



Complexity and Algorithm Analysis of Hash tables and Disjoint sets implementations.

Juan Sebastián Rueda Segura

Data structures

Universidad Nacional de Colombia Sede Bogotá

1/03/2025

Contents

1	Introduction	2
1.1	Goal	2
1.2	Motivation	2
2	Code Analysis	3
2.1	Used Libraries	3
2.2	Interfaces and Templates	3
2.3	Disjoint Set (Rank)	5
2.3.1	Constructor	5
2.3.2	Find	5
2.3.3	Find	5
2.4	Disjoint Set (Path compression)	7
2.4.1	Constructor	7
2.4.2	Find	7
2.4.3	Find	7
2.5	Fixed Size Hash Table	8
2.5.1	Constructor	8
2.5.2	Hash method	8
2.5.3	Insert	9
2.5.4	Find	9
2.5.5	Find	10
2.6	Variable Size Hash Table	11
2.6.1	Constructor	11
2.6.2	Insert	11
2.6.3	Find and Delete	12

3	Time measurements	13
3.1	Path compression vs Union by Rank	13
3.2	Fixed vs Variable size	13
3.3	Load Factor	16
3.4	Initial size	16
4	Conclusions	19

1 Introduction

To access the full code and previous implementations, you can check out this [Github repository](#).

1.1 Goal

This document's objective is to study the implementation of the **Hash tables**, and **Disjoint sets** data structures, analyzing their corresponding algorithms to study the behavior and time complexity using Big O notation. The implementation is presented in *Python*.

1.2 Motivation

These non-linear data structures arise from the need of efficient insert, find and delete times, since the linear structures such as linked lists or array lists had at least a method that required $O(N)$ time complexity. For hash tables we want to create a data structure that has $O(1)$ amortized complexity to find an element and for disjoint sets to know if there are connected components, useful in graph algorithms.

The disjoint set is a rooted-tree implementation, meaning that the leaves now point to the root. Each set is a rooted tree and the whole structure would be the forest containing them. It's array based for easy and constant time access. We have two possible implementations, using rank union and path compression.

For hashing, we will need an array, containing more arrays to handle collisions (this is called chaining) and a hash function for integers and strings. The hash maps implementation will allow insert, find and delete in almost constant time. This will depend on the longest chain but a good hash function will distribute the elements uniformly. We will show fixed size and variable hash table implementations.

2 Code Analysis

In this section we will explain details of the Tree interface in the BST and AVL trees, inspecting the corresponding code, making pertinent comments and finally we will give the time complexity in Big O notation.

2.1 Used Libraries

```
1 from abc import ABC, abstractmethod
2 import time
3 import matplotlib.pyplot as plt
4 import random
5 from typing import TypeVar, Generic, Optional, List, Any
```

We will need `abc` to create the interface of the trees, `time` and `matplotlib` to analyze the time measurements at the end of the document and `random` to generate numbers to be inserted on the trees and heap. Additionally, we will use *matplotlib* to graph the results obtained from the data analysis.

2.2 Interfaces and Templates

```
1 class DisjointSetsInterface(Generic[T], ABC):
2     @abstractmethod
3     def find(self, i: T) -> T:
4         pass
5
6     @abstractmethod
7     def union(self, i: T, j: T) -> None:
8         pass
9
10 class HashMap(Protocol, Generic[K]):
11     @abstractmethod
12     def put(self, key: K, value: Any) -> None:
13         """Insert or update a key-value pair in the hash map."""
14         pass
15
16     @abstractmethod
17     def get(self, key: K) -> Optional[Any]:
18         """Retrieve a value by key."""
19         pass
20
21     @abstractmethod
22     def remove(self, key: K) -> bool:
23         """Remove a key-value pair by key."""
24         pass
25
26     @abstractmethod
27     def size(self) -> int:
28         """Get the current size of the hash map."""
```

Python doesn't have properly defined *interfaces*, but they can be implemented through abstract classes, as the previous code shows, and then the data structures classes inherit from this class. It defines all the methods that have to be implemented and in case one or them is missing in the child class, the code won't compile and will show an error.

Templates or generics are a way to input parameters to functions or classes so that they can support multiple types of variables. This is useful in case we need to create strings or integer lists without having to change the implementation.

2.3 Disjoint Set (Rank)

When joining two trees in the forest we must ensure that we link the smallest one, so we create an array of the ranks that will be the height of the tree.

