

Analiza średniej liczby powtórzonych zapytań przez użytkownika na podstawie sumulowanego strumienia danych

Mateusz Kacpura

15 listopada 2024

Spis treści

1	Wstęp	2
2	Dane	2
3	Metody	2
3.1	Opis funkcji	2
3.1.1	load_data_from_file	2
3.1.2	simulate_data	3
3.1.3	naive_sampling	3
3.1.4	hash_sampling	3
3.1.5	compute_average_repeats	4
3.1.6	Główna część programu	4
4	Wyniki	5
5	Dyskusja	5
6	Wnioski	5
7	Bibliografia	5

1 Wstęp

Celem niniejszego sprawozdania jest analiza średniej liczby powtórzonych zapytań przez użytkownika w ciągu tego samego dnia. Ze względu na ograniczenia pamięciowe (możemy przechowywać jedynie 10% strumienia danych), zastosowano próbkowanie proporcjonalne w dwóch implementacjach:

- implementacja naiwna,
- implementacja z hashowaniem dla klucza *użytkownik*.

Przeprowadzono testy na rzeczywistych danych, które zostały wczytane z pliku tekstowego. W sprawozdaniu przedstawiono kod w języku Python, omówiono metody oraz zaprezentowano otrzymane wyniki.

2 Dane

Analizowane dane pochodzą z pliku `log.txt` i zawierają wiersze w formacie:

Adres_IP Zapytanie

Przykładowa linia z pliku:

89.113.36.135 http://www.gingiva.coop

Dla celów analizy przyjęto, że dane te reprezentują strumień krotek (*użytkownik*, *zapytanie*).

3 Metody

3.1 Opis funkcji

W celu przeprowadzenia analizy zaimplementowano następujące funkcje w języku Python:

3.1.1 `load_data_from_file`

Funkcja `load_data_from_file` wczytuje dane z pliku tekstowego i zwraca listę krotek (*użytkownik*, *zapytanie*).

```
1 import random
2 from collections import defaultdict
3
4 # Funkcja wczytująca dane z pliku
5 def load_data_from_file(filename):
6     sample_data = []
7     with open(filename, 'r') as file:
8         for line in file:
9             # Podział linii na użytkownika i zapytanie
10            parts = line.strip().split()
11            if len(parts) == 2:
12                user, query = parts
13                sample_data.append((user, query))
14    return sample_data
```

Listing 1: Funkcja `load_data_from_file`

3.1.2 simulate_data

Funkcja `simulate_data` symuluje strumień zapytań na podstawie wczytanych danych. Dla każdego użytkownika generuje losową liczbę zapytań (od 1 do 10), które mogą się powtarzać.

```
1 # Symulacja strumienia zapytań na podstawie przykładowych danych
2 def simulate_data(sample_data):
3     simulated_stream = []
4     for user, query in sample_data:
5         num_queries = random.randint(1, 10) # Liczba zapytań od 1 do 10
6         queries = [random.choice(sample_data)[1] for _ in range(num_queries)]
7         for q in queries:
8             simulated_stream.append((user, q))
9     return simulated_stream
```

Listing 2: Funkcja `simulate_data`

3.1.3 naive_sampling

Implementacja naiwnego próbkowania proporcjonalnego, gdzie każdy rekord jest wybierany z prawdopodobieństwem p .

```
1 # Naiwne próbkowanie: losowe próbkowanie par użytkownik-zapytanie z
  # prawdopodobieństwem p
2 def naive_sampling(stream, p=0.1):
3     sample = [entry for entry in stream if random.random() < p]
4     return sample
```

Listing 3: Funkcja `naive_sampling`

3.1.4 hash_sampling

Implementacja próbkowania z hashowaniem klucza *użytkownik*. Dzięki temu wszystkie zapytania wybranych użytkowników są uwzględniane.

```
1 # Hashowe próbkowanie: próbkowanie par użytkownik-zapytanie na podstawie
  # shashowanego ID użytkownika
2 def hash_sampling(stream, p=0.1):
3     sample = []
4     users = set()
5     for user, _ in stream:
6         users.add(user)
7
8     # Próbkowanie użytkowników na podstawie funkcji hash
9     sampled_users = {user for user in users if (hash(user) % 1000) < p * 1000}
10    sample = [(user, query) for user, query in stream if user in sampled_users]
11
12    return sample
```

