Research Design and pandas

Ivan Corneillet

Data Scientist



Learning Objectives

After this lesson, you should be able to:

- Define a problem and types of data
- Identify dataset types
- Define the data science workflow
- Apply the data science workflow in the pandas context
- Write an iPython notebook to import, format, and clean data using the pandas library

Outline

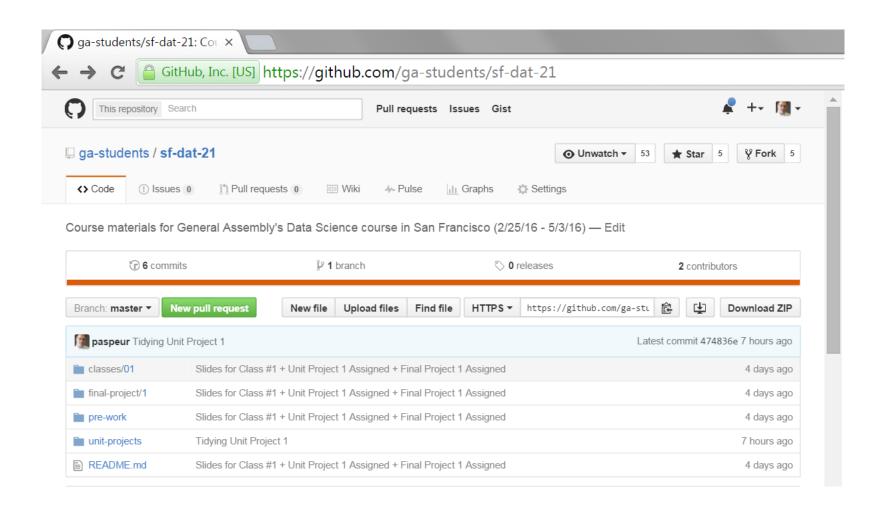
- Course GitHub Repository
- Pre-Work
- Data Science Workflow Review
- Identify the problem
 - The Why's and How's of a Good Question
 - The SMART Goals Framework
- Acquire the Data
 - Data Types
 - Logistics of Acquiring Data
 - Our "Zillow" Dataset

- Tidying Data
- File Formats
- Parse the Data
 - Documentation and Data Dictionaries
 - Codealong Introduction to pandas
 - Codealong Tidying up the Zillow dataset
 - Lab
- Review
- Unit Project 1 (due next session on 3/3)



Course GitHub Repository

SF-DAT-21's Web Interface





Pre-Work

Pre-Work

Before this lesson, you should already be able to:

- Have completed Python pre-work
- Open and create iPython notebooks



Getting a local copy of the repo via the command line

First, get a command line prompt



Then, clone the repo with the following command:

git clone https://github.com/ga-students/sf-dat-21

```
git clone https://github.com/ga-students/sf-dat-21
```

When the cloning is done ...

```
$ git clone https://github.com/ga-students/sf-dat-21
Cloning into 'sf-dat-21'...
remote: Counting objects: 50, done.
remote: Compressing objects: 100% (37/37), done.
remote: Total 50 (delta 9), reused 47 (delta 6), pack-reused 0
Unpacking objects: 100% (50/50), done.
  Checking connectivity... done.
```

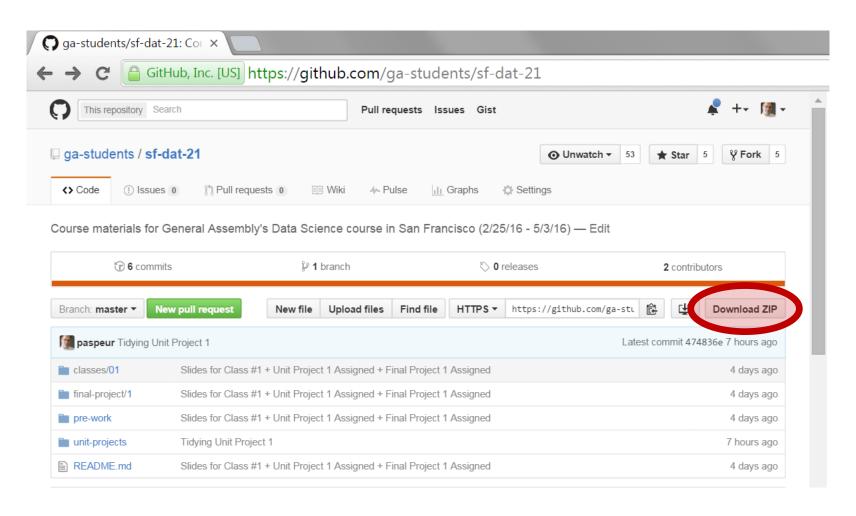
... the repo content is under sf-dat-21

```
~/sf-dat-21
$ git clone https://github.com/ga-students/sf-dat-21
Cloning into 'sf-dat-21'...
remote: Counting objects: 50, done.
remote: Compressing objects: 100% (37/37), done.
remote: Total 50 (delta 9), reused 47 (delta 6), pack-reused 0
Unpacking objects: 100% (50/50), done.
Checking connectivity... done.
$ cd sf-dat-21
  van ~/sf-dat-21
classes final-project pre-work README.md unit-projects
  van ~/sf-dat-21
```



Alternative method

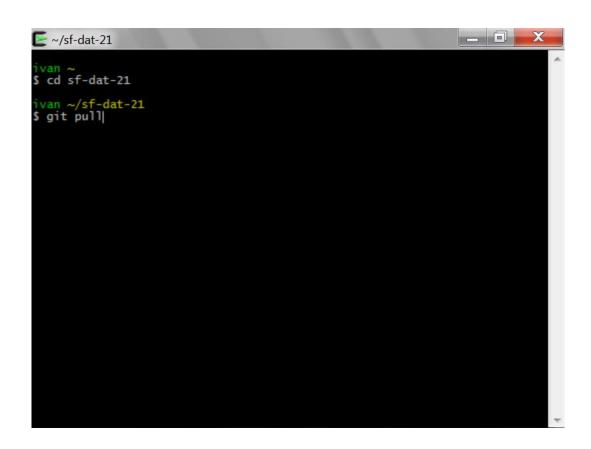
Alternatively, download the entire repository as a ZIP file





