# 2nd half overview

CSE 314A
DCDS 510

# Open Data Sets

- Quality and quantity of open data exploding
- Mature datasets are often provided with interactive dashboards and spreadsheet exports.
- API (Application Programming Interface)
  - Usually gets you more fine grained and complete data
- Web scraping: last resort but sometimes necessary

Real-world data is sourced from many sources & people

# Common Open Data Sources

- US Census
- World Bank Open Data
- Social Media: Twitter and Yelp
- Movies: OMDb,
- Sports: Sports-Reference
- FiveThirtyEight
- Yahoo! Finance API
- St Louis Regional Data Exchange
- FEC (data), FEC (API), MEC
- Science: Nature, PLoS One
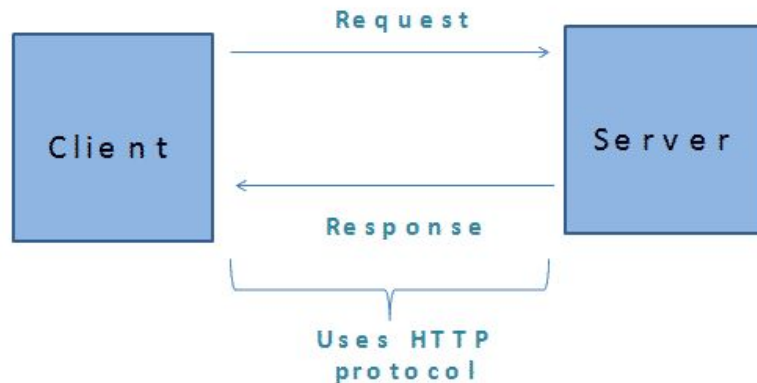- ML/AI: Kaggle, UCI ML Repository

# Other closed/personal data sources

- Business data:
  - Accounting, marketing, finance, inventory
- Organizational Data:
  - Contacts, personal information
- Science/Social Science Data:
  - Subjects, results, conditions
- Public Data
  - Freedom of Information Act (FOIA)

- Your own experiment/survey

# Anatomy of an API

- **Application Programming Interface** - Terminology
- A reminder of how the Internet works →
- Stereotypically:
  - You tell the API what you want in a URL
  - The server responds with some JSON

# API Clients in Python

- [Requests](#) package
- Provide common functionality that helps abstract away some web development knowledge (response handling, parsing, pagination, etc)
  - Prepackaged clients, e.g.: [Google API Python Client](#), [Sportsipy](#)
  - Build your own client

# Web Scraping/Crawling - when all else fails?

- Packages
  - BeautifulSoup
  - Scrapy
  - Selenium
  - Some API clients
  - Pandas

# Web Scraping Ethics and Etiquette

- Be aware of Terms of Service and proprietary information
- Be aware of your network traffic and server load
- Follow the robots.txt
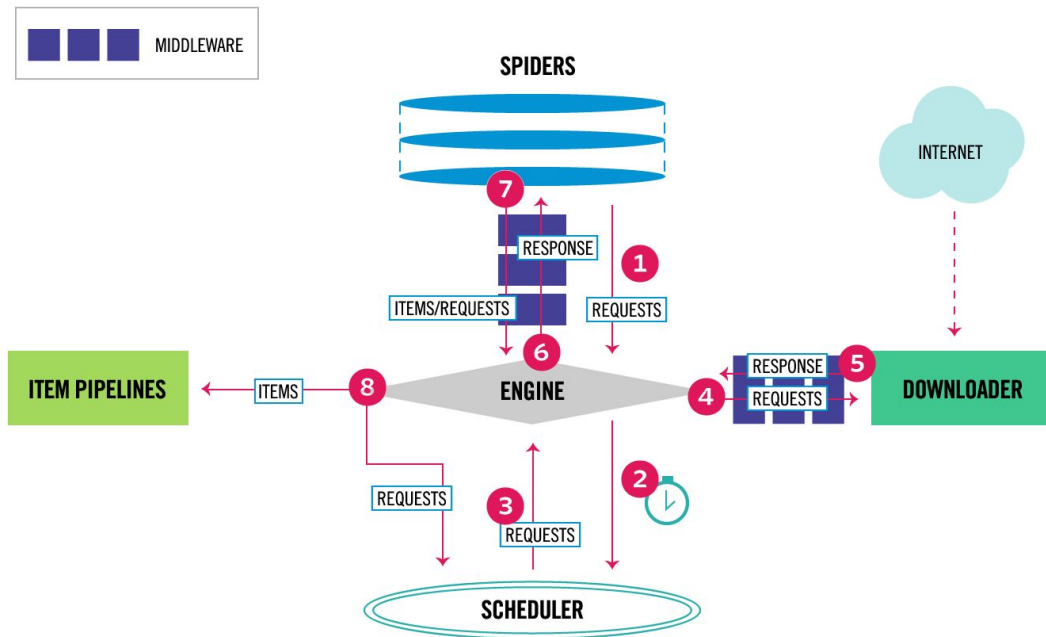
# Beautiful Soup - Parsing your HTML

# Scrapy -

[https://docs.scrapy.org/en/latest/](https://docs.scrapy.org/en/latest/)

Scrapy is an application framework for crawling web sites and extracting structured data which can be used for a wide range of useful applications, like data mining, information processing or historical archival.

