

Narzędzie do znajdowania podobnych pytań na Stack Overflow

Mizak Mateusz, Mazur Szymon

1. Wstęp

Celem napisanego przez nas narzędzia jest odnajdywanie pytań, podobnych do pytania zadanego przez użytkownika, znajdujących się na Stack Overflow. Po podaniu przez użytkownika konkretnego pytania, program daje propozycje najbardziej podobnych do niego pytań już zadanych na Stack Overflow.

2. Rozwinięcie

Z powodów problemów związanych ze zbiorem danych udostępnionych na stronie www.kaggle.com postanowiliśmy sami stworzyć prosty skrypt, który umożliwiłby nam pobranie bazy pytań Stack Overflow za pomocą otwartego API stackexchange.

```

2 import requests
3
4 api_url = "https://api.stackexchange.com/2.2/questions"
5
6 params = {
7     "order": "desc",
8     "sort": "votes",
9     "site": "stackoverflow",
10    "tagged": "python",
11    "pagesize": 100
12 }
13 num_questions = 2500
14 pagesize = 100
15 num_pages = num_questions // pagesize + (num_questions % pagesize > 0)
16 questions = []
17 for page in range(1, num_pages + 1):
18     params["page"] = page
19     print(page)
20     response = requests.get(api_url, params=params)
21     if response.status_code == 200:
22         page_questions = response.json()["items"]
23         questions += page_questions
24     else:
25         print("Error retrieving questions:", response.status_code)
26         break
27
28 with open("questions.txt", "a", encoding='utf-8') as f:
29     i = 0
30     for question in questions:
31         f.write(str(i) + ". " + question["title"] + "\n")
32     i += 1
33

```

Zbiór danych pobrany w taki sposób nie jest tak duży jak ten znajdujący się w źródle podanym przez prowadzącego. Jest to spowodowane ograniczeniami nałożonymi przez stackexchange na ilość zapytań, które API może otrzymać od pojedynczego użytkownika. Metodą na ominięcie tej blokady było by zarejestrowanie swojej aplikacji. Mimo stosunkowo małego zbioru danych jest on wystarczający, aby zaprezentować działanie naszego narzędzia.

Przykład jak wygląda wygenerowany przez nas zbiór danych:

```
1634. How to check if a float value is a whole number
1635. How do I print the model summary in PyTorch?
1636. Convert tuple to list and back
1637. Python type hinting without cyclic imports
1638. SQLAlchemy default DateTime
1639. Checking if a string can be converted to float in Python
1640. Calculating arithmetic mean (one type of average) in Python
1641. Fastest way to get the first object from a queryset in django?
1642. How to filter a dictionary according to an arbitrary condition function?
1643. Pandas percentage of total with groupby
1644. How to create a temporary directory and get its path/ file name?
1645. Does Python optimize tail recursion?
1646. multiprocessing vs multithreading vs asyncio
1647. Using property() on classmethods
1648. Pandas DataFrame Groupby two columns and get counts
1649. UnboundLocalError on local variable when reassigned after first use
```

Główny program najpierw pobiera już przygotowany zbiór danych a następnie oczekuje na wprowadzenie przez użytkownika pytania w celu sprawdzenia czy takie pytanie już istnieje lub istnieją do niego podobne również pod względem semantycznym.

Po otrzymaniu pytania program przygotowuje tablice krotek które zawierają w sobie już istniejące pytanie jak i wskaźnik podobieństwa wyrażony w postaci liczby między 0 a 1.

W celu dobrania wskaźnika podobieństw do każdego pytania wykorzystaliśmy napisaną przez nas funkcję:

```

8 def similarity(string1, string2):
9     string1_synsets = [get_synsets(*word) for word in pos_tag(word_tokenize(string1)) if get_synsets(*word)]
10    string2_synsets = [get_synsets(*word) for word in pos_tag(word_tokenize(string2)) if get_synsets(*word)]
11
12    similarity_points = 0
13    for set1 in string1_synsets:
14        max_points = 0
15        for set2 in string2_synsets:
16            points = set1.path_similarity(set2)
17            if points > max_points:
18                max_points = points
19        similarity_points += max_points
20
21    if len(string1_synsets) == 0:
22        return 0
23
24    return similarity_points / len(string1_synsets)
25
26
27 def get_synsets(word, tag):
28     wordnet_tag = None
29
30     if tag.startswith('N'):
31         wordnet_tag = 'n'
32     elif tag.startswith('V'):
33         wordnet_tag = 'v'
34     elif tag.startswith('J'):
35         wordnet_tag = 'a'
36     elif tag.startswith('R'):
37         wordnet_tag = 'r'
38     else:
39         return None
40
41     try:
42         return wordnet.synsets(word, wordnet_tag)[0]
43     except:
44         return None

```

Funkcja ta wykorzystuje wordnet czyli bazę słownikową języka angielskiego stworzoną przez Uniwersytet Princeton. Stosujemy do szukania podobieństwa semantycznego poszczególnych wyrazów zawartych w pytaniu.

