Ubuntu

Device Name

cdm-virtual-machine >

Memory	3.8 GiB
Processor	12th Gen Intel® Core™ i7-12700H × 2
Graphics	SVGA3D; build: RELEASE; LLVM;
Disk Capacity	21.5 GB

OS Name	Ubuntu 22.04.1 LTS
OS Type	64-bit
GNOME Version	42.4
Windowing System	Wayland
Virtualization	VMware
Software Updates	>

Section 1 —— 跳表的实现

跳表类及节点结构体的实现

跳表类：

- 节点结构体
 - 所存的键值对（key, val）以及高度 level
 - 一个指针指向其同层的下一个节点
 - 构造函数
 - 定义头指针和尾指针
- 构造函数
- 插入函数
- 查找函数
- 查找步长函数

- 概率
- 随机数（与概率用于生成节点高度）

```

1  class skiplist_type
2  {
3      // add something here
4  private:
5      struct Node
6      {
7          key_type key;
8          value_type val;
9          int level;
10         Node *forward[MAX_LEVEL];
11         Node(key_type key, const value_type &s, int level);
12     } *head, *tail;
13     int random_level();
14     double my_rand();
15     double Jump_Probaility;
16
17 public:
18     explicit skiplist_type(double p = 0.5);
19     void put(key_type key, const value_type &val);
20     // std::optional<value_type> get(key_type key) const;
21     std::string get(key_type key) const;
22
23     // for hw1 only
24     int query_distance(key_type key) const;
25 };

```

跳表插入的实现

- 首先要有一个指针 `cur`，初始时指向最顶层的头指针，后续就使用这个指针由上至下、由左至右移动进行查找
- 还要有一个指针 `update` 用于存放上一次向下跳跃的节点后续会使用这个指针来更新跳表
- 具体算法：
 - `cur` 指针首先从跳表“左上角”开始移动，比较当前节点的同层右侧节点的值和待插入的值，若右侧节点值更大或其为尾节点，则跳入下一层并用 `update` 记录下当前跳跃的节点，否则将当前的 `cur` 指针变更到此右侧节点
 - 当跳跃到最底层时，如果此时 `cur` 指针的节点不为尾节点或其值与待插入的值相同，则更新其 `val`，如果此时 `cur` 指针的节点不为尾节点或其值不与待插入的节点值相同，则创建一个 `n` 节点并用 `random_level()` 初始化其高度，然后将其在概率生成的高度右侧指向 `update` 指针的节点原本的右侧节点，随后将当前 `update` 指针的节点指向这个插入的 `n` 节点完成插入操作

```

1  void skiplist_type::put(key_type key, const value_type &val){
2      Node *cur = head;
3      Node *updated[MAX_LEVEL];
4      for (int i = MAX_LEVEL - 1; i >= 0; --i)
5      {
6          while (cur->forward[i] != NULL && cur->forward[i]->key < key)
7              cur = cur->forward[i];
8          updated[i] = cur;

```

```

9     }
10    if (cur->forward[0] != tail && cur->forward[0]->key == key)
11    {
12        cur->forward[0]->val = val;
13        return;
14    }
15    int rlevel = random_level();
16    Node *n = new Node(key, val, rlevel);
17    for (int i = 0; i < rlevel; ++i)
18        n->forward[i] = tail;
19    for (int i = rlevel; i >= 0; --i)
20    {
21        n->forward[i] = updated[i]->forward[i];
22        updated[i]->forward[i] = n;
23    }
24 }

```

跳表查找的实现

- 类似于插入操作，首先要有一个指针 `cur`，初始时指向最顶层的头指针，后续就使用这个指针由上至下、由左至右移动进行查找
- 比较 `cur` 右侧的节点值是否与待查找的值相同，若是则返回其 `val`，否则进行比较，当当前层 `cur` 的节点右侧节点不为空且其值小于所查找的值时，`cur` 向右移动，否则向下移动一层继续进行比较直至最后一层
- 若在最后一层仍未找到则返回一个空字符串

```

1  std::string skiplist_type::get(key_type key) const {
2      Node *cur = head;
3      for (int i = MAX_LEVEL - 1; i >= 0; --i){
4          while (cur->forward[i] != tail && cur->forward[i]->key < key)
5              cur = cur->forward[i];
6          if (cur->forward[i] != tail && cur->forward[i]->key == key)
7              return cur->forward[i]->val;
8      }
9      return "";
10 }

```

其余部分代码

```

1  //跳表构造函数
2  skiplist_type::skiplist_type(double p) {
3      Jump_Probaility = p;
4      head = new Node(0, "", MAX_LEVEL);
5      tail = new Node(0, "", MAX_LEVEL);
6      for (int i = 0; i < MAX_LEVEL; ++i) {
7          head->forward[i] = tail;
8          tail->forward[i] = tail;
9      }
10 }
11 //查找路径（类似于查找）
12 int skiplist_type::query_distance(key_type key) const {
13     Node *cur = head;

