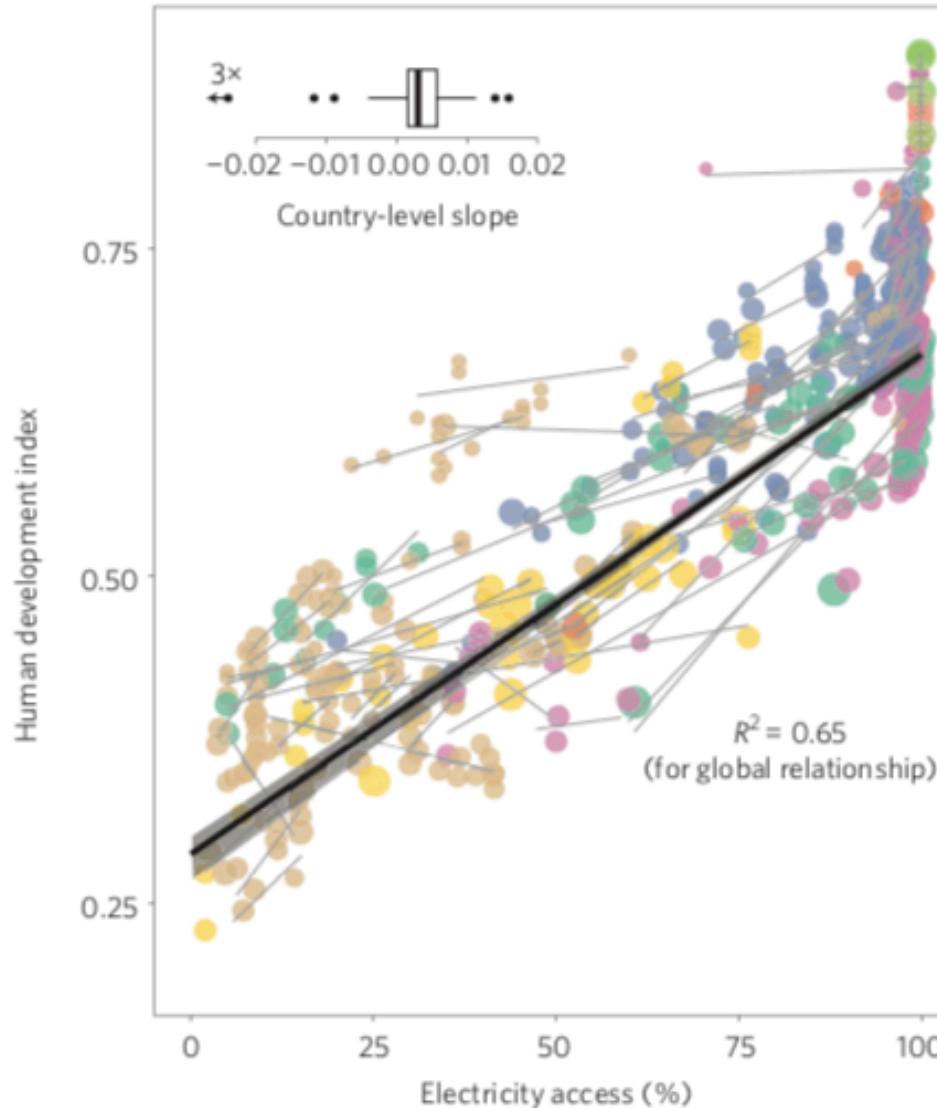
# Why data, environment and society?

