

Databases for Developers - Assignment 2

Anders Jacobsen, Dima Karaush

April 8, 2021

1 Task 1 - investigation

Produce a small writeup (around 300 words) answering the following questions.

1. What is point of NoSQL databases?
2. What is the CAP theorem?
3. What are ideal use cases of HBase?

1.1 What is the point of NoSQL databases

The point of NoSQL databases is to fill the holes where SQL databases does not. We handle data in many different ways, so it is just logical that we should have multiple tools for storing that data. Some pros and cons can be found below.

Pros

- Performance
- Scalability
- Flexibility
- Data Models

Cons

- Not mature
- Requires multiple databases
- Huge databases

But the above information depends heavily on witch NoSQL database you choose. There are many!

1.2 What is the CAP theorem

In theoretical computer science, the CAP theorem, also named Brewer's theorem after computer scientist Eric Brewer, states that it is impossible for a distributed data store to simultaneously provide more than two out of the following three guarantees [2]:

- Consistency: Every read receives the most recent write or an error
- Availability: Every request receives a (non-error) response, without the guarantee that it contains the most recent write
- Partition tolerance: The system continues to operate despite an arbitrary number of messages being dropped (or delayed) by the network between nodes

1.3 What are ideal use cases of HBase

HBase uses Google's BigTable paper as blueprint, HBase is build on Hadoop and the Hadoop Distributed File System (HDFS) and designed for scaling horizontally on clusters of commodity hardware. HBase makes strong consistency guarantees and features tables with rows and columns - with should make SQL fans feel right at home. Out-of-the-box support for versioning and compression [3].

2 Task 2 - Bloom filters

Bloom filters are used in hbase as an incredible optimization. Solve the following.

1. What is a bloom filter?
2. What is an advantage of bloom filters over hash tables?
3. What is a disadvantage of bloom filters?
4. Using your language of choice, implement a bloom filter with add and check functions. The backing bit-array can simply be a long (64 bit integer)
5. If you are to store one million ASCII strings with an average size of 10 characters in a hash set, what would be the approximate space consumption?
6. The following equation gives the required number of bits of space per inserted key, where E is the false positive rate. $b = 1.44 \log_2(1/\xi)$
7. How many bits per element are required for a 1% false positive rate?
8. How many bits per element are required for a 5% false positive rate?

9. If you are to store one million ASCII strings with an average size of 10 characters in a bloom filter, what would be the approximate space consumption, given an allowed false positive rate of 5%?.

2.1 What is a bloom filter

A Bloom filter is a data structure designed to tell you, rapidly and memory-efficiently, whether an element is present in a set. The price paid for this efficiency is that a Bloom filter is a probabilistic data structure: it tells us that the element either definitely is not in the set or may be in the set [1].

2.2 What is an advantage of bloom filters over hash tables

Bloom filters is more space efficient. Since it does not store the actual object.

2.3 What is a disadvantage of bloom filters

The bloom filter is a probabilistic data structure, this means that the result has a small risk of returning as a "false positive". This could mean that you get a positive result for an object that does not exist.

2.4 Using your language of choice, implement a bloom filter with add and check functions

We have made a very simple exapmle of a bloom filter in python, it can be found in the folder "bloom_filter".

2.5 If you are to store one million ASCII strings with an average size of 10 characters in a hash set, what would be the approximate space consumption

$$OneASCIIcharacter = 8bits$$

$$Stringlength10 = 80bits$$

$$Timesamillion = 80.000.000bits$$

$$80.000.000bits = 10Mb(Megabytes)$$

For each entry the HashSet normally takes an empty space of equal size of the object as well as a pointer to the next entry. In Java the standart HashSet implementations memory consumption can be defined as

$$32 * SIZE + 4 * CAPACITY$$

The default Capacity of a HashSet implementation is 16
Therefore

$$32 * 1000000 + 4 * 16 = 32.000.064bytes = 30.5176Mb$$

So the approximate memory consumption for a HashSet in Java for storing 1.000.000 10 letter string would be 30.5176 Mb

2.6 The following equation gives the required number of bits of space per inserted key, where E is the false positive rate

$$b = 1.44 \log_2(1/\xi)$$

2.7 How many bits per element are required for a 1% false positive rate

We insert the probability (procent as decimal) in the equation from section 2.6:

$$b = 1.44 \log_2(1/0.01) = 9.56715291328$$

We can now conclude that we need approx. 9.57 bits of space pr. inserted key, to get a 1% false positive rate.

2.8 How many bits per element are required for a 5% false positive rate

We insert the probability (procent as decimal) in the equation from section 2.6:

$$b = 1.44 \log_2(1/0.05) = 6.22357645664$$

We can now conclude that we need approx. 6.22 bits of space pr. inserted key, to get a 5% false positive rate.

2.9 If you are to store one million ASCII strings with an average size of 10 characters in a bloom filter, what would be the approximate space consumption, given an allowed false positive rate of 5%

We insert the space consumption from section 2.8:

$$6.22357645664 * 1000000 = 6223576.45664$$

The approx. space consumption would therefore be 6223576.46 bits that equals approx. 759.71 Kilobytes.

3 Task 3 - Huffman coding

HBase internally uses a compression that is a combination of LZ77 and Huffman Coding.

1. Generate Huffmann Code (and draw the Huffmann Tree) based on the following string: “beebps beebps!!!! their eerie ears hear pears”
2. How many bits is the compressed string? How many bits is the raw ASCII string?
3. Compress “pete is here” with the Huffmann tree from before.
4. Write your own 10 word sentence. Generate the Huffmann Code (a new Huffmann Tree), and write a new compressed message (ie. in binary). Swap with one of your fellow students, and decompress each other’s message.

References

- [1] Bill Mill. Bloom filters by example.
- [2] the free encyclopedia Wikipedia. Cap theorem.
- [3] Luc Perkins with Eric Redmond and Jim R. Wilson. *Seven Databases in Seven Weeks*. Andy-Hunt, 2018.