# QSS20: Modern Statistical Computing

Unit 03: Pandas for data manipulation

## Agenda

▶ Housekeeping: logistics and psets
▶ Recap of Python basics
▶ Intro to static visualization
▶ Data manipulation using pandas
  1. Intro to dataframes
  2. Aggregation
  3. Creating new columns/transforming their type
  4. Row and column filtering

# Meet Ellie, our fantastic peer tutor!

Eleanor (Ellie) Sullivan took QSS20 in Winter 2022 and was in a final project team with Ramsey, one of our TAs!

Please make use of peer tutoring sessions as a collaborative, supportive workspace for completing psets led by an experienced fellow student. Ellie will tell us more...

Peer tutoring schedule: Sun 8-9pm; Mon 7-8pm; Thurs 9-10pm

# Agenda

- **Housekeeping: logistics and psets**
- Recap of Python basics
- Intro to static visualization
- **Data manipulation using** `pandas`
  1. Intro to dataframes
  2. Aggregation
  3. Creating new columns/transforming their type
  4. Row and column filtering

# Class schedule and special guests

▶ Class moved from Monday, 01/16 at 3A time period to Tuesday, 01/17 at 3B time period
  ▶ This means class meets later than usual this coming Tuesday: 4:30-6:20 PM (OK to bring coffee/tea to that class!)

▶ Special guest on Tuesday, 01/17 (at 3:30 PM): Ashley Doolittle, Associate Director of Dartmouth's Center for Social Impact and head of the Social Impact Practicum (SIP)

▶ Special guests on Wednesday, 01/25 (at 5:00 PM): Andrea Caoili and Ann Klein of the National Center for START Services (NCSS), our project partner
  ▶ *Note:* "START" stands for "Systemic, Therapeutic, Assessment, Resources, and Treatment" and is a service model for treatment of individuals with intellectual and developmental disabilities (IDD). Here's a 1-hr documentary about their work.

# Logistics of problem set 1

▶ To submit, upload two files using the Problem Set 1 Assignment on Canvas

  ▶ Raw .ipynb file that contains your answers in response to questions; please put these answers in pset1_blank rather than starting a new .py or .ipynb file
  ▶ Compiled .html

▶ Save each with your netid—e.g.: pset1_f004bt8.ipynb and pset1_f004bt8.html

▶ Pset 1 Due Sunday 01.15 at 11:59 PM

▶ *Note:* The wording for 3.2.C and 3.3 assume two TAs, but the data in parts 1-2 has only one TA. Adjust the numbering to reflect this as I indicated on Piazza

# Logistics of problem sets in general

▶ I will release problem set 2 to Canvas and the course GitHub tomorrow; due date Sunday, 01.22 at 11:59 PM

▶ Pset 1 must be entirely your own work; we will put you in pairs for pset 2 (you will submit a single notebook between you)

▶ You will submit psets 3-5 using GitHub, which we will go over in two weeks

▶ Four late days available for use across psets; let teaching team know via Piazza if you're using a late day

## Where we are

- ▶ Housekeeping: logistics and psets
- ▶ **Recap of Python basics**
- ▶ Intro to static visualization
- ▶ Data manipulation using `pandas`
    1. Intro to dataframes
    2. Aggregation
    3. Creating new columns/transforming their type
    4. Row and column filtering

# Recap of Python basics

What do you remember from last class?