Updating your local copy of the repository via the command line

First cd to your local copy, then update the repo with git pull



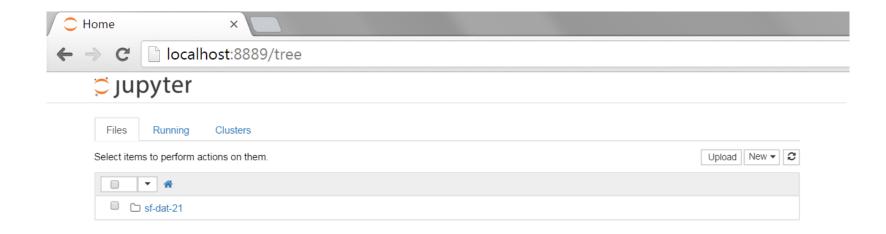
Updating the local copy (cont.)

```
E ∼/sf-dat-21
 cd sf-dat-21
ivan ~/sf-dat-21
$ git pull
Already up-to-date.
 van ~/sf-dat-21
```

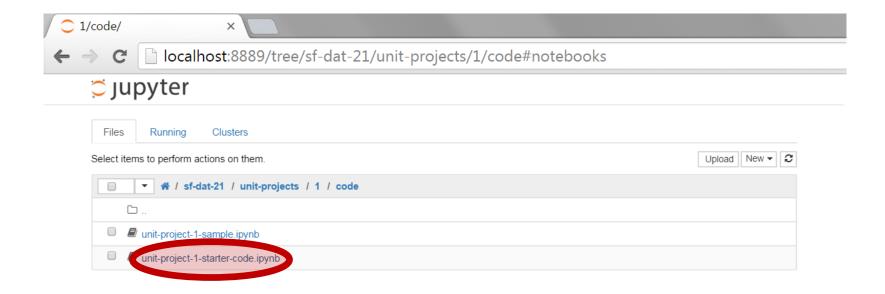


Verify that Anaconda for Python 2.7 was correctly installed

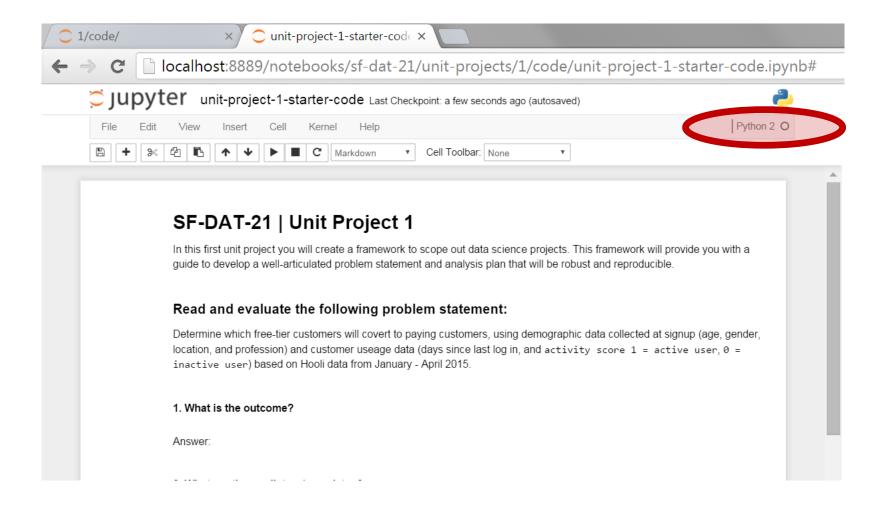
First, start IP[y]: IPython (Py 2.7) Notebook



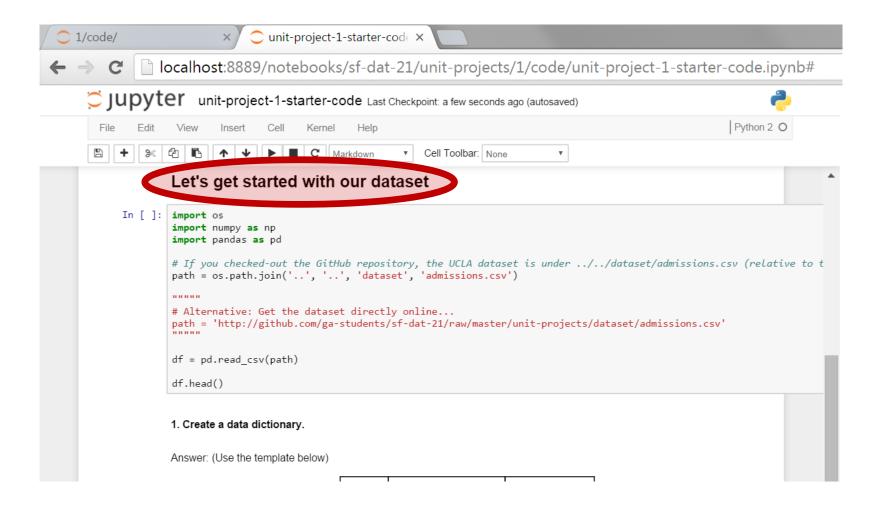
Navigate to Unit Project 1 Stater Code



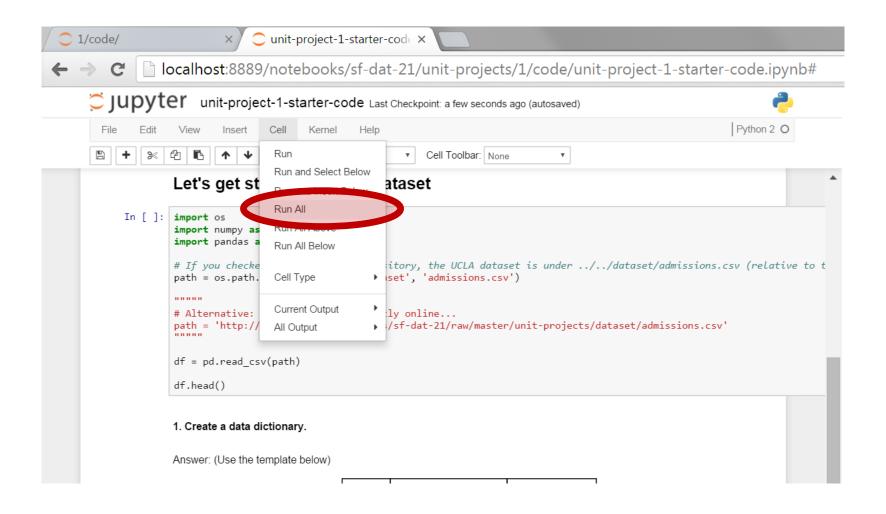
Using Python 2?