2.3.1 Constructor

```
1 def __init__(self, num: int):
2     #makeset
3     self.parents = [i for i in range(0, num)]
4     self.rank = [0 for _ in range(0, num)]
```

Note that we have implicitly created the Makeset method when creating the parents array. This simplifies the operations, but when needing another data type this must be changed. We need twice the space to store the ranks as well.

Complexity: $O(1)$

2.3.2 Find

```
1 def find(self, i):
2     while i != self.parents[i]:
3         i = self.parents[i]
4     return i
```

We iterate traversing the array until the parent of the index is equal to itself, meaning that it is contained in the same set.

Complexity: $O(\log N)$

2.3.3 Find

```
1 def union(self, i, j):
2     i_id = self.find(i)
3     j_id = self.find(j)
4
5     if i_id == j_id:
6         return
7
8     if self.rank[i_id] > self.rank[j_id]:
9         self.parents[j_id] = i_id
10    else:
11        self.parents[i_id] = j_id
12        if self.rank[i_id] == self.rank[j_id]:
13            self.rank[j_id] = self.rank[j_id] + 1
```

First we find the sets that contains i and j . If they are on the same set, we are done. Otherwise, we check the rank of the trees that represents the sets to see which is the smallest and append the other one. If they are the same size we update the rank of the one that appended the other. The complexity is $O(\log N)$ because the number of nodes in a tree is at most $\log N$ as we saw on class.

Complexity: $O(\log N)$

2.4 Disjoint Set (Path compression)

The main difference is that now we don't need the rank array because every element will now point to the element in the root, making the access faster.

2.4.1 Constructor

```
1 def __init__(self, num: int):
2     #makeset
3     self.parents = [i for i in range(0, num)]
```

As we said, now we don't need the rank array.

Complexity: $O(1)$

2.4.2 Find

```
1 def find(self, i):
2     if i != self.parents[i]:
3         self.parents[i] = self.find(self.parents[i])
4     return self.parents[i]
```

We iterate traversing the array until the parent of the index is equal to itself. The difference with the previous implementation is that now we start to update the elements that we traverse until we find the root, recursively calling the find method. This operation works in $O(\log^* N)$, the iterated logarithm as we saw on class.

Complexity: $O(\log^* N)$

2.4.3 Find

```
1 def union(self, i, j):
2     i_id = self.find(i)
3     j_id = self.find(j)
4
5     self.parents[i_id] = j_id
```

It's almost the same as the previous find but in this case we simply make the parent of i be j 's. This works because of the find method. The complexity is $O(2\log^* N)$ since we call find twice.

Complexity: $O(\log^* N)$

2.5 Fixed Size Hash Table

2.5.1 Constructor

```
1 def __init__(self, m: int):
2     self.m = m
3     self.buckets: List[List[Tuple[K, Any]]] = [[] for _ in range(m)]
4     self._count = 0
5
6     # For integer hashing (ax + b) mod p mod m
7     self.p = self._find_prime_larger_than(m * 100)
8     self.a = random.randint(1, self.p - 1)
9     self.b = random.randint(0, self.p - 1)
10
11    # For string hashing (polynomial hash)
12    self.x = random.randint(1, self.p - 1)
```

We will need m , the maximum size, buckets, the list of lists containing the hash map information, count that is the number of elements on the table, and some parameters to hash integers or strings. To find a prime number larger than m , we implement this function:

```
1 def _find_prime_larger_than(self, n: int) -> int:
2     """Find a prime number larger than n."""
3     def is_prime(num: int) -> bool:
4         if num <= 1:
5             return False
6         if num <= 3:
7             return True
8         if num % 2 == 0 or num % 3 == 0:
9             return False
10        i = 5
11        while i * i <= num:
12            if num % i == 0 or num % (i + 2) == 0:
13                return False
14            i += 6
15        return True
16
17        while not is_prime(n):
18            n += 1
19        return n
```

Complexity: $O(1)$

2.5.2 Hash method

```
1 def _hash_int(self, key: int) -> int:
2     return ((self.a * key + self.b) % self.p) % self.m
3
4     def _hash_str(self, key: str) -> int:
```



```

5     hash_value = 0
6     for char in key:
7         hash_value = (hash_value * self.x + ord(char)) % self.p
8     return hash_value % self.m
9
10    def _hash(self, key: K) -> int:
11        if isinstance(key, int):
12            return self._hash_int(key)
13        elif isinstance(key, str):
14            return self._hash_str(key)
15        else:
16            raise TypeError("Key must be an integer or string")

```

Depending if is a hash table consisting on integers or strings, we have to select the corresponding method to hash. If it's an integer the hash function will be $((ax + b) \bmod p) \bmod m$. And if it's a string it will be the polynomial hash as we saw on class.

