

Implementacja i analiza filtra *Blooma*

Mateusz Kacpura

15 listopada 2024

Streszczenie

Filtr Blooma jest strukturą danych umożliwiającą efektywne pod względem pamięciowym testowanie przynależności elementu do zbioru z dopuszczeniem kontrolowanego ryzyka fałszywych pozytywów. W niniejszym artykule przedstawiono implementację filtra Blooma z uwzględnieniem parametrów takich jak liczba funkcji mieszających, rozmiar tablicy bitowej oraz prawdopodobieństwo fałszywego pozytywu. Przeprowadzono eksperymenty dla różnych konfiguracji parametrów, analizując wpływ na liczbę fałszywych pozytywów. Wyniki potwierdzają teoretyczne zależności pomiędzy parametrami filtra a jego skutecznością.

Spis treści

1	Wprowadzenie	2
2	Teoria	2
2.1	Filtr Blooma	2
2.2	Parametry filtra	2
2.2.1	Rozmiar tablicy bitowej (m)	2
2.2.2	Liczba funkcji mieszających (k)	2
2.2.3	Prawdopodobieństwo fałszywego pozytywu (p)	2
3	Implementacja	3
3.1	Opis implementacji	3
3.2	Kod źródłowy	3
3.3	Objaśnienia	4
4	Eksperymenty i wyniki	5
4.1	Cel eksperymentów	5
4.2	Konfiguracja eksperymentów	5
4.3	Wyniki	5
4.4	Analiza wyników	5
5	Dyskusja	6
5.1	Wpływ liczby funkcji mieszających	6
5.2	Rozmiar tablicy bitowej	6
5.3	Optymalizacja parametrów	6
6	Wnioski	6
7	Bibliografia	6

1 Wprowadzenie

W erze przetwarzania dużych ilości danych kluczowe jest znalezienie struktur danych, które są zarówno efektywne pod względem pamięciowym, jak i czasowym. **Filtr Blooma** [1] jest probabilistyczną strukturą danych służącą do reprezentacji zbioru i umożliwiającą sprawdzanie, czy dany element należy do zbioru, z dopuszczeniem określonego poziomu fałszywych pozytywów. Filtr ten zużywa mniej pamięci niż struktury deterministyczne, takie jak tablice mieszające lub drzewa.

Celem niniejszej pracy jest przedstawienie implementacji filtra Blooma, uwzględnienie kluczowych parametrów wpływających na jego działanie oraz przeprowadzenie analizy liczby fałszywych pozytywów dla różnych konfiguracji.

2 Teoria

2.1 Filtr Blooma

Filtr Blooma składa się z tablicy bitów o długości m oraz k niezależnych funkcji mieszających h_1, h_2, \dots, h_k . Wszystkie bity tablicy są inicjalizowane jako 0. Przy dodawaniu elementu x do filtra, obliczane są wartości $h_i(x)$ dla $i = 1, 2, \dots, k$, a odpowiadające im pozycje w tablicy są ustawiane na 1.

Aby sprawdzić, czy element y należy do zbioru, należy sprawdzić czym są bity na pozycjach $h_i(y)$. Jeśli którykolwiek z tych bitów jest równy 0, to na pewno element nie należy do zbioru. Jeśli wszystkie bity są równe 1, to z dużym prawdopodobieństwem element jest w zbiorze, jednak istnieje ryzyko przyporządkowania fałszywego pozytywu.

2.2 Parametry filtra

2.2.1 Rozmiar tablicy bitowej (m)

Rozmiar tablicy bitowej wpływa na prawdopodobieństwo przyporządkowania fałszywego pozytywu. Większa tablica pozwala na zmniejszenie kolizji między funkcjami mieszającymi.

