



Department of Computer Science and Engineering (Data Science)

Subject: Machine Learning – IV Laboratory AY:
2024-25

Name: Jhanvi Parekh
Sap: 60009210033
Batch: D11

Experiment 3

Aim: Implement and analyze a Bloom Filter for checking the presence of elements, counting false positives, and true negatives.

Theory:

Introduction:

A Bloom Filter is a space-efficient probabilistic data structure used to test whether a given element is a member of a set or not. It is particularly useful when dealing with large datasets and is designed for fast and efficient lookups. However, a Bloom Filter may produce false positives (indicating an element is present when it's not) but never produces false negatives (if it says an element is not present, it's true). The main idea behind a Bloom Filter is to use multiple hash functions to map elements to multiple positions in a bit array.

Bloom Filter Workflow:

1. Initialization:

- Create a bit array of 'm' bits, initialized with all zeroes.
- Choose 'k' independent hash functions, each mapping to one of the 'm' bit positions.

2. Insertion:

- To add an element to the Bloom Filter, apply each of the 'k' hash functions to the element and set the corresponding bit positions in the bit array to 1.

3. Membership Check (Query):

- To check if an element is present in the Bloom Filter, apply each of the 'k' hash functions to the element.
- If all 'k' bit positions are 1, the element is in the set (may have false

positives). Counting False Positives and True Negatives:

- False Positives: These occur when the Bloom Filter incorrectly indicates that an element is in the set when it's not. False positives can happen because of hash collisions or if other elements have set the same bit positions.
- True Negatives: These occur when the Bloom Filter correctly indicates that an element is not in the set. Bloom Filters guarantee that false negatives do not occur.



Shri Vile Parle Kelavani Mandal's

DWARKADAS J. SANGHVI COLLEGE OF ENGINEERING

(Autonomous College Affiliated to the University of Mumbai)

NAAC Accredited with "A" Grade (CGPA : 3.18)



Department of Computer Science and Engineering (Data Science)

Sample problem:

Q] Insert elements 10 and 7. The Bloom's filter size is given as 5.

$h_1(x) = x \bmod 5$
 $h_2(x) = (2x+6) \bmod 5$

Comment whether 14 and 15 are present

Elements to be inserted	$h_1(x)$	$h_2(x)$	Bloom Filter										
10	$x \bmod 5$ $10 \bmod 5$ $= 0$	$(2x+6) \bmod 5$ $26 \bmod 5$ $= 1$	<table border="1"><tr><td>1</td><td>1</td><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td></tr></table>	1	1	0	0	0	0	1	2	3	4
1	1	0	0	0									
0	1	2	3	4									
7	$x \bmod 5$ $= 2$	$= 0$	<table border="1"><tr><td>1</td><td>1</td><td>1</td><td>0</td><td>0</td></tr></table>	1	1	1	0	0					
1	1	1	0	0									
Elements to check													
14	$14 \bmod 5$ $= 4$	$34 \bmod 5$ $= 4$	<p>Not present as per Bloom filter</p>										
15	$= 0$	$= 1$	<p>Acc. to Bloom's filter 15 is present (FALSE POSITIVE)</p>										

Lab Exercise:

Step 1: Initialize a Bloom Filter with a specified size 'm' and the number of hash functions 'k'.

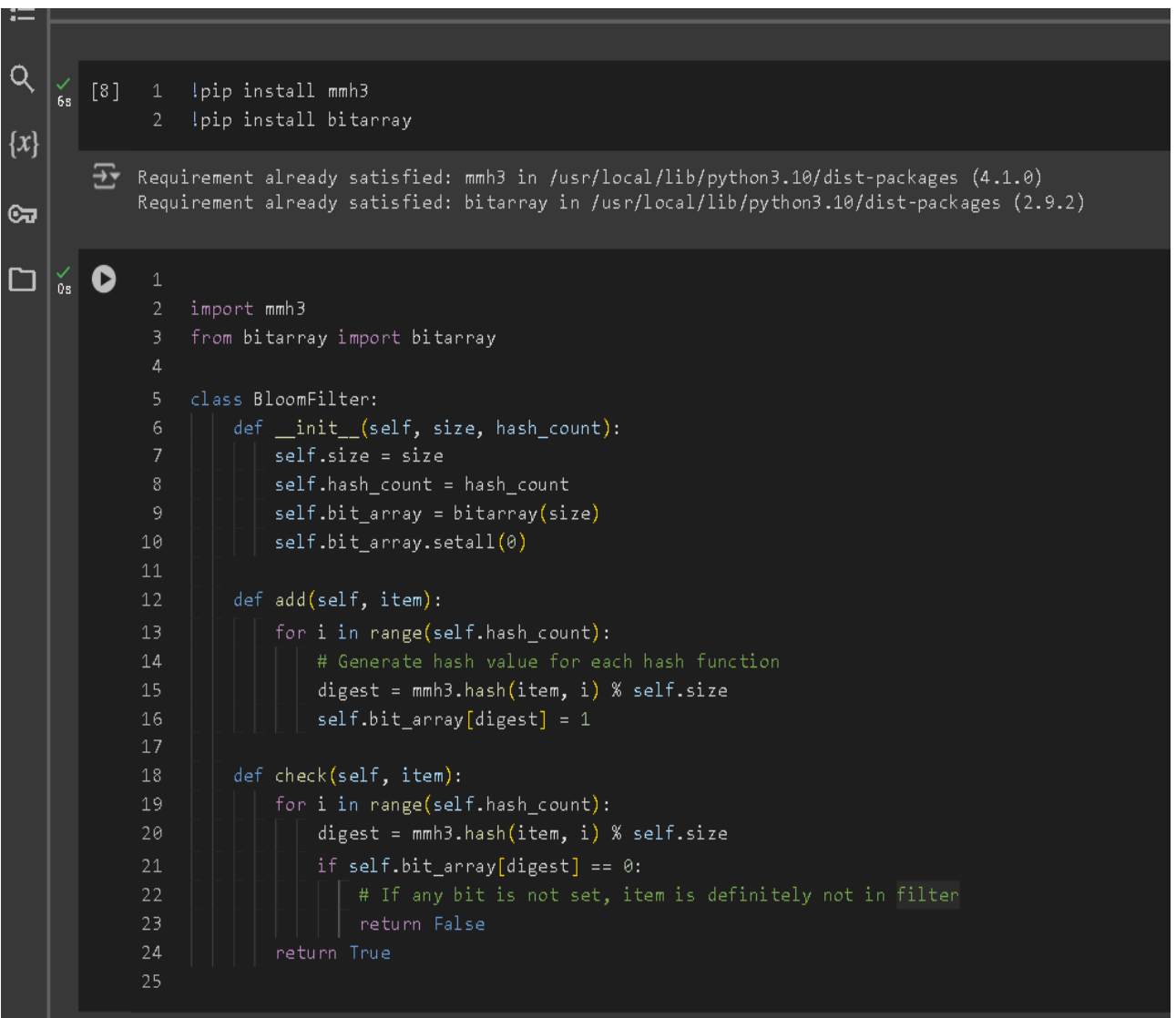
Step 2: Insert a set of elements into the Bloom Filter using the 'k' hash functions.

