



亞洲大學
ASIA UNIVERSITY

Midterm Project Report

Advanced Computer Programming

Student Name : RADIAN TRY DARMAWAN
Student ID : 113021220
Teacher : DINH-TRUNG VU

2025-04

Chapter 1 Introduction

1.1 Github

- 1) **Personal Github Account:** <https://github.com/radiandrmwn>
- 2) **Group Project Repository:**
<https://github.com/113021198/ACP/tree/main/midterm/radian-try>

1.2 Overview

This Midterm project implements a GitHub repository scraper using the Scrapy framework in Python. The scraper extracts detailed information about all repositories belonging to a specified GitHub user. The implementation leverages several advanced Python language features and libraries:

- **Scrapy Framework:** Used in the process of creating the web crawler that navigates through GitHub pages
- **XML Processing:** Utilized to format and output the scraped data in a structured XML format
- **CSS Selectors:** Used for precise targeting of HTML elements containing repository data

The program successfully extracts repository information including URL, description, and last updated date. Data is exported as a well-structured XML file that can be easily parsed and analyzed further.

Chapter 2 Implementation

2.1 Class 1

2.1.1 The main spider class that inherits from scrapy(github_repos.py) and handles the crawling and data extraction process.

2.1.2 Fields

This snippet of the code is where we put our GitHub repositories address

```
1 class GitHubSpider(scrapy.Spider):
2     name = 'github_repos'
3     github_username = 'radiandrmwn'
4
5     start_urls = [f'https://github.com/{github_username}?tab=repositories']
6     custom_settings = {
7         'USER_AGENT': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/120 Safari/537.36'
8     }
```

- name: Identifier for the spider, set to 'github_repos'
- github_username: GitHub username to scrape repositories from, set to 'radiandrmwn'
- start_urls: Initial URL to begin scraping from, constructed using the username

This is important as it is pointing the GitHub address to start the scraping process

2.1.3 Methods

init(self, *args, **kwargs)

Initializes the spider by:

- Calling the parent class constructor
- Setting up an empty list to store repository data

```
1 def __init__(self, *args, **kwargs):
2     super().__init__(*args, **kwargs)
3     self.repositories = []
```

parse(self, response)

The entry point method that processes the repository list page:

- Extracts basic information from each repository listing
- Determines if repositories are empty
- Creates follow-up requests for each repository to gather detailed information
- Handles pagination to crawl through all repository pages

```
1 class GithubSpider(scrapy.Spider):
2     name = 'github_repos'
3     github_username = 'radiandrmmn'
4
5     start_urls = [f'https://github.com/{github_username}?tab=repositories']
6     custom_settings = {
7         'USER_AGENT': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/120 Safari/537.36'
8     }
9
10    def __init__(self, *args, **kwargs):
11        super().__init__(*args, **kwargs)
12        self.repositories = []
13
14    def parse(self, response):
15        repositories = response.css('li[itemprop="owns"]')
16        self.logger.info(f'Found {len(repositories)} repositories')
17
18        for repo in repositories:
19            repo_url = response.urljoin(repo.css('a[itemprop="name codeRepository"]::attr(href)').get())
20            repo_name = repo.css('a[itemprop="name codeRepository"]::text').get().strip()
21
22            about = repo.css('p[itemprop="description"]::text').get()
23            if about:
24                about = about.strip()
25
26            last_updated = repo.css('relative-time::attr(datetime)').get()
27            if last_updated:
28                last_updated = datetime.datetime.fromisoformat(
29                    last_updated.replace('Z', '+00:00')
30                ).strftime('%Y-%m-%d')
31
32            is_empty = 'This repository is empty.' in repo.get()
33
34            if not about and not is_empty:
35                about = repo_name
36
37            repo_data = {
38                'url': repo_url,
39                'name': repo_name,
40                'about': about,
41                'last_updated': last_updated,
42                'is_empty': is_empty,
43            }
```

```
1     if not is_empty:
2         yield scrapy.Request(
3             url=repo_url,
4             callback=self.parse_repository_details,
5             meta={'repo_data': repo_data}
6         )
7     else:
8         repo_data['languages'] = None
9         repo_data['commits'] = None
10        self.repositories.append(repo_data)
11
12    next_page = response.css('a.next_page::attr(href)').get()
13    if next_page:
14        yield response.follow(next_page, callback=self.parse)
```

2.1.4 Functions

parse_repository_details(self, response)

This method handles the processing of individual repository pages:

