

# Data Collection Methods



**SDAIA**

الهيئة السعودية للبيانات  
والذكاء الاصطناعي  
Saudi Data & AI Authority



# Agenda



Intro to Data Collection Methods



Primary versus Secondary Approaches



Direct Data Collection Approach



Survey & Experiments



In-Direct Data Collection Approach



Web Scraping



API Access



# ► Introduction to Data Collection Methods

**Data Collection** is a foundational step in the field of data science and decision-making.

It is gathering information from various sources to create datasets that accurately represent the phenomenon or subject of interest.

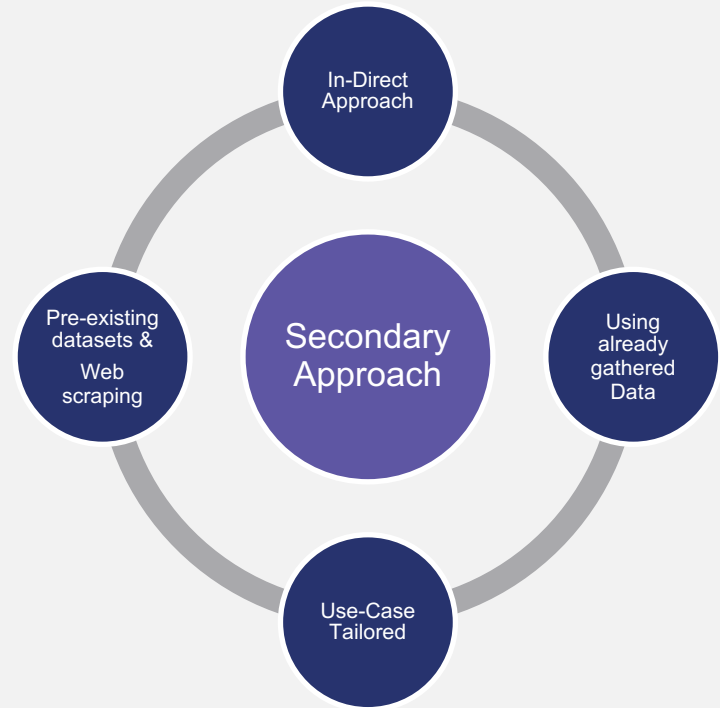
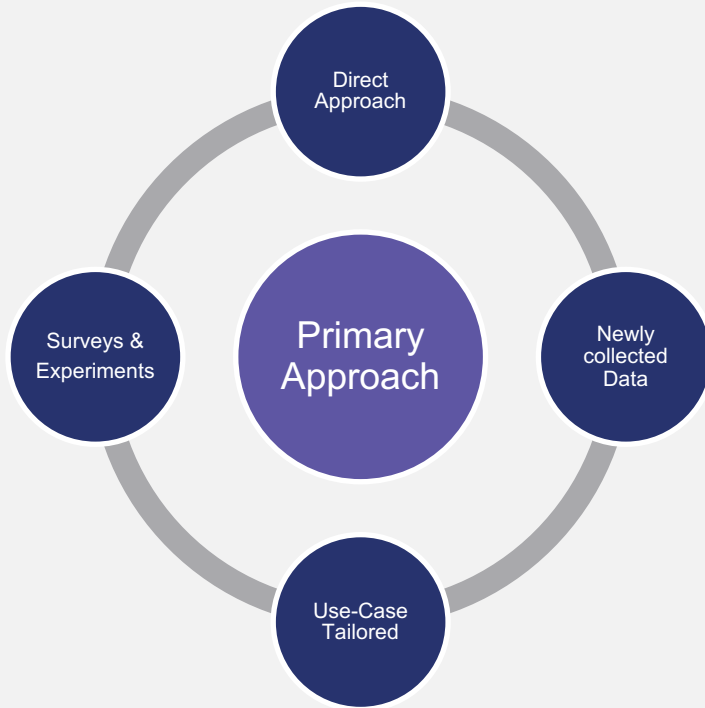
This process can be performed using *surveys*, *experiments*, *web scraping*, and *accessing public datasets or APIs*.

## **Data Collection Considerations:**

The goal of data collection in data science is to amass data that is **relevant**, **accurate**, and of **high quality**.



# Primary versus Secondary Approaches



# Direct Data Collection Approach

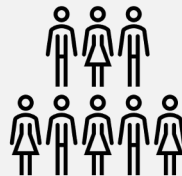
## Primary Approach

**The Direct Data Collection Approach** refers to a method of gathering information straightforwardly from the source for the first time.

Surveys and experiments are two fundamental methods for data collection in various fields, including *social sciences*, *marketing*, *health research*, and *many areas of data science*.



**Surveys**



**Experiments**





# Survey

## Direct Data Collection Approach

Surveys gather self-reported data through formats like online questionnaires and interviews. They're useful for collecting opinions and behaviors.

### Their advantages include:

**Scalability:** Efficiently reach many respondents.

**Versatility:** Collect diverse data, from demographics to opinions.

**Comparability:** Standardized questions facilitate cross-group analysis.

Limitations include biases from question phrasing, respondent interpretation, and accuracy of responses





# Experiments

## Direct Data Collection Approach

Experiments manipulate variables to study effects and infer causality, often under controlled settings like labs, though they can also be in the field or online.

### Key features include:

**Control:** Ability to manage conditions and isolate variables.

**Randomization:** Random assignment to groups to limit bias and support causal conclusions.

**Repeatability:** Can be replicated to confirm findings.

Valuable for exploring cause-and-effect

Limitations include high costs, time demands, and practical or ethical limitations.



