

Coding Sample

Bishmay Barik

January 1, 2026

1 Project Overview: `lipi-swap`

This Python script was developed to scrape data from an Indian administrative database, where names of entities like Panchayats and Blocks are often written in Hindi (Devanagari script). The script uses web scraping techniques to collect data and leverages the `unidecode` library to transliterate Hindi names into English, ensuring compatibility when saved to a CSV file.

1.1 Purpose and Functionality

This script is built to scrape hierarchical data from Indian datasets (e.g., states, districts, blocks, and panchayats) using `requests` and `BeautifulSoup`. It handles non-Roman scripts like Devanagari by converting them to English equivalents with `unidecode`. The scraped data, including monthly figures, is organized into a structured CSV file using `pandas`. English names remain unchanged, while Hindi names are transliterated for readability and compatibility.

1.2 Why Transliteration?

When dealing with datasets containing Devanagari script (e.g., Hindi names), saving them directly to a CSV can result in encoding issues or unreadable characters on systems not configured for non-Roman scripts. Transliteration solves this by converting Hindi characters to their closest English (Roman) equivalents, making the data portable and usable across different platforms.

1.3 Installation

To run `lipi-swap`, you'll need the following Python libraries:

- `requests` - For making HTTP requests
- `beautifulsoup4` - For parsing HTML content
- `pandas` - For handling data and CSV output
- `unidecode` - For transliterating Hindi to English

Install them using pip:

```
pip install requests beautifulsoup4 pandas unidecode
```

2 Refactored Code Following Rigorous Standards

I have refactored this code to adhere to rigorous coding standards, including:

- **Extensive commenting:** One comment for every line of code (minimum 1:1 ratio)

- **Proper naming conventions:** Using `df_` prefix for DataFrames, `_list` suffix for lists, `_path` suffix for file paths
- **Clear function documentation:** Comprehensive docstrings with Args, Returns, and Example sections
- **Verbose parameter:** Added to control debug output
- **Clear error messages:** Informative messages for debugging
- **Interpretability over brevity:** Code prioritizes clarity and maintainability

2.1 lipi-swap.py (Refactored)

```

1 """
2 lipi-swap: Web scraper for Indian administrative database with Hindi
3   transliteration
4
5 This script scrapes hierarchical data (states, districts, blocks, panchayats)
6   from
7 Indian administrative databases and transliterates Devanagari (Hindi) names to
8   Roman
9 script for CSV compatibility.
10 """
11
12
13 import requests
14 from bs4 import BeautifulSoup
15 import pandas as pd
16 from urllib.parse import urljoin, urlparse, parse_qs
17 import time
18 from unidecode import unidecode
19
20 # define base URL for constructing full URLs
21 BASE_URL_PATH = "put/your/base/url/here"
22
23
24 def get_soup(url, verbose=True):
25     """
26         Fetch and parse a webpage into a BeautifulSoup object.
27
28     Args:
29         url (str): The URL to fetch
30         verbose (bool): If True, print status messages. Default is True.
31
32     Returns:
33         BeautifulSoup: Parsed HTML content, or None if request fails
34
35     Example:
36         >>> soup = get_soup("https://example.com")
37         >>> soup is not None
38         True
39     """
40
41     try:
42         # send GET request to the URL
43         response = requests.get(url)
44
45         # raise exception if status code indicates error

