

Dokumentacja: Przetwarzanie języka naturalnego w systemach sztucznej inteligencji - Projekt I

Autor: Aleksandra Ignacyk

1 Przegląd

Celem tego projektu jest przetworzenie i analiza danych tekstowych pozyskanych z artykułów Wikipedii. W szczególności projekt obejmuje następujące zadania:

1. Pobieranie danych tekstowych z artykułów Wikipedii,
2. Przetwarzanie danych (czyszczenie i normalizacja tekstu),
3. Analiza rozkładu częstotliwości słów, w tym wizualizacja przy użyciu Prawa Zipfa,
4. Budowanie grafów zależności między słowami i analiza najważniejszych węzłów.

W niniejszej dokumentacji omówiono poszczególne kroki wraz z wyjaśnieniem kluczowych funkcji i fragmentów kodu.

2 Importowanie bibliotek

Pierwszym krokiem jest zaimportowanie bibliotek, które są niezbędne do wykonania wszystkich zadań związanych z przetwarzaniem tekstu i analizą danych.

- `requests` - służy do pobierania danych z URL,
- `BeautifulSoup` (z `bs4`) - do parsowania kodu HTML i wyodrębniania treści z artykułów,
- `Counter` (z `collections`) - do zliczania częstotliwości występowania słów,
- `matplotlib.pyplot` i `numpy` - do wizualizacji danych,
- `re` - do czyszczenia tekstu (usuwanie liczb i znaków specjalnych),
- `json` - do przetwarzania odpowiedzi z API w formacie JSON,
- `networkx` - do budowania grafów i analizy zależności między słowami.

Kod importujący te biblioteki:

```
import requests
from bs4 import BeautifulSoup
from collections import Counter
import matplotlib.pyplot as plt
import numpy as np
import re
import json
import networkx as nx
```

3 Pobieranie treści z URL

Aby pobrać treść artykułu, implementujemy funkcję `get_page_content`, która:

- Wysyła żądanie HTTP do podanego URL,
- Wyodrębnia treści z wszystkich paragrafów `<p>` w kodzie HTML przy pomocy `BeautifulSoup`,
- Zwraca zebrany tekst lub pusty ciąg w przypadku niepowodzenia żądania.

Kod funkcji:

```
def get_page_content(url):
    response = requests.get(url)
    if response.status_code == 200:
        soup = BeautifulSoup(response.content, 'html.parser')
        paragraphs = soup.find_all('p')
        text = ' '.join([p.get_text() for p in paragraphs])
        return text
    return ""
```

4 Przetwarzanie tekstu

Kolejnym krokiem jest przetworzenie pobranego tekstu, aby był gotowy do analizy. Funkcja `preprocess_text` wykonuje następujące operacje:

- Usuwa liczby za pomocą wyrażeń regularnych,
- Usuwa znaki specjalne, pozostawiając jedynie litery i spacje,
- Zamienia tekst na małe litery i dzieli go na listę słów.

Kod funkcji:

```
def preprocess_text(text):
    text = re.sub(r'\d+', '', text) # Usunięcie liczb
    text = re.sub(r'[^a-zA-Z\s]', '', text) # Usunięcie znaków specjalnych
    return text.lower().split()
```

5 Pobieranie artykułów z API Wikipedii

Korzystając z API Wikipedii, pobieramy do 100 artykułów. Proces obejmuje:

- Wybór języka wyszukiwania (w tym przypadku polski, czyli `pl`),
- Wysyłanie żądania do API Wikipedii, przetwarzanie odpowiedzi w formacie JSON,
- Zapisanie linków do pełnych wersji artykułów w liście `urls`.

Kod:

```
language_code = 'pl'
search_query = 'the'
number_of_results = 100
headers = {
    'User-Agent': 'YOUR_APP_NAME (YOUR_EMAIL_OR_CONTACT_PAGE)'
}

base_url = 'https://api.wikimedia.org/core/v1/wikipedia/'
endpoint = '/search/page'
url = base_url + language_code + endpoint
```

```

parameters = {'q': search_query, 'limit': number_of_results}

response = requests.get(url, headers=headers, params=parameters)
response = json.loads(response.text)

urls = []
for page in response['pages']:
    article_url = 'https://' + language_code + '.wikipedia.org/wiki/' +
        page['key']
    urls.append(article_url)

print(urls)

```

6 Analiza czestotliwości słów

Na podstawie pierwszego artykułu pobranego z listy `urls`, przeprowadzamy analize rozkładu słów:

- Pobieramy treść artykułu i przetwarzamy tekst na liste słów,
- Zliczamy wystąpienia słów przy użyciu `Counter`,
- Tworzymy wykres rozkładu czestotliwości słów zgodnie z Prawem Zipfa.

Kod:

```

url = urls[0]
text = get_page_content(url)

if text:
    words = preprocess_text(text)
    word_count = Counter(words)
    sorted_words = sorted(word_count.items(), key=lambda x: x[1], reverse=
        True)

# Prawo Zipfa – przygotowanie danych
ranks = np.arange(1, len(sorted_words) + 1)
freq = np.array([count for _, count in sorted_words])

# Wykres log-log (prawo Zipfa)
plt.loglog(ranks, freq, marker="o")
plt.title("Prawo Zipfa – zależność rangi a czestotliwosci")
plt.xlabel("Ranga")
plt.ylabel("Czestotliwosc")
plt.grid(True)
plt.show()

```

7 Tworzenie grafu słów

Na końcu budujemy graf zależności między słowami na podstawie ich sąsiedztwa w tekście. Tworzymy graf, w którym krawędź występuje, jeśli dwa słowa są sąsiadujące w artykule.

- Dodajemy krawędzie między sąsiadującymi słowami w każdym artykule,
- Wyodrębniamy 10% węzłów o największej liczbie sąsiadów i wizualizujemy powstały graf.

Kod:

```

G = nx.Graph()
G.add_nodes_from(all_words)

# Dodanie krawędzi między słowami występującymi obok siebie

```

```

for url in urls:
    text = get_page_content(url)
    if text:
        words = preprocess_text(text)
        for i in range(len(words) - 1):
            word1 = words[i]
            word2 = words[i + 1]
            G.add_edge(word1, word2)

# Wyodrębnienie 10% w z w o najwi kszej liczbie s siad w
num_nodes_to_extract = int(0.1 * len(G.nodes()))
sorted_nodes_by_degree = sorted(G.degree, key=lambda x: x[1], reverse=True)
top_nodes = [node for node, degree in sorted_nodes_by_degree[:
    num_nodes_to_extract]]

# Tworzenie podgrafu
H = G.subgraph(top_nodes)

# Wizualizacja grafu
plt.figure(figsize=(12, 12))
pos = nx.spring_layout(H, k=0.5)
nx.draw(H, pos, with_labels=True, node_size=1000, node_color="lightblue",
    font_size=8, font_weight="bold")
plt.title("Graf 10% w z w z najwi ksz liczb s siad w (100
    artyku w)")
plt.show()

```

8 Wnioski

Projekt przedstawia sposób pobierania, przetwarzania i analizy tekstu z artykułów Wikipedii. Zastosowano Prawo Zipfa do zrozumienia rozkładu słów oraz narzędzia grafowe do wizualizacji zależności między słowami. Możliwości analizy tekstu w naturalnych językach pozwalają na odkrywanie interesujących zależności, co ma szerokie zastosowanie w systemach sztucznej inteligencji.

9 Link do Githuba

https://github.com/olaignacyk/przetwarzanie_jzyka/blob/main/zadanie1-3.ipynb