6. There is an incredible growing toolkit available to anyone with a computer...
  - a. New Algorithms
  - b. Improving computing infrastructure
  - c. New analysis tools: Jupyter, Python libraries

In sum: *Problems and data are abundant, the algorithms are maturing, and everything is accessible to someone with relatively little training*

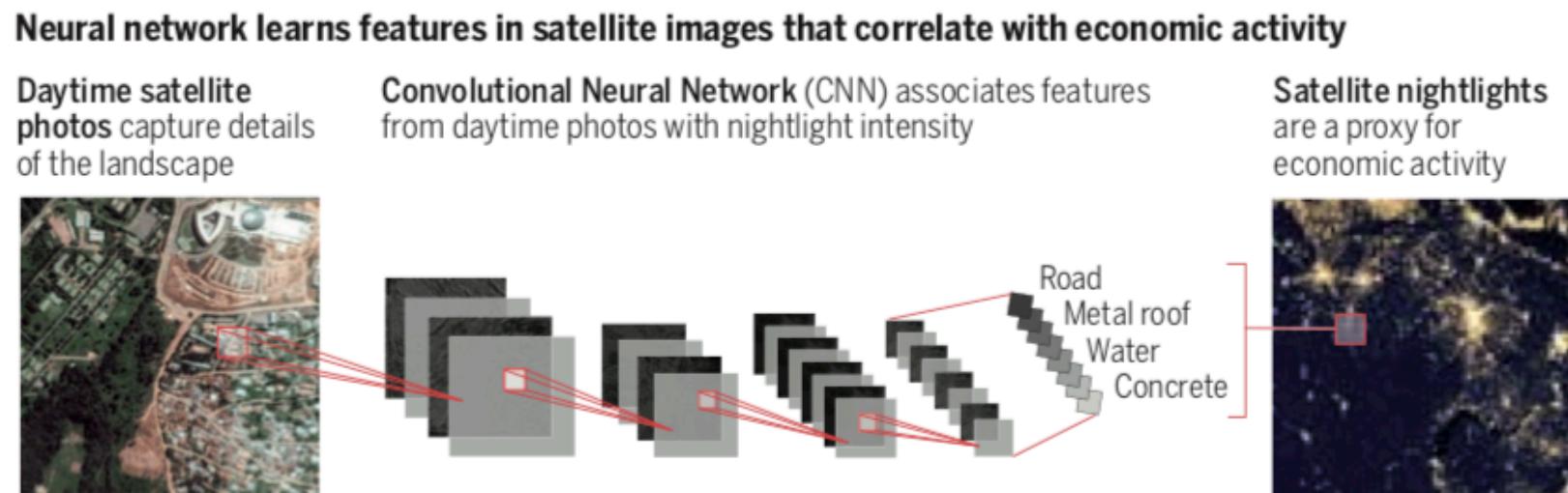
# Example: Human development and energy

- From Alstone, Gershenson and Kammen, Nature Climate Change (2015)

