# QSS20: Modern Statistical Computing

## Unit 14: Recap & work session

## Goals for today

- ▶ Recap of web-scraping
- ▶ Recap of SQL
- ▶ Upcoming deadlines
- ▶ Final project work session

## Goals for today

- **Recap of web-scraping**
- Recap of SQL
- Upcoming deadlines
- Final project work session

## Recap of web-scraping

**Tips**:

- ▶ Grabbing information from web pages is *web-scraping*. Systematically finding what pages to scrape is *web-crawling*, which usually uses *link extraction* to find & follow links in certain domains

- ▶ Most people do *narrow crawling*—grabbing specific information from a few pages with similar structure. To do *broad crawling* over a range of websites, use scrapy for flexible, non-blocking (i.e., fast) operations

- ▶ To maximize *extensibility* or resilience to new sites/HTML layouts, use an HTML exclusion list to scrape visible text (rather than relying on specific HTML/CSS; BS code, spider code)

**Useful general commands:**

```
scrapy startproject schools # create scrapy project
scrapy shell 'http://quotes.toscrape.com' # shell to test scrapy
scrapy genspider -t crawl broad site.com # create broad spider
    from CrawlSpider template
scrapy crawl broad -o output.json # run broad spider, save JSON
```

# Code for broad spider in `broad.py`

```python
import scrapy
from scrapy.linkextractors import LinkExtractor
from scrapy.spiders import CrawlSpider, Rule
from schools.items import SchoolsItem

class BroadSpider(CrawlSpider):
    name = 'broad'
    allowed_domains = ['quotes.toscrape.com']
    start_urls = ['http://quotes.toscrape.com/']

    rules = (Rule(
        LinkExtractor(), callback='parse_item', follow=True))

    def parse_item(self, response):
      for quote in response.css('div.quote'):
        item = SchoolsItem() # initialize

        item['text'] = quote.css('span.text::text').get(),
        item['author'] = quote.css('small.author::text').get(),
        item['tags'] = quote.css('div.tags a.tag::text').getall(),

        yield item
```

## Where we are

- Recap of web-scraping
- **Recap of SQL**
- Upcoming deadlines
- Final project work session

What do you remember?

## Recap of SQL: explanation

**When to use**:

▶ You walk into an internship/job and they say, "Here's how to authenticate to our database, it's in SQL. You can use that, right?" (very common)

▶ You have a huge dataset on which to run complex queries → can put in SQL and work from your laptop (requires little memory)

▶ You want to access a secured, external dataset without downloading → can put in SQL and work from your laptop (like an API)

**Tips:**

▶ SQL like condensed version of Pandas; similar but diff. query syntax

▶ Usually easy to install Python MySQL connector (database setup less easy)

▶ Use creds .yaml file to access connector—or copy into code (less secure)

▶ Subqueries in () are like functions as in df.apply(func)

▶ Use count(*) with group by to aggregate each slice (e.g., by race x crime)

# Recap of SQL: simple examples

**Simple ex. w/ basic syntax:**

```
1 select CASE_ID, AGE_AT_INCIDENT    —select cols
2 from caseinit    —tablename
3 —table ops like joins or subqueries would go here
4 where AGE_AT_INCIDENT > 40    —condition(s)
```

**Ex. creating new col w/ case/when:**

```
1 select *,    —take all cols
2 CASE
3     WHEN OFFENSE_CATEGORY = UPDATED_OFFENSE_CATEGORY
4     THEN 'Same offense'
5     ELSE 'Diff offense'
6 END as charge_update
7 from caseinit
```

## Recap of SQL: subquery to create col, filter w/ inner join

```sql
1  select *,
2  from caseinit
3  inner join
4      (select CASE_ID as cid, CASE_PARTICIPANT_ID as cpid,
5      CASE
6          WHEN OFFENSE_CATEGORY = UPDATED_OFFENSE_CATEGORY
7          THEN 'Same offense' ELSE 'Diff offense'
8      END as charge_update
9      from caseinit) as tmp
10      on tmp.cid = caseinit.case_ID and
11      tmp.cpid = caseinit.CASE_PARTICIPANT_ID
12 where charge_update = "Diff offense"
```

**Equivalent in pandas:**

```python
1 tmp = caseinit.copy()
2 tmp['charge_update'] = np.where(
3      tmp.OFFENSE_CATEGORY == tmp.UPDATED_OFFENSE_CATEGORY,
4      'Same offense', 'Diff offense')
5 merged = pd.merge(left=caseinit, right=tmp, how='inner',
6      left_on=('CASE_ID', 'CASE_PARTICIPANT_ID'),
7      right_on=('cid','cpid'))
8 merged = merged[merged.charge_update == 'Diff offense']
```

## Recap of SQL: group by offense, add top 5 w/ inner join

```
1  select *
2  from caseinit
3  inner join (
4      select UPDATED_OFFENSE_CATEGORY as tmp_oc,
5      count(*) as count_offense
6      from caseinit
7      where RACE in ("Black", "White")
8      group by UPDATED_OFFENSE_CATEGORY
9      order by count_offense desc
10     limit 5) as top5
11     on caseinit.UPDATED_OFFENSE_CATEGORY = top5.tmp_oc
```

**Equivalent in pandas:**

```
1  top5 = caseinit[caseinit.RACE.isin(['Black', 'White'])].
2      groupby('UPDATED_OFFENSE_CATEGORY').
3      agg({'CASE_ID': 'nunique'}).reset_index().
4      rename(columns = {'CASE_ID': 'count_offense',
5          'UPDATED_OFFENSE_CATEGORY': 'tmp_oc'}).
6      sort_values(by='count_offense', ascending=False)[:5]
7  merged = pd.merge(left=caseinit, right=top5, how='inner',
8      left_on='UPDATED_OFFENSE_CATEGORY',
9      right_on='tmp_oc')
```

## Where we are

- Recap of web-scraping
- Recap of SQL
- **Upcoming deadlines**
- Final project work session

## Upcoming deadlines

- ▶ Problem set five
  - ▶ **Due next Friday, 11-18**
  - ▶ Submit in usual way (assign Prof. & TA a GitHub issue)
  - ▶ Reminder: Your lowest pset grade will be dropped
- ▶ Final project presentation
  - ▶ **Delivered in class this coming Monday, 11-14**
  - ▶ Please also submit slides before that class: PDF $\rightarrow$ Canvas, share w/ Prof. on Overleaf
- ▶ Final paper
  - ▶ **Due Tuesday, 11-22** (last day of finals)
  - ▶ How to submit: PDF $\rightarrow$ Canvas, share w/ Prof. on Overleaf

## Goals for today

▶ Recap of web-scraping
▶ Recap of SQL
▶ Upcoming deadlines
▶ **Final project work session**

## Final project work session

**Things to work on:**

- ▶ Slides for final presentation
- ▶ Analysis, codebase development, documentation
- ▶ Final paper

**Places to work:**

- ▶ This classroom
- ▶ Outside in hall
- ▶ Rockefeller Center atrium (downstairs)—**please mind your volume**

**Feel free to consult/ask questions:**

- ▶ Prof. Haber
- ▶ Other groups/research teams
- ▶ The whole class (consult Prof. first)