2.2.2 Liczba funkcji mieszających (k)

Optymalna liczba funkcji mieszających k zależy od rozmiaru tablicy m oraz liczby wstawionych elementów n , i jest dana wzorem:

$$k = \frac{m}{n} \ln 2 \quad (1)$$

2.2.3 Prawdopodobieństwo fałszywego pozytywu (p)

Prawdopodobieństwo przyporządkowania fałszywego pozytywu można oszacować wzorem:

$$p = \left(1 - e^{-kn/m}\right)^k \quad (2)$$

3 Implementacja

3.1 Opis implementacji

Implementacja została napisana w języku Python z wykorzystaniem bibliotek `numpy`, `math`, `bitarray` oraz `mmh3` (MurmurHash3). Kod został podzielony na klasę `BloomFilter`. Na końcu został przedstawiony przykład użycia zdefiniowanej class-y `BloomFilter`.

3.2 Kod źródłowy

```
1 import numpy as np
2 import math
3 import mmh3
4 from bitarray import bitarray
5
6 class BloomFilter():
7     def __init__(self, n, p):
8         self.n = n # Liczba elementów, które planujemy dodać
9         self.p = p # Akceptowalne prawdopodobieństwo fałszywego pozytywu
10
11         # Oblicz rozmiar tablicy bitowej (m) i liczbę funkcji mieszających (k)
12         self.m = self.get_size(n, p)
13         self.k = self.get_hash_count(self.m, n)
14
15         # Inicjalizacja tablicy bitowej zerami
16         self.bit_array = bitarray(self.m)
17         self.bit_array.setall(0)
18
19     def get_size(self, n, p):
20         m = -(n * math.log(p)) / (math.log(2) ** 2)
21         return int(m)
22
23     def get_hash_count(self, m, n):
24         k = (m / n) * math.log(2)
25         return int(k)
26
27     def add(self, item):
28         for i in range(self.k):
29             digest = mmh3.hash(item, i) % self.m
30             self.bit_array[digest] = True
31
32     def check(self, item):
33         for i in range(self.k):
34             digest = mmh3.hash(item, i) % self.m
35             if self.bit_array[digest] == False:
36                 return False
37         return True
38
39 # Wartości prawdopodobieństwa do przetestowania
40 p_values = [0.01, 0.05, 0.1, 0.2]
41
42 # Wyniki testów
43 print(f"{'Próba':<10} {'p':<10} {'m':<10} {'k':<10} {'Fałszywe pozytywy':<20}
44       {'Oczekiwane FP':<15}")
45 print("-" * 80)
46
47 for i, p in enumerate(p_values):
48     bloom = BloomFilter(n, p)
```

```

48     # Dodawanie elementów do filtru Blooma
49     for item in animals:
50         bloom.add(item)
51
52
53     # Sprawdzanie obecności elementów niewstawionych
54     false_positives = 0
55     for item in other_animals:
56         if bloom.check(item):
57             false_positives += 1
58
59     # Obliczanie oczekiwanych fałszywych pozytywów
60     expected_fp = len(other_animals) * ((1 - math.exp(-bloom.k * n / bloom.m)) **
61         bloom.k)
62
63     # Wyświetlanie wyników
64     print(f"{i + 1:<10} {p:<10} {bloom.m:<10} {bloom.k:<10} {false_positives:<20}
65         {expected_fp:<15.2f}")

```

Listing 1: Implementacja filtru Blooma

3.3 Objaśnienia

- `__init__`: Konstruktor klasy oblicza rozmiar tablicy bitowej `m` oraz liczbę funkcji mieszających `k` na podstawie podanych parametrów `n` i `p`.
- `get_size`: Metoda oblicza rozmiar tablicy bitowej za pomocą wzoru:

$$m = -\frac{n \ln p}{(\ln 2)^2} \quad (3)$$

- `get_hash_count`: Metoda oblicza optymalną liczbę funkcji mieszających:

$$k = \frac{m}{n} \ln 2 \quad (4)$$

- `add`: Metoda dodaje element do filtru, ustawiając odpowiednie bity na 1.
- `check`: Metoda sprawdza, czy element może znajdować się w filtrze.