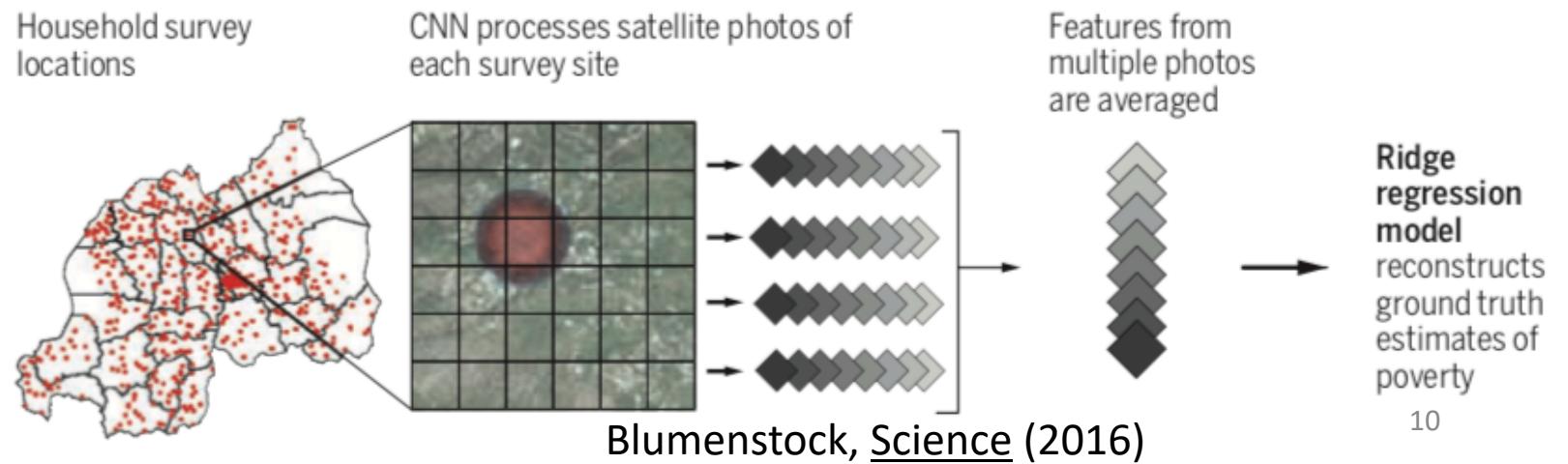


# Example: Predicting poverty

Using day and night satellite data, survey data to predict poverty in locations without surveys (Jean *et al* Science 2018)



**Daytime satellite images can be used to predict regional wealth**



# What is data science?

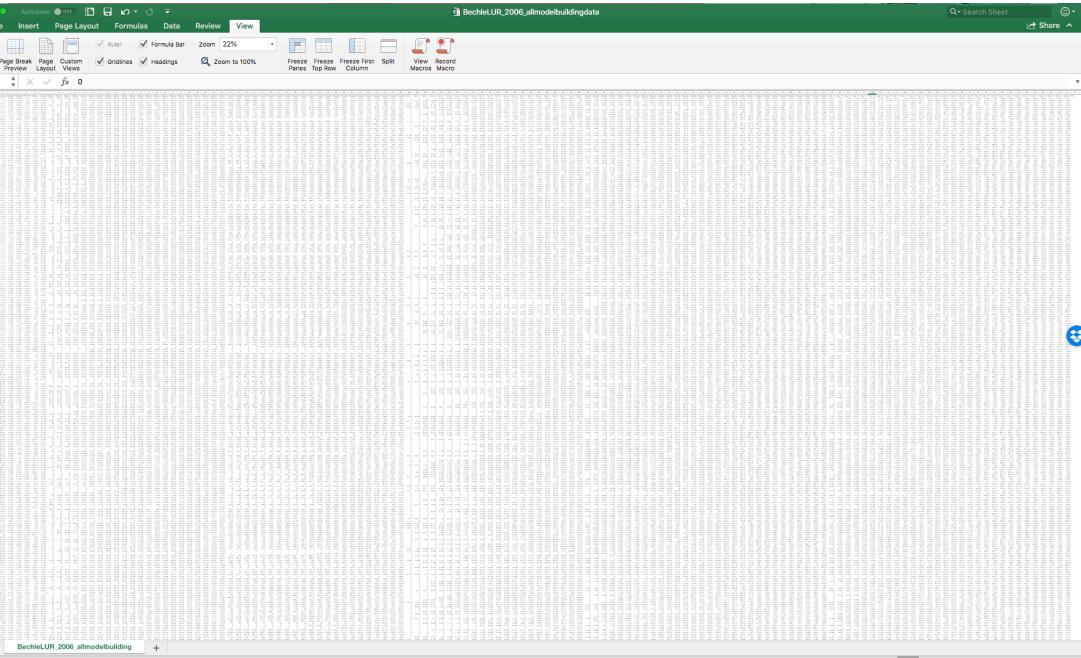
*“Data science focuses on exploiting the modern deluge of data for prediction, exploration, understanding, and intervention”*

--Blei and Smith (2017)

- It is a mix of data, computing and statistics.
- It is not **only** statistics: the algorithms have deep origins in computing and computer science
- It is not **only** computing: one must understand the fundamentals of the origins and characteristics of data

# What does “big data” mean?

- Two characteristics
  - Data sets are large (many data points in total, or cells in a spreadsheet)
  - Data sets are *wide* (many different categories of observations, or columns in a spreadsheet)
- Many of the “data science” things we do in Python are things one could do in excel with small data sets.



