

Task 1

key considerations

1. Dynamic Content Loading (click_load_more): Use 12 in the script to avoid taking too much time.
 - It repeatedly searches for and clicks buttons or links containing "load" (case-insensitive) to reveal more faculty profiles that might be hidden behind "load more" or pagination mechanisms.
2. Identifying Profile Containers (inspect_frequent_combos):
 - It counts the frequency of each unique class string (combination of classes). It then takes the most frequent class combinations .
 - For each of these frequent class combinations, it calculates a "hit ratio": the proportion of elements with that class combination that also contain keywords like "professor" or "lecturer" in their text content. (The container for the profile should include title information.)
 - The class combination with the highest hit ratio is selected as the most likely container for individual faculty profiles.
3. LLM-Powered Rule Inference: extract the structure of the container.
 - Once the main profile container class is identified, the script takes a sample HTML snippet from one of these containers.
 - It constructs a prompt for the DeepSeek LLM, providing the sample HTML and asking the LLM to return a JSON object.
 - This JSON object should map predefined field names (e.g., "name", "title", "email", "research interest") to:
 - "selector": A concise CSS selector to find the element containing the field's value within the profile card.
 - "tag_pattern": The opening and closing HTML tags (including class attributes) that usually enclose the field's value. (Though the script primarily uses the selector).
4. Data Extraction Loop:
 - The script finds all elements matching the best profile container class identified earlier.
 - It iterates through each of these "profile card" elements.
 - For each card, and for each field defined in the LLM-generated rules (name, title, etc.): It attempts to find the specific data element using the CSS selector provided by the LLM for that field. It extracts the text content. If the text is empty, it tries a fallback by checking the data-value attribute of the parent element (adjustment for fitting more websites).

Task 2

key considerations

1. Product List Identification (find_products):

- It analyzes the page for frequently repeated CSS class combinations on major layout elements since the containers of the products are often duplicated.
- It filters out class combinations that appear at least min_products_for_list times (e.g., 10)

2. Then, for each container: Main Link Extraction from Containers

- Within each potential product container, it tries to find the primary clickable link leading to the product page.

If there is none, return. Select the second container.

- It handles lazy loading by scrolling if a link isn't immediately found after a successful previous extraction (very likely the right container). If multiple consecutive attempts fail to find a link in containers, it stops processing that class combination.

3. Within the main loop, process the corresponding container of the selector each time by its index.

4. Product Detail Extraction (extract_product_info):

- Attempts to extract the title and price using a predefined list of common CSS selectors and XPath. (Handle most cases) (fast and robust: no extra time needed for api calling)

corner cases: No price is displayed within the price area.

The screenshot displays a Walmart product page for the Sony PlayStation 4 Pro 1TB Gaming Console - Wireless Game Pad - Black. The page layout includes a top navigation bar with categories like Departments, Services, and various delivery options. The main content area features a large product image on the left, a title and rating section in the center, and a detailed 'About this item' section on the right. The 'About this item' section lists specifications such as Processor Type, Memory, and Display. Below this, there is a 'View more' link and a 'At a glance' section with key features like Brand, Resolution, Is portable, Smart tech, Optical drive, and Condition. On the far right, there are sections for 'Walmart Protection Plan by Allstate' and 'More seller options (2)'.

Departments Services | Get it Fast My Items Memorial Day Dinner Solutions Pharmacy Delivery Father's Day Pet Month New Arrivals Auto Service Only at Walmart Registry

Sony
Sony PlayStation 4 Pro 1TB Gaming Console - Wireless Game Pad - Black
★★★★☆ (4.4) | 599 ratings Everyone

About this item

- Product Type Gaming Console
- Processor & Chipset Processor Manufacturer AMD Processor Type Jaguar Processor Core Octa-core (8 Core)
- Memory Standard Memory 8GB Memory Technology GDDR5
- Display & Graphics Controller Manufacturer AMD Graphics Controller Model Radeon Graphics Memory Technology GDDR5 Maximum Resolution 3840 x 2160
- Video Aspect Ratio 16:9 Scan Format 2160p...

[View more](#)

At a glance

Brand Sony	Resolution 4K UHD (2160p)	Is portable N
Smart tech Y	Optical drive Blu-ray Disc Player	Condition New

Free shipping Free 30-day returns
[See more seller options](#)

Walmart Protection Plan by Allstate
What's covered
(Only one option can be selected at a time)
☐ 3-Year Plan - \$59.00
☐ 4-Year Plan - \$75.00

[Add to list](#) [Add to registry](#)

More seller options (2)
Starting from \$594.99 [Compare all sellers](#)

- If the price is not found, doesn't contain a '\$' and digits, or is exactly '\$0.00', it also uses extract_with_deepseek for the price.
- The extract_with_deepseek function sends the HTML content and a targeted prompt (for 'price' or 'title') to the DeepSeek API.

5. Pagination:

- After processing products on a page, it searches for a "next page" button using various common selectors.
- It attempts to click the button and verifies if the URL changes since the selector I defined might unintentionally capture a common "next" button. It also has a mechanism to try alternative next page buttons if the first one fails.

6. Error Handling

Things to optimize:

For different pages, the selector for the product container should remain consistent. Currently, it's being re-selected each time a new page is loaded. I can simplify the logic in the future, which I believe will reduce the execution time.