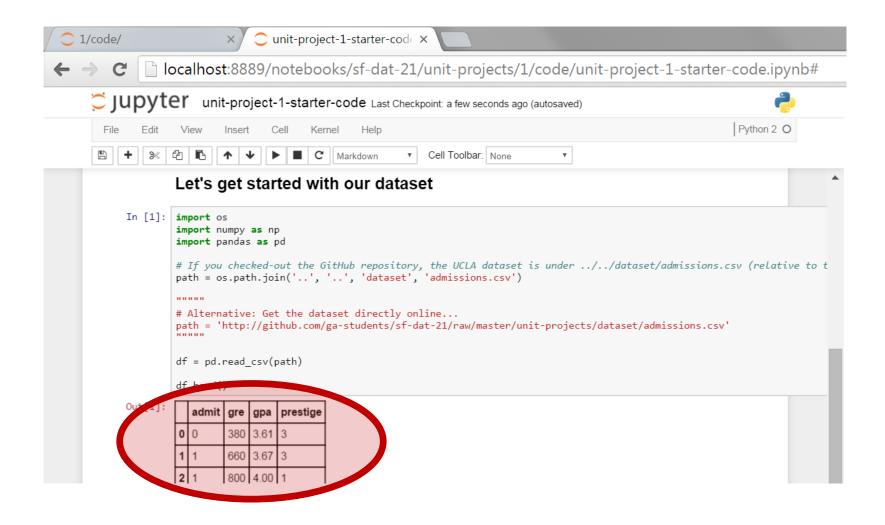
Scroll down to "Let's get started with our dataset"



Cell > Run All



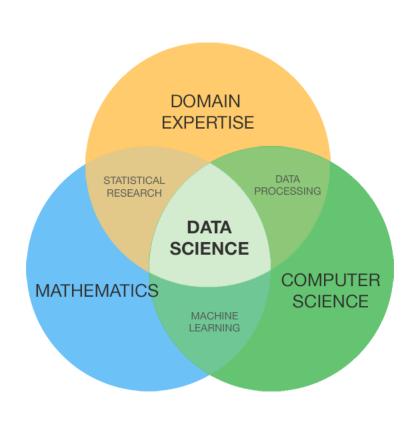
If all goes well, the dataset's head is shown

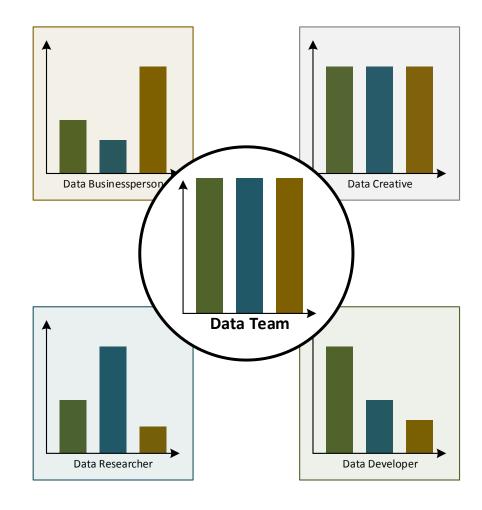




Review What is Data Science? Who are Data Scientists?

What is Data Science? Who are Data Scientists?

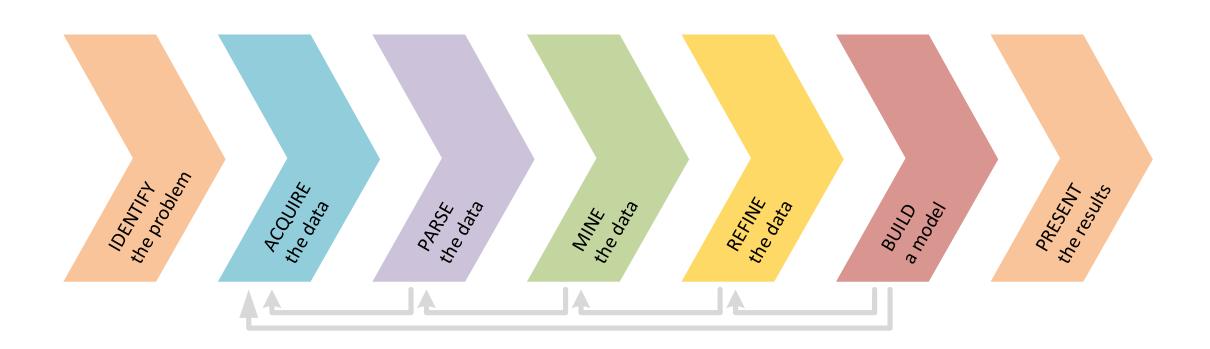




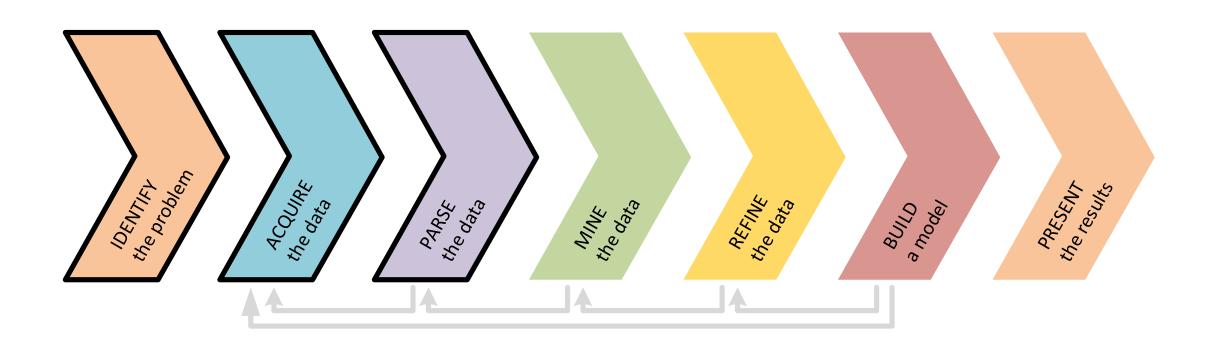


Review and Exercise Data Science Workflow

The Data Science Workflow (again...)