The image shows two screenshots of Microsoft Excel. The top screenshot displays a very wide dataset with many columns (labeled A through N) and approximately 100 rows of data. The bottom screenshot shows a similar dataset, but with only 28 rows of data. Both datasets include columns for Monitor\_ID, State, Latitude, Longitude, Observed\_N, Predicted\_N, WRF+DOMIN, Distance\_to\_River, Elevation, and various Imperviousness metrics. The Excel interface at the top includes the ribbon, search bar, and various view and formula options.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
1	Monitor_ID	State	Latitude	Longitude	Observed_N	Predicted_N	WRF+DOMIN	Distance_to_River	Elevation	Impervious_1	Impervious_2	Impervious_3	Impervious_4	Impervious_5
2	04-013-0019-AZ		33.48385	-112.14257	23.8847062	20.9866433	11.6152229	313	0.304	59.4431	59.4715	59.481	59.4572	59.3039
3	04-013-3002-AZ		33.45793	-112.04601	25.0898863	20.9900961	11.4726772	323.8	0.304	72	72	72	72	71.9109
4	04-013-3003-AZ		33.47968	-111.91721	19.2819688	18.0881533	8.9903717	308.4	0.304	53	53	52.7818	52.6609	52.7508
5	04-013-3010-AZ		33.46093	-112.11748	30.6451381	20.3580027	11.9192677	309	0.304	61.3099	62.2828	62.6643	62.8645	62.9876
6	04-013-4011-AZ		33.37005	-112.6207	11.0704122	20.9428562	11.7318916	269.5	0.293	12	12	11.5703	10.7575	10.0111
7	04-013-9997-AZ		33.503731	-112.09581	22.3935655	20.9428562	11.7318916	314.7	0.304	66	66	66	66	65.9252
8	04-019-1011-AZ		32.204411	-110.87867	15.7498554	13.3944963	2.03146195	275.1	0.304	57	57	57	57	56.4274
9	04-019-1028-AZ		32.29515	-110.9823	14.7733132	12.2451094	2.13054681	271.6	0.304	48	48.0235	47.0218	46.5397	46.2783
10	04-027-0006-AZ		32.677855	-114.47586	10.412472	5.90995768	0.02611787	112.4	0.063	16	15.7191	15.3588	15.1553	15.6506
11	05-035-0005-AR		35.197288	-90.193141	11.158921	8.15179949	3.7115612	629.7	0.065	24.0641	24.7609	24.9682	24.5339	24.1025
12	05-119-0007-AR		34.756189	-92.281296	11.8618027	10.0012435	2.8131434	619.7	0.07	40	41.6501	42.6964	43.2522	43.5951
13	06-001-0007-CA		37.6875	-121.7842	14.1472496	11.6085415	5.31101799	41.4	0.138	45.0138	45.2184	45.3083	45.3552	44.6595
14	06-001-1001-CA		37.5358	-121.9619	15.2519112	12.9785652	7.42330694	16.3	0.018	45.1205	45.729	45.9772	45.8182	45.3849
15	06-007-0002-CA		39.7575	-121.84222	9.29446239	7.96905035	0.26094782	215.5	0.061	47	47	47	47	47.0146
16	06-013-0002-CA		37.936	-122.0262	10.8010778	12.2288549	6.42227507	30.2	0.06	43.3642	42.2195	41.8174	41.6142	41.5619
17	06-013-1002-CA		38.010556	-121.64139	7.85950413	6.89874529	4.34174776	68.9	0.001	5.1292	5.6994	6.3623	6.7508	6.9942
18	06-013-1004-CA		37.96028	-122.35667	12.6705497	11.3483599	6.07488012	3.3	0.014	47	47	47	46.9084	44.7375
19	06-013-3001-CA		38.029167	-121.90222	10.896315	9.92793557	5.90302849	39.6	0.003	38.7771	36.2271	34.401	33.5606	33.0885
20	06-019-0007-CA		36.705556	-119.74139	17.1590768	12.0576727	6.70139551	215.4	0.081	38	38	38	37.624	36.173
21	06-019-0008-CA		36.781389	-119.77222	17.105538	12.9702012	6.15611124	220.2	0.091	49.3561	49.1798	49.1201	49.0901	49.0721
22	06-019-0242-CA		36.841389	-119.87444	11.3524899	9.72084333	4.80730915	213	0.091	22	22.3427	22.5483	22.6578	22.3872
23	06-019-4001-CA		36.5975	-119.50361	11.4508176	8.82198205	4.21073055	221.6	0.098	6	6	6.099	6.1838	6.2005
24	06-019-5001-CA		36.819167	-119.71639	14.4078723	11.705453	5.51833105	228.2	0.097	35.5158	35.5079	35.5053	35.504	36.0172
25	06-025-0005-CA		32.676111	-115.48333	13.9211227	6.85386558	1.11637271	115.4	0.003	26.0578	28.0004	29.2236	30.0497	30.5913
26	06-025-0006-CA		32.677778	-115.38972	11.9643335	4.38844285	1.04167461	110.7	0.011	6	6	6	6	6.0321
27	06-025-1003-CA		32.791667	-115.56167	10.9193461	8.60121463	0.89209485	132.9	0	48	48.0144	48.1148	47.8191	47.4461
28	06-029-0007-CA		35.346111	-118.85111	10.5846284	9.48948971	5.43211842	147.7	0.193	13	12.1803	10.877	10.1911	9.7835