## Recap of Python basics

**Tips:**

▶ Python data types: string, int, float, boolean, list, datetime, objects
▶ Another basic data type for course: numpy array
  ▶ Like lists except only store a single data type and can be visualized as matrices/tables. Great for math (ML, NLP, etc)
▶ To transform each element of a list, use a list comprehension
▶ Can use if or if/else with list comprehension; these have slightly different syntax (see exs below)

**Useful code snippets:**

```
type(var) # get data type of var
int(var), float(var), str(var) # coerce type
len(somelist) # get length of a list
somearray.shape # get num rows/columns of arr
somelist[1:3] # get 2nd & 3rd elems of list
[el for el in somelist if len(el)>0] # filter list
[el if len(el)>0 else 0 for el in somelist] # filter/except
```

# Where we are

- ▶ Housekeeping: logistics and psets
- ▶ Recap of Python basics
- ▶ **Intro to static visualization**
- ▶ Data manipulation using `pandas`
    1. Intro to dataframes
    2. Aggregation
    3. Creating new columns/transforming their type
    4. Row and column filtering

# Walk through notebook with plotting example code

Link to static viz notebook using plotnine

Can use any plotting syntax for problem sets — popular ones are matplotlib (covered by DataCamp last chapter of introduction to pandas); seaborn; plotnine

## Where we are

▶ Housekeeping: logistics and psets

▶ Recap of Python basics

▶ Intro to static visualization

▶ **Data manipulation using** `pandas`
1. Intro to dataframes
2. Aggregation
3. Creating new columns/transforming their type
4. Row and column filtering

## Policy background for pset 2 sentencing data

▶ **Data**: deidentified felony sentencing data from Cook County State's Attorney's Office (SAO)

▶ Released as part of push towards transparency with election of a new prosecutor in 2016

Opinion
**EDITORIAL**

**Unequal Sentences for Blacks and Whites**

By The Editorial Board

Dec. 17, 2016

f ⊙ 𝕏 ✉ ↗ 🔖

> *Earlier this month Kim Foxx, the state's attorney for Cook County, Illinois, which covers Chicago, released six years' worth of raw data regarding felony prosecutions in her office. It was a simple yet profound act of good governance, and one that is all too rare among the nation's elected prosecutors. Foxx asserted that "for too long, the work of the criminal justice system has been largely a mystery. That lack of openness undermines the legitimacy of the criminal justice system." Source*

## Concepts in problem set 2

> Creating new columns
> Aggregating using groupby and agg
> Row filtering
> Pandas operations like quantile, pd.to_datetime()
> value_counts() and sort_values()
> Loops and functions
> Visualizing results

**Now let's practice these concepts using crime reports from DC!**

# Where we are

▶ Housekeeping: logistics and psets

▶ Recap of Python basics

▶ Intro to static visualization

▶ **Data manipulation using** pandas
   1. **Intro to dataframes**
   2. Aggregation
   3. Creating new columns/transforming their type
   4. Row and column filtering

# Intro to data wrangling: how do dataframes differ from lists and arrays?

- ▶ In last lecture, we covered two structures for storing information in python:
    - ▶ Lists: structure built into python; 1-dimensional storage of information that can deal with information of different types in the same list
    - ▶ Arrays: requires the numpy package; n-dimensional (can be $> 2$) storage of information - usually use to store numeric information for efficient math calculations/model estimation
- ▶ DataFrames: 2-dimensional with rows (first dimension) and columns (second dimension) — sometimes called tabular data structure
    - ▶ pandas package (usually aliased as pd)
    - ▶ Each column can contain a different type of information
    - ▶ Each row references some unit of analysis (person; nation; city; 911 call; etc)

# Two ways of interacting with dataframes

1. Creating our own
   - ▶ Less common
   - ▶ **Main use**: when we're creating data from some non-tabular source (e.g., a text string of a short politician bio; extract their name, party, and religious affiliation)
2. Reading in data stored in different formats
   - ▶ Focus for now: csv
   - ▶ Others we'll get to: excel; json; pkl or other serialized format; txt; spatial data stored as shapefiles

# Creating our own dataframe: dictionary approach

**Keys:** column names; **values:** lists or arrays containing information

```python
## create own df
### approach 1: dictionary where keys
## are names of columns
### and items are lists with information
name_list = ['Rebecca', 'Yifan', 'Sonali']
role_list = ['Instructor', 'TA', 'TA']
ht_list = [63, 70, 63.5]
my_df = pd.DataFrame({'names': name_list,
                      'role': role_list,
                      'fictional_height': ht_list})
my_df
my_df.dtypes
```

|   | names | role | fictional_height |
|---|-------|------|------------------|
| **0** | Rebecca | Instructor | 63.0 |
| **1** | Yifan | TA | 70.0 |
| **2** | Sonali | TA | 63.5 |

```
names                object
role                 object
fictional_height    float64
dtype: object
```