```

```
46     response.raise_for_status()
47
48     # add delay to prevent server overload
49     time.sleep(1)
50
51     # parse response text into BeautifulSoup object
52     return BeautifulSoup(response.text, 'html.parser')
53
54 except requests.RequestException as e:
55     # print error message if request fails
56     if verbose:
57         print(f"Error fetching {url}: {e}")
58     return None
59
60
61 def extract_links(soup):
62     """
63     Extract links and names from a table on the webpage.
64
65     Args:
66         soup (BeautifulSoup): Parsed HTML content
67
68     Returns:
69         list: List of tuples containing (name, url) pairs. Names are
70             transliterated
71             from Hindi to English using unidecode.
72
73     Example:
74         >>> from bs4 import BeautifulSoup
75         >>> html = '<table id="t1"><tr><td></td><td><a href="/test">
76             </a></td></tr></table>'
77         >>> soup = BeautifulSoup(html, 'html.parser')
78         >>> links = extract_links(soup)
79         >>> len(links) > 0
80             False
81
82     """
83
84     # return empty list if soup is None
85     if soup is None:
86         return []
87
88     # find the table with id 't1'
89     table = soup.find('table', id='t1')
90
91     # return empty list if table not found
92     if not table:
93         return []
94
95     # get all rows except first 4 header rows
96     rows = table.find_all('tr')[4:]
97
98     # initialize empty list for storing links
99     link_list = []
100
101    # loop over each row in the table
102    for row in rows:
103        # extract all cells from the row
104        cells = row.find_all('td')
105
106        # check if row has at least 2 cells
107        if len(cells) >= 2:
108            # find the link tag in the second cell
109            link_tag = cells[1].find('a')
```

```
107     # process link if found
108     if link_tag:
109         # extract and transliterate name from Hindi to English
110         name = unidecode(link_tag.text.strip())
111
112         # get href attribute
113         href = link_tag['href']
114
115         # construct full URL from base and relative URL
116         full_url = urljoin(BASE_URL_PATH, href)
117
118         # append tuple of name and URL to list
119         link_list.append((name, full_url))
120
121     # return list of extracted links
122     return link_list
123
124
125 def scrape_table_data(soup, state, state_id, fy, district, block):
126     """
127     Scrape panchayat-level table data from a block webpage.
128
129     Args:
130         soup (BeautifulSoup): Parsed HTML content
131         state (str): State name
132         state_id (str): State code/ID
133         fy (str): Financial year
134         district (str): District name
135         block (str): Block name
136
137     Returns:
138         list: List of dictionaries containing panchayat data with monthly
139             figures
140
141     Example:
142         >>> soup = None
143         >>> data = scrape_table_data(soup, "State", "01", "2023", "District", " "
144             Block")
145         >>> len(data)
146         0
147     """
148
149     # return empty list if soup is None
150     if soup is None:
151         return []
152
153     # find the table with id 't1'
154     table = soup.find('table', id='t1')
155
156     # return empty list if table not found
157     if not table:
158         return []
159
160     # get all rows except first 4 header rows
161     rows = table.find_all('tr')[4:]
162
163     # initialize empty list for storing data
164     data_list = []
165
166     # loop over each row in the table
167     for row in rows:
168         # extract all cells from the row
169         cells = row.find_all('td')
```

```

168     # check if row has at least 14 cells (name + 12 months)
169     if len(cells) >= 14:
170         # extract and transliterate panchayat name from Hindi to English
171         panchayat = unidecode(cells[1].text.strip())
172
173         # extract monthly data from columns 2-13 (12 months)
174         monthly_data_list = [cell.text.strip() for cell in cells[2:14]]
175
176         # create dictionary with all data for this panchayat
177         data_list.append({
178             'State': state,
179             'State ID': state_id,
180             'FY': fy,
181             'District': district,
182             'Block': block,
183             'Panchayat': panchayat,
184             'April': monthly_data_list[0],
185             'May': monthly_data_list[1],
186             'June': monthly_data_list[2],
187             'July': monthly_data_list[3],
188             'August': monthly_data_list[4],
189             'September': monthly_data_list[5],
190             'October': monthly_data_list[6],
191             'November': monthly_data_list[7],
192             'December': monthly_data_list[8],
193             'January': monthly_data_list[9],
194             'February': monthly_data_list[10],
195             'March': monthly_data_list[11]
196         })
197
198     # return list of scraped panchayat data
199     return data_list
200
201
202 def main(verbose=True):
203     """
204     Main function to orchestrate the hierarchical scraping of administrative
205     data.
206
207     Scraps data at multiple levels: states -> districts -> blocks ->
208     panchayats,
209     and saves all collected data to a CSV file.
210
211     Args:
212         verbose (bool): If True, print status messages during scraping. Default
213             is True.
214
215     Returns:
216         None. Saves output to CSV file.
217
218     Example:
219         >>> # This would run the full scraping process
220         >>> # main(verbose=False)
221     """
222
223     # initialize empty list to store all scraped data
224     all_data_list = []
225
226
227     # loop over each state URL in the list
228     for state_url in state_url_list:
229         # parse URL to extract query parameters
230         parsed_url = urlparse(state_url)
231
232         # extract query parameters from URL

