WL

I have the following code for a web scraper for the Daily Pennslyvanian, a newspaper at the University of Pennsylvania. Can you help me understand how it works?

""" Scrapes a headline from The Daily Pennslyvanian website and saves it to a JSON file that tracks headlines over time. """

import os import sys

import daily_event_monitor

import bs4 import requests import loguru

def scrape_data_point(): """ Scrapes the main headline from The Daily Pennsylvanian home page.

Returns: str: The headline text if found, otherwise an empty string. """ headers = { "User-Agent": "cis3500-scraper" }

req = requests.get("https://www.thedp.com", headers=headers) loguru.logger.info(f"Request URL: {req.url}") loguru.logger.info(f"Request status code: {req.status_code}")

```
if req.ok: soup = bs4.BeautifulSoup(req.text, "html.parser") target_element = soup.find("a", class_="frontpage-link") data_point = "" if target_element is None else target_element.text loguru.logger.info(f"Data point: {data_point}") return data_point

if __name__ == "__main__":

    # Setup logger to track runtime loguru.logger.add("scrape.log", rotation="1 day")

    # Create data dir if needed loguru.logger.info("Creating data directory if it does not exist") try: os.makedirs("data", exist_ok=True) except Exception as e: loguru.logger.error(f"Failed to create data directory: {e}") sys.exit(1)

    # Load daily event monitor loguru.logger.info("Loading daily event monitor") dem = daily_event_monitor.DailyEventMonitor( "data/daily_pennsylvanian_headlines.json" )

    # Run scrape loguru.logger.info("Starting scrape") try: data_point = scrape_data_point() except Exception as e: loguru.logger.error(f"Failed to scrape data point: {e}") data_point = None

    # Save data if data_point is not None: dem.add_today(data_point) dem.save() loguru.logger.info("Saved daily event monitor")

    def print_tree(directory, ignore_dirs=[".git", "__pycache__"]): loguru.logger.info(f"Printing tree of files/dirs at {directory}") for root, dirs, files in os.walk(directory): dirs[:] = [d for d in dirs if d not in ignore_dirs] level = root.replace(directory, "").count(os.sep) indent = " " * 4 * (level) loguru.logger.info(f"{indent}+--{os.path.basename(root)}/") sub_indent = " " * 4 * (level + 1) for file in files: loguru.logger.info(f"{sub_indent}+--{file}")

    print_tree(os.getcwd())

    loguru.logger.info("Printing contents of data file {}".format(dem.file_path)) with open(dem.file_path, "r") as f: loguru.logger.info(f.read())

    # Finish loguru.logger.info("Scrape complete") loguru.logger.info("Exiting")
```

I'd be happy to break down this web scraper code for you! Here's how it works:

## Overview

This script scrapes the main headline from The Daily Pennsylvanian (Penn's student newspaper) website and saves it to a JSON file, tracking headlines over time.

## Key Components

### Imports

- Standard libraries: `os`, `sys` for file operations and system functions
- Custom module: `daily_event_monitor` (not shown here, but used for data tracking)
- Third-party libraries: `bs4` (BeautifulSoup), `requests` for web scraping, and `loguru` for logging

## `scrape_data_point()` Function

This function does the actual web scraping:

1. Sets up a custom User-Agent header to identify the scraper
2. Makes a GET request to "https://www.thedp.com"
3. Logs the request URL and status code
4. If the request is successful:
   - Creates a BeautifulSoup object to parse the HTML
   - Finds the first element with tag `<a>` that has the class "frontpage-link"
   - Extracts and returns the text from this element (the headline)
   - Returns an empty string if the element isn't found

## Main Block

When the script runs directly:

1. Configures `loguru` to save logs to "scrape.log" with daily rotation
2. Creates a "data" directory if it doesn't exist
3. Initializes the `DailyEventMonitor` with a specific JSON file path
4. Calls `scrape_data_point()` to get the headline
5. Adds today's headline to the monitor and saves it
6. Prints a directory tree of the current directory (excluding Git and Python cache)
7. Prints the contents of the data file
8. Logs completion

## Error Handling

The script has error handling for:

- Directory creation
- Web scraping process
- Any other exceptions that might occur

## `daily_event_monitor` Module

While not shown in your code, this appears to be a custom module that handles:

- Loading/saving the JSON file

- Adding new data points with timestamps
- Tracking headlines over time

## Overall Workflow

1. Setup logging and create directories
2. Load existing headline data
3. Scrape today's headline from The Daily Pennsylvanian
4. Add the new headline to the tracking system
5. Save the updated data
6. Log the file structure and data content
7. Complete the process

This is a well-structured scraper that not only collects data but also maintains good logging practices and error handling for reliability.

ⓘ
Claude does not have internet access. Links provided may not be accurate or up to date.

※