# Creating our own dataframe: nested lists

```python
## approach 2: list of lists
### each person's information is one list
rj_info = ['Rebecca', 'Instructor', 63]
yl_info = ['Yifan', 'TA', 70]
ss_info = ['Sonali', 'TA', 63.5]

### together they're a nested list
nested_info = [rj_info, yl_info, ss_info]
nested_info

### we can then make that into a dataframe
### need to specify column names
my_df_2 = pd.DataFrame(nested_info,
                       columns = ['names',
                                  'role',
                                  'fictional_height'])
my_df_2
```

```
[['Rebecca', 'Instructor', 63], ['Yifan', 'TA', 7
0], ['Sonali', 'TA', 63.5]]
```

|   | names | role | fictional_height |
|---|-------|------|------------------|
| 0 | Rebecca | Instructor | 63.0 |
| 1 | Yifan | TA | 70.0 |
| 2 | Sonali | TA | 63.5 |

# More common way of interacting with dataframes: reading in data

- ▶ os package is important for finding the path of the file
  - ▶ os.getcwd() tells you the working directory you're in
- ▶ Two ways to structure path names (will return to these when we cover command line + GitHub in a couple weeks)
- ▶ **Way one (avoid if possible) absolute paths:**
  '/Users/rebeccajohnson/Dropbox/ppol564_prepwork/prep_activities/f22_materials'
- ▶ **Better way: relative paths to .py or .ipynb:** my data is stored two levels up from where my notebook is; can provide abbreviated pathname:
  '../../data/example_data.csv'
- ▶ Structure of read command pd.read_csv('path to file')

Once we have the data in Python, can summarize using built-in attributes or functions

- ▶ `dfname.dtypes`: returns a pandas series where the index is the name of the column; the value is the type of data it contains (eg str; int; float)
- ▶ `dfname.shape`: returns a length-2 tuple with dimensions of dataframe (number of rows, number of columns)
- ▶ `dfname.head()`: prints first n rows (defaults to 5)
- ▶ `dfname.tail()`: prints last n rows (defaults to 5)

## Manipulations of data

**Examples:** finding mean height across the TAs; recoding heights into different categories; subsetting to a dataframe only containing TAs

| names | role | fictional_height |
|-------|------|------------------|
| Jaren | Instructor | 68.0 |
| Ramsey | TA | 70.0 |
| Eunice | TA | 63.5 |
| Ellie | Peer Tutor | 64.5 |

## Where we are

► Housekeeping: logistics and psets

► Recap of Python basics

► Intro to static visualization

► **Data manipulation using** `pandas`
  1. Intro to dataframes
  2. **Aggregation**
  3. Creating new columns/transforming their type
  4. Row and column filtering

# Basic aggregation syntax: one grouping variable, summarizing one column)

```
1  grouping_result = df.groupby('grouping_varname')
2                    ['varname_imsummarizing'].agg(
3                    {functiontosummarize
4                    }).reset_index()
```

▶ **Why might we use reset_index()?** This helps us treat the output as a DataFrame with clear, one-level columns

More aggregation syntax: one grouping variable, custom function to summarize one column

```
1 grouping_result = df.groupby('grouping_varname').agg(
2                    {'varname_imsummarizing': functiontosummarize
3                    }).reset_index(drop = True)
```

- ▶ **Why is there a dictionary inside of agg?** Lets us name the specific columns to summarize by and what functions to use; keys in the dictionary are the variables to summarize by, values are what function to use
- ▶ **Why might we use reset_index(drop = True)?** If we don't want to keep the old index (e.g., to avoid cluttering our DataFrame)