# What are those other buzzwords?

**Artificial intelligence**

**Machine Learning**

**Statistical Learning**

These refer to processes to build models that use data to learn from and make decisions about the world based on available observations.

We'll learn from a statistical learning textbook, and many of the tools we use are also known as machine learning.

Central theme of the course: using *prediction*  
methods in support of *resource allocation*  
decisions

Let's unpack this a little further...

# What can we do with data science?

## Rain dances vs. Umbrellas

(adapted from Kleinberg AER 2017)

- Consider two policy makers studying rain
  - One seeks to understand if a rain dance will make it rain
  - One seeks to decide whether or not to invest in umbrellas
- Both require models and data, the analyses are qualitatively different.  
Why?
  - “Rain dance” problems are *impact analyses*. They seek to identify a *causal* effect.
  - We call “umbrella” problems *resource allocation* problems. They seek to *predict* the future so people can decide where to allocate resources and effort.



(amazon.com)

# Is prediction enough for resource allocation?

(Adapted from Susan Athey's 2017 Science article.)

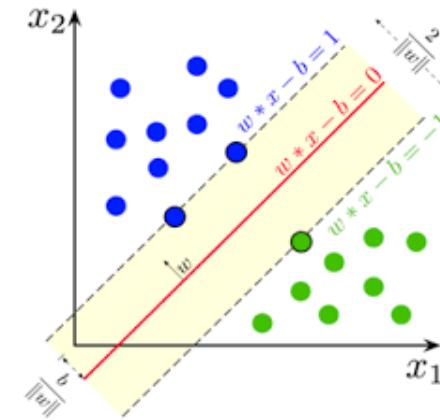
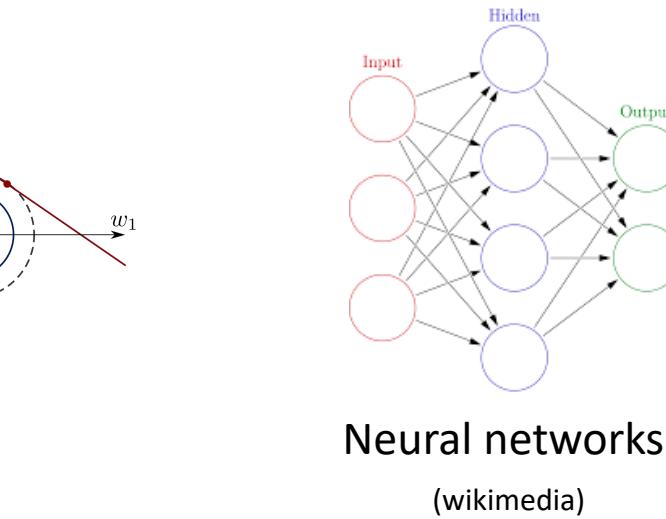
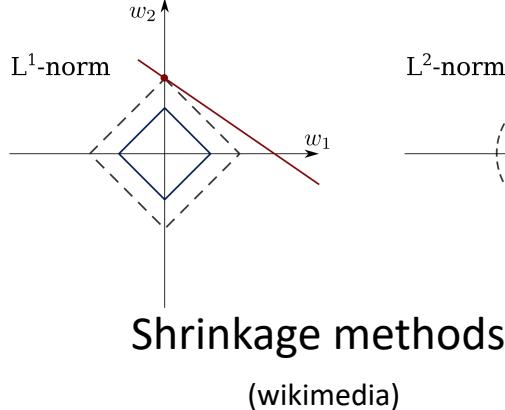
- Consider a scarce resource: Environmental health and safety inspectors at CalEPA.
- If one could predict which chemical facilities are *likely* to have an unplanned toxic release, one could dispatch inspectors to those locations
- Why might this be a bad strategy?
  - The behavior of the facilities being inspected might change during an inspection.
  - The cost to remedy identified problems might be different at different locations, which might make some problems less likely to be solved.