Przy użyciu pierwszych linijek tworzymy dwie tablice zawierające już przetworzone słowa z pytania użytkownika oraz pytania z zestawu danych z pominięciem słów, które nie wpływają na znaczenie semantyczne pytania (np. “the”, “a”).

```

8 def similarity(string1, string2):
9     string1_synsets = [get_synsets(*word) for word in pos_tag(word_tokenize(string1)) if get_synsets(*word)]
10    string2_synsets = [get_synsets(*word) for word in pos_tag(word_tokenize(string2)) if get_synsets(*word)]

```

Po przygotowaniu takich tablic za pomocą funkcji path_similarity() porównujemy ze sobą po kolei każdy wyraz z pytania użytkownika z każdym wyrazem pytania z zestawu i szukamy wyniku najwyższego po czym dodajemy wynik każdego wyrazu i uśredniamy.

```

12     similarity_points = 0
13     for set1 in string1_synsets:
14         max_points = 0
15         for set2 in string2_synsets:
16             points = set1.path_similarity(set2)
17             if points > max_points:
18                 max_points = points
19             similarity_points += max_points
20
21     if len(string1_synsets) == 0:
22         return 0
23
24     return similarity_points / len(string1_synsets)

```

Mając już tablice krotek z pytaniami oraz ich wskaźnikiem podobieństwa sortujemy ją od wartości najwyższych do najniższych. Następnie bierzemy 4 najwyższe wyniki o ile posiadają wskaźnik podobieństwa na poziomie przynajmniej 0.75 w celu otrzymania tylko faktycznie podobnych pytań i zwracamy te pytania w formie listy.

```

55     similarity_list.sort(key=lambda x: x[1], reverse=True)
56
57     for j in range(0, 4):
58         if similarity_list[j][1] > 0.75:
59             result_questions.append(similarity_list[j][0])

```

Posiadając już tak przygotowaną tablice, jesteśmy w stanie wyświetlić te najlepsze wyniki dla użytkownika.

3. Podsumowanie i sprawdzenie działania

Poniżej przykład przy zadaniu pytania “How to use ternary operator in Python?”

```

Enter question that you want to ask and i will show you if there are already similar questions on stack overflow: how to use ternary operator in Python?
List of similar questions:
1) Does Python have a ternary conditional operator?

```

Jak widać, narzędzie to odnalazło pytanie o bardzo podobnym znaczeniu, mimo że zostało ono zapisane w widocznie różny sposób (Dodatkowo

pojawiła się w pytaniu literówka która udowodnia, że program radzi sobie nawet z drobnymi literówkami). Wypisane zostało tylko jedno proponowane pytanie co oznacza, że w naszym zbiorze tylko to jedno pytanie miało wskaźnik podobieństwa na wymaganym poziomie.

Kolejny przykład tym razem ze zmianą języka programowania, którego dotyczy pytanie z Python na Java:

```
Enter question that you want to ask and i will show you if there are already similar questions on stack overflow:how to use ternary operator in Java
I can't find any questions similar to yours
```

Na tym przykładzie widać, że mimo iż pytanie było sformułowane w bardzo podobny sposób co poprzednie to różnica wywołana zmianą nazwy języka programowania była wystarczająco duża, aby program nie zaproponował ponownie pytania co w poprzednim przykładzie.

Poniżej widzimy przykład dla pytania “How to get object size in Python”

```
Enter question that you want to ask and i will show you if there are already similar questions on stack overflow:how to get object size in Python
List of similar questions:
1) Getting file size in Python?

2) How do I check file size in Python?
|

3) How do I determine the size of an object in Python?

4) Create an empty list with certain size in Python
```

Jak widać, tym razem otrzymaliśmy 4, czyli maksymalną ilość pytań podobnych, aczkolwiek nie wszystkie są już aż tak podobne. Pierwsze jak i drugie zaproponowane przez program pytania są skonstruowane bardzo podobnie, ale dotyczą różnych rzeczy (użytkownik zapytał się o rozmiar obiektu a otrzymał rozmiar pliku). Trzecie pytanie za to jest już bardzo dobrze dopasowane pod względem semantycznym. Ostatnie pytanie zaproponowane przez nasze narzędzie nie jest już aż tak podobne do zadanego pytania jak poprzednie, rozpoznał podobieństwo związane z językiem Python i tym, że chodzi nam o wielkość czegoś, jednakże użytkownikowi chodziło o rozmiar obiektu a nie listy i nie o to jak stworzyć go z określonym rozmiarem tylko jak ten rozmiar otrzymać.

Przykład ten pokazuje, że nasze narzędzie pomimo dość dużej skuteczności, wciąż jest dalekie od perfekcji. W celu osiągnięcia lepszych wyników, należałoby między innymi zarejestrować aplikacje w stock exchange w celu rozbudowy zbioru danych oraz dopracowanie funkcji przypisującej punkty, aby poprawić jej skuteczność.

4. Źródła:

- <https://api.stackexchange.com/docs>
- <https://www.kaggle.com/c/predict-closed-questions-on-stack-overflow/data>
- <https://www.nltk.org/howto/wordnet.html>
- <https://wordnet.princeton.edu/>
- <https://www.nltk.org/api/nltk.tag.html>