Step 3: Test the presence of various elements in the Bloom Filter and observe if there are any false positives.

Step 4: Count and record the number of false positives and true negatives.

[Carry out application of a Bloom's filter, like checking if a word is present or not in a document provided time permits]

Code with Output:



The image shows a terminal window and a code editor. The terminal window at the top shows the installation of mmh3 and bitarray. The code editor below shows the implementation of a BloomFilter class.

```
[8] 1 |pip install mmh3
    2 |pip install bitarray

Requirement already satisfied: mmh3 in /usr/local/lib/python3.10/dist-packages (4.1.0)
Requirement already satisfied: bitarray in /usr/local/lib/python3.10/dist-packages (2.9.2)

1
2 import mmh3
3 from bitarray import bitarray
4
5 class BloomFilter:
6     def __init__(self, size, hash_count):
7         self.size = size
8         self.hash_count = hash_count
9         self.bit_array = bitarray(size)
10        self.bit_array.setall(0)
11
12    def add(self, item):
13        for i in range(self.hash_count):
14            # Generate hash value for each hash function
15            digest = mmh3.hash(item, i) % self.size
16            self.bit_array[digest] = 1
17
18    def check(self, item):
19        for i in range(self.hash_count):
20            digest = mmh3.hash(item, i) % self.size
21            if self.bit_array[digest] == 0:
22                # If any bit is not set, item is definitely not in filter
23                return False
24        return True
25
```

```

1 # Parameters
2 bloom_filter_size = 10000 # Size of the bit array
3 hash_count = 10 # Number of hash functions
4
5 # Initialize Bloom Filter
6 bloom = BloomFilter(size=bloom_filter_size, hash_count=hash_count)
7
8 # Insert a set of elements
9 elements_to_add = ["apple", "banana", "cherry", "date", "elderberry"]
10
11 for element in elements_to_add:
12     bloom.add(element)
13
14 print("Elements added to Bloom Filter.")
15

```

Elements added to Bloom Filter.

```

1 # Test elements
2 test_elements = ["apple", "grape", "watermelon", "banana", "fig"]
3
4 for element in test_elements:
5     result = bloom.check(element)
6     print(f'{element} is in Bloom Filter: {result}')
7

```

'apple' is in Bloom Filter: True
 'grape' is in Bloom Filter: False
 'watermelon' is in Bloom Filter: False
 'banana' is in Bloom Filter: True
 'fig' is in Bloom Filter: False

```

[12] 1 # Example document
2 document = """A Bloom Filter is a space-efficient probabilistic data structure used to test whether a given element is a
3 member of a set or not. It is particularly useful when dealing with large datasets and is designed for fast
4 and efficient lookups. However, a Bloom Filter may produce false positives (indicating an element is
5 present when it's not) but never produces false negatives (if it says an element is not present, it's true).
6 The main idea behind a Bloom Filter is to use multiple hash functions to map elements to multiple
7 positions in a bit array. """
8
9

```

```

[13] 1 # Add words to Bloom Filter
2 # Parameters
3 bloom_filter_size = 10000 # Size of the bit array
4 hash_count = 10 # Number of hash functions
5 bloom = BloomFilter(size=bloom_filter_size, hash_count=hash_count)
6 words = document.lower().split()
7 for word in words:
8     bloom.add(word)
9
10 # Check for a word in the document
11 word_to_check = "bloom"
12 if bloom.check(word_to_check):
13     print(f"The word '{word_to_check}' might be in the document.")
14 else:
15     print(f"The word '{word_to_check}' is definitely not in the document.")
16

```

The word 'bloom' might be in the document.

```

1
2 word_to_check = "cat"
3 if bloom.check(word_to_check):
4     print(f"The word '{word_to_check}' might be in the document.")
5 else:
6     print(f"The word '{word_to_check}' is definitely not in the document.")

```

The word 'cat' is definitely not in the document.

Conclusion: Thus, gained hands-on experience estimating cardinality, counting false positives, and analysing true negatives using Bloom's Filter.