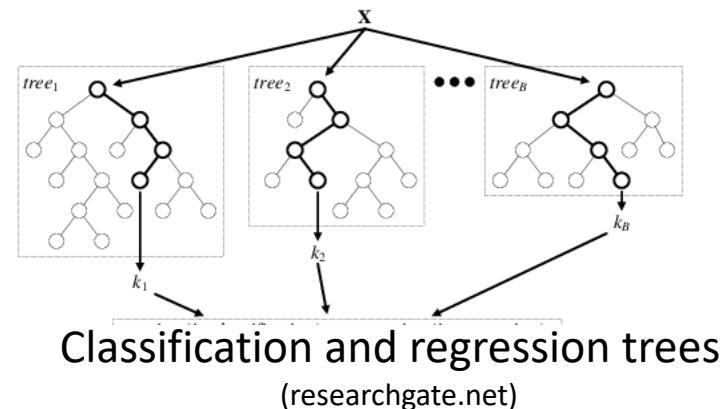
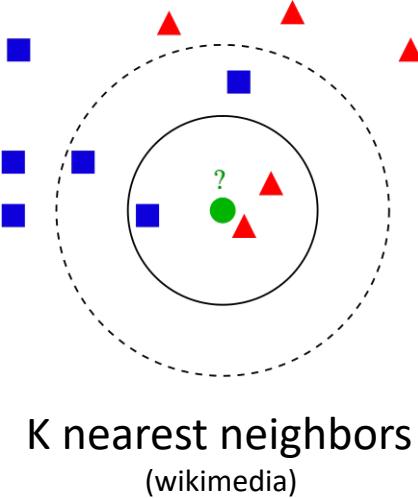
The models we study in this course are good for  
*prediction*

We'll focus on using them to answer resource  
allocation problems

# Some of the specific tools we'll learn



Support vector machines  
(Wikimedia)



- Visualization
- Exploratory data analysis
- Gradient descent
- Pandas and scikit-learn

# The importance of context and narratives

From Pastor's 2007 article in Reclaiming Nature on Environmental Justice (EJ):

*"Traditional environmentalists tend to favor 'rational' processes of debate, objective scientific research on hazards and their risks...such frameworks tend to produce negotiation between businesses and their hired experts..."*

*"EJ activists favor 'democratic' epistemologies in which community participation facilitates story-telling about lived experiences...in the minds of many EJ advocates, this allows for community empowerment."*

Which do you prefer?

# What's the awareness gap?

- For many in the room, talking about these issues is abstract, but for others it may not be
- We will discuss ideas about intent vs impact, for example in identifying environmental racism
- It is also important to think about intent vs impact in the words we use.
- ...it matters little if you have the best of intentions if your actions still have a negative or unfair impact on others



# There is a human side to the algorithms themselves...

TheUpshot

## *Who's to Blame When Algorithms Discriminate?*

A proposed rule from HUD would make it harder to hold people accountable for subtler forms of discrimination.

