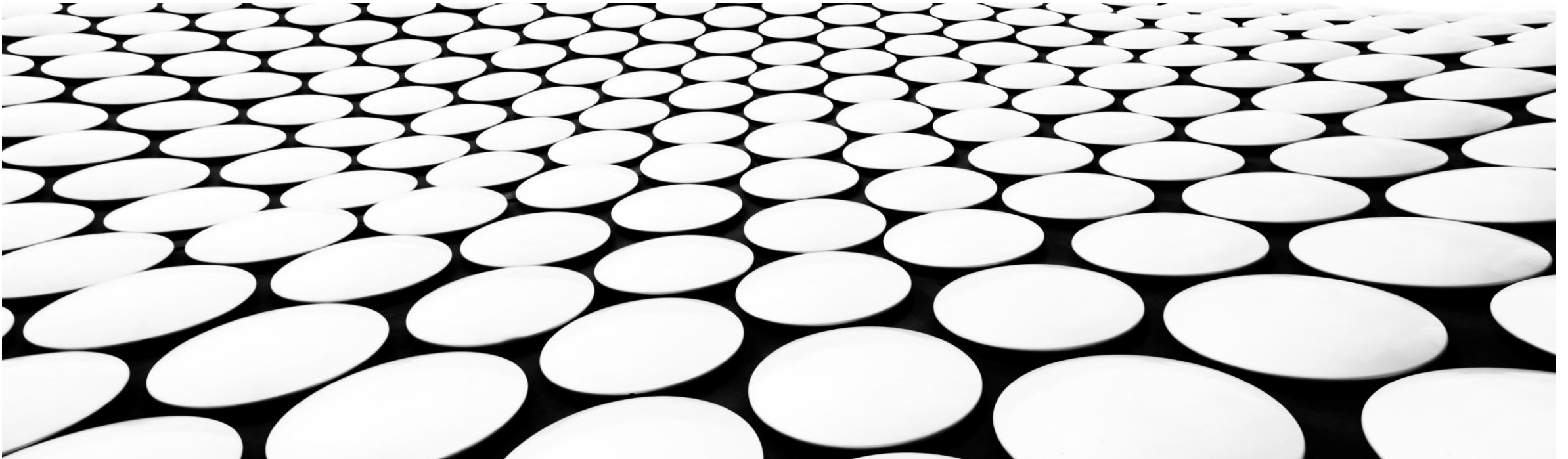
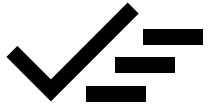


DATA COLLECTION

Exploration of diverse approaches





DATA COLLECTION

- Data Science Process
- Types of questions (insights)
- Overview of data collection methods

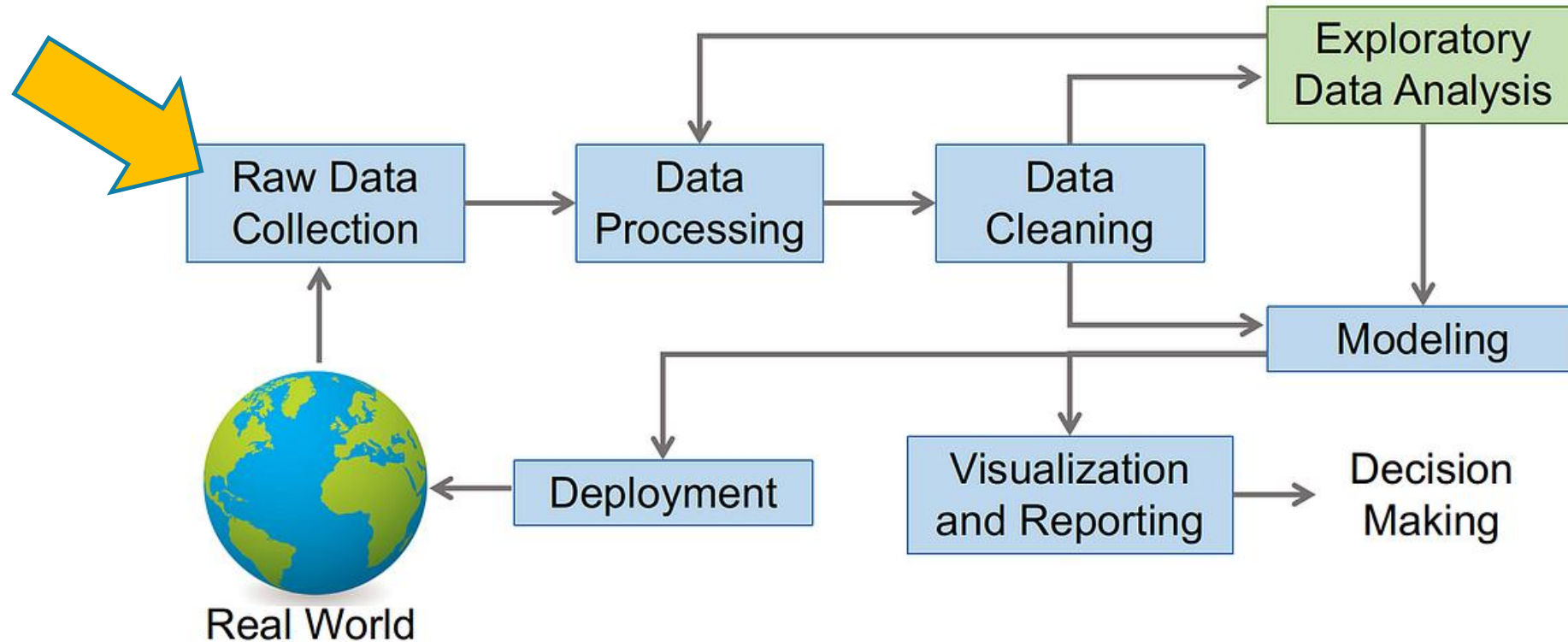
Data Science Process

Data Science Process: A Comprehensive Guide



Abhijit · Follow
6 min read · Jan 15, 2024

Data Science Process



7 Steps to a Successful Data Science Project

Beginners Guide on Completing a Data Science Project from Scratch

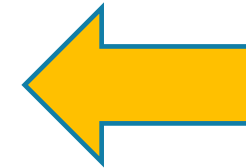


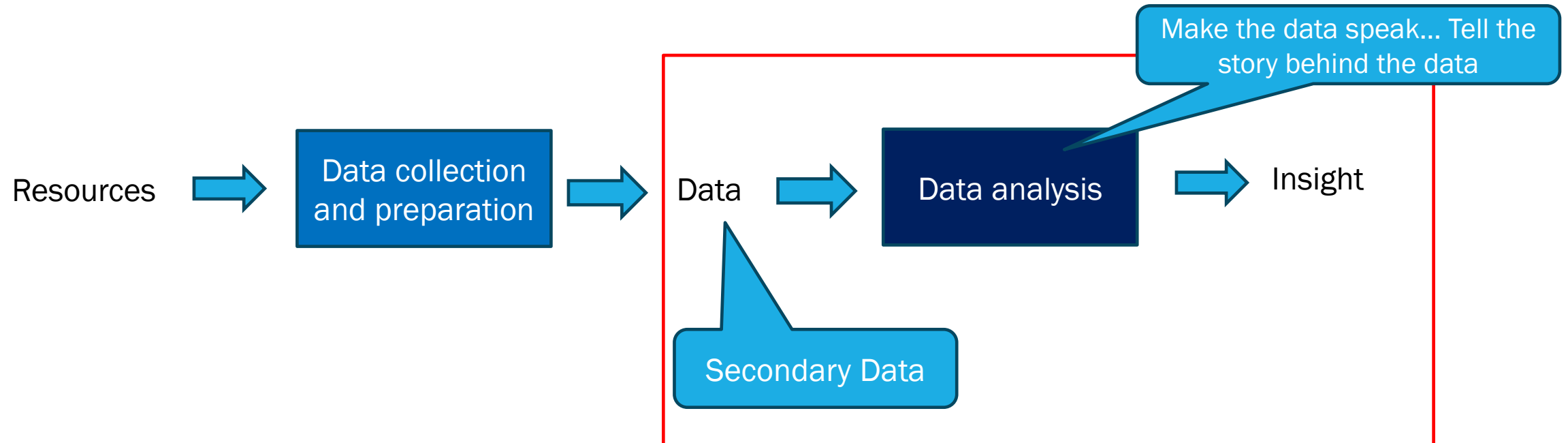
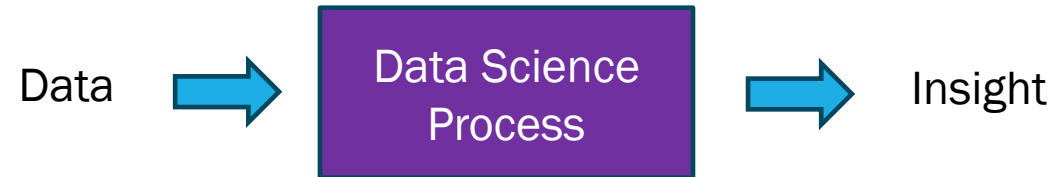
Amit Bharadwa · Follow