Complexity: $O(1)$

2.5.3 Insert

```

1 def put(self, key: K, value: Any) -> None:
2     hash_value = self._hash(key)
3     bucket = self.buckets[hash_value]
4
5     # Check if key already exists
6     for i, (k, v) in enumerate(bucket):
7         if k == key:
8             bucket[i] = (key, value) # Update value
9             return
10
11    # Key not found, append new key-value pair
12    bucket.append((key, value))
13    self._count += 1

```

We hash the key and search in it's bucket. If a value is found, we update it to the new one. And in case it's not in the bucket, we append it and increase the count attribute by one. The complexity depends on the longest chain c .

Complexity: $O(c + 1)$

2.5.4 Find

```

1 def get(self, key: K) -> Optional[Any]:
2     hash_value = self._hash(key)
3     bucket = self.buckets[hash_value]
4

```

```

5         for k, v in bucket:
6             if k == key:
7                 return v
8
9         return None

```

We repeat the process in insert but if no key in the bucket matches the key to find, we return None.

Complexity: $O(c + 1)$

2.5.5 Find

```

1 def remove(self, key: K) -> bool:
2     hash_value = self._hash(key)
3     bucket = self.buckets[hash_value]
4
5     for i, (k, v) in enumerate(bucket):
6         if k == key:
7             bucket.pop(i)
8             self._count -= 1
9             return True
10
11     return False

```

We repeat the process in find. When the key is found we delete from the bucket the tuple value, key, subtract 1 to count and return True. If no key in the bucket matches the key to find, we return False.

Complexity: $O(c + 1)$

2.6 Variable Size Hash Table

We will use a fixed size table to use some operations, but will rebalance after an insertion.

2.6.1 Constructor

```
1 def __init__(self, initial_size: int, load_factor_threshold: float)
  :
2     self.load_factor_threshold = load_factor_threshold
3     self.count = 0
4     self.capacity = initial_size
5     self.map = FixedSizeHashMap[K](initial_size)
```

We set a load factor that represents $\frac{n}{m}$, the number of elements in the table and the maximum size. If we cross that balance factor, we rebalance.

Complexity: $O(1)$

2.6.2 Insert

```
1 def put(self, key: K, value: Any) -> None:
2     # Check if key already exists
3     old_value = self.map.get(key)
4     if old_value is not None:
5         self.map.put(key, value)
6         return
7
8     # New key, check load factor
9     current_load_factor = (self.count + 1) / self.capacity
10    if current_load_factor > self.load_factor_threshold:
11        self._resize()
12
13    self.map.put(key, value)
14    self.count += 1
```

We insert it in the fixed table size and then we check if the load factor is still satisfied. Otherwise we resize:

```
1 def _resize(self) -> None:
2     """Resize the hash map by doubling its size."""
3     old_map = self.map
4     new_capacity = self.capacity * 2
5     new_map = FixedSizeHashMap[K](new_capacity)
6
7     # Rehash all entries
8     for bucket in old_map.buckets:
9         for key, value in bucket:
10            new_map.put(key, value)
```

```

11
12     self.map = new_map
13     self.capacity = new_capacity

```

We double the capacity and create a new map. We must then rehash all the elements present in the old table so we iterate through all the list accessing each list element and inserting it in the new table. Finally we update the pointer to new map and new capacity. Similarly to the array variable size, resizing takes lineal time but amortized constant time.

Complexity: $O(c + 2)$

2.6.3 Find and Delete

```

1  def get(self, key: K) -> Optional[Any]:
2      return self.map.get(key)
3
4  def remove(self, key: K) -> bool:
5      if self.map.remove(key):
6          self.count -= 1
7          return True
8      return False

```

We simply call the methods in the map attribute. In the delete we decrease the count.