Listing 4: Funkcja `hash_sampling`

3.1.5 compute_average_repeats

Funkcja oblicza średnią liczbę powtórzonych zapytań na użytkownika.

```
1  # Obliczanie średniej liczby powtórzonych zapytań na użytkownika
2  def compute_average_repeats(sample):
3      user_queries = defaultdict(list)
4
5      # Grupowanie zapytań według użytkownika
6      for user, query in sample:
7          user_queries[user].append(query)
8
9      total_repeats = 0
10     total_users = len(user_queries)
11
12     # Zliczanie powtórzonych zapytań dla każdego użytkownika
13     for queries in user_queries.values():
14         total_queries = len(queries)
15         unique_queries = len(set(queries))
16         repeats = total_queries - unique_queries
17         total_repeats += repeats
18
19     # Obliczenie średniej
20     average_repeats_per_user = total_repeats / total_users if total_users > 0
21     else 0
22     return average_repeats_per_user
```

Listing 5: Funkcja compute_average_repeats

3.1.6 Główna część programu

Poniższy kod stanowi główną część programu, który wykorzystuje powyższe funkcje do analizy danych z pliku log.txt.

```
1  # Wczytanie danych z pliku log.txt
2  sample_data = load_data_from_file('log.txt')
3
4  # Symulacja strumienia zapytań
5  simulated_stream = simulate_data(sample_data)
6
7  # Naiwne próbkowanie
8  naive_sample = naive_sampling(simulated_stream, p=0.1)
9  naive_average_repeats = compute_average_repeats(naive_sample)
10
11 # Hashowe próbkowanie
12 hash_sample = hash_sampling(simulated_stream, p=0.1)
13 hash_average_repeats = compute_average_repeats(hash_sample)
14
15 # Wyświetlenie wyników
16 print(f" rednia  liczba powtórzonych zapytań na użytkownika (naiwne próbkowanie):
17       {naive_average_repeats}")
17 print(f" rednia  liczba powtórzonych zapytań na użytkownika (hashowe
18       próbkowanie): {hash_average_repeats}")
```

Listing 6: Główna część programu

4 Wyniki

Uruchomienie powyższego kodu na rzeczywistych danych z pliku `log.txt` dało następujące wyniki:

Średnia liczba powtórzonych zapytań na użytkownika (naiwne próbkowanie): 142.761
Średnia liczba powtórzonych zapytań na użytkownika (hashowe próbkowanie): 9077.836

5 Dyskusja

Z otrzymanych wyników wynika, że średnia liczba powtórzonych zapytań na użytkownika jest znacznie wyższa w przypadku hashowego próbkowania niż w przypadku naiwnego próbkowania. Przyczyną tego jest sposób, w jaki obie metody traktują dane użytkowników:

- **Naiwne próbkowanie:** Każda para (użytkownik, zapytanie) jest próbkowana niezależnie z prawdopodobieństwem $p = 0.1$. W rezultacie, dla danego użytkownika w próbie mogą znaleźć się tylko niektóre z jego zapytań, co prowadzi do mniejszej liczby powtórzeń i potencjalnie zaniżonej średniej liczby powtórzonych zapytań.
- **Hashowe próbkowanie:** Użytkownicy są próbkowani na podstawie wartości funkcji hash ich identyfikatora. Jeśli użytkownik zostanie wybrany, to wszystkie jego zapytania zostają uwzględnione w próbie. Dzięki temu zachowana jest kompletność danych dla wybranych użytkowników, co prowadzi do większej liczby powtórzonych zapytań i wyższej średniej.

Istotne jest również zwrócenie uwagi na generowanie danych w funkcji `simulate_data`. W tej funkcji dla każdego użytkownika przypisujemy losową liczbę zapytań, które są losowo wybierane z puli wszystkich zapytań. Może to prowadzić do powtórzeń, zwłaszcza przy ograniczonej liczbie unikalnych zapytań.

6 Wnioski

Przeprowadzona analiza wykazała, że metoda próbkowania z hashowaniem klucza *użytkownik* dostarcza dokładniejsze oszacowanie średniej liczby powtórzonych zapytań na użytkownika w ciągu dnia w porównaniu z naiwnym próbkowaniem. Dzieje się tak, ponieważ hashowe próbkowanie zachowuje pełne informacje o aktywności wybranych użytkowników, co jest kluczowe przy analizie wzorców zachowań.

W praktycznych zastosowaniach, gdzie istotne jest dokładne modelowanie aktywności użytkowników przy ograniczonych zasobach pamięciowych, zaleca się wykorzystanie metod opartych na hashowaniu kluczy identyfikujących użytkowników.

7 Bibliografia

Literatura

- [1] Muthukrishnan, S. (2005). *Data streams: Algorithms and applications*. Now Publishers Inc.
- [2] Leskovec, J., Rajaraman, A., & Ullman, J. D. (2014). *Mining of Massive Datasets*. Cambridge University Press.