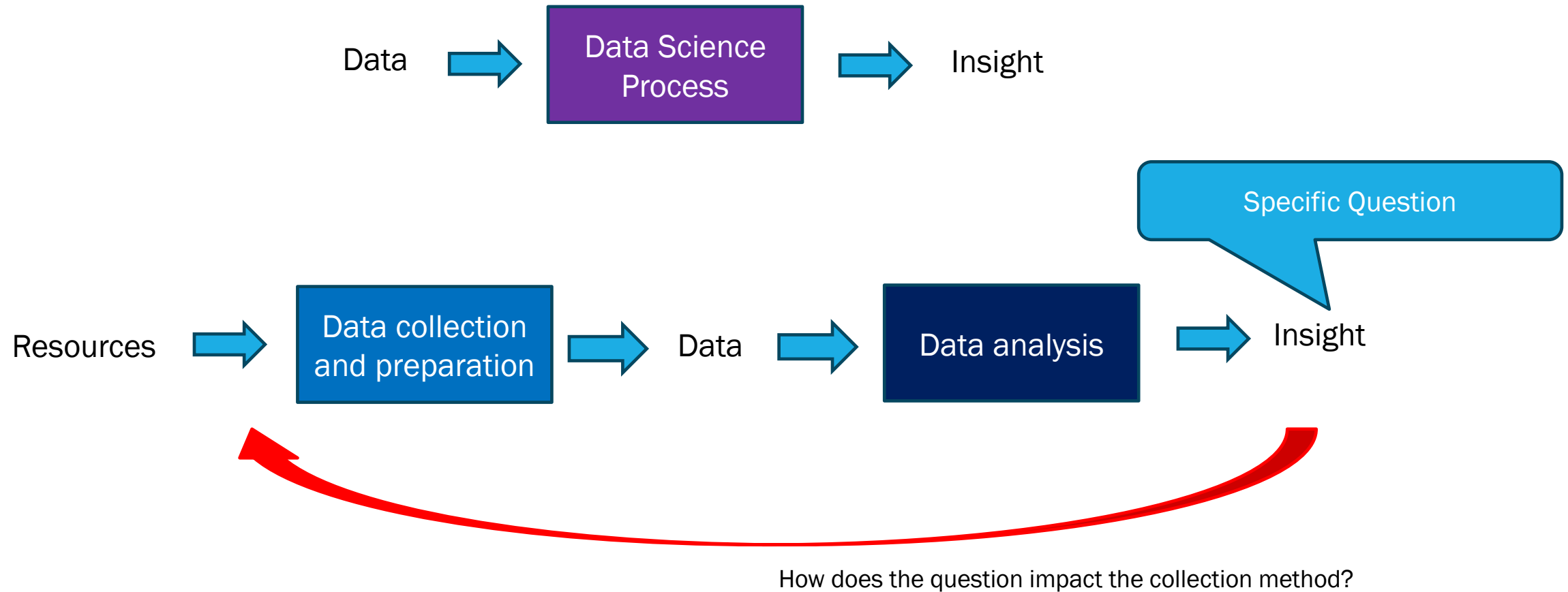
Published in Towards Data Science · 7 min read · Feb 6, 2021

The method is as follows:

1. Problem Statement
2. Data Collection
3. Data Cleaning
4. Exploratory Data Analysis (EDA)
5. Feature Engineering
6. Modelling
7. Communication







Questions

Example 1 – Retail Domain

How much money are younger people spending on clothing compared to older people?

Example 2 – Fashion Domain

How do consumers feel about the ethical sourcing of products in the fashion industry?

Example 3 – Digital Marketing Domain

What is the expected bounce rate if the company adopts the suggested website redesign?

Example 4 – Environment Domain

How do air quality levels vary across different regions of Ottawa during peak traffic hours?

Example 5 – Urban Planning Domain

What are the most common complaints about public transportation services in Gatineau?

Example 6 – Public Health Domain

How do social media users perceive the new health policy announced by the government?

Example 7 – E-commerce Domain

What is the average number of hours a user spends on an e-commerce site per week?

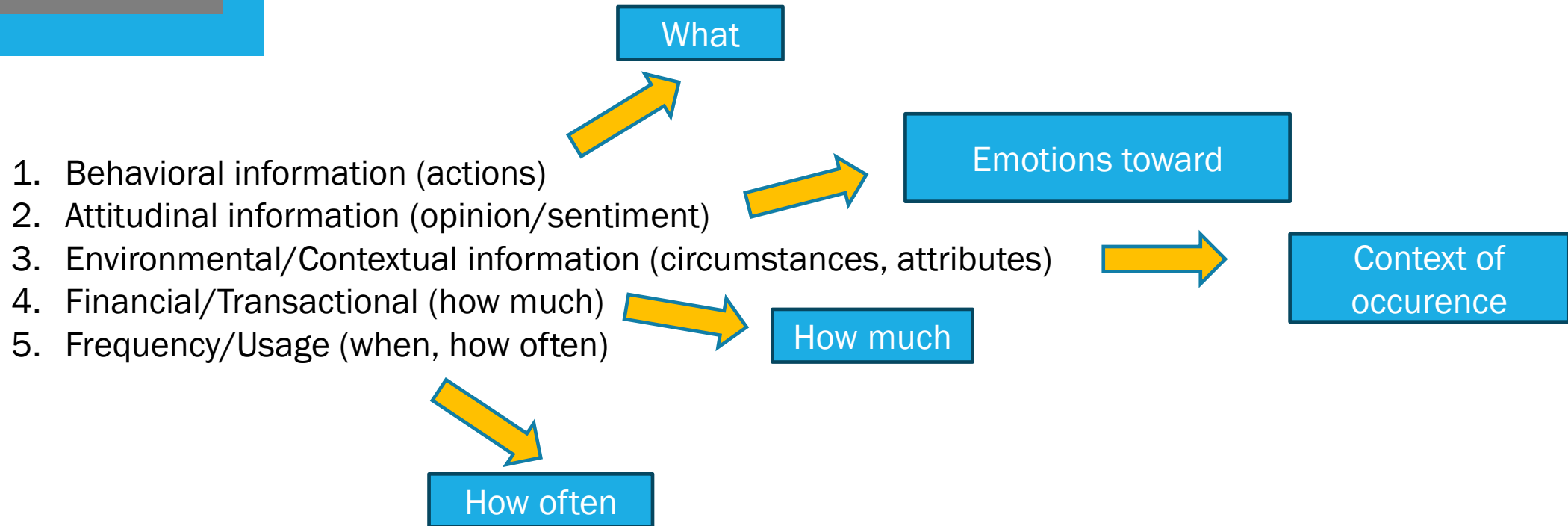
Example 8 – Human Resources Domain

What are the primary motivations of employees for choosing remote work over on-site work in a company?

Example 9 – Healthcare Domain

What are the effects of a new drug treatment on the recovery rate of patients with a specific condition?

What are we trying to capture?