DEPT. OF TECHNOLOGY NOVEMBER 12, 2018 ISSUE

## IN THE AGE OF A.I., IS SEEING STILL BELIEVING?

*Advances in digital imagery could deepen the fake-news crisis—or help us get out of it.*



By Joshua Rothman November 5, 2018

f t m



# What are the course objectives?

- Teach students to build, estimate and interpret models that describe phenomena in the areas of energy & environmental decision-making.
- We will
  - Focus on analysis and prediction
  - Learn a suite of data-driven modeling approaches
  - Build the programming and computing skills to use those models (in Python and Jupyter notebooks)
  - Develop the expertise to formulate questions that are appropriate for available data and models.
- Students will leave the course as both critical consumers and responsible producers of data driven analysis.

Now let's go into the nuts and bolts of the course

# Prerequisites

- (required) Foundations of Data Science (CS/ INFO/ STAT C8)
- (recommended) Computing: An introductory programming course (CS61A or CS88).
- Math:
  - (required) High school or college calculus.
  - (recommended) Linear Algebra (Math 54, EE 16a, or Stat89a).

# How does this course fit into Berkeley's curriculum?

- The course is designed to fit into Berkeley's emerging data science curriculum. It is similar to Data 100.
  - But we will place a stronger emphasis on how to use *prediction* methods as decision-making tools in energy and environment contexts
  - It has less emphasis on web technologies, working with text, databases and statistical *inference*.
- The course satisfies:
  - Upper division domain emphasis for Data Science major
  - Engineering Elective for Energy Engineering
  - Upper division requirement for Energy and Resources Group minor



# What resources are required?

- Hardware: You'll need a laptop.
  - Mac, Windows, Linux, Chromebook all ok.
  - If you don't have a laptop please see me or Salma after class.
- Software:
  - All work can be done in the cloud, on Berkeley's datahub
  - If you wish, you can install Python locally with the Anaconda distribution.
  - We'll use Python 3.x
- Internet connection: you'll be computing the cloud.
- Read-ware:
  - Introduction to Statistical Learning, available as PDF or for sale from online retailers.
  - The DS100 online textbook
  - Other reading will be distributed through github.

# Note taking in lecture

- I'll post slides ahead of time in github. They'll be numbered.
- Prior students report that taking notes on paper during lecture works best
  - ...then review slides after.
- If you wish you can take notes on your computer
- There will be a few lectures where you'll need your computer
- (...and you'll always need your computer in lab sections)

# Github and bCourses

- bCourses
  - We'll use it for distributing homework (mainly hyperlinks to jupyter notebooks).
  - You'll upload your homework there too.
- Github: [https://github.com/duncancallaway/ER131\\_2019](https://github.com/duncancallaway/ER131_2019)
  - Readings posted here
  - Lecture notes here
  - Homework and lab files
  - All other materials

# Datahub

- Berkeley's super cool jupyter notebook cloud server.
- You'll get there by using links in posted labs and homeworks.
- For example, here is [hw01](#)
  - You can do all your work here.
  - Your personal copy will get saved here.
  - When you're done, create a pdf, save it and upload it to bCourses for grading.

# Datahub vs Anaconda

- If you wish, you can install a local version of Python and a starter package of data science libraries with the “Anaconda” distribution
  - See <http://www.ds100.org/sp18/setup> for instructions on how to do this.
- If you encounter problems running a notebook, the first thing we’ll ask is “is it working on datahub?”
  - We can’t provide support if the answer is yes...
- In other words: You don’t *\*have\** to use Datahub. But we encourage it.

# How are students evaluated?

1. Homework (20%). Submission by submitting pdf of Jupyter notebook. 10 assignments, we drop the lowest grade.
2. Labs (15%). Submission by submitting pdf of Jupyter notebook. 9 assignments, we drop the lowest grade.
3. Mid term (25%). In class on November 19
4. Final project:
  - a. Poster session presentation (10%). December 17
  - b. Jupyter notebook (30%). Due morning of December 18

# Late policy

- You may request up to two extensions of two days over the course of the semester.
- The poster must be presented during the poster session to receive credit.
- For the final project, we drop 10 points out of 100 for each day late, or roughly a full letter grade. Projects submitted after 11:59am on December 20, 2019 will not receive credit.