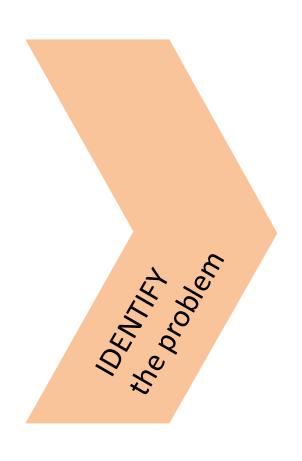
Today we'll focus on the first three (IDENTIFY the problem, ACQUIRE the data, and PARSE the data)





1 IDENTIFY the Problem

• Identify the Problem



- Identify the Problem
 - Identify business/product objectives
 - Identify and hypothesize goals and criteria for success
 - Create a set of questions for identifying correct dataset

• Identify the Problem (cont.)

- Identify the Problem
 - Identify business/product objectives
 - Identify and hypothesize goals and criteria for success
 - Create a set of questions for identifying correct dataset

- The Why's and How's of a GoodQuestion
- → The SMART Goals Framework



The Why's and How's of a Good Question

Why do we need a good question?

 "The scientist is not a person who gives the right answers, he's one who asks the right questions." – Claude Lévi-Strauss





- "If they can get you asking the wrong questions, they don't have to worry about answers." – Thomas Pynchon
- "Judge a man by his questions rather than by his answers." –
 Voltaire



By asking a good question and setting a clear aim:

- You set yourself up for success
 - "A problem well stated is half solved" Charles Kettering
- You help other data scientists learn from and reproduce your work
 - You establish the basis for making your analysis reproducible
- You also help them expand on your work in the future



The SMART Goals Framework

The SMART Goals Framework provides a good foundation to set a clear aim

Specific	Who, What, Where, When, Why, Which	Define the goal as much as possible, with no ambiguous language. WHO is involved, WHAT do I want to accomplish, WHERE will it be done, WHY am I doing this — reasons, purpose, WHICH constraints and/or requirements do I have?	
MEASURABLE (MEANINGFUL)	From and To	Can you track the progress and measure the outcome? How much, how many, how will I know when my goal is accomplished?	
ATTAINABLE (ACTION ORIENTED)	How	Is the goal reasonable enough to be accomplished? How so? Make sure the goal is not out reach or below standard performance	
Relevant (REALISTIC)	Worthwhile	Is the goal worthwhile and will it meet your needs? Is each goal consistent with other goals you have established and fits with your immediate and long term plans?	
TIMELY (TIME-BOUND)	When	Your objective should include a time limit. "I will complete this step by month/day/year." It will establish a sense of urgency and prompt you to have better time management	

The SMART Framework can be tuned up for DS

Specific	The dataset and key variables are clearly defined		
MEASURABLE	The type of analysis and major assumptions are articulated		
ATTAINABLE	The question you are asking is feasible for your dataset and is not likely to be biased		
Reproducible	Another person (or you in 6 months!) can read your state and understand exactly how your analysis is performed		
T _{IME-BOUND}	You clearly state the time period and population for which this analysis will pertain		

Trends often change over time and vary by the population of source of your data. It is important to clearly define who/what you included in your analysis as well as the time period for the analysis

Determine the association of foods in the home with child dietary intake. Is this Aim SMART?



 Using one 24-hour recall from the cross-sectional NHANES 2007-2010 we will determine the factors associated with food available in the homes of American children and adolescents. We will test if self-reported availability of fruits, dark green vegetables, low fat milk or sugar sweetened beverages available in the home increases the likelihood that children and adolescents will meet their USDA recommended dietary intake for that food. Our hypothesis is that children will be more likely to meet the USDA recommended intake level when food is always available in their home compared to rarely or never

The aim is Specific

SPECIFIC

The dataset and key variables are clearly defined

Determine the association of foods in the home with child dietary intake. Using one 24-hour recall from the cross-sectional NHANES 2007-2010 we will determine the factors associated with food available in the homes of **American** children and adolescents. We will test if self-reported availability of fruits, dark green vegetables, low fat milk or sugar sweetened beverages available in the home increases the likelihood that children and adolescents will meet their **USDA recommended dietary intake** for that food. Our hypothesis is that children will be *more likely* to meet the USDA recommended intake level when food is always available in their home compared to rarely or never

The aim is Specific (cont.)

☐ How data was collected: □ 24-hour recall, self-reported ☐ What data was collected: ☐ Fruits, dark green vegetables, low fat milk or sugar sweetened beverages, always vs. rarely available The dataset and key ☐ How data will be analyzed: variables are ☐ Using USDA recommendations as a gold-standard to measure the clearly association defined ☐ The specific hypothesis and direction of the expected associations: ☐ American children and adolescents will be more likely to meet their recommended intake level

The aim is Measurable

MEASURABLE

The type of analysis and major assumptions are articulated

- ☐ Determine the association of foods in the home with child dietary intake
- ☐ We will test if the self-reported availability of certain foods increases the likelihood that children and adolescents will meet their USDA recommended dietary intake for food

The aim is Attainable

ATTAINABLE

The question
you are
asking is
feasible for
your dataset
and is not
likely to be
biased

- ☐ Association, not causation
 - ☐ We are determining the association between two items (food availability in homes and children meeting their dietary recommendations)
 - ☐ Because we are using cross-sectional data, we cannot say that having fruits and vegetables at home actually causes children to meet their dietary requirements
- Note that cross-sectional data has inherent limitations, one being that causal inference is typically not possible

The aim is Reproducible

REPRODUCIBLE

Another person can read your state and understand exactly how your analysis is performed

☐ With all the specifics, it would be straightforward to pull the data from NHANES and reproduce the analysis

The aim is Time-Bound

TIME-BOUND

You clearly
state the
time period
and
population
for which this
analysis will
pertain

☐ Using one 24-hour recall from NHANES 2007-2010, we will determine the factors associated with food available in the homes of American children and adolescents

Context is important

- The previous example laid out research goals but in a business setting, you will need to articulate business objectives
 - Example: Success for the Netflix recommendation engine may be if 70% of customers over the age of 18 select a movie from the recommended queue during Q3 of 2015
- Regardless of setting, start your question with the SMART framework to help achieve your objectives



Activity: Knowledge Check

Activity: Knowledge Check



ANSWER THE FOLLOWING QUESTIONS (5 minutes)

- 1. Which of the following uses the SMART framework? Why? What is missing?
 - a. I am looking to determine the value of 2 bedroom homes in San Francisco
 - b. Determine what real estate factors besides number of bedrooms (e.g., bathrooms, home size, lot size, amenities, etc.) explains the difference in values between 1 bedroom and 2 bedroom homes in San Francisco using sale data from Zillow between November 2015 and February 2016
 - c. When finished, split into pairs and share your answers with each other