Example 1 – Retail Domain

How much money are younger people spending on clothing compared to older people?

transactional

context

Example 2 – Fashion Domain

How do consumers feel about the ethical sourcing of products in the fashion industry?

opinion

context

Example 3 – Digital Marketing Domain

What is the expected bounce rate if the company adopts the suggested website redesign?

behavior

context

Collection Methods

Types of Data Sources: A Comprehensive Guide to Understanding Different Data Sources



Sumana Dotnettricks · Follow

10 min read · Jun 13, 2023

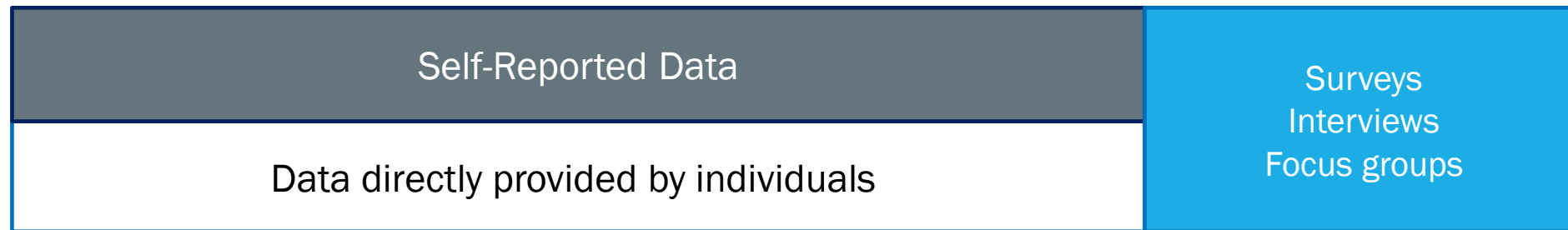
[Source](#)

Really nice resource for listing types of data sources. I'll just organize them a bit differently

Types of Data Collection Methods

1. Self-Reported Data
2. Behavioral and Observational Data
3. Automated Data
4. Digital and Platform-based Data

1. Self-Reported Data
2. Behavioral and Observational Data
3. Automated Data
4. Digital and Platform-based Data





Surveys

Purpose:

To systematically collect quantifiable data from a sample of individuals to understand their opinions, preferences and satisfaction levels.

Example:

A tourism board wants to understand the preferences of international tourists visiting their country.

2023 Survey on Consumer Perceptions of Food

Wave VI • Methodology

March 16 to 28, 2023

3,343 Canadian adults (18+) who have at least shared (50% or greater) responsibility for grocery shopping for the household.



Statistics Canada performs many surveys

70%

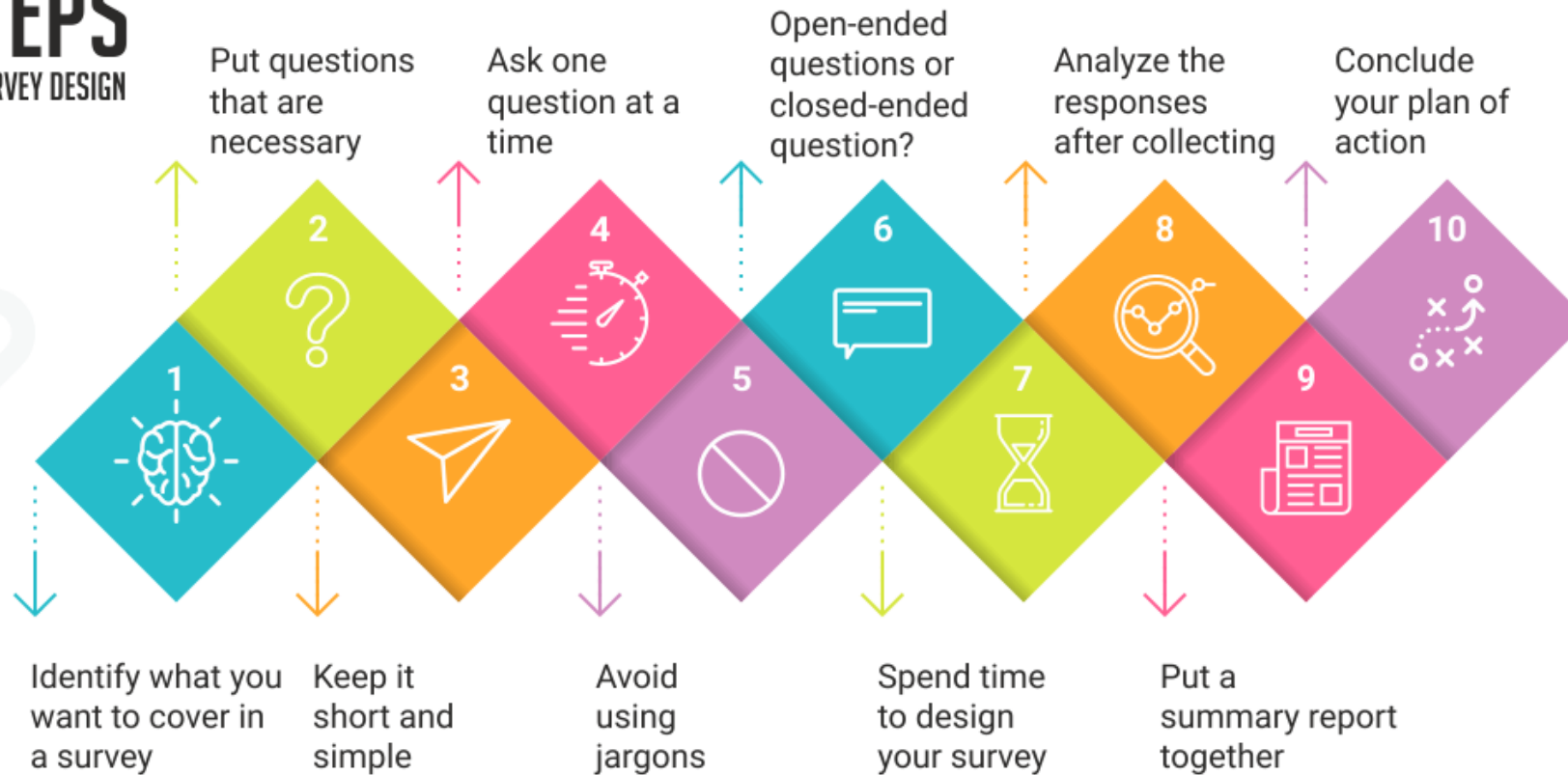
of consumers say they have changed their food purchasing habits in the last year because of increasing food prices



