# Web Scraping for Football Data

Building skills in extracting, cleaning, and analyzing football statistics from the web.

By: Sahil Gidwani

# Objectives

## Master Web Scraping Fundamentals

Understand methodologies, legal frameworks, and ethical considerations for responsible data extraction.

## Identify Data Sources

Discover football websites and analyze the types of statistics and metrics available for analysis.

## Apply Python Tools

Use BeautifulSoup, Requests, Selenium, and Pandas to scrape, clean, and store football data.

## Build Complete Pipelines

Design end-to-end workflows from website discovery to analysis-ready datasets.

# What We Will Cover

## 01

### The Role of Web Scraping in Football Analytics

Understanding why automated data extraction is essential for modern football analysis.

## 02

### Legal and Ethical Considerations

Navigating terms of service, respecting website policies, and maintaining ethical practices.

## 03

### Data Source Identification

Exploring major football websites and understanding available data types.

## 04

### Website Structure and Planning

Analyzing HTML structure and planning effective scraping strategies.

## 05

### Python Tools Setup

Setting up and understanding key libraries for web scraping and data processing.

## 06

### Practical Exercises

Hands-on projects covering real-world scraping scenarios and challenges.

# The Role of Web Scraping in Football Analytics

## Why Web Scraping Matters

Automated extraction of data from websites enables comprehensive football analysis. Most football sites don't provide easy data downloads, making scraping essential for analysts.

- Performance tracking across seasons
- Predictive modeling for matches
- Player scouting and recruitment
- Real-time statistics collection



Web scraping allows analysts to collect large, up-to-date statistics that power modern football analytics and decision-making processes.

# Legal and Ethical Considerations

## Check Terms of Service

Always review website terms and copyright policies before scraping sports data. Understand what's allowed and what's prohibited.

## Respect Technical Guidelines

Honor robots.txt files, implement proper rate limiting, and use appropriate user agent headers to identify your scraper.

## Practice Fair Use

Avoid aggressive scraping patterns and ensure responsible use of publicly available data to maintain access.

- **Bot/Scraping/Crawler Traffic on Sports-Reference.com Sites**
- **Hudl Acceptable Use Policy**

Made with GAMMA

# Data Source Identification
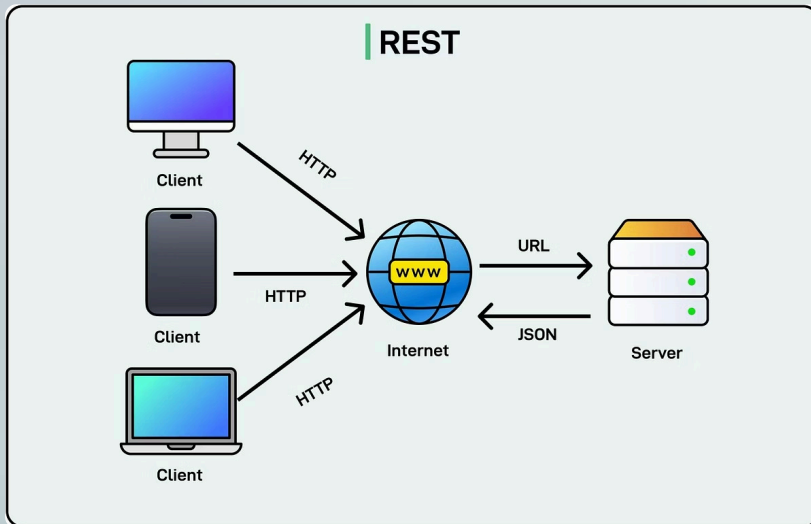
## Common Football Websites

1. FBref
2. WhoScored
3. Sofascore
4. FotMob
5. Understat
6. Transfermarkt

## Available Data Types

- Match scores and results
- Individual player statistics
- Live match events
- League tables and standings
- Advanced performance metrics
- Historical data archives

Each source offers unique data types and structures, requiring different scraping approaches and techniques.

# Website Structure & Planning



**REST**

Client — HTTP →
Client — HTTP →
Client — HTTP →
Internet
URL → Server
← JSON

### 1. Inspect HTML Source

Examine the website's HTML structure using browser developer tools to understand data organization.

### 2. Identify Data Loading Method

Determine if data appears as HTML tables or loads dynamically via JavaScript API calls.

### 3. Analyze Content Rendering

Distinguish between static HTML content and dynamically generated JavaScript-based data.

### 4. Plan Navigation Strategy

Study pagination patterns, URL structures for seasons/teams, and identify potential anti-scraping measures.

# Python Tools Setup

### Requests

HTTP library for fetching web pages and making API calls. Essential for accessing website content programmatically.

### BeautifulSoup

HTML parsing library for extracting specific elements from web pages. Perfect for navigating and searching HTML structures.

### Pandas

Data manipulation and analysis library for cleaning, transforming, and organizing scraped tabular data into usable formats.

### Selenium

Browser automation tool for handling JavaScript-heavy pages and interactive content that requires user simulation.

# Practical Exercises

- **Scraping player data from FBref with pandas**

  Leverage pandas' powerful data structures for efficient player data extraction.

- **Collecting team data from FBref with requests and BeautifulSoup**

  Utilize HTTP requests and HTML parsing to gather comprehensive team statistics.

- **Extracting FBref player data with Selenium for interactive content**

  Automate browser interactions to capture dynamic player information from interactive pages.

- **Gathering player shot data from Understat using requests and BeautifulSoup**

  Combine HTTP requests and HTML parsing to collect detailed shooting statistics.

