

# Structured and Hyper Text Model

## Information Retrieval



**Submitted by:**

*2021-CS-46 Muhammad Usman Asghar*

**Submitted to:**

*Dr. Syed Khaldoon Khurshid*

# Contents

<b>1</b>	<b>Structured Guided Model</b>	<b>2</b>
<b>2</b>	<b>Flask Application Code</b>	<b>2</b>
2.1	Function: load_content_map . . . . .	2
2.2	Route: / . . . . .	3
2.3	Route: /content/<section> . . . . .	3
<b>3</b>	<b>Web Scraping Code</b>	<b>3</b>
3.1	Function: get_used_cars_structure . . . . .	3
3.2	Function: scrape_car_data . . . . .	5
<b>4</b>	<b>Saving Data to JSON</b>	<b>5</b>
<b>5</b>	<b>Hyper-Text Model</b>	<b>6</b>
<b>6</b>	<b>Flask Application Code</b>	<b>6</b>
6.1	Route: / . . . . .	6
6.2	Route: /section . . . . .	6
6.3	Function: add_hyperlinks . . . . .	7
<b>7</b>	<b>Web Scraping Code</b>	<b>7</b>
7.1	URL: <a href="https://en.wikipedia.org/wiki/Pakistan">https://en.wikipedia.org/wiki/Pakistan</a> . . . . .	7
7.2	Scraping Content . . . . .	7
7.3	Function: format_content . . . . .	7
7.4	Saving Content to JSON . . . . .	8
<b>8</b>	<b>Sample JSON Structure</b>	<b>8</b>

# 1 Structured Guided Model

This project involves the creation of a web application using Flask, combined with a web scraping component to fetch and display data about cars and bikes from PakWheels. The application dynamically generates and serves content from JSON files.

## 2 Flask Application Code

### 2.1 Function: load\_content\_map

**Purpose:** This function loads data from JSON files and combines them into a content map.

**Code:**

Listing 1: load\_content\_map

```
1 def load_content_map():
2     used_cars_data = []
3     bikes_data = []
4     new_cars_data = []
5
6     # Load used cars data
7     with open('used_cars.json') as file:
8         used_cars_data = json.load(file)
9
10    # Load bikes data
11    with open('bikes.json') as file:
12        bikes_data = json.load(file)
13
14    # Load new cars data
15    with open('new_cars.json') as file:
16        new_cars_data = json.load(file)
17
18    # Combine the data into a single content map
19    content_map = {
20        'used_cars': used_cars_data,
21        'bikes': bikes_data,
22        'new_cars': new_cars_data
23    }
24    return content_map
```

**Explanation:** This function reads three JSON files (used\_cars.json, bikes.json, and new\_cars.json) and combines their data into a dictionary.

## 2.2 Route: /

**Purpose:** Serves the homepage displaying the full content map.

**Code:**

Listing 2: Route for index

```
1 @app.route('/')
2 def index():
3     content_map = load_content_map()
4     return render_template('index.html', content_map=content_map)
```

**Explanation:** This route uses the `load_content_map` function to fetch content and passes it to the `index.html` template for rendering.

## 2.3 Route: /content/<section>

**Purpose:** Serves content for a specific section.

**Code:**

Listing 3: Route for specific section

```
1 @app.route('/content/<section>')
2 def get_content(section):
3     content_map = load_content_map()
4     if section in content_map:
5         return render_template('content.html', section=section, items=
6             content_map[section])
7     else:
8         return jsonify({'error': 'Section not found'}), 404
```

**Explanation:** This route dynamically serves content for a requested section (e.g., `used_cars`, `bikes`, etc.) or returns an error if the section doesn't exist.

# 3 Web Scraping Code

## 3.1 Function: `get_used_cars_structure`

**Purpose:** Scrapes data about used cars from the given URL.

**Code:**

Listing 4: Scraping used cars

```
1 def get_used_cars_structure(url):
2     response = requests.get(url)
3     soup = BeautifulSoup(response.content, 'html.parser')
4
5     cars = []
6     car_elements = soup.find_all('div', class_='cards-content')
7
8     for car in car_elements:
9         title_tag = car.find('h3', class_='nomargin truncate')
10        price_tag = car.find('div', class_='generic-green')
```

```
11     location_tag = car.find('div', class_='generic-gray')
12
13     if title_tag and price_tag and location_tag:
14         title = title_tag.text.strip()
15         price = price_tag.text.strip()
16         location = location_tag.text.strip()
17
18         cars.append({
19             'title': title,
20             'price': price,
21             'location': location
22         })
23
24     return cars
```

**Explanation:** This function sends an HTTP request to the specified URL, parses the response using BeautifulSoup, and extracts data about cars.