COST OF FOOD



10 STEPS TO A GOOD SURVEY DESIGN





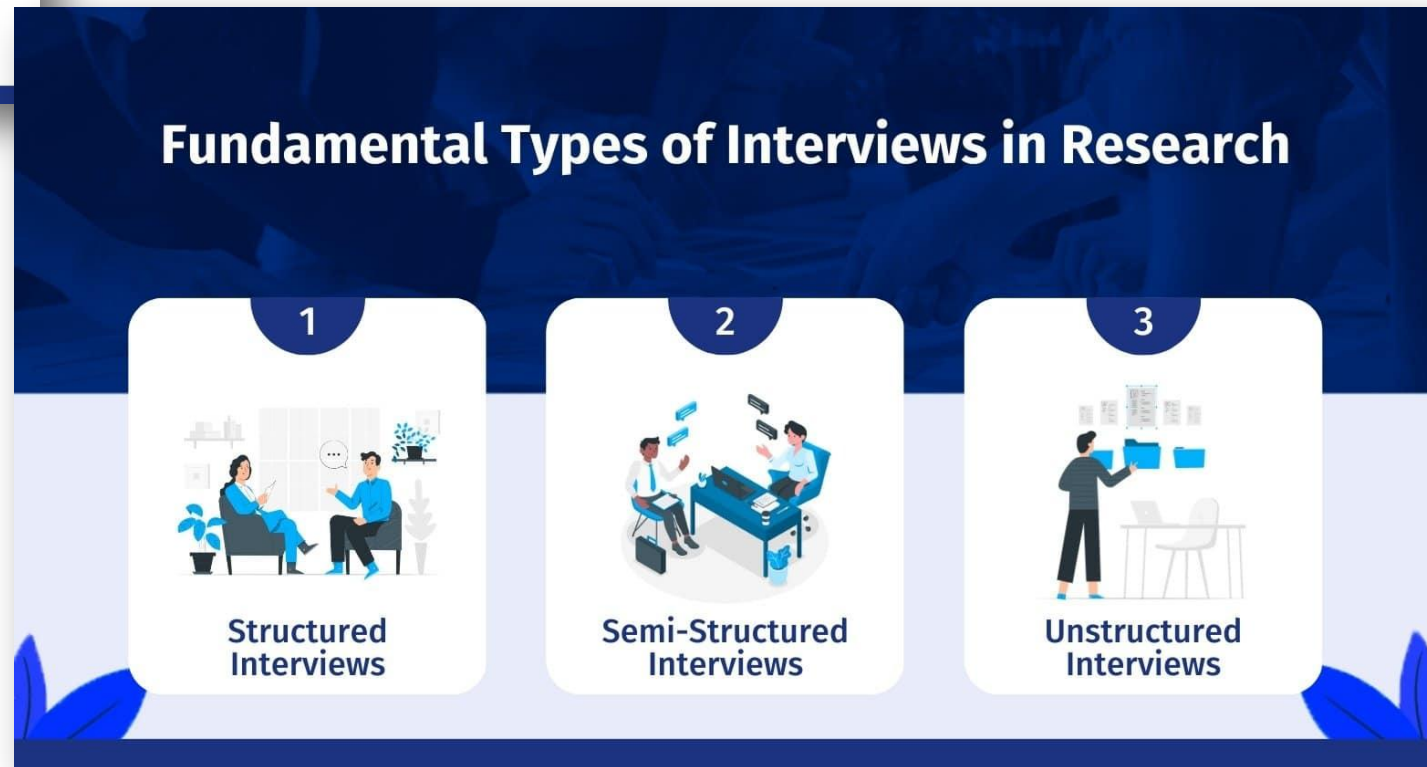
Interviews

Purpose:

To gather detailed, qualitative insights by exploring an individual's experiences, perspectives, or knowledge through one-on-one conversation.

Example:

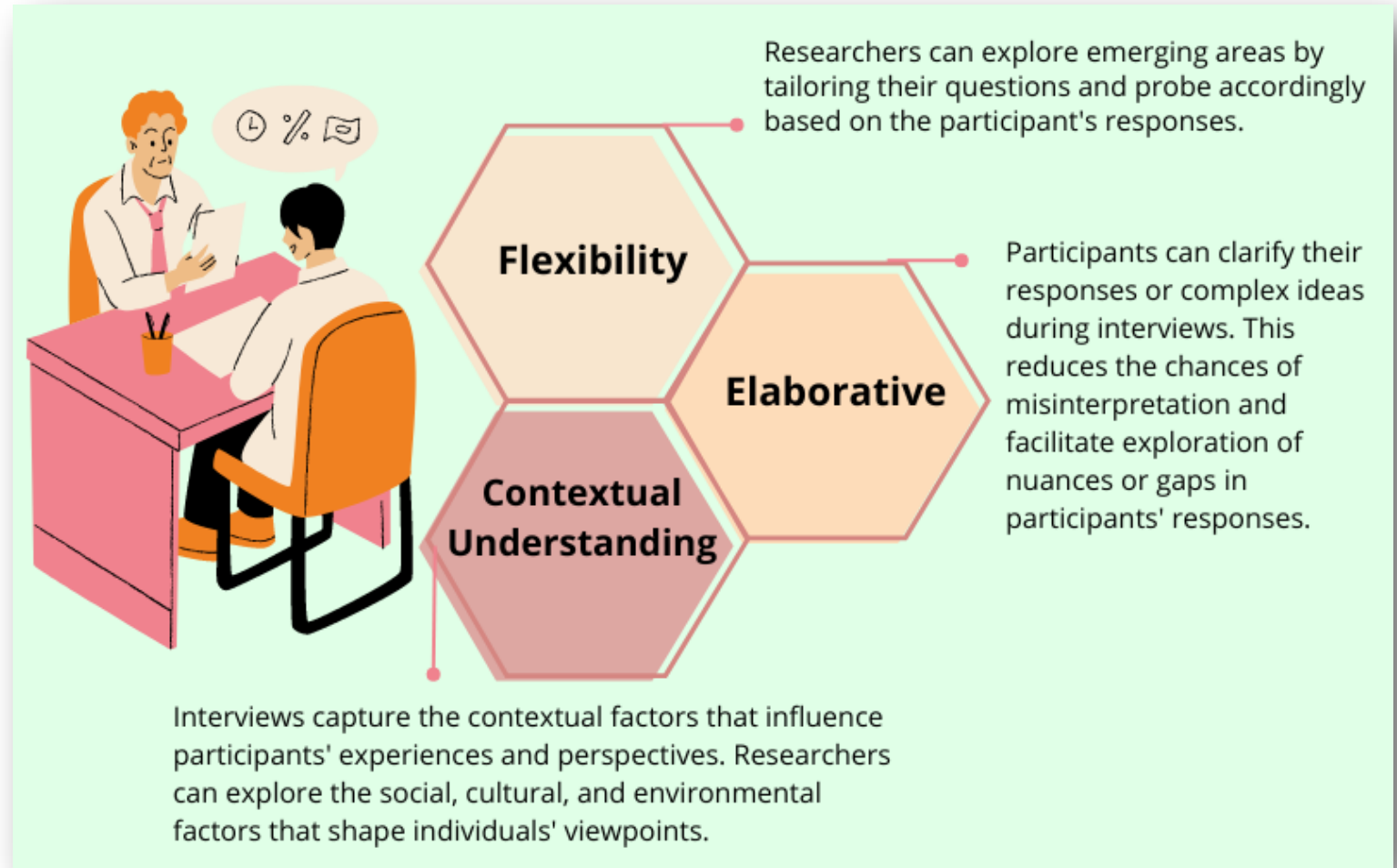
Conducting interviews with patients to understand their personal experiences with a new medical treatment.