- **Pulling player profiles and stats from Transfermarkt with requests and BeautifulSoup**

  Scrape market values, transfer histories, and player profiles from Transfermarkt.

- **Capturing match event data from WhoScored via Selenium**

  Use browser automation to extract real-time event data from live match pages.

- **Accessing player data from Sofascore through API response JSON**

  Parse JSON responses from Sofascore's API to retrieve structured player information.

- **Reproducing Sofascore API calls to retrieve structured data**

  Understand and replicate API requests to programmatically fetch data from Sofascore.

- **Querying Sofascore API endpoints for targeted information**

  Construct specific API queries to retrieve precise and targeted football data.

- **Intercepting Sofascore API calls with Selenium for data extraction**

  Employ Selenium to monitor and extract data from hidden API calls made by the browser.

- **Addressing Anti-Scraping Measures**

  Implement strategies like delays, user-agent rotation, and rate limit handling to ensure robust scraping.

- **Exploring the soccerdata library for pre-built football datasets**

  Utilize a specialized Python library for easy access to pre-compiled football statistics.

# Resources

### GitHub Repository

Complete code examples and project files for all exercises covered in this course.

**https://github.com/sahil-gidwani/football-data-webscraping**

### Video Tutorial Series

Comprehensive YouTube playlist covering sports analytics and web scraping techniques.
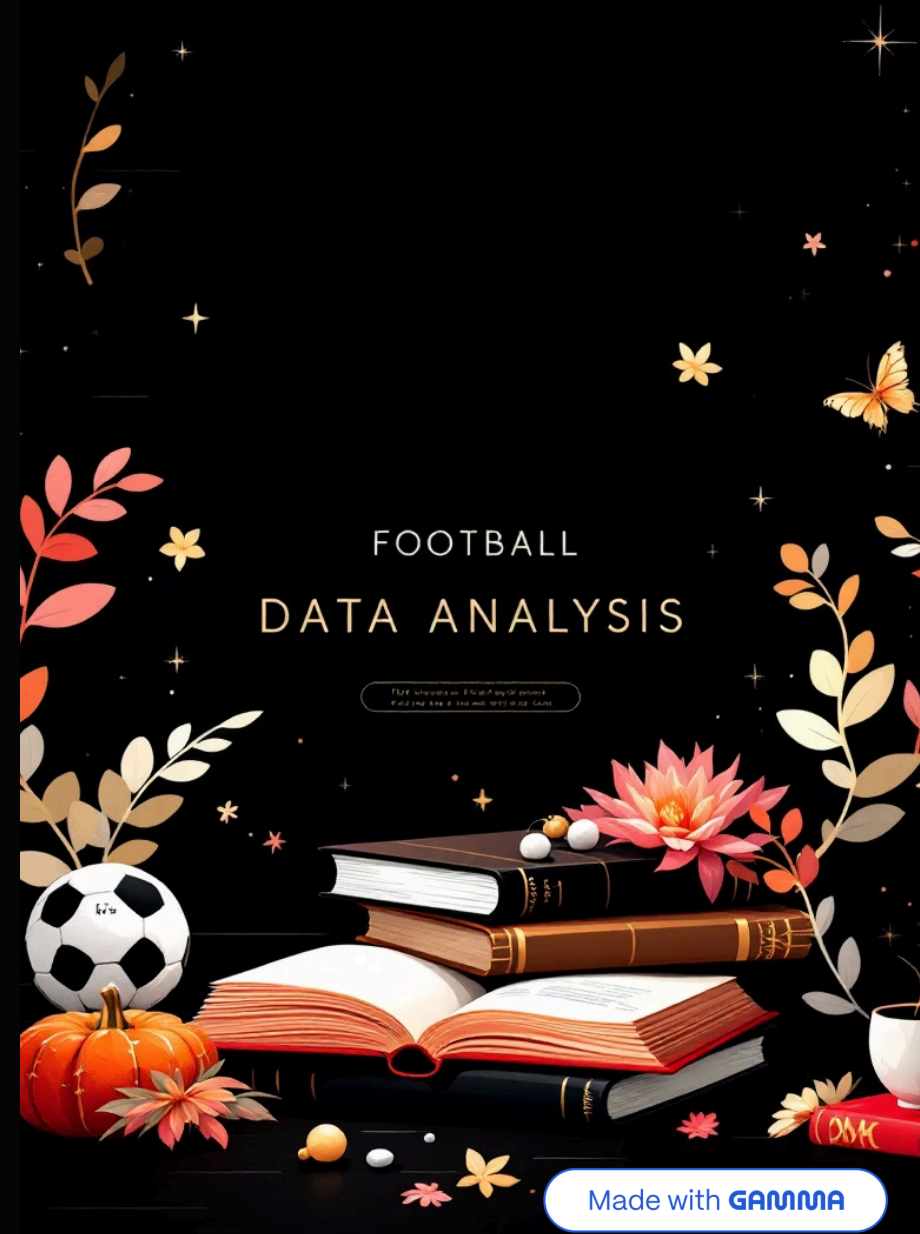
**Sports Analytics by McKay Johns - YouTube**

### SoccerData Library

Pre-built Python library for accessing football datasets without manual scraping.

**SoccerData's Documentation**

These resources provide ongoing support for your football data scraping journey, from beginner tutorials to advanced techniques.

# Thank You!

Kudos to you for hanging in and tolerating me this long.

## Any questions?

Feel free to ask anything about the presentation or topic.

## Any thoughts?

Share your insights, feedback, or ideas with us.