# Module 1 - Data Acquisition & Programming
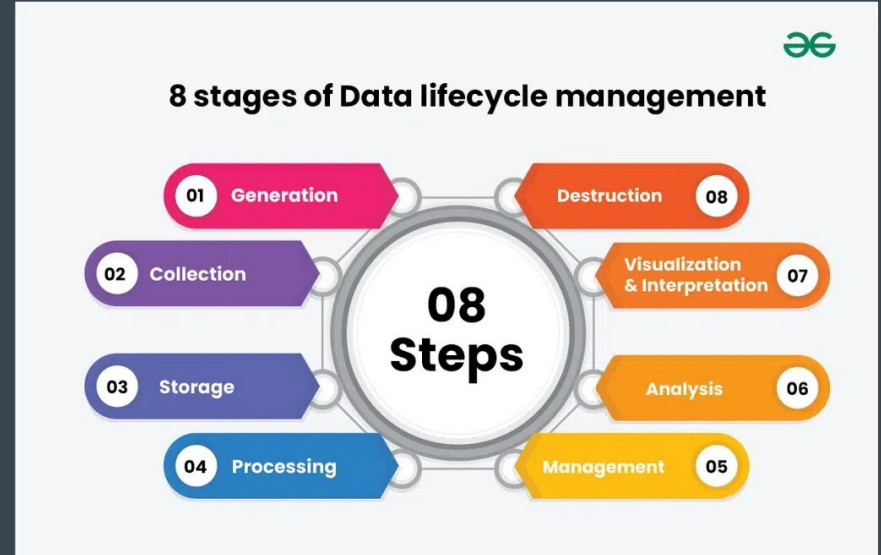
● ● ●

CSc 46000 - Professor Sheng

Alex Abraham, Amir Nabiyev, Azim Rahat, Md Mamun, Wen Jie Long

# What is Data Acquisition?

- Data acquisition is the process of gathering data from various sources for analysis

- Data can come from sensors, databases, the web, APIs, or real-time streams

- Critical for accurate insights and modeling



8 stages of Data lifecycle management

08 Steps

- 01 Generation
- 02 Collection
- 03 Storage
- 04 Processing
- 05 Management
- 06 Analysis
- 07 Visualization & Interpretation
- 08 Destruction

# Data Science Programming Languages

- Languages used in data science are designed for efficiently handling large datasets, machine learning, and data visualization

- Most commonly used languages:

    - Python - Most widely known language

    - R - Designed for statistics and data visualization

    - MATLAB - Built for mathematical modeling and engineering

# Python

- Python is an open-source, high-level language known for its readability and versatility

- Used in every step of the data science workflow

- Why it's Popular

  - Extensive Libraries: NumPy, Pandas, Matplotlib, Scikit-learn, TensorFlow, BeautifulSoup

  - Easy Integration: Works easily with web frameworks, databases, cloud services

  - Versatility: Used in analytics, AI research, automation, and engineering



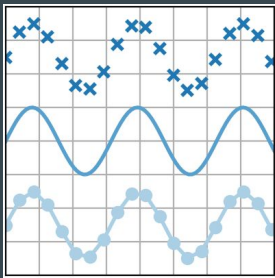Pull data from API    Clean data    Visualize data    Train model

# R

- R is an open source language built for statistical computing and graphs

- Useful in data analysis, hypothesis testing, and creating visualizations

- Why it's Popular

    - Easy Integration: Works with RStudio, able to interact with Python or SQL

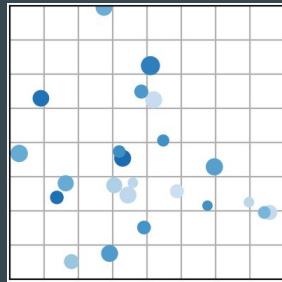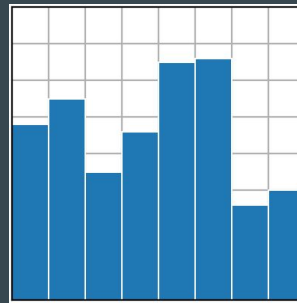    - Visualization Libraries: ggplot2 and plotly

