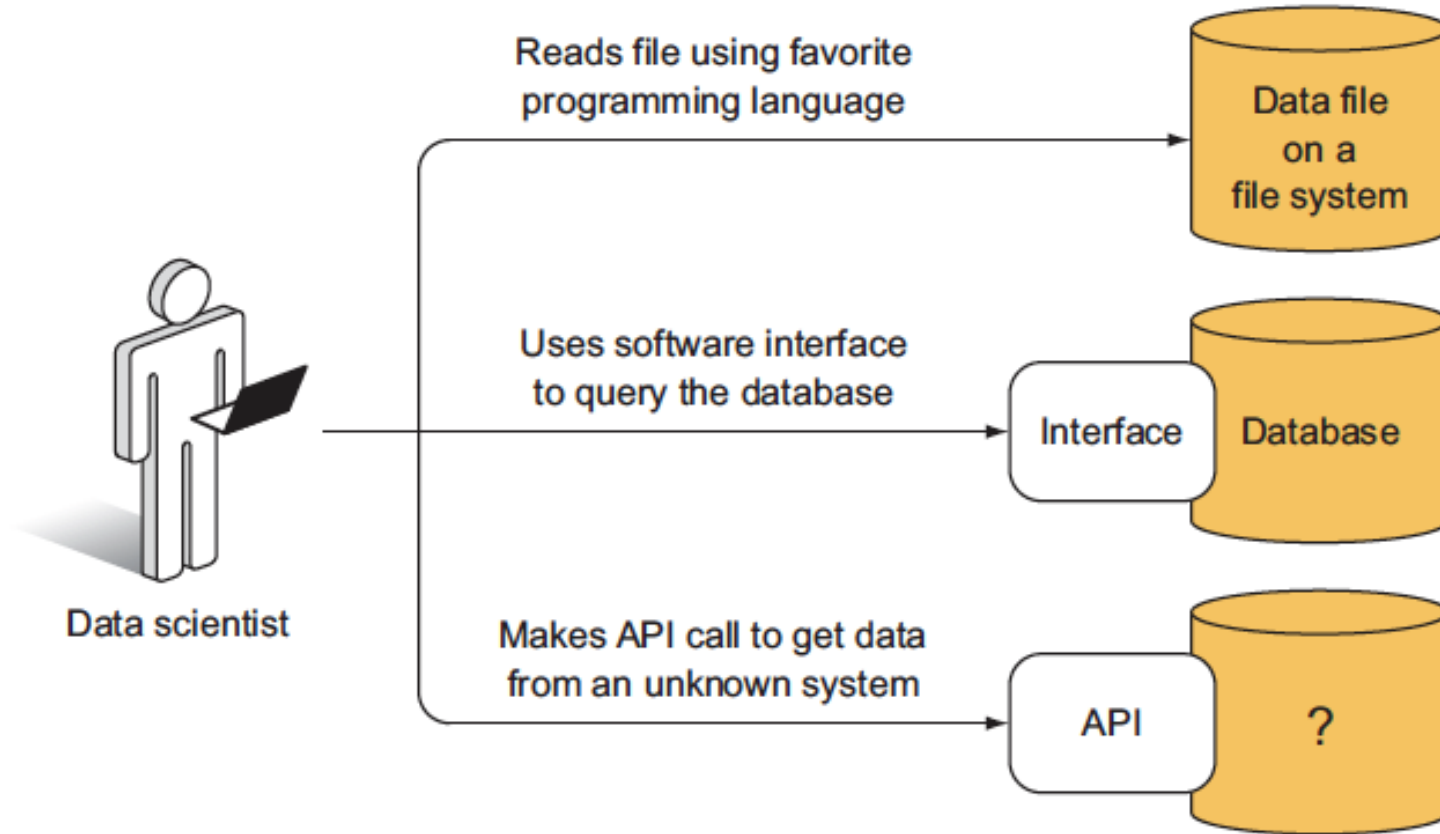# Data Search

## CDS-492 | Dr. Slamani

(original slides by Dr. Ron Mahabir)
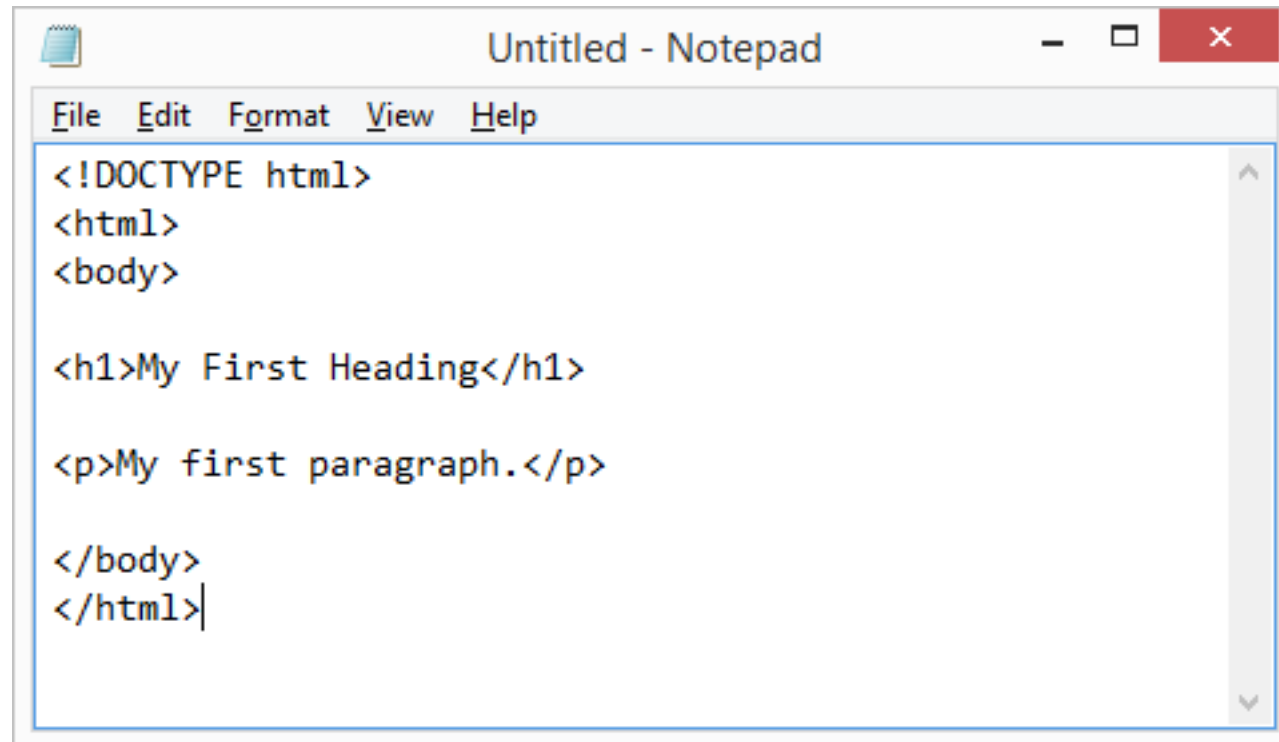
# Common ways to access data

# Data formats - Flat files

- Just text (e.g. csv or tsv)

- Minimal

- Word documents are not flat files
  - They contain additional overhead such as style etc.

- Many popular programming languages have packages to read these

- Can be slow to read if very large
  - Database may be an option here

- Some files may be compressed

# Data formats - HTML

- Plain text marked up with tags or specially denoted instructions for how the text should be interpreted.