2



Semi-Structured Interviews





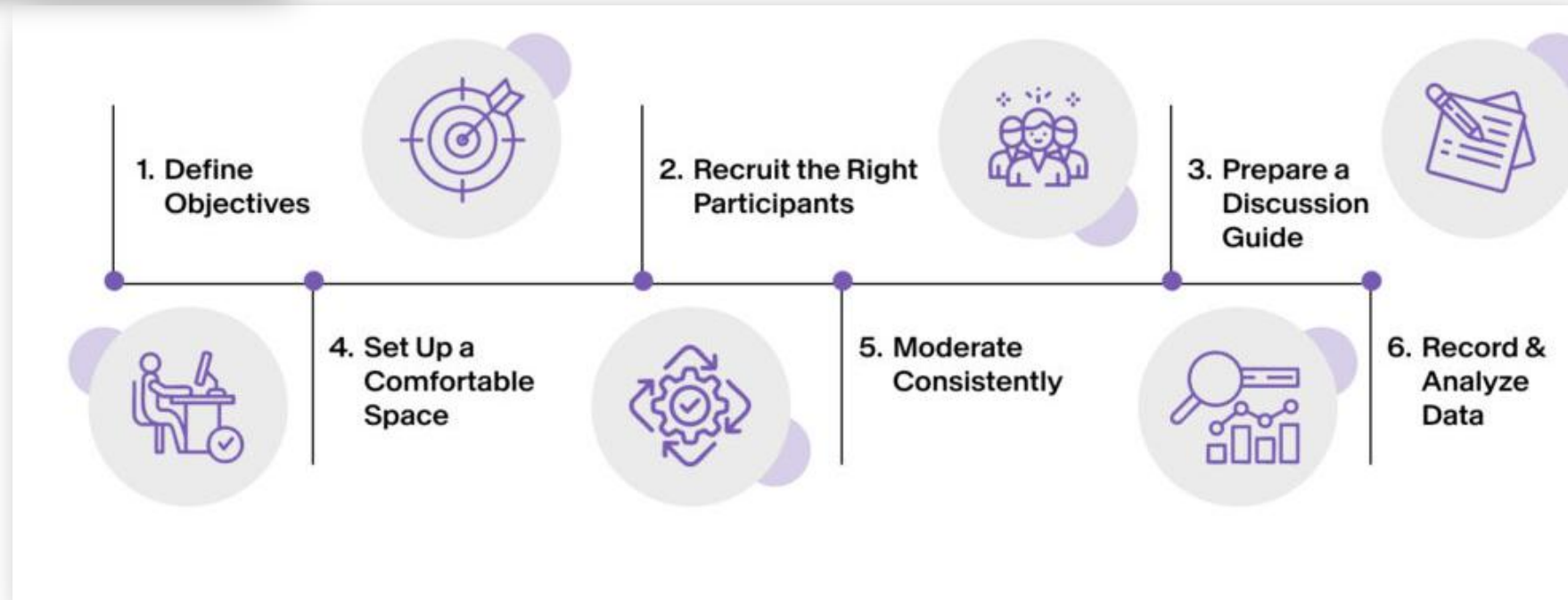
Focus Groups

Purpose:

To gather diverse perspectives, ideas, or feedback through guided discussions among a small group of participants on a specific topic.

Example:

Organizing a focus group of parents to explore opinions on proposed changes to school nutrition programs.



1. Self-Reported Data
2. Behavioral and Observational Data
3. Automated Data
4. Digital and Platform-based Data

Behavioral and Observational Data

Data collected through monitoring or observing human behavior.

Observations
Experiments



Observations

Purpose:

To collect data by systematically watching and recording behaviors, actions, or events in their natural or controlled settings.

Examples:

Observing customers in a retail store to analyze shopping behaviors and product preferences.

Observing traffic flow at an intersection to study the effects of a newly implemented traffic light system.



Observational Research Methods



Have a Clear Objective



Get Permission



Unbiased Observation



Hide Your Observers



Documentation



Data Analysis



Experiments

Purpose:

Data collected by manipulating variables in a controlled setup (e.g., clinical trials, A/B testing).

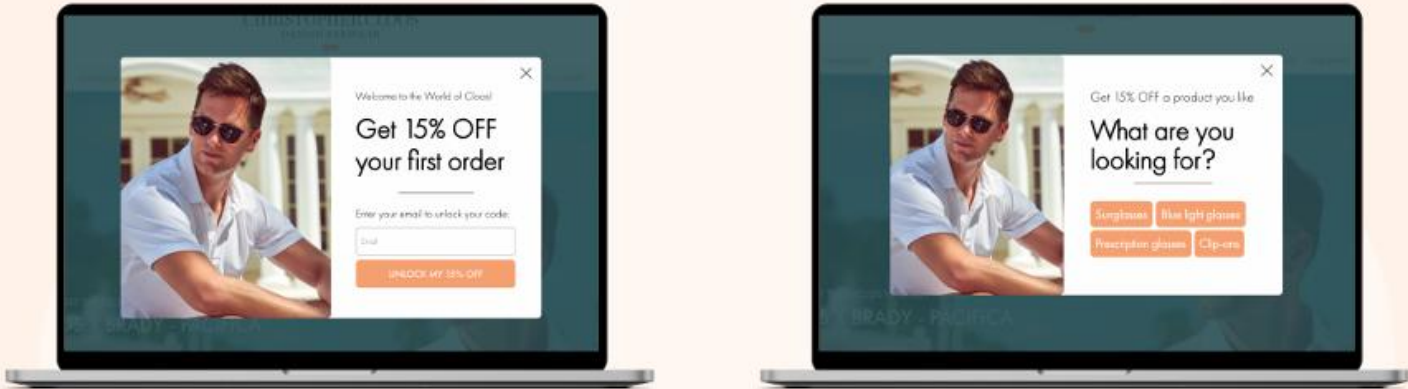
Examples:

A travel agency wants to test whether offering a limited-time discount on adventure packages increases bookings.

Conducting an experiment to test the effectiveness of a new drug by comparing patient recovery rates between a treatment group and a control group.

A conversion rate records the percentage of users who have completed a desired action. Conversion rates are calculated by taking the total number of users who 'convert' (for example, by clicking on an [advertisement](#)), dividing it by the overall size of the audience and converting that figure into a percentage.

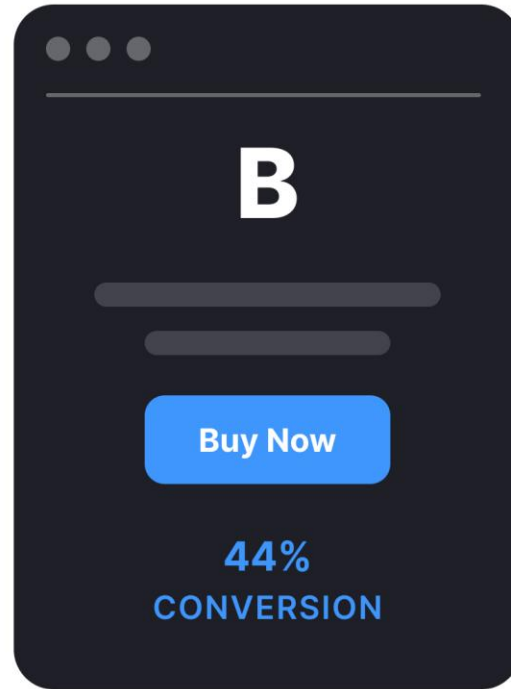
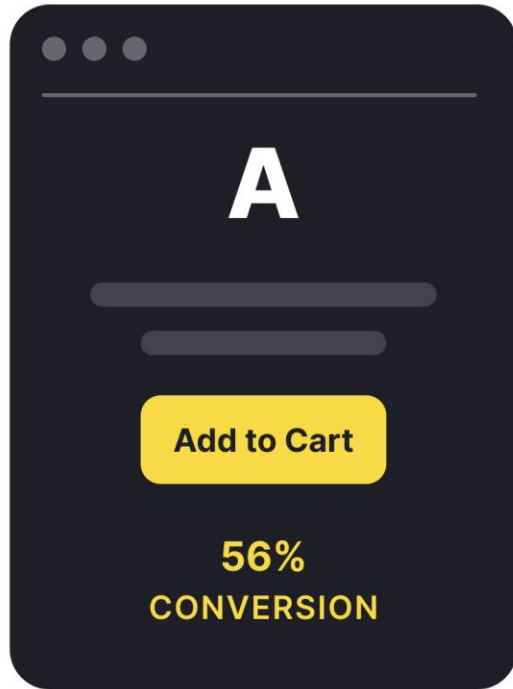
+15.38% conversion rate



Classic welcome popup

Conversational popup

A/B testing



Contact Dropdown - Default Static top Fixed top

Welcome to our website

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat.

Learn more

Click rate: 52 %



Project name Home About Contact Dropdown - Default Static top Fixed top

Welcome to our website

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat.

→ Learn more

72 %

1. Self-Reported Data
2. Behavioral and Observational Data
3. Automated Data
4. Digital and Platform-based Data

Automated Data	Logs Sensor Transactions
Data generated obtained through an automated system.	



Logs

Purpose:

Analyse automatically generated records of system or user activities (e.g., app usage, website clicks). Log files capture information about system errors, user interactions, and more.

Example:

Analyzing website log files to identify popular pages and navigation patterns.