More aggregation syntax: one grouping variable, custom function to summarize multiple columns

```
1  grouping_result = df.groupby('grouping_varname').agg(
2                    {'varname_imsummarizing': functiontosummarize,
3                     'othervarname_imsummarizing': functiontosummarize
4                    }).reset_index()
```

- ▶ **When might this be useful?** Can summarize different columns in different ways

# More aggregation syntax: two grouping variables

```
1  grouping_result = df.groupby(['grouping_varname1',
2                                'grouping_varname2']).agg(
3                                {'varname_imsummarizing': 'functiontosummarize'
4                                }).reset_index()
```

- ▶ **When might this be useful?** things like "how does this vary by time and category x?"

# How do we structure the function inside the aggregation?

Three common ways of calling the function:

1. Functions that operate on pandas series, e.g.:

   ```
   df.groupby('month').agg({'offense': ['nunique', 'first']})
   ```

2. Functions from numpy (aliased here as np), e.g.:

   ```
   df.groupby('month').agg({'offense': [np.mean, np.median]})
   ```

3. "Lambda" functions we write ourself that take an argument

   ```
   df.groupby('month').agg({'offense':
                           lambda x: len(x.unique())})
   ```

# Summarizing over multiple rows or columns (without aggregation)

```
1  mean_threecols =        df [[ 'colA ', 'colB ',
2                          'colC ']]. apply ( np . mean ,
3                          axis = 0)
```

▶ Pandas apply function: https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.apply.html

▶ axis argument tells us whether to apply the function over columns (axis 0) or rows (axis 1)

▶ Can see in activity what happens

# Before we code, let's group!

Count off to 8...

Then find your number to form a group!

Pause for practice

Aggregation section (section 1) of `01_pandas_blank.ipynb`

## Where we are

- ▶ Housekeeping: logistics and psets
- ▶ Recap of Python basics
- ▶ Intro to static visualization
- ▶ **Data manipulation using** pandas
    1. Intro to dataframes
    2. Aggregation
    3. **Creating new columns/transforming their type**
    4. Row and column filtering

## First type of column creation: binary indicators

Two general approaches that are "vectorized," or they work across all rows automatically without you needing to do a for loop:

1. `np.where`: similar to `ifelse` in R; useful if there's only 1-2 True/False conditions; can be used in conjunction with things like `df.varname.str.contains(``some pattern'')` if the column is string/char type

2. `np.select`: similar to `case_when` in R; useful for when there's either (1) several True/False conditions or (2) you're coding one set of categories into a different set of categories (this may come up in the problem sets)

# Different types of `np.where`

```
1
2 ## indicator for after 2020 christmas or not (make sure to
3 ## format date in same way)
4 df['is_after_christmas'] = np.where(
5                             df.nameofdatecol > "2020-12-25",
6                             True, False)
7
8 ## indicator for whether month is in spring quarter (april, may,
      june)
9 df['is_spring_q'] = np.where(
10                      df.monthname.isin(["April", "May", "June"]),
11                      True, False)
12
13 ## indicator for whether someone's name contains johnson
14 df['is_johnson'] = np.where(
15                      df.fullname.str.contains("Johnson"),
16                      True, False)
17
18 ## strip string of all instances of johnson
19 df['no_johnson'] = df.fullname.str.replace("Johnson", "")
```

# Then, if we created binary indicator, can use for subsetting rows

```
1
2  ## subset to after christmas
3  df_afterchristmas = df[df.is_after_christmas].copy()
4
5  ## subset to after christmas AND spring quarter
6  ## note parentheses around each
7  df_postc_spring = df[(df.is_after_christmas) &
8                       (df.is_spring_q)].copy()
9
10 ## subset to after christmas BUT NOT spring quarter
11 ## note tilde ~ for negation
12 df_postc_notspring = df[(df.is_after_christmas) &
13                         (~df.is_spring_q)].copy()
```