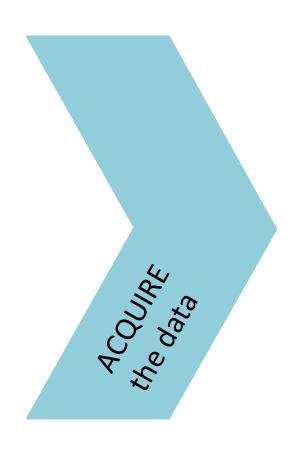
DELIVERABLE

Answers to the above questions



2 ACQUIRE the Data

2 Acquire the Data



- Acquire the Data
 - Identify the "right" dataset(s)
 - Import data and set up local or remote data structure
 - Determine most appropriate tools to work with data

2 Acquire the Data (cont.)

- Acquire the Data
 - Identify the "right" dataset(s)
 - Import data and set up local or remote data structure
 - Determine most appropriate tools to work with data

- Data Types
- Logistics of Acquiring Data
- Our "Zillow" Dataset
- Tidying Data
- File Formats

2 Acquire the Data (cont.)

- Acquire the Data
 - Identify the "right" dataset(s)
 - Import data and set up local or remote data structure
 - Determine most appropriate tools to work with data

- Questions to ask:
 - What type of data is it, cross-sectional or longitudinal?
 - How well was the data collected?
 - Is there much missing data?
 - Was the data collection instrument calibrated?
 - Is the dataset aggregated?
 - Do we need pre-aggregated data?

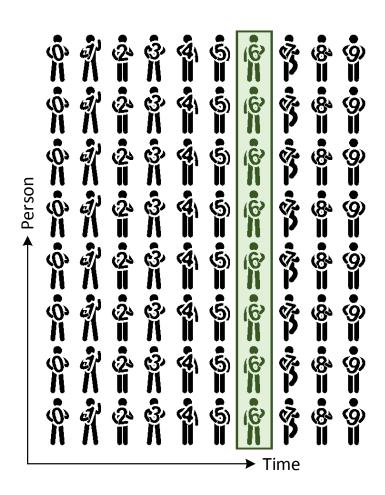


Data Types

Why Data Types Matter

- Different data types have different limitations and strengths
- Certain types of analyses aren't possible with certain data types
- There are 2 types of data which we may might use for analysis

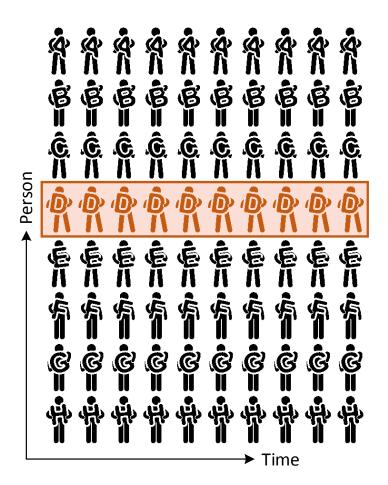
Cross-Sectional Data



- All information is determined at the same time; all data comes from the same time period
- Note: There is no distinctionbetween exposure and outcome

Time Series/Longitudinal Data

Information is collected over a period of time



Data Types: Strengths and Weaknesses

	Strengths	Weaknesses
Cross-Sectional Data	 Often population-based Generalizable Less expensive compared to other types of data collection methods 	 Separation of cause and effect may be difficult (or impossible) Variables/cases with long duration are over-represented
Time Series/Longitudinal Data	 Unambiguous temporal sequence; exposure precedes outcome Multiple outcomes can be measured 	 Takes a long time to collect data Vulnerable to missing data More expense compared to other types of data collection methods



Activity: Knowledge Check

Activity: Knowledge Check



ANSWER THE FOLLOWING QUESTIONS (5 minutes)

- 1. What type of data is shown by Zillow? (http://www.zillow.com/san-francisco-ca/sold/)
- 2. Can you create a cross-sectional analysis from a longitudinal data collection? How? Is this applicable from the data above?
- 3. When finished, split into pairs and share your answers with each other

DELIVERABLE

Answers to the above questions



Activity: Write a Research Question with Raw Data

Activity: Write a Research Question with Raw Data



DIRECTIONS (10 minutes)

- 1. Individually, look at the data from <u>Kaggle's Titanic competition</u> (https://www.kaggle.com/c/titanic/data) (also in the course repo) and write a high quality research question
- 2. Make sure you answer the following questions:
 - a. What type of data is this, cross-sectional or longitudinal?
 - b. What will we be measuring?
 - c. What is the SMART aim for this data?
- 3. When finished, split into pairs and share your answers with each other

DELIVERABLE

Research Question



Logistics of Acquiring Data

Logistics of Acquiring Data

- Data can be acquired through a variety of sources
 - Web (e.g., Google Analytics, HTML)
 - Databases
 - SQL (Structured Query Language)
 - NoSQL ("Not only SQL")

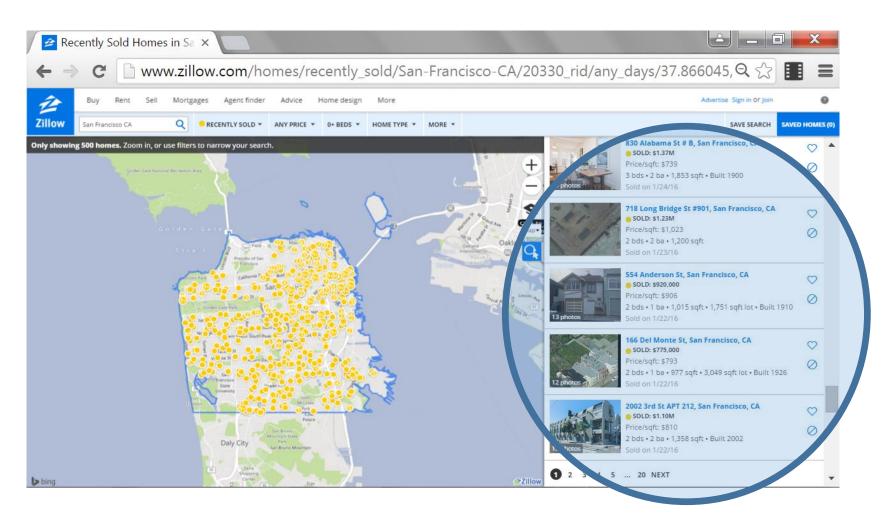
- Files
 - CSV (Comma-Separated Values)
 - TSV/TXT (Tab-Separated Values)
 - JSON (JavaScript Object Notation)
 - XML (eXtensible Markup Language)

Zillow dataset: a dataset we will use consistently use across this course



- Recently Sold Homes (Source: Zillow)
 - 1,000 homes sold in San Francisco between 11/10/2015and 2/12/2106

Raw data was scrapped from the Zillow website (20 pages, each listing 50 homes = 1,000 homes)



Raw data is Messy™...

