

Python Web Scraping Midterm Project – Documentation

Avtandil Beradze

Project Overview

This project is a modular Python application built to scrape, model, analyze, and visualize quotes from the website <https://quotes.toscrape.com>. It demonstrates the full pipeline of a web scraping task, from data acquisition and parsing to statistical analysis and GUI-based presentation.

Project Structure

The codebase is divided into logical modules for clear separation of concerns:

- `scraper/`: Handles HTML fetching and parsing
- `models/`: Contains classes representing structured data (quotes and authors)
- `utils/`: Includes tools for file handling and data analysis
- `main.py`: The entry point that orchestrates the scraping and analysis
- `gui.py`: A Tkinter-based GUI for visualizing quote data

Implementation Layers

1. Scraper Engine (`scraper/collector.py`)

This layer handles web page requests and navigation.

- **Session Setup**: A `requests.Session` object is configured with custom headers to simulate a real browser. The scraper respects the site's `robots.txt` and includes a delay (`time.sleep`) between requests to avoid overwhelming the server.
- **HTML Fetching**: Pages are fetched via `fetch_page()`, and paginated navigation is implemented through `fetch_all_pages()` using a while-loop until no more pages are found.

- **Error Handling:** Uses try/except blocks and logs issues using Python's logging module.

2. Data Parsing (scraper/parser.py)

Responsible for converting raw HTML into structured Python data.

- **Quote Parsing:** Utilizes BeautifulSoup to extract quote text, author, and tags. CSS selectors like .quote, .author, and .tag are used for precise selection.
- **Author Links:** If present, the author link is extracted to allow for future expansion (e.g., collecting author bios).

3. Data Modeling (models/data_models.py)

Defines Quote and Author classes to encapsulate structured data.

- **Quote Class:** Contains text, author, tags, and author_link, with a to_dict() method for easy export.
- **Author Class:** Stores metadata like name, bio, born_date, and born_location. Not all fields are used yet but are included for extensibility.

4. Data Storage & Analysis (utils/)

- **File Handling:** file_handler.py supports saving and loading data in JSON and CSV formats. Files are saved in an output/ folder, with automatic folder creation.
- **Analysis** (analyzer.py):
 - Calculates the most quoted author
 - Finds the most common tag
 - Computes the average number of tags per quote

Execution Flow (main.py)

The main script glues all components together:

1. Initializes the scraper and fetches all paginated pages
2. Parses quotes from each page using QuoteParser
3. Converts quote data into Quote objects and saves them as JSON/CSV
4. Runs analysis functions and displays stats in the console

GUI Interface (gui.py)

The GUI is built using Tkinter and matplotlib for visualization.

- **Table Display:** A TreeView widget lists all quotes with columns for text, author, and tags.
- **Statistics Section:** Displays the most common tag, most quoted author, and average tags per quote.
- **Tag Plot:** Shows a bar chart of the top 10 tags using matplotlib embedded in Tkinter via FigureCanvasTkAgg.

Key Design Principles

- **Modularity:** Each file or class has a clearly defined responsibility.
- **Reusability:** Models and parsers are designed to be extended (e.g., for adding new data fields or new output formats).
- **Robustness:** All critical operations (requests, parsing, file I/O) are wrapped with exception handling.
- **Ethical Scraping:** The scraper obeys robots.txt and uses request delays to minimize server load.

Challenges Faced

1. **Pagination Handling**
 - a. Solution: Used a while-loop and dynamic URL formatting to navigate through all pages until the end.
2. **Missing Data**
 - a. Solution: Used try/except blocks and fallback values to skip or log problematic entries.
3. **Author Data Redundancy**
 - a. Solution: Implemented structure to allow future caching of author metadata to avoid repeated HTTP requests.
4. **GUI Visualization**
 - a. Challenge: Displaying structured data cleanly.
 - b. Solution: Used ttk.TreeView for tabular data and matplotlib for charts.

Data Summary

- **Unique Authors:** ~50
- **Unique Tags:** ~100

Key Findings:

- **Most Quoted Author:** Albert Einstein
- **Top Tags:** love, life, inspirational, humor
- **Avg Tags per Quote:** Approximately 2.5

Possible Improvements

- **Concurrency:** Current scraping is sequential; using `asyncio` or `concurrent.futures` could significantly speed it up.
- **Author Metadata:** Extend the parser to retrieve and store author bios and birth locations.
- **GUI Enhancements:** Implement tabbed views, search functionality, or export buttons using more advanced GUI frameworks like `PyQt`.
- **Logging:** Expand logging with timestamps and log levels to better track scraping events and failures.

Conclusion

This project demonstrates a complete scraping pipeline: web requests, HTML parsing, data modeling, storage, analysis, and GUI visualization. Through a clean modular architecture, error handling, and ethical design choices, it is both maintainable and extensible. It also lays a strong foundation for more complex scraping or analytics projects in the future.