

# Report — Spotify Playlist Scraper Automation

(Task No. 1)

Project Title: Spotify Featured Playlist Data  
Author: Rohan Jadhav  
Date: 22, July 2025

## Objective:

This project automates the process of extracting metadata and track-level details from any public Spotify playlist. The extracted data is saved to a formatted Excel workbook for further analysis or archival.

## Tools/Technologies used:

Requirement Type	Name / Tool	Version (Recommended)	Purpose
Programming Language	Python	3.9 or above	Core language for scripting and automation
Browser Driver	ChromeDriver	Compatible with Chrome	Required for Selenium to control the Chrome browser
Web Browser	Google Chrome	Latest (curr. version 138.0.7204.158 )	Used for loading and scraping playlist content
Automation Library	Selenium	4.x	For web scraping and DOM interaction
Excel Library	openpyxl	3.x	To create and modify Excel (.xlsx) workbooks
Text Editor / IDE	VSCode / PyCharm / Sublime	Latest	Recommended for editing and running the Python scripts

Table 1: Software Requirements

## Methodology:

- The script takes a Spotify playlist URL from the user via a command-line interface.
- It uses Selenium WebDriver to open the playlist page in a Chrome browser and simulate user scrolling to load all track data.
- Playlist metadata and detailed track information (track name, artist, album, duration, date added) are extracted by targeting specific DOM elements.
- The openpyxl library is used to export this data into a structured Excel workbook with formatted sheets for both metadata and tracks.
- Within Excel, simple visualizations like pivot tables, bar charts, and duration bins are applied for basic data analysis.
- The entire process is automated and modular, enabling reuse for any public Spotify playlist URL with minimal changes.

## Summarized Approach:

- Selenium WebDriver is used to automate browser interaction and simulate scrolling to load tracks.
- Playlist metadata (title, description, saves, etc.) and track-level information (track name, artists, album, date added, duration) are parsed from the DOM.
- Extracted data is organized and exported using openpyxl into an Excel workbook with bold headers and separate sheets.

## Project Directory Structure:

The project is modular and organized into three main scripts:

**project-1/**

|— **main.py**

|— **scraper.py**

|— **save\_excel.py**

|— **outputs/**

|     |— **Spotify\_Playlist\_Export.xlsx**

## └─ README.pdf

1. *scraper.py*: Handles browser setup and scraping logic
2. *save\_excel.py*: Responsible for Excel export using openpyxl
3. *main.py*: CLI runner script that connects components
4. *outputs/*: Auto-generated folder for exported Excel files
5. *README.pdf*: Report Document

## Execution:

The script is executed via the CLI using the command:

**python main.py "url"**

e.g. python main.py "https://open.spotify.com/playlist/37i9dQZF1DWUAOn5dYbrDa"

It opens the playlist, extracts relevant data, and saves it to path - *outputs/Spotify\_Playlist\_Export.xlsx*.

## Screenshots:

```
def get_track_data(driver, wait):
    scroll_to_load(driver)
    # Find the first 21 rows of tracks
    rows = driver.find_elements(By.XPATH, "//*[div[@role='row' and .//div[@aria-colindex='2']]]")
    data = []
    for idx, row in enumerate(rows):
        try:
            #track and artist
            track_col = row.find_element(By.XPATH, ".//div[@aria-colindex='2']")
            track_name = track_col.find_element(By.XPATH, ".//div[@dir='auto' and contains(@class, 'Text')]")
            artist_links = track_col.find_elements(By.XPATH, ".//span[contains(@class, 'UudGc')]")
            artists = ", ".join([a.text.strip() for a in artist_links])
            #album and date added, duration
            album_elem = row.find_element(By.XPATH, ".//div[@aria-colindex='3']//a")
            album_name = album_elem.text.strip()
            date_added = row.find_element(By.XPATH, ".//div[@aria-colindex='4']//span").text.strip()
            duration = row.find_element(By.XPATH, ".//div[@aria-colindex='5']").text.strip()
            # Append the data to the list
            data.append({
                "Track Name": track_name,
                "Artist(s)": artists,
                "Album Name": album_name,
                "Date Added": date_added,
                "Duration (mm:ss)": duration
            })
        except Exception as e:
            print(f"[Row {idx}] Skipped due to error: {e}")
    return data
```

Fig 1. Code Snippet



## **Insights:**

1. Track durations revealed that most songs in the playlist fall between 3–4 minutes, aligning with commercial music norms and listener preferences.
2. By extracting the "Date Added" field, we observed how the playlist evolved over time — showing when most tracks were added (e.g., in bursts or steadily).
3. Sorting tracks by album name or artist allowed us to spot repeated contributors, helping highlight if the playlist heavily features a specific album or artist.

## **Notes:**

- Headless execution is supported via config
- Currently limited to top 21 tracks (can be customized)
- Playlist must be public