... and needs to be parsed and tidied up (a.k.a., organized)

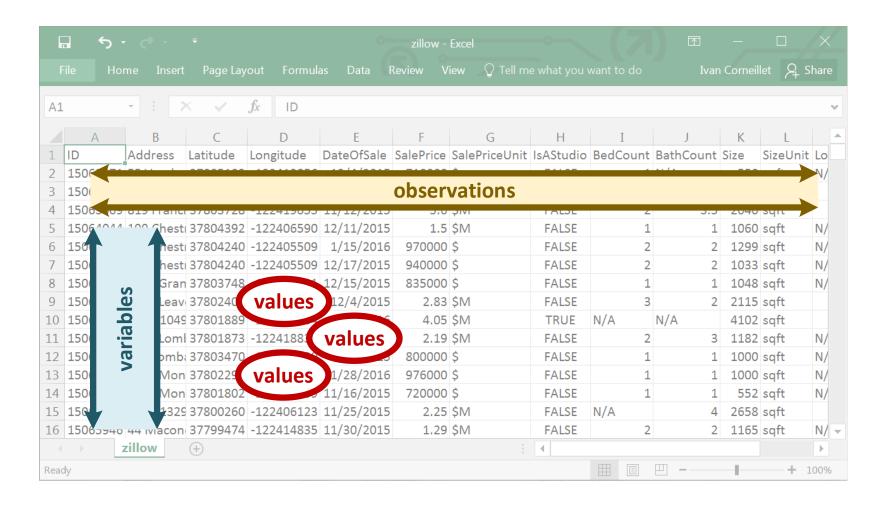
Tidying up your Data

- Tidying data is the most fruitful skill you can learn as a data scientist
 - It will save you hours of time and make your data much easier to visualize, manipulate, and model
- Many data science tools follow a set of conventions that makes one layout of tabular data much easier to work with than others. Your data will be easier to work with if you follow three rules:
 - Each observation is placed in its own row
 - Each variable in the dataset is placed in its own column
 - Each value is placed in its own cell

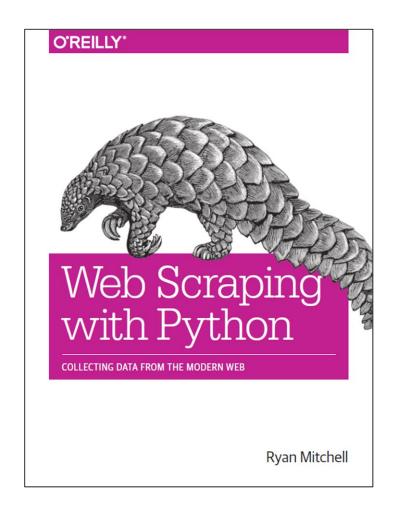
Really, data can be incredibly raw

- Trouble tickets inspect and maintain manholes in New Year City
 - "Service box," a common piece of infrastructure, had at least 38 variants, including SB, S, S/B, S.B, S?B, S.B., SBX, S/BX, SB/X, S/XB, /SBX, S.BX, S &BX, S?BX, S BX, S/B/X, S BOX, SVBX, SERV BX, SERV-BOX, SERV/BOX, and SERVICE BOX