## 3.2 Function: `scrape_car_data`

**Purpose:** Scrapes detailed data about new cars, including ratings and reviews.

**Code:**

Listing 5: Scraping new cars

```
1 def scrape_car_data(url):
2     car_data = []
3     response = requests.get(url)
4
5     if response.status_code == 200:
6         soup = BeautifulSoup(response.text, 'html.parser')
7         car_listings = soup.find_all('li', class_='col-md-3')
8
9         for car in car_listings:
10             title = car.find('h3', class_='nomargin truncate').text.strip()
11                 if car.find('h3', class_='nomargin truncate') else "N/A"
12
13             price_element = car.find('div', class_='generic-green truncate fs14')
14             price = price_element.text.strip() if price_element else "N/A"
15
16             rating_element = car.find('span', class_='rating')
17             rating_value = 0
18             if rating_element:
19                 rating = rating_element.find_all('i')
20                 rating_value = sum(1 if 'fa-star' in str(r) else 0 for r
21                                     in rating)
22
23             reviews_count_element = car.find('span', class_='fs14 generic
24                 -gray ml5 dib')
25             reviews_count = reviews_count_element.text.strip() if
26                 reviews_count_element else "N/A"
27
28             car_data.append({
29                 'Title': title,
30                 'Price': price,
31                 'Rating': rating_value,
32                 'Reviews Count': reviews_count
33             })
34
35     return car_data
```

**Explanation:** This function extracts more detailed information from car listings, including reviews and ratings.

## 4 Saving Data to JSON

**Code:**

Listing 6: Saving data to JSON

```

1 with open('used_cars.json', 'w') as f:
2     json.dump(used_cars_data, f, indent=4)
3
4 with open('new_cars.json', 'w') as f:
5     json.dump(new_cars_data, f, indent=4)
6
7 with open('bikes.json', 'w') as f:
8     json.dump(bikes_data, f, indent=4)

```

**Explanation:** The scraped data is saved into separate JSON files for later use in the Flask application.

## 5 Hyper-Text Model

This project implements a hypertext model using Flask and web scraping. The application scrapes structured content from Wikipedia, processes it into sections, and displays it with hyperlinks between related terms.

## 6 Flask Application Code

### 6.1 Route: /

**Purpose:** Serves the homepage with a list of headings and their content.

**Code:**

Listing 7: Home Route

```

1 @app.route('/')
2 def home():
3     # Prepare headings for sidebar and content rendering
4     headings = list(wiki_content.keys())
5     return render_template('home.html', headings=headings, content=
        wiki_content)

```

**Explanation:** This route loads all the section headings and their respective content from the `pakistan.content` file and renders them on the homepage.

### 6.2 Route: /section

**Purpose:** Displays the content of a specific section with clickable links to other sections.

**Code:**

Listing 8: Section Route

```

1 @app.route('/<section>')
2 def show_section(section):
3     content = wiki_content.get(section)
4     if content:
5         linked_content = add_hyperlinks(" ".join(content), wiki_content.
            keys())
6         return render_template('section.html', section=section, content=
            linked_content)
7     return "Section not found", 404

```

**Explanation:** This route dynamically renders the content for a requested section, adding hyperlinks for all referenced headings within the content.

### 6.3 Function: add\_hyperlinks

**Purpose:** Adds hyperlinks to terms in the content that correspond to other headings.

**Code:**

Listing 9: Hyperlinking Function

```
1 def add_hyperlinks(content, headings):
2     for heading in headings:
3         content = re.sub(rf"\b{heading}\b", f'<a href="#{heading}">{
4             heading}</a>', content)
5     return content
```

**Explanation:** The function scans the content for matches with the provided headings and converts these matches into clickable hyperlinks.

## 7 Web Scraping Code

### 7.1 URL: <https://en.wikipedia.org/wiki/Pakistan>

**Purpose:** Scrapes headings and paragraphs from the Wikipedia page on Pakistan.

### 7.2 Scraping Content

**Purpose:** Extracts sections and paragraphs from the Wikipedia page.

**Code:**

Listing 10: Scraping Content

```
1 content = {}
2 current_section = None
3
4 for element in soup.find_all(['h2', 'h3', 'p']):
5     if element.name in ['h2', 'h3']:
6         current_section = element.text.strip()
7         content[current_section] = []
8     elif element.name == 'p' and current_section:
9         paragraph = ''.join(format_content(child) for child in element.
10                             children)
11         content[current_section].append(paragraph)
```

**Explanation:** The code parses the Wikipedia page for headings and their corresponding paragraphs, organizing them into sections.

### 7.3 Function: format\_content

**Purpose:** Formats text with hyperlinks and citations.

**Code:**



Listing 11: Formatting Content

```

1 def format_content(element):
2     if element.name == 'a':
3         href = element.get('href')
4         if href and href.startswith('/'):
5             href = base_url + href
6         link_text = element.text
7         return f'<a href="{href}" target="_blank">{link_text}</a>'
8     elif element.name == 'sup':
9         return f'<sup>{element.text}</sup>'
10    return element.text

```

**Explanation:** The function identifies anchor tags and superscripts in the HTML, converts them into clickable links, and retains citations.

## 7.4 Saving Content to JSON

**Purpose:** Saves the structured content into a JSON file.

**Code:**

Listing 12: Saving JSON

```

1 with open('pakistan_content.json', 'w') as f:
2     json.dump(content, f, indent=4)

```

**Explanation:** This code saves the scraped and structured content into `pakistan_content.json` for use in the Flask application.

## 8 Sample JSON Structure

**Example:**

Listing 13: Example JSON Content

```

1 {
2     "History": [
3         "Pakistan gained independence in 1947.",
4         "It was created as a separate state for Muslims."
5     ],
6     "Geography": [
7         "Pakistan is located in South Asia.",
8         "It shares borders with India, Afghanistan, and Iran."
9     ]
10 }

```

**Explanation:** Each section contains a list of paragraphs, with headings as keys.