13 Best Log Analysis Tools of 2025. Top Paid, Free & Open-Source Log Analyzers Reviewed

TOOLS & COMPARISONS

Updated on: September 29, 2025

Best Log Analyzers

1. **Sematext Logs**

2. Elastic Stack

3. Graylog

5. Loggly

5. Splunk

6. Logz.io

7. Sumo Logic

8. SolarWinds

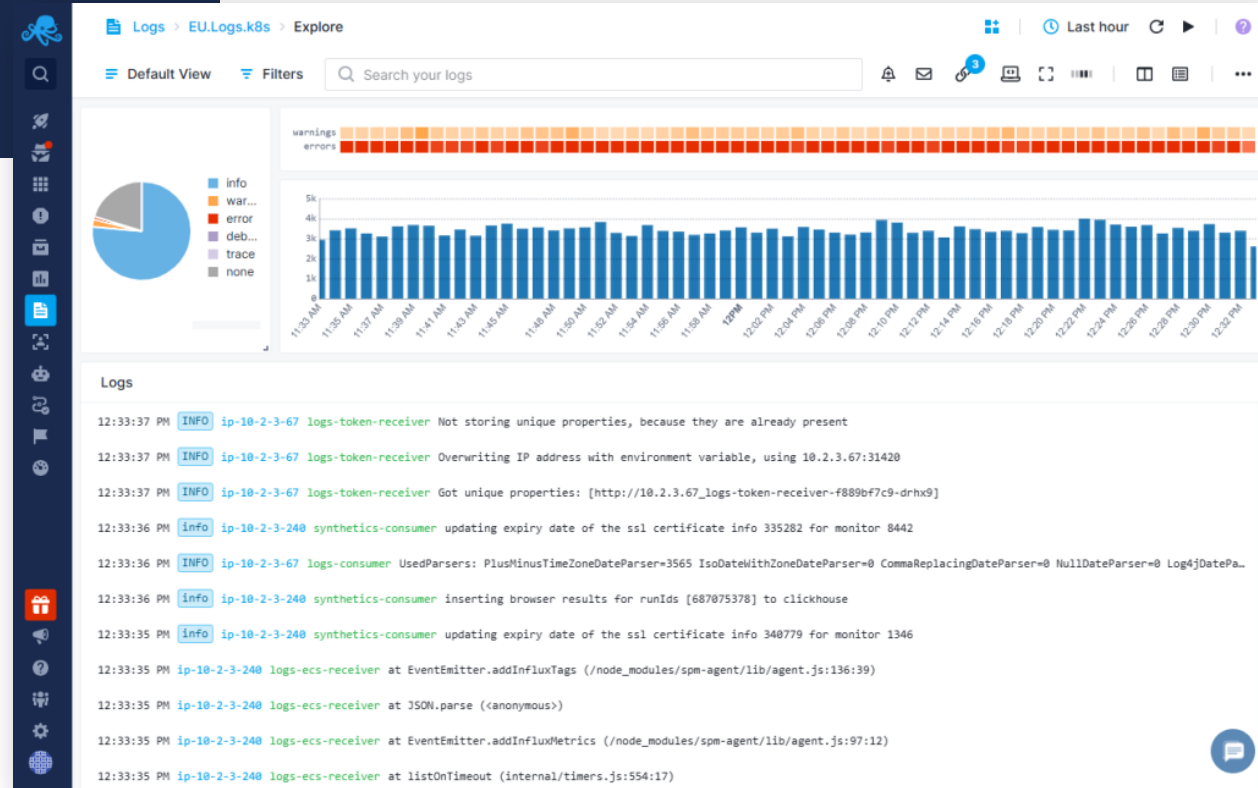
9. ManageEngine

10. Papertrail

11. Mezmio

12. Datadog

13. LogicMonitor





Sensors

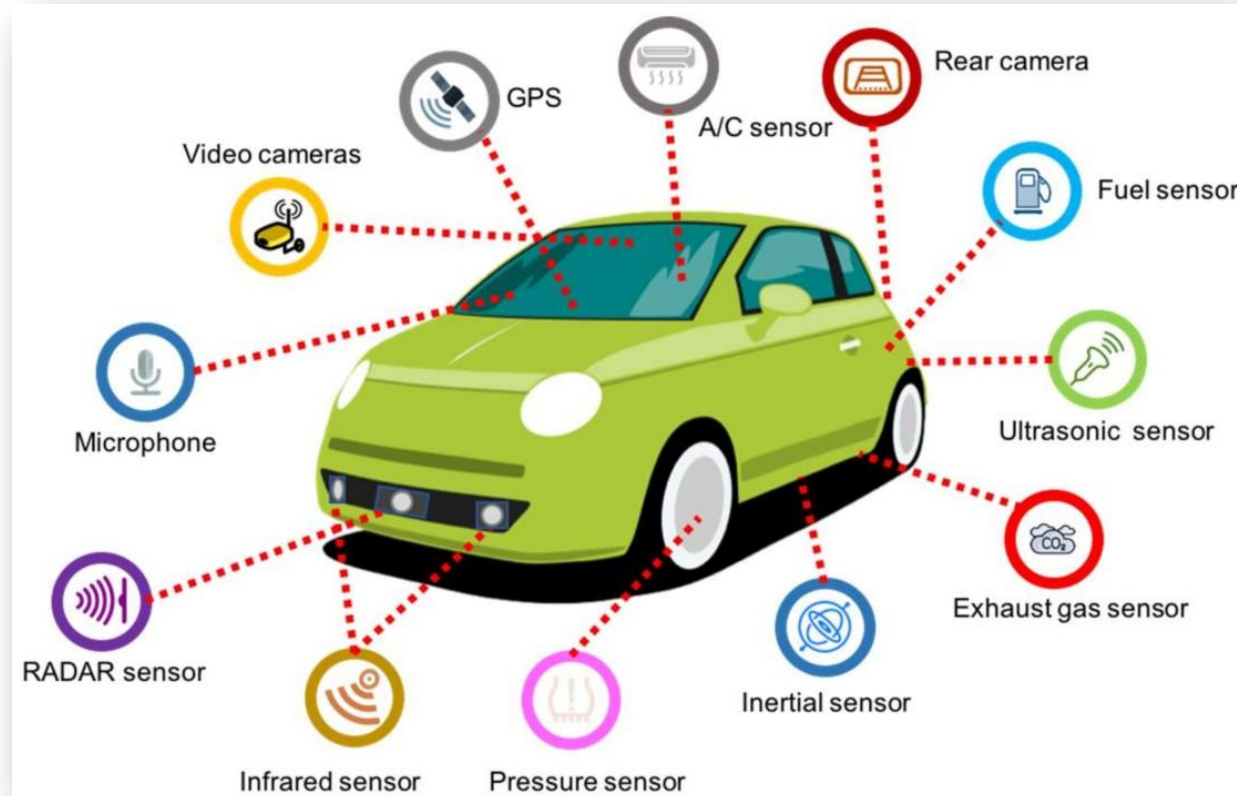
Purpose:

Analyse measurements coming from devices monitoring physical or environmental conditions (e.g., motion, temperature, air quality). Sensors are ubiquitous in today's world, collecting data on various physical parameters such as temperature, pressure, motion, and more.