Our tidy Zillow dataset



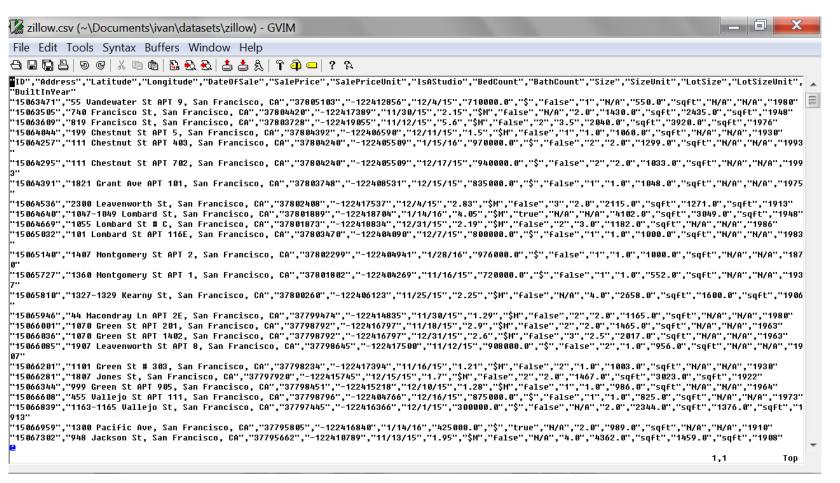
Going further



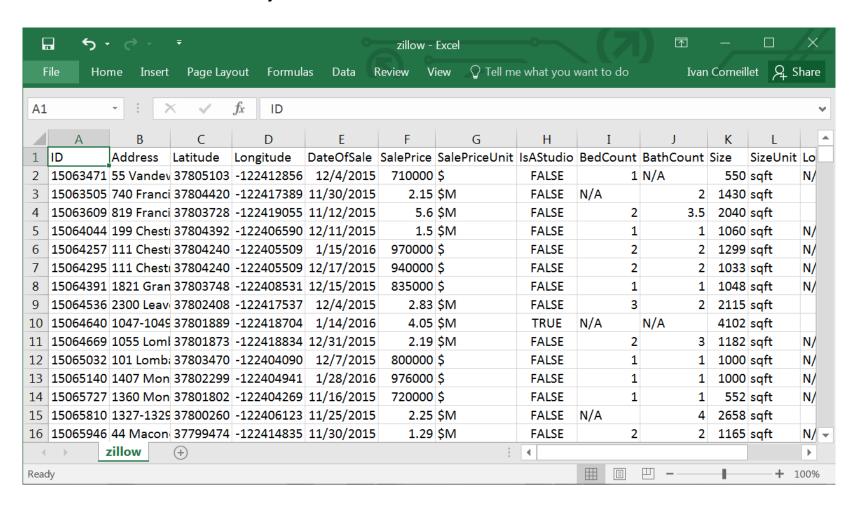


File Formats

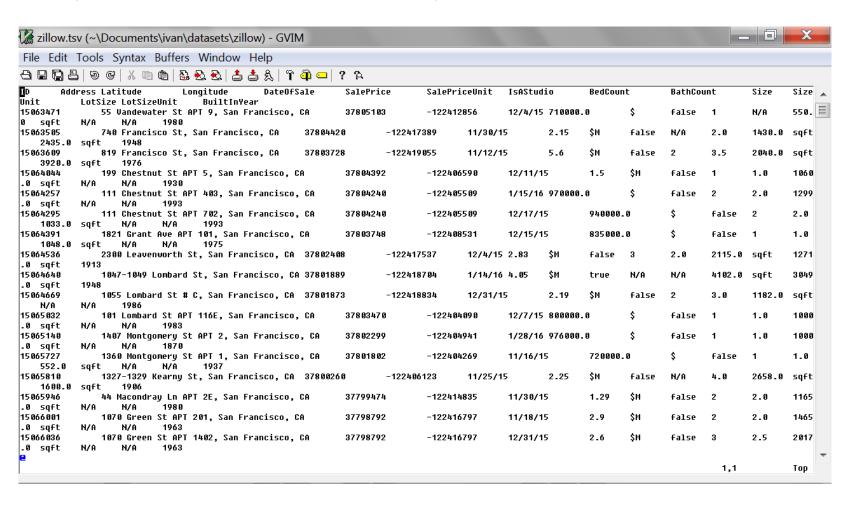
Our tidy Zillow dataset in CSV format: each observation is in one line; within each line, variables are separated with commas (and here delimited with double-quotes)



Excel reads CSV files natively (and our Python toolkit will too...)



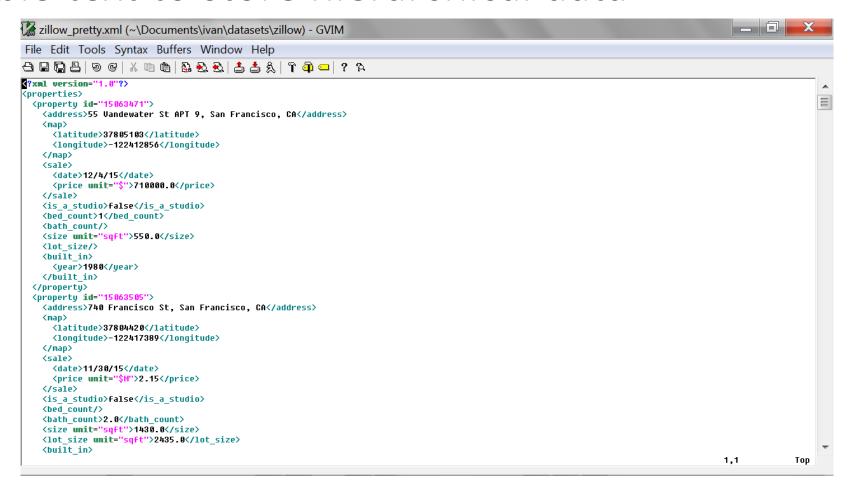
TSV is another simple text format for storing data in a tabular structure: each observation in the table is one line of the text file and each variable is separated from the next by a tab character



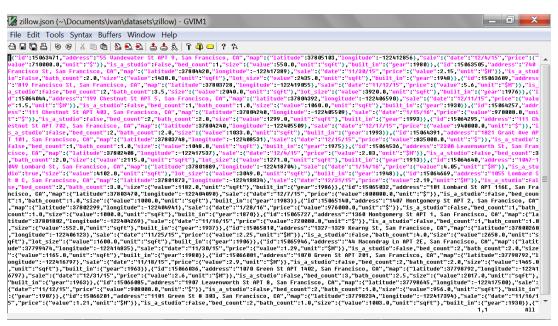
JSON, the most common open standard data format used for asynchronous browser/server communication, uses human-readable text to store data using key—value pairs and lists. Unlike CVS/TSV format, data can be represented hierarchically

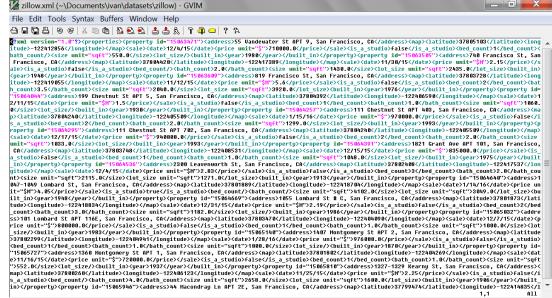
```
zillow pretty.json (~\Documents\ivan\datasets\zillow) - GVIM
File Edit Tools Syntax Buffers Window Help
"id": 15063471,
   "address": "55 Vandewater St APT 9, San Francisco, CA",
     "latitude": 37805103,
     "longitude": -122412856
   "sale": {
     "date": "12/4/15";
     "price": {
      "value": 710000.0,
      "unit": "$"
   "is a studio": false,
   "bed count": 1,
   "size": {
     "value": 550.0,
     "unit": "sqft'
   "built in": {
     "year": 1980
   "id": 15063505,
   "address": "740 Francisco St, San Francisco, CA",
   "map": {
    "latitude": 37804420,
     "longitude": -122417389
   },
   "sale": {
     "date": "11/30/15",
     "price": {
       "value": 2.15,
       "unit": "$M"
                                                                                                                        1,1
                                                                                                                                    Top
```

XML, another data format used for asynchronous browser/server communication, also uses human-readable text to store hierarchical data



JSON and XML are harder to read by humans when indentation is removed (usually the default) although it is still straightforward for machines...

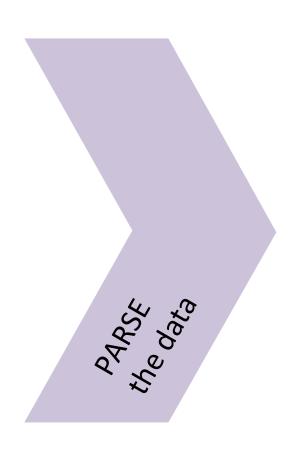






3 PARSE the Data

Parse the Data



Parse the Data

- Read any documentation provided with the data (session 2)
- Perform exploratory data analysis (session 3)
- Verify the quality of the data(sessions 2/3)

Parse the Data (cont.)