# Data formats - XML

- Extension of HTML
- More suitable for storing and transmitting documents and data other than web pages

# Data formats - JSON

- JavaScript Object Notation

- Typically describes something more like a data structure
  - such as a list, map, or dictionary in many popular programming languages
  - Leaner in number of characters compared to XML

# Relational databases

- Storage systems used to optimize retrieval, writing and management of data
  - Indexing
- Queries
- Relationship between tables
- Some analytical capabilities

# Non-relational databases

- Not only SQL (NoSQL)
- Don't conform to the norm
- Greater flexibility
- Examples include:
  - MongoDB
  - Elasticsearch – document oriented database
    - Full-text search engine with an HTTP web interface and schema-free JSON documents.
  - Graph database
    - Graph structures for semantic queries with nodes, edges, and properties to represent and store data instead of tables, or documents.

# API – Application Programming Interface

- Set of rules for communicating with a piece of software to access data

- The gatekeeper

- Most programming language have packages to assist

# Find a work around

- MS Word
- MS Excel
- PDF
- .mbox
- netcdf – Maybe a software (network Common Data Form - file format for storing multidimensional scientific data (variables) such as temperature, humidity, pressure)
- Then again
  - For a quick peek
  - For very small files

# Finding the DATA

- Google
  - General vs Specific search terms
- Repositories
  - Github
  - Kaggle
- Web scraping
  - Maybe against the terms of service for many websites
    - REMEMBER IF YOU DO WEBSCRAPE – Humans are erratic at some things! ☺
- Download an entire website
  - YIKES!!!
- Measure or collect it yourself

# Copyright and licensing

- Data may have licensing, copyright, or other restrictions that can make it illegal to use the data for certain purposes.
  - University – research
  - Twitter – don't distribute
- Without confirming that your use case is legal, you remain at risk of losing access to the data or, even worse, a lawsuit
- Other ethical considerations
  - Hand, D.J., 2018. Aspects of data ethics in a changing world: Where are we now?. *Big data*, *6*(3), pp.176-190.

# Got the data, is it ENOUGH?

- On the surface things can be all bright and shiny

- Below the surface, no happy faces

- What to do
  - Dive in
  - Work on a sample
  - Combine another dataset?

# Links to data

- https://www.dataquest.io/blog/free-datasets-for-projects/
- https://flowingdata.com/2009/10/01/30-resources-to-find-the-data-you-need/
- https://dsc.gmu.edu
- Uber
  - https://data.cityofnewyork.us/Transportation/uber-Data/3jeu-mn7j
  - https://help.uber.com/riders/article/download-your-data?nodeId=2c86900d-8408-4bac-b92a-956d793acd11
  - https://data.world/datasets/uber
  - https://www.kaggle.com/fivethirtyeight/uber-pickups-in-new-york-city
  - https://www.reddit.com/r/uber/comments/74c2mn/export_uber_trips_to_csv/

# Links to data

- Medicare
  - https://data.medicare.gov
  - https://www.medicare.gov/download/downloaddb.asp
  - https://www.cms.gov/newsroom/data
  - https://www.kaggle.com/cms/cms-medicare
  - https://catalog.data.gov/dataset?q=medicare
  - https://www.ama-assn.org/practice-management/medicare/medicare-claims-data-release
  - https://www.cms.gov/OpenPayments/Explore-the-Data/Dataset-Downloads.html
  - https://data.medicare.gov/download-data
  - https://data.medicaid.gov
  - https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5433516/
  - https://fcw.com/articles/2018/03/06/cms-blue-button-api.aspx

# Links to data

- https://data.humdata.org/group/bhs
- https://catalog.data.gov/dataset?tags=bahamas
- https://github.com/GBPA/gis-datafiles
- https://www.digitalglobe.com/ecosystem/open-data
- https://www.nga.mil/MediaRoom/PressReleases/Pages/dorian.aspx
- Where else?
  - FEMA
  - NOAA
  - USAID
  - USGS EarthExplorer

# Auxiliary

- UCI repositories: archive.ics.uci.edu/ml/datasets.php

- reddit.com/r/datasets

- Health
  - ibi.gmu.edu/faculty-directory/niloofar-ramezani