4 Eksperymenty i wyniki

4.1 Cel eksperymentów

Celem eksperymentów jest analiza wpływu liczby funkcji mieszających k oraz rozmiaru tablicy bitowej m na liczbę fałszywych pozytywów dla ustalonej liczby elementów n .

4.2 Konfiguracja eksperymentów

Przeprowadzono testy dla różnych wartości akceptowalnego prawdopodobieństwa fałszywego pozytywu p . Dla każdego testu:

- Ustalono liczbę elementów wstawianych do filtru $n = 20$ (liczba zwierząt w liście `animals`).
- Obliczono rozmiar tablicy bitowej m oraz liczbę funkcji mieszających k na podstawie n i p .
- Dodano elementy z listy `animals` do filtru Bloom.
- Sprawdzone obecność elementów z listy `other_animals`.
- Zanotowano liczbę fałszywych pozytywów.
- Parametry p zostały ustalone na następujące wartości: 0.01, 0.05, 0.1, 0.2.

4.3 Wyniki

Próba	p	m	k	Fałszywe pozytywy	Oczekiwane FP
1	0.01	182	6	0	0.20
2	0.05	118	4	1	1.02
3	0.1	91	3	4	2.02
4	0.2	63	2	3	2.10

Tabela 1: Wyniki eksperymentów dla różnych wartości p

Objaśnienia do tabeli

- **p :** Akceptowalne prawdopodobieństwo fałszywego pozytywu.
- **m :** Rozmiar tablicy bitowej obliczony dla podanego n i p .
- **k :** Liczba funkcji mieszających.
- **Fałszywe pozytywy:** Liczba fałszywych pozytywów na 20 testów (elementów z `other_animals`).
- **Oczekiwane FP:** Oczekiwana liczba fałszywych pozytywów, obliczona jako $p \times$ liczba testów.

4.4 Analiza wyników

Zaobserwowano, że liczba fałszywych pozytywów rośnie wraz ze wzrostem akceptowalnego prawdopodobieństwa p . Wyniki eksperymentów są zbliżone do wartości oczekiwanych, co potwierdza poprawność implementacji oraz zgodność z teoretycznymi założeniami.

5 Dyskusja

5.1 Wpływ liczby funkcji mieszających

Liczba funkcji mieszających k wpływa na prawdopodobieństwo fałszywego pozytywu. Zbyt mała liczba funkcji może prowadzić do zwiększenia liczby fałszywych pozytywów, natomiast zbyt duża liczba funkcji zwiększa czas potrzebny na dodawanie i sprawdzanie elementów.

5.2 Rozmiar tablicy bitowej

Większy rozmiar tablicy bitowej m pozwala na zmniejszenie prawdopodobieństwa kolizji hashy różnych elementów, co z kolei zmniejsza liczbę fałszywych pozytywów. Jednak zwiększenie m wiąże się z większym zużyciem pamięci.

5.3 Optymalizacja parametrów

Dobór optymalnych wartości m i k jest kluczowy dla efektywnego działania filtru Blooma. Parametry te powinny być dostosowane do oczekiwanego prawdopodobieństwa fałszywego pozytywu oraz liczby elementów, które planujemy wstawić do filtru.

6 Wnioski

Filtr Blooma jest efektywną strukturą danych pozwalającą na oszczędne reprezentowanie zbiorów z dopuszczeniem kontrolowanego ryzyka fałszywych pozytywów. Implementacja i eksperymenty przeprowadzone w ramach pracy potwierdzają teoretyczne zależności między parametrami filtru a jego skutecznością.

7 Bibliografia

Literatura

- [1] Bloom, B. H. (1970). *Space/time trade-offs in hash coding with allowable errors*. Communications of the ACM, 13(7), 422-426.
- [2] Mitzenmacher, M. (2002). *Compressed Bloom Filters*. IEEE/ACM Transactions on Networking, 10(5), 604-612.
- [3] Broder, A., & Mitzenmacher, M. (2004). *Network applications of bloom filters: A survey*. Internet Mathematics, 1(4), 485-509.