

# Moving Closer Towards Web Data Integration

An Interactive Lecture

Manpreet Singh Juneja

Former FinTech Intern, Global Research, JP Morgan Chase & Co.

# What is WDI?

WDI is an extension and specialization of data integration that views the web as a collection of heterogeneous databases.

Example Website: <https://www.covid19india.org/>



# The Need for Web Data Integration

1. At the beginning of 2020, the digital universe was estimated to consist of 44 zettabytes of data.
2. Modern SMEs (Small and Medium Enterprises) have access to an unprecedented amount of data. Unfortunately, most of it goes unused in terms of analytics, meaning business owners are failing to capitalize on the data boom.
3. With AI-powered services growing in number, the greed for data is exploding.



# The Stages of a WDI Life Cycle



1. **Identify:** Find the page, and where on the page, web data is located.
2. **Extract:** Collect either displayed or hidden content, even if it's behind a login or require user interactions.
3. **Prepare:** Cleanse and normalize web data using functions and formulas.
4. **Integrate:** Bring prepared web data into other business applications for faster insights.
5. **Consume:** Use graphs, charts, and other visualizations to see how web data changes over time.

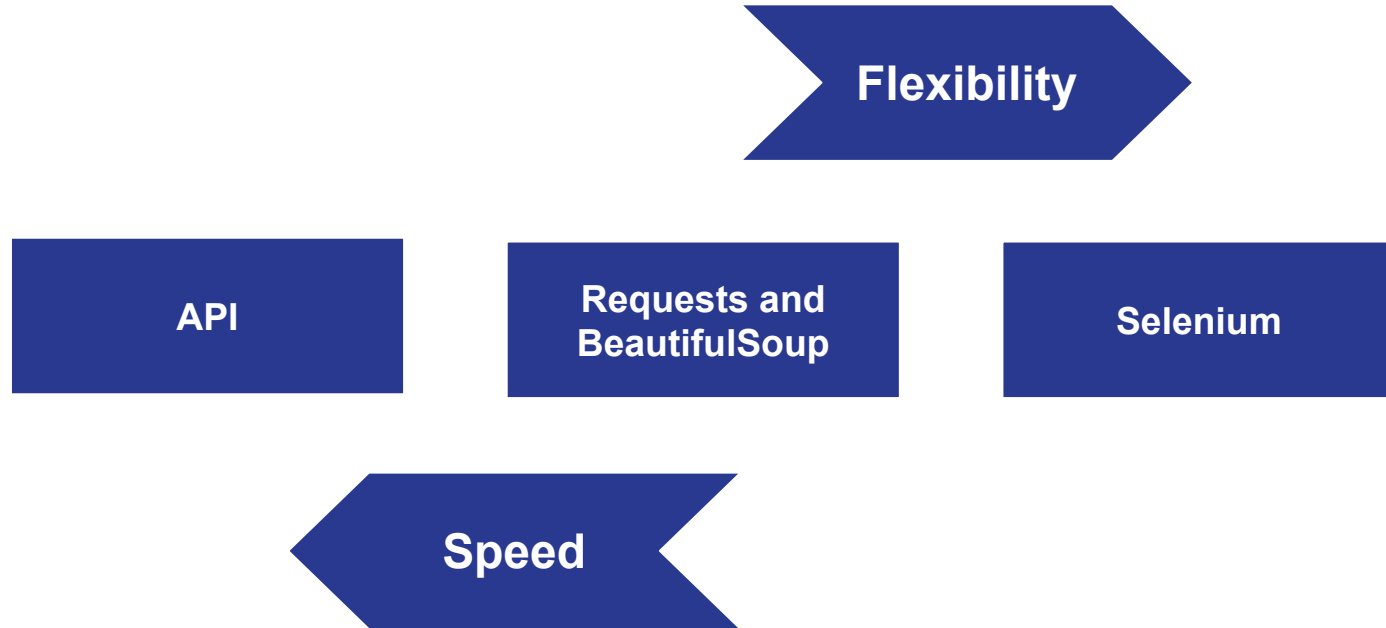
## Web Scraping

## Web Data Integration

<b>Data quality</b>	Poor quality, unreliable and dirty data	High quality, reliable and clean data
<b>Resource Needed</b>	Resource intensive, requiring specialized skills	Little or no internal resources required
<b>Legal Risk</b>	High risk	Low or little risk
<b>Integration</b>	Stand alone files—needs to be integrated manually or with other tools	Integrated directly into business process and apps, analytics and AI platforms, and data warehouses—ready for consumption
<b>Time to Value</b>	Time consuming and difficult to maintain	Rapid onboarding and data delivery, and automatically maintained
<b>Components</b>	Scrape only, a subset of Extract	Extract, Prepare and Integrate
<b>Types of data extracted</b>	Displayed only	Displayed, Hidden and Derived
<b>Scaling</b>	Poor and fragile	Web scale

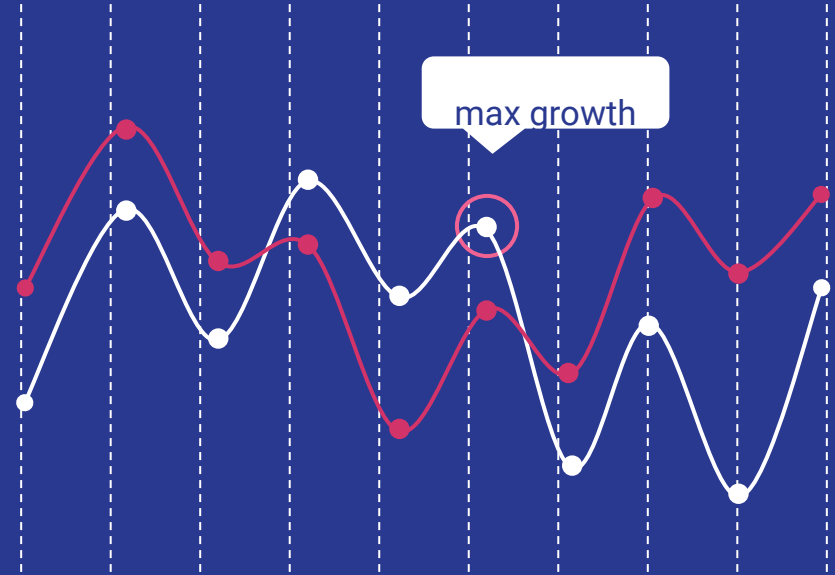
Source: import.io

# Methods of Extracting Web Data



# Live Demo

Demo of the three web data extraction methods.



# Career Opportunities

Since this is a very upcoming field, not many companies offer a dedicated job profile for the same.

However, freelancing is quite popular and some startups do offer A dedicated role for the same.

1. <https://www.upwork.com/freelance-jobs/web-scraping/>
  2. <https://www.import.io/careers/>
  3. <https://www.octoparse.com/blog/top-30-free-web-scraping-software?qu=>
-



# References

1. <https://www.import.io/post/web-data-integration-vs-web-scraping/>
2. <https://pupuweb.com/quality-control-web-data-enters-enterprise-mainstream/>
3. <https://realpython.com/learning-paths/python-web-scraping/>