# MATLAB

- A high-level language built for numerical computation, visualization, and algorithm development

- Why it's Popular

    - Matrix and math focus: Optimized for signal processing, control systems, and modeling

    - Built-in Toolkits: Specialized toolkits for image processing, machine learning, and statistics

    - Easy Integration: Easily connects with hardware sensors and Python for data collection
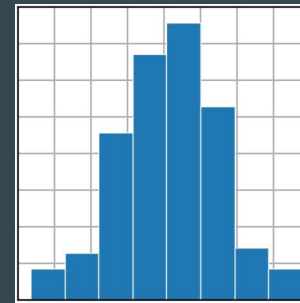
Plot

Scatter

Bar

Histogram
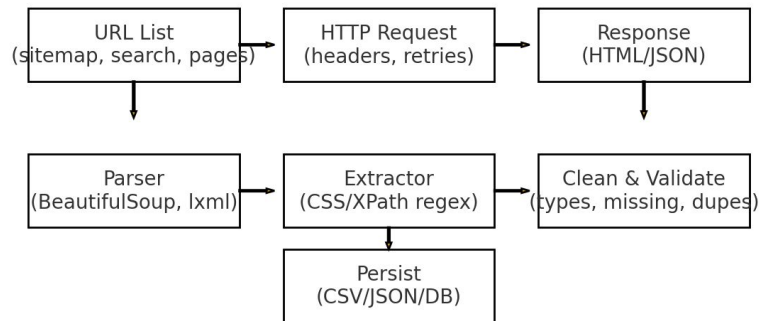
# What is Data Extraction in Data Science?

- Data extraction refers to to retrieving information from different sources including

    - Websites, databases, APIs, sensors, etc.

- Main techniques include

    - Web Scraping

    - APIs (Application Programming Interfaces)

    - cURL

    - WebSockets

# Web Scraping

- The process of extracting information from webpages by parsing their HTML structure

- How it works

  - Send a request to the website's page

  - Retrieve the HTML source code

  - Parse and extract specific elements

  - Store the data in a set format (CSV, JSON, or database)

### Web Scraping: Data Pipeline

| URL List (sitemap, search, pages) | → | HTTP Request (headers, retries) | → | Response (HTML/JSON) |
| --- | --- | --- | --- | --- |

| Parser (BeautifulSoup, lxml) | → | Extractor (CSS/XPath regex) | → | Clean & Validate (types, missing, dupes) |
| --- | --- | --- | --- | --- |

Persist (CSV/JSON/DB)

# API - Application Programming Interface (Part 1)



What is an API?

Database → Web Server → API → application, with Internet connected to API and application

# API - Application Programming Interface (Part 2)

- An official method provided by a website to retrieve data from them
- Different sites have different APIs
- Can be free or paid
- Can have a rate limit
- Can require an API key as an identifier

# API - Application Programming Interface (Part 3)

Websites with APIs will tell you how to use them.

## NYC OpenData

| Download file | **API endpoint** |
| --- | --- |

**Data format**
JSON

○ All data (41.3M rows)

Version ● SODA3 ○ SODA2

ⓘ Note: The SODA3 API requires authentication. For more details see this article

**API endpoint**
https://data.cityofnewyork.us/api/v3/views/erm2-nw ...

**API documentation** ☐    **Developer portal** ☐

[ Cancel ]    [ **Copy to clipboard** ]

## Columns (41)

| Column Name | Description | API Field Name | Data Type |
| --- | --- | --- | --- |
| Tᴛ Unique Key | Unique identifier of a Service Request (SR) in the open data set | unique_key | Text |
| 🗓 Created Date | Date SR was created | created_date | Floating Timestamp |
| 🗓 Closed Date | Date SR was closed by responding agency | closed_date | Floating Timestamp |
| Tᴛ Agency | Acronym of responding City Government Agency | agency | Text |
| Tᴛ Agency Name | Full Agency name of responding City Government Agency | agency_name | Text |
| Tᴛ Complaint Type | This is the first level of a hierarchy identifying the topic of the incident or condition. Complaint Type may have a corresponding Descriptor (below) or may stand alone. | complaint_type | Text |