Complexity: $O(c + 1)$ each

3 Time measurements

3.1 Path compression vs Union by Rank

```
1 if __name__ == "__main__":
2     sizes = [1000, 10000, 100000, 1000000]
3     times = []
4     times2 = []
5     for size in sizes:
6         start = time.time()
7         ds = DisjointSets(size)
8         for i in range(0, size - 1):
9             ds.union(random.randint(0, size - 1), random.randint(0,
10                size - 1))
11         end = time.time()
12         times.append(end - start)
13
14         start = time.time()
15         ds2 = DisjointSets2(size)
16         for i in range(0, size - 1):
17             ds2.union(random.randint(0, size - 1), random.randint
18                (0, size - 1))
19         end = time.time()
20         times2.append(end - start)
21
22     plt.plot(sizes, times2, marker='o', linestyle='-', label="Union-
23 -Find with Path Compression")
24
25     plt.plot(sizes, times, marker='o', linestyle='-', label="Union-
26 Find with Rank")
27     plt.xscale("log")
28     plt.yscale("log")
29     plt.grid(True, which="both", linestyle="--", alpha=0.7)
30     plt.xlabel("Size (n)")
31     plt.ylabel("Total Time (s)")
32     plt.legend()
33     plt.show()
```

We compare different sizes and create two disjoint sets, one with path compression and the other with rank implemented. We see that they behave very similarly with the only exception being that the rank works better for small sizes.

3.2 Fixed vs Variable size

```
1 if __name__ == "__main__":
2     sizes = [1000, 10000, 100000, 1000000]
3     trials = 5
4     times = []
5     times2 = []
```

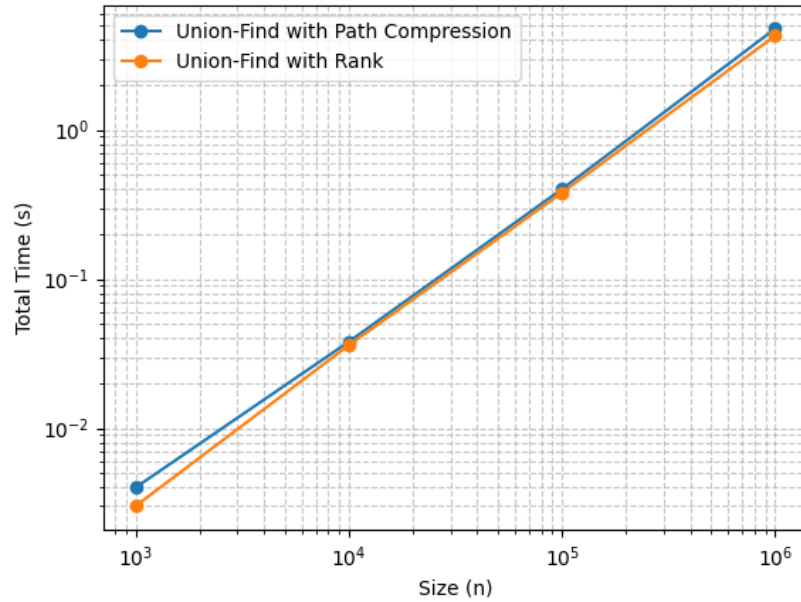


Figure 1: Comparison between Rank and Path compression

```

6
7   for size in sizes:
8       fhm_times = []
9       rhm_times = []
10
11      for _ in range(trials):
12          random.seed(42) # Ensure the same random values for
13                          # both hash maps
14
15          # FixedSizeHashMap timing
16          start = time.time()
17          fhm = FixedSizeHashMap(size)
18          for i in range(size):
19              fhm.put(random.randint(0, size - 1), random.randint
20                      (0, size - 1))
21          end = time.time()
22          fhm_times.append(end - start)
23
24          random.seed(42)
25          # ResizableHashMap timing
26          start = time.time()
27          rhm = ResizableHashMap()
28          for i in range(size):
29              rhm.put(random.randint(0, size - 1), random.randint
30                      (0, size - 1))
31          end = time.time()

```

```

29         rhm_times.append(end - start)
30
31         times.append(sum(fhm_times) / trials)
32         times2.append(sum(rhm_times) / trials)
33
34     plt.plot(sizes, times2, marker='o', linestyle='-', label="
Resizable Hash Map")
35     plt.plot(sizes, times, marker='o', linestyle='-', label="Fixed
Size Hash Map")
36     plt.xscale("log")
37     plt.yscale("log")
38     plt.grid(True, which="both", linestyle="--", alpha=0.7)
39     plt.xlabel("Size (n), Load Factor = 0.75")
40     plt.ylabel("Total Time (s)")
41     plt.legend()
42     plt.show()

```

We compare different sizes and create two tables, one with path fixed size and the other with variable size. We set the initial size to be the default, 16 and the load factor to be 0.75 and we initialize the fixed size with the number of elements. We ran 5 tests using the same random seed, to ensure that the items are the same for each test. We see that the resizable hash map performs better. I thought that the static had the advantage to get the same size as the total of elements, but it seems that the resize outperforms this.

