

Jan Kaczmarek - Algorytmy Tekstowe

Laboratorium 1

parse_publications.py

```
58 import re
59 from typing import Optional
60
61 def parse_publication(reference: str) -> Optional[dict]:
62     """
63     Parse academic publication reference and extract structured information.
64
65     Expected reference format:
66     Lastname, I., Lastname2, I2. (Year). Title. Journal, Volume(Issue), StartPage-EndPage.
67
68     Example:
69     Kowalski, J., Nowak, A. (2023). Analiza algorytmów tekstowych. Journal of Computer Science, 45(2), 123-145.
70
71     Args:
72     reference (str): Publication reference string
73
74     Returns:
75     Optional[dict]: A dictionary containing parsed publication data or None if the reference doesn't match expected format
76     """
77     authors_year_pattern = r"(?P<authors>(?:[A-ZÀÇÈÌÑÓŠŽŽ][a-zàçèìñóšžž]+, [A-ZÀÇÈÌÑÓŠŽŽ]\.(?:, )?) \((?P<year>\d{4})\)"
78     title_journal_pattern = (
79         r"\. (?P<title>.+?)\.( (?P<journal>[A-ZÀÇÈÌÑÓŠŽŽa-zàçèìñóšžž\s]+), "
80     )
81     volume_issue_pages_pattern = (
82         r"(?P<volume>\d+)\((?P<issue>\d+)?\)?, (?P<start_page>\d+)-(?P<end_page>\d+)\. "
83     )
84     full_pattern = (
85         authors_year_pattern + title_journal_pattern + volume_issue_pages_pattern
86     )
87
88     match = re.match(full_pattern, reference)
89     if not match:
90         return None
91
92     authors_raw = match.group("authors")
93     authors_list = []
94     author_pattern = r"([A-ZÀÇÈÌÑÓŠŽŽ][a-zàçèìñóšžž]+), ([A-ZÀÇÈÌÑÓŠŽŽ])\."
95     for author_match in re.finditer(author_pattern, authors_raw):
96         authors_list.append(
97             {"last_name": author_match.group(1), "initial": author_match.group(2)}
98         )
99
100     return {
101         "authors": authors_list,
102         "year": int(match.group("year")),
103         "title": match.group("title"),
104         "journal": match.group("journal"),
105         "volume": int(match.group("volume")),
106         "issue": int(match.group("issue")) if match.group("issue") else None,
107         "pages": {
108             "start": int(match.group("start_page")),
109             "end": int(match.group("end_page")),
110         },
111     }
```

Do testowania czy regex jest poprawny korzystałem ze strony <https://regex101.com/>. Moim pomysłem jest podzielenie jednego dużego regexa na 3 mniejsze części, tak jak zostało to zaproponowane w szkielecie zadania. Głównym zabiegiem które ułatwiło mi rozwiązanie tego zadania jest użycie **named capturing groups** które ułatwiają późniejszy dostęp do konkretnych pól w sparsowanym regex'ie. Dzięki zastosowaniu tego mechnizmu wykorzystuję `re.match.group` co można zobaczyć w kodzie źródłowym.

Wyniki testów

```

❯ **$ pytest tests/test_parse_publication.py
===== test session starts =====
platform darwin -- Python 3.12.9, pytest-8.3.5, pluggy-1.5.0
rootdir: /Users/jkaczmarwski/private/text-algorithms-course/python-labs/lab_1
plugins: anyio-4.9.0
collected 4 items

tests/test_parse_publication.py .... [100%]

===== 4 passed in 0.01s =====

```

extract_links.py

```

42 import re
43
44
45
46
47
48 def extract_links(html: str) -> list[dict[str, str]]:
49     """
50     Extract all links from the given HTML string.
51
52     Args:
53         html (str): HTML content to analyze
54
55     Returns:
56         list[dict]: A list of dictionaries where each dictionary contains:
57             - 'url': the href attribute value
58             - 'title': the title attribute value (or None if not present)
59             - 'text': the text between <a> and </a> tags
60     """
61
62     # TODO Implement a regular expression pattern to extract links from HTML.
63     # The pattern should capture three groups:
64     # 1. The URL (href attribute value)
65     # 2. The title attribute (which might not exist)
66     # 3. The link text (content between <a> and </a> tags)
67     pattern = r'<a\s+[>]*?href=["\'](?:P<url>[^\']*+)[^\']*+)<?>(?<s+title=["\'](?:P<title>[^\']*+)[^\']*+)<?>(?<text>.*?</a>'
```

Regex w tym przypadku jest stosunkowo prosty. Znowu bardzo pomaga tu użycie name capturing groups jak w przypadku zadania `parse_publications`. Jest to najważniejszy aspekt tego rozwiązania. Z rozwiązań zastosowanych w tym regexie wymienię 2 jako “ciekawsze”, są to:

- `[^>]` matchowanie jednego znaku który NIE należy do danego zbioru (w tym przypadku, wszystkie znaki oprócz `>`)
- `\s+` matchowanie jednego lub więcej białych znaków. Sprawia to, że regex jest "czystszy"

Wyniki testów

```
▶(venv) jkaczmariski@JAKACZMA-M-RT4J:python-labs/lab_1 (k8s: db-operations) 'main*$ pytest tests/test_extract_links.py
===== test session starts =====
platform darwin -- Python 3.12.9, pytest-8.3.5, pluggy-1.5.0
rootdir: /Users/jkaczmariski/private/text-algorithms-course/python-labs/lab_1
plugins: anyio-4.9.0
collected 3 items

tests/test_extract_links.py ... [100%]

===== 3 passed in 0.01s =====
```

analzye_test_files.py

```
13
12 # TODO Implement word extraction using regex
11 # Find all words in the content (lowercase for consistency)
10 words = re.findall(r"\b\w+\b", content)
9 word_count = len(words)
8
7 # TODO Implement sentence splitting using regex
6 # A sentence typically ends with ., !, or ? followed by a space
5 # Be careful about abbreviations (e.g., "Dr.", "U.S.A.")
4
3 sentence_pattern = r"s*([^.!?]*?(?:Dr|Mr|Mrs|Ms|Prof|Jr|Sr|U\.S\.A|etc|e\.g|i\.e)\.|\s*([^.!?]+)[.!?]"
2 sentences = [m.group().strip() for m in re.finditer(sentence_pattern, content)]
1 sentence_count = len([s for s in sentences if s.strip()])

73
1 # TODO Implement email extraction using regex
2 # Extract all valid email addresses from the content
3 email_pattern = r"[A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Za-z]{2,}"
4 emails = re.findall(email_pattern, content)
5
6 # TODO Calculate word frequencies
7 # Count occurrences of each word, excluding stop words and short words
8 # Use the Counter class from collections
9 frequent_words = Counter([word for word in words if len(word) > 1])
10 for word in stop_words:
11     if word in frequent_words:
12         del frequent_words[word]
13 frequent_words = dict(frequent_words.most_common(10))
14
15 # TODO Implement date extraction with multiple formats
16 # Detect dates in various formats: YYYY-MM-DD, DD.MM.YYYY, MM/DD/YYYY, etc.
17 # Create multiple regex patterns for different date formats
18 date_patterns = [
19     r"\b\d{4}-\d{2}-\d{2}\b", # YYYY-MM-DD
20     r"\b\d{2}\.\d{2}\.\d{4}\b", # DD.MM.YYYY
21     r"\b\d{2}/\d{2}/\d{4}\b", # MM/DD/YYYY or DD/MM/YYYY
22     r"\b\d{2}-\d{2}-\d{4}\b", # MM-DD-YYYY or DD-MM-YYYY
23     r"(?<!\d)\b\d{1,2} [A-Za-z]+ \d{4}\b", # 12 March 2025 (Avoids matching "23")
24     r"\b[A-Za-z]+ \d{1,2}, \d{4}\b", # March 12, 2025
25 ]
26 date_regex = re.compile("|".join(date_patterns))
27 dates = date_regex.findall(content)
28
29 # TODO Analyze paragraphs
30 # Split the content into paragraphs and count words in each
31 # Paragraphs are typically separated by one or more blank lines
32 paragraphs = re.split(r"\n\s*\n", content)
33 paragraph_sizes = {
34     i + 1: len(re.findall(r"\b\w+\b", para)) for i, para in enumerate(paragraphs)
35 }
36
```

W moim rozwiązaniu wykorzystuję proste regexy. Ciekawe rozwiązanie zastosowałem w przypadku `sentence_pattern` w celu wykluczenia skrótów, potrzebuję bazy danych

skrótów, a przynajmniej nie jestem w stanie znaleźć poprawnego warunku, na odróżnienie zdania od skrótu. Też wykorzystałem nieoczywiste rozwiązanie w przypadku wyboru słów. Nie używam zwykłego `split` z powodu specyfikacji zadania. Na przykład przy użyciu `.split` słowa “Analyze” i “Analyze,” co powodowało problemy przy kolejnych zdaniach. Dlatego zdecydowałem się na regex.

Wyniki testów

```
▶ pytest tests/test_analyze_text_file.py
===== test session starts =====
platform darwin -- Python 3.12.9, pytest-8.3.5, pluggy-1.5.0
rootdir: /Users/jkaczmarek/private/text-algorithms-course/python-labs/lab_1
plugins: anyio-4.9.0
collected 8 items

tests/test_analyze_text_file.py ..... [100%]

===== 8 passed in 0.01s =====
```