- Parse the Data
 - Read any documentation provided with the data
 - Perform exploratory data analysis
 - Verify the quality of the data

- Documentation and DataDictionary
- Codealongs
 - Introduction to pandas
 - Tidying up (more) the Zillow dataset
- Lab

Parse the Data (cont.)

- Parse the Data
 - Read any documentation provided with the data
 - Perform exploratory data analysis
 - Verify the quality of the data

- You need to understand what you're working with
- To better understand your data
 - Create or review the data dictionary
 - Perform exploratory surface analysis
 - Describe data structure and information being collected
 - Explore variables and data types

Documentation and Data Dictionary

- Data dictionaries
 - Help you judge the quality of the data
 - Also help understand how it's coded
 - Does "gender = 1" mean female or male?
 - Is the currency dollars or euros?
 - Help identify any requirements, assumptions, and constraints of the data
 - Make it easier to share data

Kaggle's Titanic Data Dictionary



VARIABLE DESCRIPTIONS: survival Survival

urvival Survival

(0 = No; 1 = Yes)

pclass Passenger Class

(1 = 1st; 2 = 2nd; 3 = 3rd)

name Name

sex Sex age Age

sibsp Number of Siblings/Spouses Aboard

parch Number of Parents/Children Aboard

ticket Ticket Number fare Passenger Fare

cabin Cabin

embarked Port of Embarkation

(C = Cherbourg; Q = Queenstown;

S = Southampton)

SPECIAL NOTES:

Pclass is a proxy for socio-economic status (SES)

1st ~ Upper; 2nd ~ Middle; 3rd ~ Lower

Age is in Years; Fractional if Age less than One (1) If the Age is Estimated, it is in the form xx.5

With respect to the family relation variables (i.e. sibsp and parch) some relations were ignored. The following are the definitions used for sibsp and parch.

Sibling: Brother, Sister, Stepbrother, or

Stepsister of Passenger Aboard Titanic

Spouse: Husband or Wife of Passenger Aboard

Titanic (Mistresses and Fiancés Ignored)

Parent: Mother or Father of Passenger Aboard

Titanic

Child: Son, Daughter, Stepson, or Stepdaughter of

Passenger Aboard Titanic

Other family relatives excluded from this study include cousins, nephews/nieces, aunts/uncles, and in-laws. Some children travelled only with a nanny, therefore parch=0 for them. As well, some travelled with very close friends or neighbors in a village, however, the definitions do not support such relations.



Codealong: Introduction to pandas

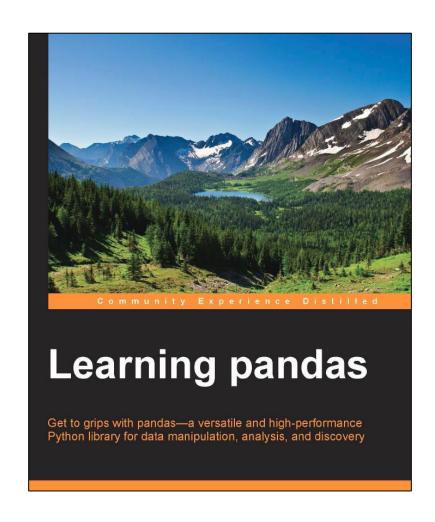
pandas is a Python library to manipulate and perform statistical and mathematical analysis on tabular and multidimensional datasets

- pandas provides the ability to index, retrieve, tidy, reshape, combine, slice, and perform various analyses on both single and multidimensional data
- It also includes loading and saving data from local and Internet-based resources
- We will use *pandas* to explore and manipulate the Zillow dataset

Resources

http://pandas.pydata.org/pandas- docs/stable/tutorials.html	Guide to many pandas tutorials, geared mainly for new users.
https://github.com/jvns/pandas-cookbook	Great resource with examples from weather, bikes, and 311 calls from Julia Evans
https://bitbucket.org/hrojas/learn-pandas	Great series of Pandas tutorials from Dave Rojas
https://github.com/ResearchComputing/Meetup-Fall-2013/tree/master/python	Awesome set of python notebooks from a meetup- based course exclusively devoted to pandas

Going further





Codealong: Tidying up the Zillow dataset



Activity: Write a Data Dictionary for the Zillow dataset

Activity: Write a Data Dictionary for the Zillow dataset



DIRECTIONS (15 minutes)

- 1. Divide into 4 groups, each located at a whiteboard
- 2. Data Dictionary: Modeling the previous example, each group should develop a data dictionary for the Zillow dataset (10 minutes)
 - a. Highlight how the variables are coded as well as any requirements, assumptions, and constraints of the data
- 3. Present: Present your dictionary to the class (5 minutes)
 - a. Choose one student to present for the group

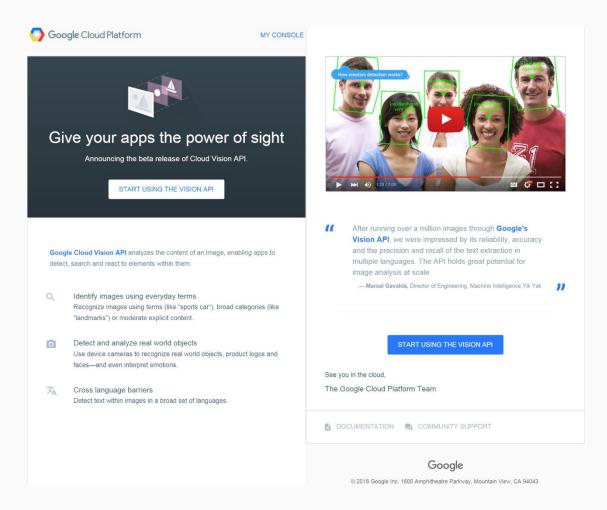
DELIVERABLE

Data Dictionary



Lab

Today's Closing Thought





Review

Review

You should now be able to:

- Define a problem and types of data
- Identify dataset types
- Define the data science workflow
- Apply the data science workflow in the pandas context
- Write an iPython notebook to import, format, and clean data using the Pandas library



Q & A



Before Next Class

Before Next Class

- Reminder on the in-flight projects
 - Unit Project 1 (due next time on 3/3)
 - Final Project 1 (due 3 weeks from now on 3/22)



Exit Ticket

Don't forget to fill out your exit ticket here