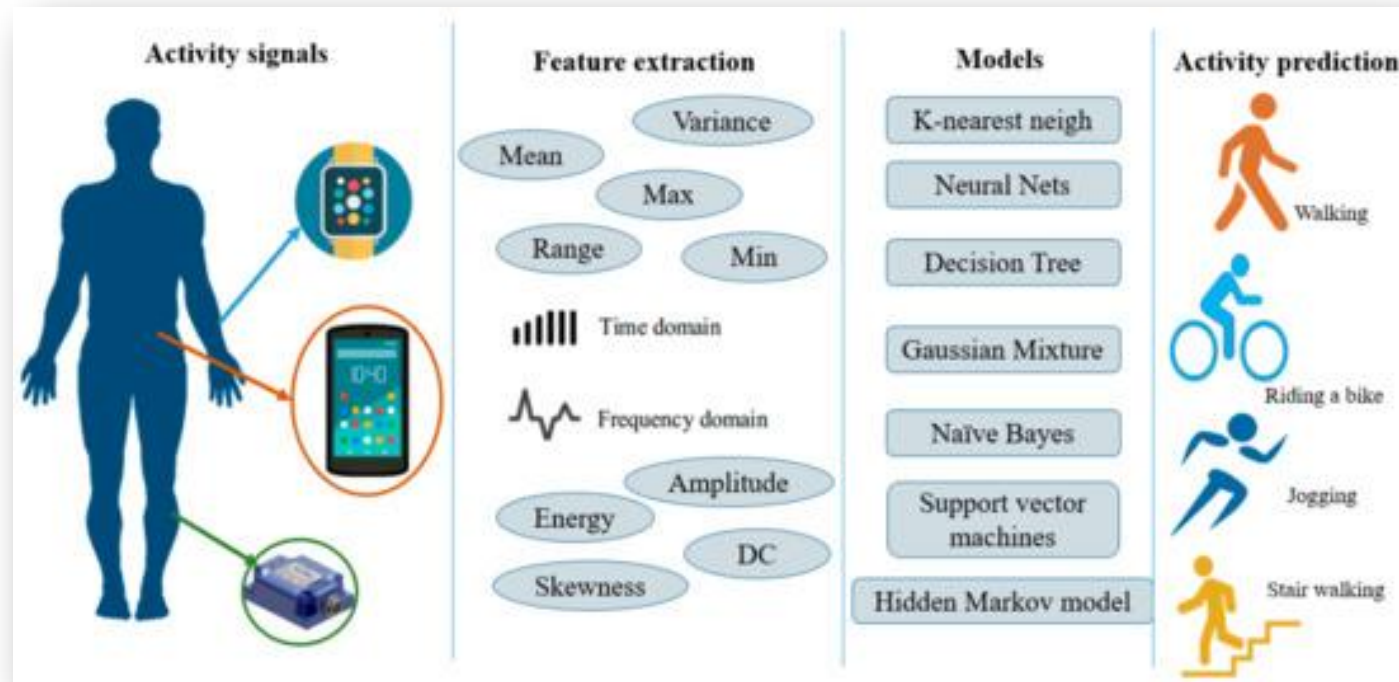
Example:

Using air quality sensors to monitor pollution levels in urban areas over time.

Sensor Data Analytics - Automobile



New Sensor Data Structuring for Deeper Feature Extraction in Human Activity Recognition †



Industry-Specific Uses of Sensor Data Analytics



Manufacturing

- Monitoring equipment & machinery (e.g., temperature, vibration, lubrication, pressure), environmental conditions, resource utilization, and employees' activity.



Agriculture

- Monitoring weather conditions, crop growth and health.
- Monitoring the chemical composition (acidity, moisture, nutrient levels) of soil and agricultural waste.
- Real-time livestock monitoring (e.g., location, body temperature).
- Monitoring and automatically adjusting greenhouse conditions based on pre-set parameters and real-time sensor data.

Industry-Specific Uses of Sensor Data Analytics



Smart cities

- Monitoring air quality, energy consumption, driving patterns, vehicle and foot traffic flow, etc.
- Enabling predictive city infrastructure maintenance.
- Forecasting energy and water consumption.
- Root cause analysis and intelligent suggestions to improve traffic management, resource consumption, waste management, and public health.



Healthcare

- Remote patient monitoring and IoT for medical devices.
- Identifying trends in patient symptoms to determine the factors influencing patient health and help improve care outcomes.



Transactions

Purpose:

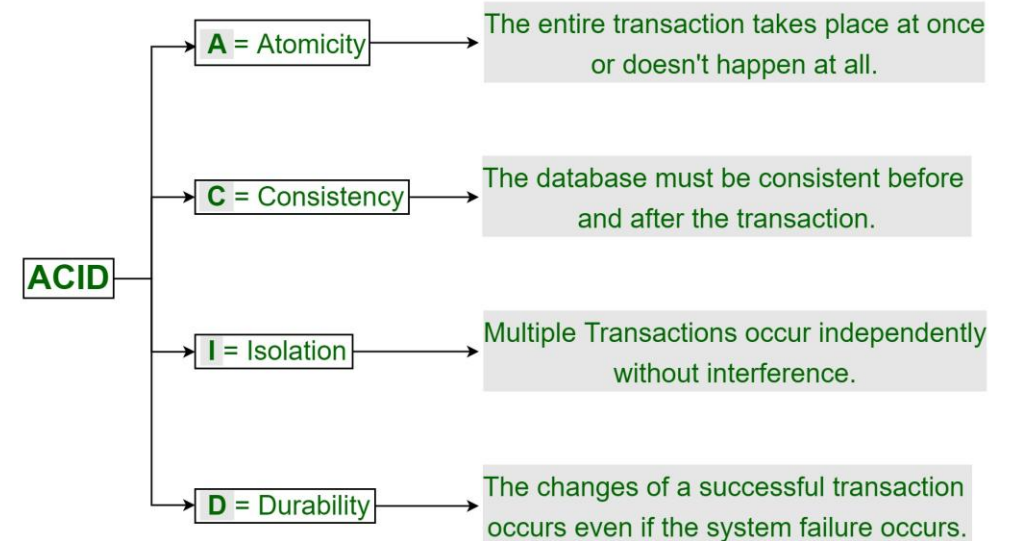
To collect detailed records of transactions, such as purchases or exchanges, to analyze patterns, behaviors, and financial activity. Transactional data typically captures the details of exchanges or interactions between entities (e.g., customers, organizations, or systems) and provides a record of a specific event or action.

Examples:

Review purchase records of renewable energy equipment from manufacturers, distributors, or financial institutions providing green loans.

	A	B	C	D	E	F
1	Date	Transaction ID	Customer Name	Category	Amount	Payment Method
2	2023-01-01 00:00:00	T00001	John Doe	Forex	54697.34	Check
3	2023-01-01 11:40:32	T00002	Alice Brown	Equity	-6230.63	Cash
4	2023-01-01 23:21:05	T00003	Jane Smith	Derivative	99837.92	Wire Transfer
5	2023-01-02 11:01:37	T00004	Jane Smith	Commodity	4473.36	Credit Card
6	2023-01-02 22:42:10	T00005	Jane Smith	Fixed Income	71438.19	Cash
7	2023-01-03 10:22:42	T00006	Bob Johnson	Forex	80311.67	Credit Card
8	2023-01-03 22:03:15	T00007	Jane Smith	Forex	31036	Check
9	2023-01-04 09:43:47	T00008	Charlie Green	Forex	11653.73	Wire Transfer
10	2023-01-04 21:24:19	T00009	John Doe	Derivative	868.59	Wire Transfer
11	2023-01-05 09:04:52	T00010	Alice Brown	Commodity	72346.66	Credit Card
12	2023-01-05 20:45:24	T00011	Alice Brown	Forex	39791.89	Credit Card
13	2023-01-06 08:25:57	T00012	John Doe	Commodity	68508.95	Wire Transfer
14	2023-01-06 20:06:29	T00013	Charlie Green	Fixed Income	90694.84	Credit Card
15	2023-01-07 07:47:02	T00014	Bob Johnson	Equity	6124.21	Wire Transfer
16	2023-01-07 19:27:34	T00015	Alice Brown	Equity	91108.81	Online Banking
17	2023-01-08 07:08:06	T00016	Bob Johnson	Fixed Income	35278.91	Check
18	2023-01-08 18:48:39	T00017	Bob Johnson	Equity	23579.37	Wire Transfer
19	2023-01-09 06:29:11	T00018	Alice Brown	Forex	93736.85	Check
20	2023-01-09 18:09:44	T00019	Bob Johnson	Derivative	98971.69	Check
21	2023-01-10 05:50:16	T00020	John Doe	Equity	11878.14	Online Banking
22	2023-01-10 17:30:49	T00021	Jane Smith	Fixed Income	62252.22	Wire Transfer
23	2023-01-11 05:11:21	T00022	Alice Brown	Derivative	1714.48	Credit Card
24	2023-01-11 16:51:54	T00023	Jane Smith	Fixed Income	61600.54	Credit Card
25	2023-01-12 04:32:26	T00024	Alice Brown	Forex	81004.46	Wire Transfer
26	2023-01-12 16:12:58	T00025	Charlie Green	Fixed Income	65294.84	Wire Transfer
27	2023-01-13 03:53:31	T00026	Jane Smith	Forex	35906.65	Check
28	2023-01-13 15:34:03	T00027	Alice Brown	Equity	32137.3	Check
29	2023-01-14 03:14:36	T00028	Bob Johnson	Derivative	33243.47	Cash
30	2023-01-14 14:55:08	T00029	Jane Smith	Commodity	54868.3	Check

ACID Properties in DBMS



1. Self-Reported Data
2. Behavioral and Observational Data
3. Automated Data
4. Digital and Platform-based Data

Digital and Platform-Based Data

Data originating from digital platforms, systems, or online activity.