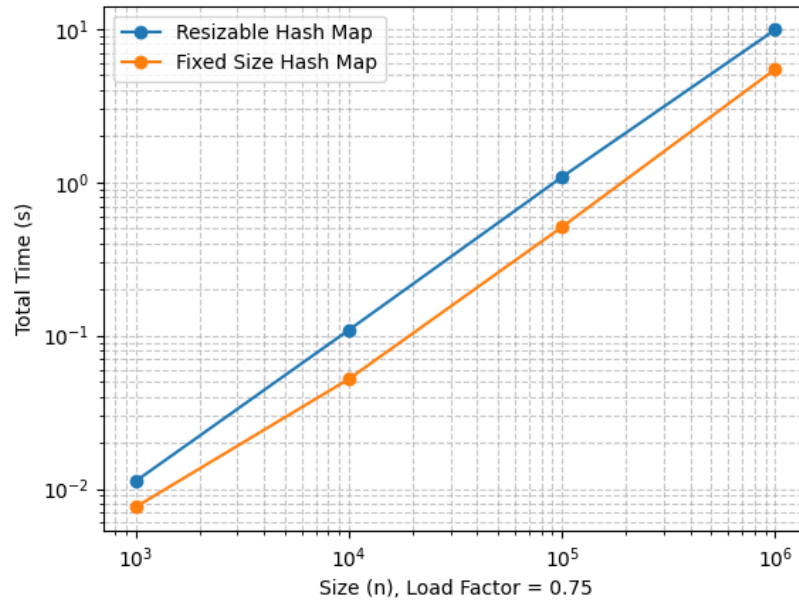


Figure 2: Comparison between Rank and Path compression

3.3 Load Factor

```
1 if __name__ == "__main__":
2     sizes = [1000, 10000, 100000, 1000000]
3     trials = 1
4     load_factors = [0.1, 0.25, 0.5, 0.75, 0.9, 1.2, 1.5]
5     times = {lf: [] for lf in load_factors}
6
7     for size in sizes:
8         rhm_times = {lf: [] for lf in load_factors}
9
10        for _ in range(trials):
11            for lf in load_factors:
12                random.seed(42) # Reset the seed for consistent
13                                values
14
15                                # ResizableHashMap timing
16                                start = time.time()
17                                rhm = ResizableHashMap(load_factor_threshold=lf)
18                                for i in range(size):
19                                    rhm.put(random.randint(0, size - 1), random.
20                                        randint(0, size - 1))
21                                end = time.time()
22                                rhm_times[lf].append(end - start)
23
24            for lf in load_factors:
25                times[lf].append(sum(rhm_times[lf]) / trials)
26
27        for lf in load_factors:
28            plt.plot(sizes, times[lf], marker='o', linestyle='--', label
29                =f"Resizable Hash Map (Load Factor = {lf})")
30
31        plt.xscale("log")
32        plt.yscale("log")
33        plt.grid(True, which="both", linestyle="--", alpha=0.7)
34        plt.xlabel("Size (n)")
35        plt.ylabel("Total Time (s)")
36        plt.legend()
37        plt.show()
```

We compare different sizes and different load factors, following the same considerations of the last test. We can see that the smallest values, 0.1 and 0.2 don't behave stably and fluctuate depending the size. The loading factors from 0.5 to 1 behave the same and the 1.5 load factor works better with larger input sizes. I conclude that the loading factor is very situational and choosing it from 0.5 to 1.5 is better.