## `np.where` is useful for single conditions, but what about multiple conditions?

▶ **Example**: code to fall q if September, October, November, or December; code to winter q if January, February, or March; code to spring q if April, May, or June; code to summer q if otherwise

▶ Gets pretty ugly if nested `np.where`

```
1
2 ## quarter ind
3 df["quarter_type"] = np.where(df.monthname.isin(["Sept",
4                       "Oct", "Nov", "Dec"]), "fall_q",
5                       np.where(df.monthname.isin(["Jan",
6                       "Feb", "March"]), "winter_q",
7                       np.where(df.monthname.isin(["April",
8                       "May", "June"]), "spring_q", "summer_q")))
```

# Better for recoding with multiple categories: `np.select`

```
 1
 2  ## step one: create a list of conditions/categories
 3  ## i can omit last category if i want or specify it
 4  quarter_criteria = [df.monthname.isin(["Sept", "Oct",
 5                       "Nov", "Dec"]),
 6                      df.monthname.isin(["Jan", "Feb", "March"]),
 7                      df.monthname.isin(["April", "May", "June"])]
 8
 9  ## step two: create a list of what to code each category to
10  quarter_codeto = ["fall_q", "winter_q", "spring_q"]
11
12  ## step three: apply and add as a col
13  ## note i can use default to set to the residual category
14  ## and here that's a fixed value; could also retain
15  ## original value in data by setting default to:
16  ## df["monthname"] in this case
17  df["quarter_type"] = np.select(quarter_criteria,
18                                 quarter_codeto,
19                                 default = "summer_q")
```

Pause for practice

Recoding section (section 2) of `01_pandas_blank.ipynb`

## Where we are

- ▶ Housekeeping: logistics and psets
- ▶ Recap of Python basics
- ▶ Intro to static visualization
- ▶ **Data manipulation using** pandas
    1. Intro to dataframes
    2. Aggregation
    3. Creating new columns/transforming their type
    4. **Row and column filtering**

## Row filtering: combining multiple conditions

Say we have a column in our data that contains Dartmouth courses (e.g., QSS17, QSS20, ECON20, GOV10, COSC1) and we want to select all rows where the course code is "30" and the prefix is "QSS".

```
1  qss_30_courses = df[(df.coursename.startswith("QSS")) &
2                      (df.coursename.str.contains("30")].copy()
```

Two notes

▶ Using pandas built in methods (startswith and str accessor)- what would happen with latter if the variable was not a string?

▶ Use parentheses around each when combining multiple conditions (weird for R users)

# A useful tool for row filtering before we learn regex: pandas built-in functions

▶ `contains` and `startswith` are functions built into `pandas` that help us work with string/character variables

▶ General syntax:
nameofdata.nameofstringcol.str.nameoffunction(argument if relevent)

▶ Examples:
  ▶ df.stringvar.str.upper() and lower()
  ▶ df.stringvar.str.replace()
  ▶ df.stringvar.str.split() — defaults to splitting on spaces; can feed it other delimiters like ;

# Column filtering: combining with list comprehension

Say we have a huge DataFrame listing all Dartmouth courses ever, and we want to find all the QSS courses plus COSC1. We could filter to "COSC1" plus any columns that have "QSS" in the name.

```
1  QSS_any2 = df[["COSC1"] +
2              [col for col in
3              df.columns
4              if "QSS" in col]].copy()
```

Notes:

▶ Use .copy() to tell python that we're assigning a copy of the original dataframe (df) to the new object QSS_any2; otherwise, gives us SettingwithCopy warning because unsure whether we want our changes to propagate back to original DF (we almost never want this)

▶ Put "COSC1" in brackets to combine two lists

Pause for practice

Filtering section (section 3) of `01_pandas_blank.ipynb`

# We covered

▶ Intro to static visualization
▶ Data manipulation using `pandas`
  1. Intro to dataframes
  2. Aggregation
  3. Creating new columns/transforming their type
  4. Row and column filtering