# Working in groups

- Homework and labs
  - Brainstorm together!
  - However the work you submit must clearly be your own. Be sure to finish assignments on your own. Comments and markdown cells must clearly be your own.
- Final project
  - You must work in groups of 2-3 for the final project

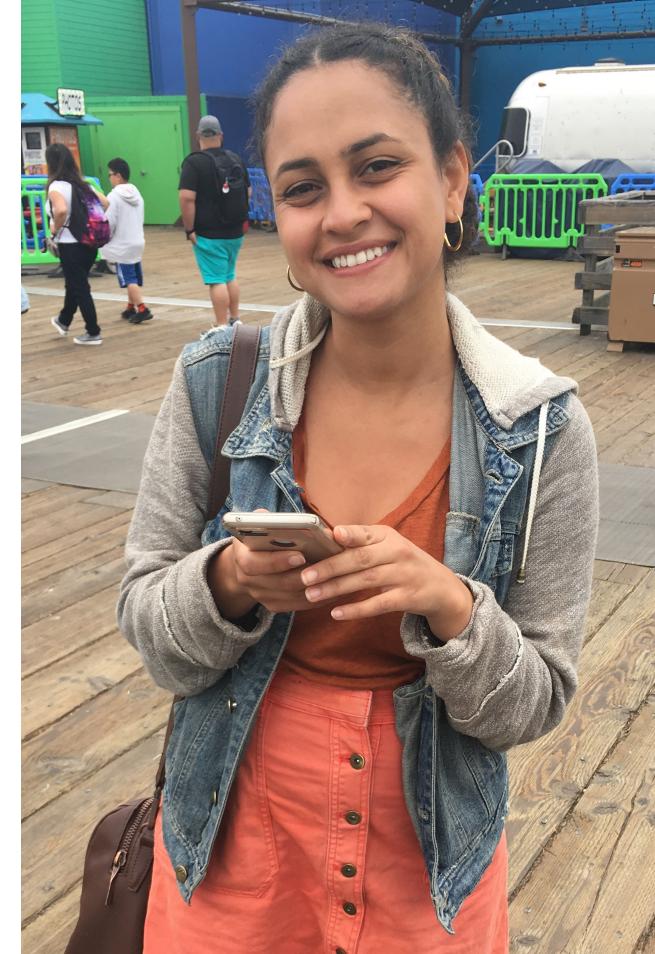
# Final project

- We'll get started in late September
- We'll ask you to do small steps toward the final goal on your homework assignments
- The final project will be structured to include
  - A clear motivating resource allocation question
  - Data cleaning
  - Exploratory data analysis
  - A deep exploration of several prediction methods
- We'll give a few project ideas, but the goal is for you to work independently

# Who are the instructors?



Duncan Callaway



Salma Elmallah

# When and where?

- Lectures: Tu/Th 9:30-11, Wheeler 202
- Labs:
  - Mo 10-12, Barrows 166
  - Tu 2-4, VLSB 2032
- Office hours
  - Duncan (Barrows 325): We 9:30-10:30, by appointment;  
Th 11:15-12:15, open door.
  - Salma (Barrows 399): We 3-4, open door
- Labs due Mondays, pdf upload to bCourses
- HWs due Thursdays before class, pdf upload to bCourses

# Next week

- No lab meetings (due to Labor Day)
  - Salma will hold extended office hours during Tuesday lab time, in VLSB 2032
- Lab 1 notebook will be distributed for you to complete on Monday
- Tuesday 9/3:
  - Reading: Science and Data Science, available on Github
- Thursday 9/5:
  - Reading: Chapter 3 from the DS100 textbook – see link in Github
- HW 1 will be released Thursday 9/5.