Social Media Data
Web Scraping
APIs



Social Media Data

Purpose:

To collect data from social media platforms to analyze user sentiment, trends, or interactions.

Examples:

Analyzing Twitter data to understand public sentiment regarding a political candidate during an election campaign.

Analyze social media post to understand people awareness and sentiment about climate change are increasing in response to global weather events.

What People Share On Social Networks



What People Like On Social Networks

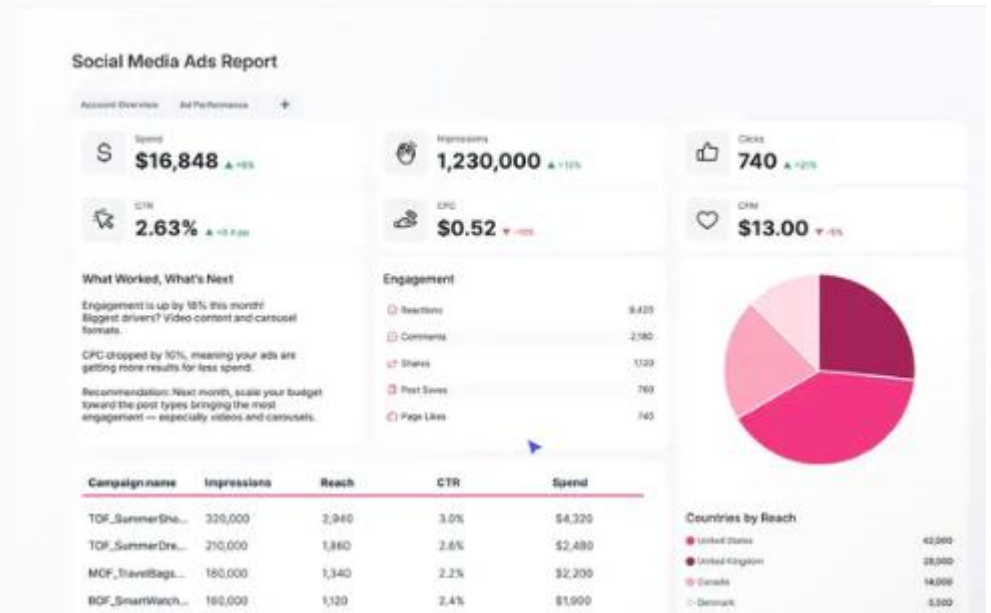
Each element needs to be analyzed differently



12 Best Social Media Analytics Tools in 2026

Here's a summary of the top social media analytics tools we'll be reviewing:

1. Whatagraph
2. Metrics Watch
3. Keyhole
4. Google Analytics 4 (GA4)
5. Sprout Social
6. Buffer
7. Hootsuite
8. Brandwatch
9. Rival IQ
10. Zoho Social
11. Mentionlytics
12. BuzzSumo





Web Scraping

Purpose:

To extract data from websites automatically, often for collecting large amounts of structured information from online sources. Web scraping is particularly useful for gathering information from websites that do not offer structured APIs or downloadable datasets.

Example:

Using web scraping to gather product pricing data from multiple e-commerce websites for competitive analysis.

How web scraping works



[Source](#)

Top 10 Best Web Scraping Frameworks for Data Extraction



[Source](#)



APIs

Purpose:

To collect data by accessing and interacting with external platforms or services through predefined interfaces for specific data requests. APIs and web services enable developers to access and interact with data from various online platforms and databases.

Examples:

Using a weather API to collect real-time temperature and humidity data for a mobile app.

Open Weather Map
Get weather and weather forecasts for multiple cities.
Verified ✓
9.9 414ms 100%

US Weather By Zip Co...
Provides current weather information for US cities by zip code
Verified ✓
8.8 37ms 100%

AerisWeather
An advanced weather API to p... all of your business needs, from...
Verified ✓
9.4 1024ms 100%

ClimaCell
MicroWeather API – An all-in-one weather API: Get the most accurat...
Verified ✓
9.6 512ms 100%

Weather
Current weather data API, and Weather forecast API - Basic acces...
Verified ✓
9.9 216ms 100%

Visual Crossing Wea...
Visual Crossing Weather API provides instant access to both ...
Verified ✓
9.4 699ms 100%

Dark Sky
The easiest, most advanced, weather API on the web.
Verified ✓
9.6 215ms 100%

US Weather by City
Provides current weather information for US cities by city an...
Verified ✓
6.5 87ms 100%

Simple Weather
Simple tool for get current weather
0.4 430ms 100%










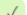












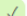












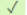















National Weather Ser...
National Weather Service API (api.weather.gov) NOAA (National Oceanic and Atmospheric ...
9.2 804ms 96%

UK Air Quality
Provides Air Quality Data across the UK, updated hourly and provided by DEFRA
6.1 1367ms 100%

AccuWeather
AccuWeather provides hourly and Minute by Minute™ forecasts with ...
Verified ✓
5 108778ms 32%

RapidAPI
Top Weather APIs
List of the Best Weather APIs to provide historical and trending weather forecasts.

The « problem » is which to choose?

 Morning Star Morning Star API helps to query for all information about finance summary, stocks, quotes, movers a...  9.5  1527ms  100%	 Bloomberg Market an... These APIs helps to query for all information about Indices, Commodities, Currencies, Futures,...  9.8  1556ms  100%	 Alpha Vantage The simplest and most effective way to receive stock, ETF, forex, ... Verified   9.8  1077ms  100%
 BraveNewCoin Latest and historic cryptocurrency market data  9.8  61ms  100%	 Morningstar Financial data for over 50,000 stocks on over 50 exchanges. Requests generally require a MIC ...  9.7  1703ms  100%	 Currency Converter Provides exchange rates based on the official banks data. Verified   9.9  61ms  100%
 Investors Exchange (IE... The IEX API is a free, web-based API supplying IEX quoting and trading data.  9.8  267ms  100%	 Yahoo Finance Yahoo Finance API helps to query for all information about finance summary, stocks, quotes, movers, ...  9.9  1304ms  100%	 Fixer Currency Powered by 15+ exchange rate data sources, the Fixer API is capable of ... Verified   9.8  301ms  100%
 Seeking Alpha Query for news, market moving, price quotes, chart, indices, analysis from investors and experts, etc...  9.8  1527ms  100%	 Yahoo Finance Low La... Real time low latency Yahoo Finance API for stock market  9.8  1556ms  100%	 Yahoo Finance Yahoo Finance official API for stocks, options, ETFs, mutual funds and news.  9.8  1077ms  100%

Not just for weather, for finance and other areas.

DECISION FACTORS

Self-Reported Data	Surveys Interviews Focus groups
Data directly provided by individuals	
Behavioral and Observational Data	Observations Experiments
Data collected through monitoring or observing human behavior.	

Automated Data	Logs Sensor Transactions
Data generated obtained through an automated system.	
Digital and Platform-Based Data	Social Media Data Web Scraping APIs
Data originating from digital platforms, systems, or online activity.	



Quite overwhelming... How to choose/decide on data collection method?

Resource constraints
(Effort/Cost)

Time constraints

Goal of the study

Availability of already collected data

Domain expertise

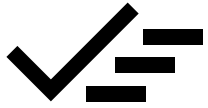
Technical expertise

Data storage possibilities

Ease of use and integration

Reliability

Ethical and Privacy concerns



DATA COLLECTION

- Data Science Process
- Types of questions (insights)
- Overview of data collection methods

Continue with case studies on Friday