```

```
228     query_params = parse_qs(parsed_url.query)
229
230     # get state name from query parameters (default to 'Unknown')
231     state = query_params.get('state_name', ['Unknown'])[0]
232
233     # get state ID from query parameters (default to 'Unknown')
234     state_id = query_params.get('state_code', ['Unknown'])[0]
235
236     # get financial year from query parameters (default to 'Unknown')
237     fy = query_params.get('fin_year', ['Unknown'])[0]
238
239     # print status message for current state
240     if verbose:
241         print(f"Processing state: {state}")
242
243     # fetch and parse state page
244     state_soup = get_soup(state_url, verbose=verbose)
245
246     # skip to next state if page fetch failed
247     if not state_soup:
248         if verbose:
249             print(f"Failed to fetch state page for {state}. Skipping.")
250         continue
251
252     # extract district links from state page
253     district_link_list = extract_links(state_soup)
254
255     # print number of districts found
256     if verbose:
257         print(f"Found {len(district_link_list)} districts in {state}.")
258
259     # loop over each district in the state
260     for district_name, district_url in district_link_list:
261         # print status message for current district
262         if verbose:
263             print(f"Processing district: {district_name} in {state}")
264
265         # fetch and parse district page
266         district_soup = get_soup(district_url, verbose=verbose)
267
268         # skip to next district if page fetch failed
269         if not district_soup:
270             if verbose:
271                 print(f"Skipping {district_name} due to fetch error.")
272             continue
273
274         # extract block links from district page
275         block_link_list = extract_links(district_soup)
276
277         # print number of blocks found
278         if verbose:
279             print(f"Found {len(block_link_list)} blocks in {district_name}.")
280
281         # loop over each block in the district
282         for block_name, block_url in block_link_list:
283             # print status message for current block
284             if verbose:
285                 print(f"Processing block: {block_name} in {district_name}, {state}")
286
287             # fetch and parse block page
288             block_soup = get_soup(block_url, verbose=verbose)
```

```
289     # skip to next block if page fetch failed
290     if not block_soup:
291         if verbose:
292             print(f"Skipping {block_name} due to fetch error.")
293         continue
294
295     # scrape panchayat data from block page
296     panchayat_data_list = scrape_table_data(
297         block_soup, state, state_id, fy, district_name, block_name
298     )
299
300     # add panchayat data to overall data list
301     all_data_list.extend(panchayat_data_list)
302
303     # print number of panchayats collected
304     if verbose:
305         print(f"Collected data for {len(panchayat_data_list)}"
306               f" panchayats in {block_name}.")
307
308     # check if any data was collected
309     if all_data_list:
310         # convert list of dictionaries to DataFrame
311         df_panchayat = pd.DataFrame(all_data_list)
312
313         # define output file path
314         output_path = "put/your/output/path/here.csv"
315
316         # save DataFrame to CSV file
317         df_panchayat.to_csv(output_path, index=False, encoding='utf-8')
318
319         # print success message with row count
320         if verbose:
321             print(f"Data saved to {output_path} with {len(df_panchayat)} rows."
322         )
323     else:
324         # print message if no data was collected
325         if verbose:
326             print("No data collected.")
327
328 if __name__ == "__main__":
329     # run main function when script is executed
330     main()
```