- Extracts programming languages used in the repository
- Extracts the number of commits using regular expressions
- Adds the detailed information to the repository data dictionary
- Stores the complete repository data for final output

```
1 def parse_repository_details(self, response):
2     repo_data = response.meta['repo_data']
3
4     # Languages
5     lang_data = response.css('.repository-content .Layout-sidebar span[itemprop="programmingLanguage"]::text').getall()
6     if not lang_data:
7         lang_data = ['None']
8     repo_data['languages'] = lang_data
9
10    # Commits
11    commits_text = response.css('li.Commits .d-none.d-sm-inline::text').get()
12    if commits_text:
13        commits = re.search(r'(\d+)', commits_text)
14        repo_data['commits'] = commits.group(1) if commits else '0'
15    else:
16        repo_data['commits'] = '0'
17
18    self.repositories.append(repo_data)
```

closed(self, reason)

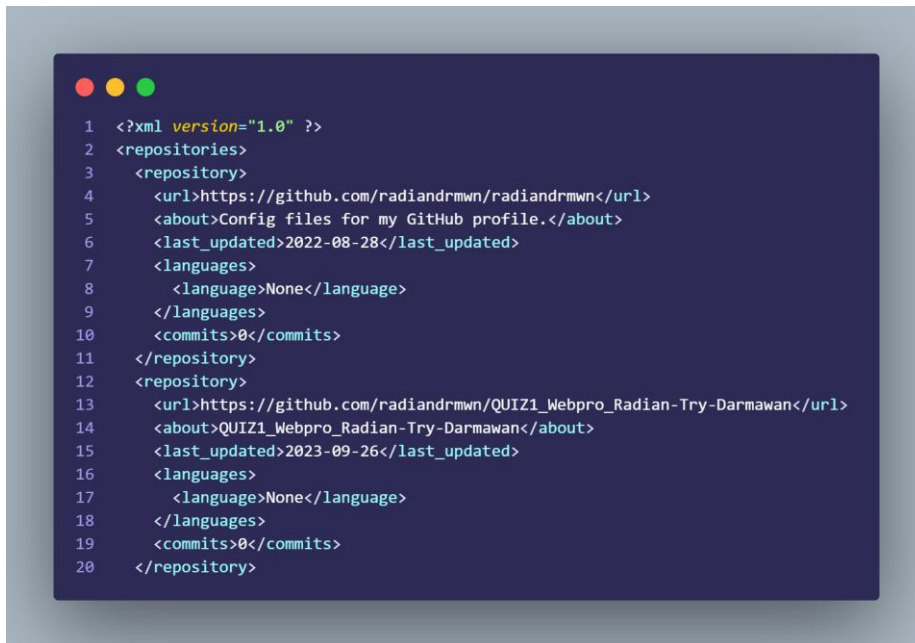
A special method that is called when the spider finishes crawling:

- Creates an XML structure to store all repository data
- Formats the data with proper indentation for readability
- Writes the formatted XML to a file named after the GitHub username
- Logs the completion of the data saving process

```
1 def closed(self, reason):
2     self.logger.info(f"Collected {len(self.repositories)} repositories")
3     root = Element('repositories')
4
5     for repo in self.repositories:
6         repo_elem = SubElement(root, 'repository')
7
8         SubElement(repo_elem, 'url').text = repo['url']
9         SubElement(repo_elem, 'about').text = repo['about'] or 'None'
10        SubElement(repo_elem, 'last_updated').text = repo['last_updated'] or 'None'
11
12        langs_elem = SubElement(repo_elem, 'languages')
13        if repo['languages']:
14            for lang in repo['languages']:
15                lang_el = SubElement(langs_elem, 'language')
16                lang_el.text = lang
17
18        SubElement(repo_elem, 'commits').text = repo['commits'] or 'None'
19
20    xml_str = tostring(root, encoding='utf-8')
21    pretty_xml = xml.dom.minidom.parseString(xml_str).toprettyxml(indent="  ")
22
23    with open(f"{self.github_username}_repositories.xml", "w", encoding="utf-8") as f:
24        f.write(pretty_xml)
25
26    self.logger.info(f"Saved data to {self.github_username}_repositories.xml")
```

Chapter 3 Results

The spider generates an XML file named after the GitHub username (e.g., radiandrmwn_repositories.xml). The XML file contains a structured representation of all repositories and their details



For each repository, the following information is successfully extracted:

- **URL:** Direct link to the repository
- **About:** The description of the repository or its name if no description exists
- **Last Updated:** The date when the repository was last modified

This data provides a comprehensive overview of a GitHub user's repositories, allowing for analysis of their programming preferences, activity level, and project topics.

Chapter 4 Conclusions

This project successfully implements a web scraper for GitHub repositories using the Scrapy framework. The scraper is designed to be robust, handling edge cases such as empty repositories and missing information. It also implements polite scraping practices by using a realistic user agent string to avoid being blocked by GitHub's anti-bot measures.

The project demonstrates the power of Python's ecosystem for web scraping tasks, showcasing how relatively few lines of code can extract significant amounts of structured data from complex websites. The resulting data can be used for various purposes, such as portfolio analysis, tracking project progress, or studying programming language trends among GitHub users.