3.4 Initial size

```
1 if __name__ == "__main__":
2     sizes = [1000, 10000, 100000, 1000000]
3     trials = 1
```

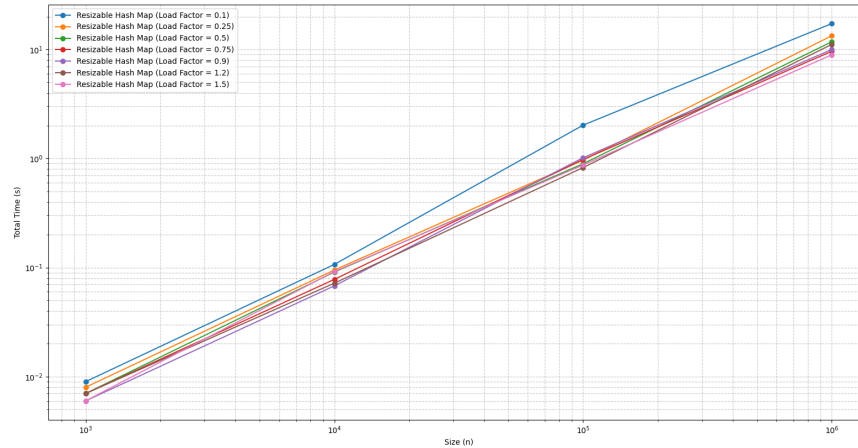



Figure 3: Comparison between Rank and Path compression

```

4  initial_sizes = [0.5, 0.75, 1.0, 1.5] # Multipliers for the
   initial size
5  times = {isize: [] for isize in initial_sizes}
6
7  for size in sizes:
8      fhm_times = {isize: [] for isize in initial_sizes}
9
10     for _ in range(trials):
11         random.seed(42) # Ensure the same random values for
   all hash maps
12
13         for isize in initial_sizes:
14             m = int(size * isize)
15             start = time.time()
16             fhm = FixedSizeHashMap(m)
17             for i in range(size):
18                 fhm.put(random.randint(0, size - 1), random.
   randint(0, size - 1))
19             end = time.time()
20             fhm_times[isize].append(end - start)
21
22         for isize in initial_sizes:
23             times[isize].append(sum(fhm_times[isize]) / trials)
24
25     for isize in initial_sizes:
26         plt.plot(sizes, times[isize], marker='o', linestyle='-',
   label=f"Fixed Size Hash Map (m = {isize} * size)")
27
28     plt.xscale("log")
29     plt.yscale("log")
30     plt.grid(True, which="both", linestyle="--", alpha=0.7)

```

```

31 plt.xlabel("Size (n)")
32 plt.ylabel("Total Time (s)")
33 plt.legend()
34 plt.show()

```

We compare different sizes and different initial sizes for the fixed hash table, following the same considerations of the previous test. Here we see that if we assign the size of the list as m we have faster insertion in small cases but then it becomes the worst. The most stable is 1.5 times the size of the inputs and when the size is big the smallest size, 0.5, is good, which means that the hash function distributes uniformly.

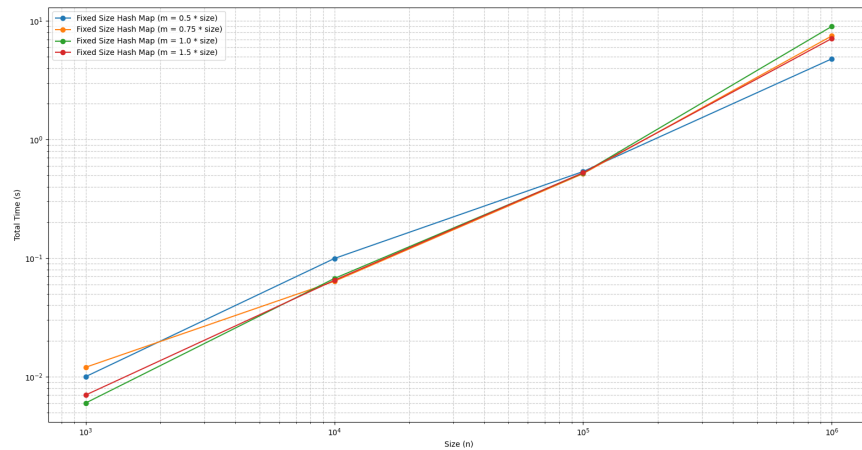


Figure 4: Comparison between Rank and Path compression

4 Conclusions

This document summarized two non-linear data structures implementations, analyzing the code, time complexity and obtained results given by testing the code and plotting the time vs size graphics. One thing to note is the faster results given by implementing these data structures. Note that Python is an interpreted language and therefore is slower than if we had implemented these data structures on a language like C or C++.

The results obtained were as I expected and as we saw on class: the best loading factor is between 0.5 and 0.9, if we assign a good hash function then the time is linear, which is represented on the graphs since we call N times the insert function, the resize is $O(1)$ amortized time or else it would have showed worse results and finally the disjoint sets implementation using path compression is as efficient as using the rank but better in memory usage.