# Scrapy -

[https://docs.scrapy.org/en/latest/](https://docs.scrapy.org/en/latest/)

Even though Scrapy was originally designed for web scraping, it can also be used to extract data using APIs or as a general purpose web crawler.

MIDDLEWARE

SPIDERS

INTERNET

RESPONSE

ITEMS/REQUESTS

REQUESTS

ITEM PIPELINES

ITEMS

ENGINE

RESPONSE
REQUESTS

DOWNLOADER

REQUESTS

REQUESTS

SCHEDULER

# Selenium

- Designed for automated UI/UX testing
- Most useful for Javascript-focused web browsing
  - Hidden/Dynamic Link navigation
  - Manual clicking processes
- Actually uses browser engines
- Requires some OS-level setup/drivers

# Pandas

- Extracts tables directly from HTML
- Easy, convenient
- Still may require some manual pruning
- Won't help for data outside of tables

# More on storing/caching data

- No go-to answer
- Pandas provides flexible options
- Common approach - store a response and work from that response
  - VCRpy records "cassettes" - designed for testing
- Keep memory limitations in mind, cache as you go
- Focus on what helps you in your development process
  - Don't repeat yourself
  - Choose good example cases

# Historical context - ggplot2 (R Tidyverse)

# Python Visualization in 2022

 +  Dash



 (HoloViz)

- Libraries can be used by Jupyter notebooks and Pandas API

# Stack Overflow Trends

## **Viz-first data analysis**

- When presenting findings
    - Keeps you focused on outcomes
    - Evaluate hypotheses
    - Keep your supervisors happy
- When exploring data
    - Narrows down possibilities
    - Fosters ongoing thinking about data structure and domain particularities
    - Familiarizes you with the data if it is external
- Simply more fun and less discouraging
- Answers a lot of philosophical questions up front
    - Resource: [Fundamentals of Data Visualization](#) by Claus O Wilke

# Feeding your plot - Quality in, Quality Out

- Modern plotting libraries can do much of the data processing for common visualizations
- Organized data allows for single-command plotting
  - More on this next week - for now use preorganized data in hw/demos
  - Fake your data if your data is unorganized or unavailable!

# Basic plotting (Plotly)

- Raw Figure objects - Just fancy dictionaries!
  - Use when customizing, examining structure, or testing
- Plotly Express
  - Works best (but not only)with "tidy" data – but it essentially means that you have one row per observation.
  - "Wide" data support added in last year or two.
- Use templates
- Minimize tinkering
- Base level interactivity

# Dashboard basics (Dash)

- Works with Flask under the hood (easily deployable – relatively)
- Uses HTML to arrange components
  - Dash Bootstrap Components can make your app look pretty and provide some supplemental functionality
- Uses "callbacks" to update the dashboard on user actions
  - Just functions – reads attributes from HTML components and updates them.
- Callbacks are chainable – you can use use them to update other input components as well as your figure data

# Exporting your images and deploying your dashboards

- Images
  - PNG button on dashboard
  - I/O methods on Figure objects
    - Always use vector formats if possible – .svg, .pdf, etc
  - Post-processing
    - Inkscape - free and open source
    - Adobe Illustrator
- Deploying Dashboards – tricky but doable
  - [The Magical Guide](The Magical Guide)

# The future: 3D dataviz, VR/AR, and data storytelling

https://www.tiktok.com/@the.data.guy/video/7043589231066860846

https://www.tiktok.com/@the.data.guy/video/7028635623833488646

# SQL from 10,000 feet -

```
1   SELECT
2   DISTINCT
3       SUM(cont.Amount) AS Total
4   FROM mec.Contributions AS cont
5       LEFT JOIN mec.Committees AS cmt
6   ON cont.MECID = cmt.MECID
7   WHERE cmt.Name = "Cori Bush for Congress"
8   GROUP BY cont.City
9   HAVING Total > 1000
10  ORDER BY Total
```

SQL order of execution - SELECT

1. FROM clause                    6. HAVING clause

2. ON clause                      7. SELECT clause

3. OUTER clause                   8. DISTINCT clause

4. WHERE clause                   9. ORDER BY clause

5. GROUP BY clause                10. TOP clause

# CRUD and ACID

- CRUD operations
  - Create
  - Read
  - Update
  - Delete
- ACID transactions
  - Atomicity - All or nothing
  - Consistency - Data is correct
  - Isolation - Independent of other transactions
  - Durability - changes are persistent

**RDBMS systems**

(SQL Standard)

Differences:
- User Permissions/Security
- Programmatic features
- Some architecture/indexing decisions
- Secret Sauce

# NoSQL - Think "Nonrelational"

- Key-Value Databases (common for caching [Redis])
- Document Databases (MongoDB) - think JSON/XML
- Graph Databases

Key Considerations:

- "Schema-less"
- Horizontal Scalability
- Performance gains at cost of ACID compliance