# ▶ In Direct Data Collection Approach

## Secondary Approach

**Indirect(Secondary) Approach:** Involves using data that has already been collected by someone else for a different purpose and leveraging existing resources to gather information that can be applied to the current research.

Web scraping and API Calling are two fundamental secondary methods for data collection:



**API**



**Web Scraping**







# Web Scraping for Data Collection

## In-Direct Approach

**Web scraping** is the process of extracting data from websites, automating the collection of information available online. It serves as a powerful tool in data collection, enabling analysts and scientists to gather vast amounts of data quickly, which is essential for analysis, research, and decision-making processes.

**Key tools and technologies for web scraping include:**



**Beautiful Soup:** A Python library for parsing HTML and XML documents. It's widely used for simple projects and tasks that require quick data extraction from websites.



**Selenium:** Originally a tool for testing web applications, Selenium can automate web browser interaction, making it suitable for scraping dynamic content that requires interaction with the webpage.



**Scrapy:** An open-source and collaborative framework for extracting the data you need from websites. It's designed for web scraping. Scrapy is highly efficient, scalable, and versatile, making it suitable for large-scale web scraping projects.





# Beautiful Soup for Data Collection

## In-Direct Approach

### Typical Steps to handle a website in Beautiful Soup

Fetching the web page content using requests.

Parsing the content with Beautiful Soup to create a parse tree.

Using Beautiful Soup's searching and navigation methods to find relevant data.

Extracting and processing the data you need from the elements found.

Iteratively refining your approach based on the specific requirements of your web scraping project and the structure of the web pages you're working with.



# ▶ Selenium for Data Collection

## ▶ In-Direct Approach

### What is Selenium?

An open-source automation tool primarily used for automating testing web applications.

Allows for browser automation, enabling tasks to be performed as if a real user is navigating the site so it can also render websites Dynamically.

### Why Use Selenium for Web Scraping?

**Dynamic Content:** Selenium can interact with webpages that load content dynamically, making it ideal for scraping modern sites.

- **Real Browser Interaction:** Performs operations in a real browser environment, allowing for actions like clicking buttons, filling forms, and scrolling.



# ▶ API Access for Data Collection

## ▶ In-Direct Approach

**APIs (Application Programming Interfaces) are software** are tools that allow different software applications to communicate with each other. They acts as intermediaries allowing different software applications to communicate, simplifying the process of data collection by providing structured ways to request and receive data.

### Advantages of Using APIs:

**Efficiency:** Streamlines data access and functionality.

**Real-Time Data:** Offers access to live data, crucial for up-to-date application needs.

**Scalability:** Eases handling of growing data or demand with minimal infrastructure adjustments.

**Cost-Effectiveness:** More affordable than developing custom data collection systems.



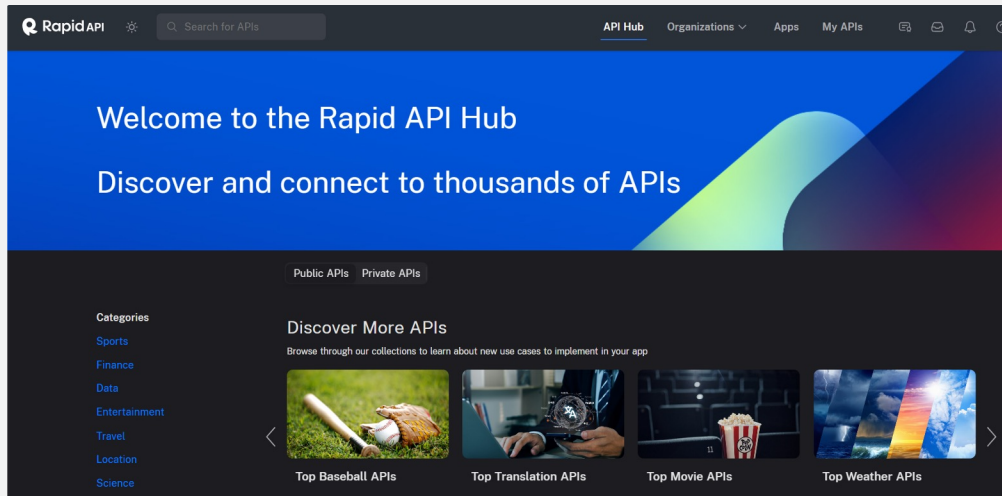
# RapidAPI

## In-Direct Approach

**RapidAPI** [\[link\]](#) is a comprehensive platform that aggregates thousands of APIs across various domains

It presents a unified platform for developers to discover, connect, and manage APIs through a single, standardized interface.

It offers access to diverse data sources across various categories, including finance, sports, entertainment, weather, and more.



## RapidAPI Considerations

### In-Direct Approach

**API Limits:** Be aware of rate limits and quotas to avoid service interruptions.

**Costs:** Understand the pricing model of the API and usage charges.

**Security:** Keep your API key confidential to prevent unauthorized usage.

**Performance:** Test response times and reliability.

**Documentation:** Read the API documentation thoroughly.

**Updates:** Stay informed about any changes or updates to the API.

**Support:** Check the support options and community forums for help for Q&A.



# Thank you



**SDAIA**  
الهيئة السعودية للبيانات  
والذكاء الاصطناعي  
Saudi Data & AI Authority