# cURL (Client URL)

- Command-line utility for transferring data to and from a server
  - HTTP, HTTPS, FTP, SCP, and SFTP
- Usages:
  - Testing REST APIs (GET, POST, PUT, or DELETE)
  - Downloading files
  - Verify a website is running
  - Retrieve HTML content

# cURL Syntax

- curl [options] [URL]
- curl.exe [options] [URL]

PS C:\Users\alexa> curl -X GET https://api.sampleapis.com/coffee/hot
Invoke-WebRequest : A parameter cannot be found that matches parameter name 'X'.
At line:1 char:6
+ curl -X GET https://api.sampleapis.com/coffee/hot
+      ~~
    + CategoryInfo          : InvalidArgument: (:) [Invoke-WebRequest], ParameterBindingException
    + FullyQualifiedErrorId : NamedParameterNotFound,Microsoft.PowerShell.Commands.InvokeWebRequestC
   ommand

PS C:\Users\alexa> curl.exe -X GET https://api.sampleapis.com/coffee/hot
[{"title":"Svart Te","description":"Svart te föddes i Kina. Det är tillverkat av blad från en växt som
 kallas Camellia och kan smaksättas olika med frukter till exempel. En trevlig, varm, smakfull och aro
matisk dryck som passar till vardagen.","ingredients":["Te"],"image":"https://images.unsplash.com/phot
o-1576092768241-dec231879fc3?auto=format&fit=crop&q=60&w=800&ixlib=rb-4.0.3&ixid=M3wxMjA3fDB8MHxzZWFyY
2h8MjB8fHRlYXxlbnwwfHwwfHx8MA%3D%3D","id":13},{"title":"Apelsinjuice","description":"Vi har inget att
säga om vår nypressade apelsinjuice. Du måste prova den själv.","ingredients":["Färska Apelsiner","Is"
],"image":"https://images.unsplash.com/photo-1600271886742-f049cd451bba?auto=format&fit=crop&q=60&w=80
0&ixlib=rb-4.0.3&ixid=M3wxMjA3fDB8MHxzZWFyY2h8NzF8fG9yYW5nZSUyMGp1aWNlfGVufDB8fDB8fHww","id":18},{"tit
le":"Lemonad","description":"Var känd i Paris först och blev sedan mycket populär i hela Europa. Denna
 söta, färglösa, kolsyrade dryck görs genom att blanda citronsaft och kolsyrat vatten.","ingredients":
["Citronsaft","Kolsyrat vatten","Honung"],"image":"https://images.unsplash.com/photo-1621263764928-df1
444c5e859?auto=format&fit=crop&q=60&w=800&ixlib=rb-4.0.3&ixid=M3wxMjA3fDB8MHxzZWFyY2h8Nnx8bGVtb25hZGV8
ZW58MHx8MHx8fDA%3D","id":20},{"title":"title","description":"desc","ingredients":"[t,e,s,t]","image":"
gerald","id":"313513"}]
PS C:\Users\alexa>

curl -X GET https://api.sampleapis.com/coffee/hot

curl https://google.com

# Websockets

- Communication protocol
  - HTTP is short-lives and simple
  - WebSocket real-time data exchange without repeated requests.

- Usage
  - Chat applications (WhatsApp, Messenger)
  - Collaborative apps (Google Docs)
  - Live dashboards or stock tickers (real-time updates)
  - Multiplayer games



- Connection Start – HTTP request with Upgrade
- Handshake – Server gives 101 Switching Protocols
- Persistent TCP Link – Switches from HTTP to WebSocket protocol, staying open.
- Close Handshake – One side sends a close frame; the other replies, then the TCP link ends.

# Q & A