```

```

14     int path_length = 1;
15     for (int i = MAX_LEVEL - 1; i >= 0; --i){
16         while (cur->forward[i] != tail && cur->forward[i]->key < key) {
17             cur = cur->forward[i];
18             path_length++;
19         }
20         if (cur->forward[i] != tail && cur->forward[i]->key == key) {
21             return path_length++;
22         }
23         path_length++;
24     }
25     return path_length++;
26 }

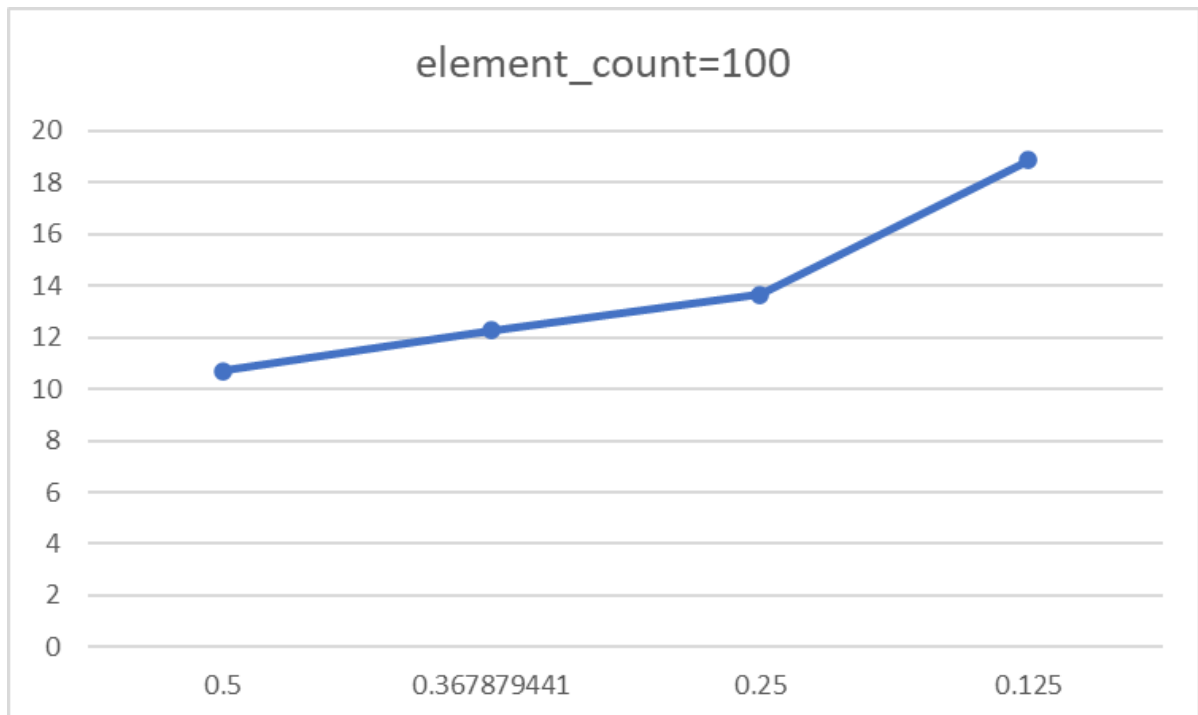
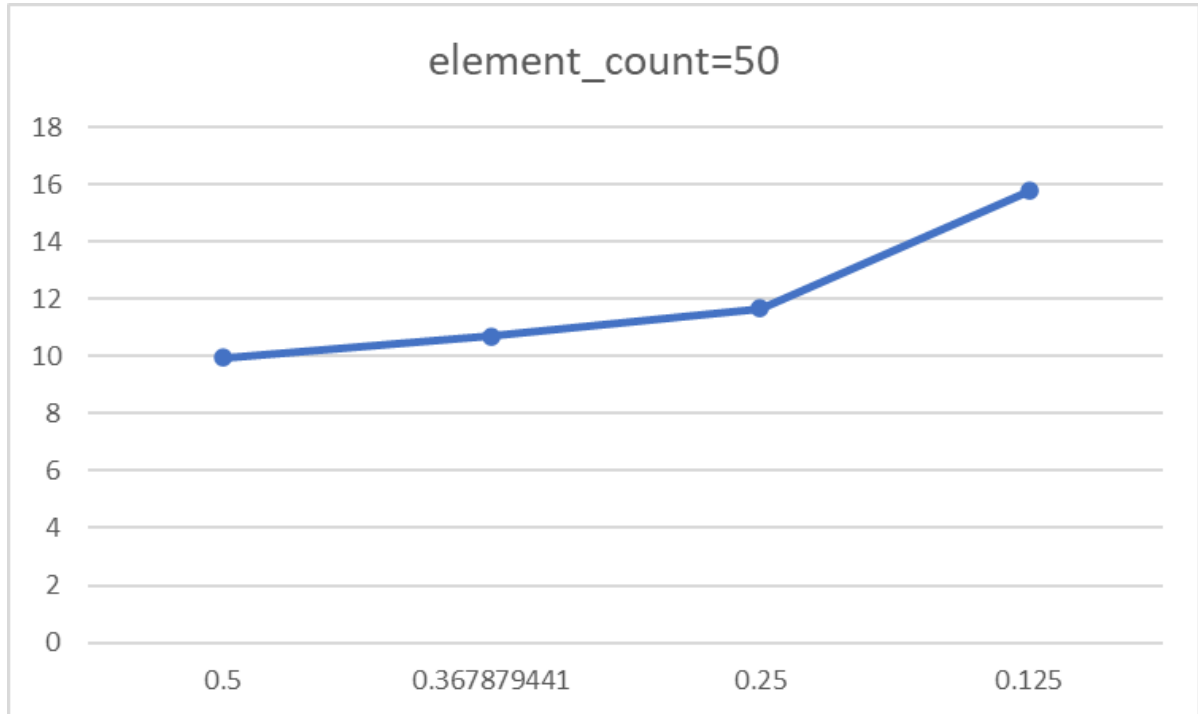
```

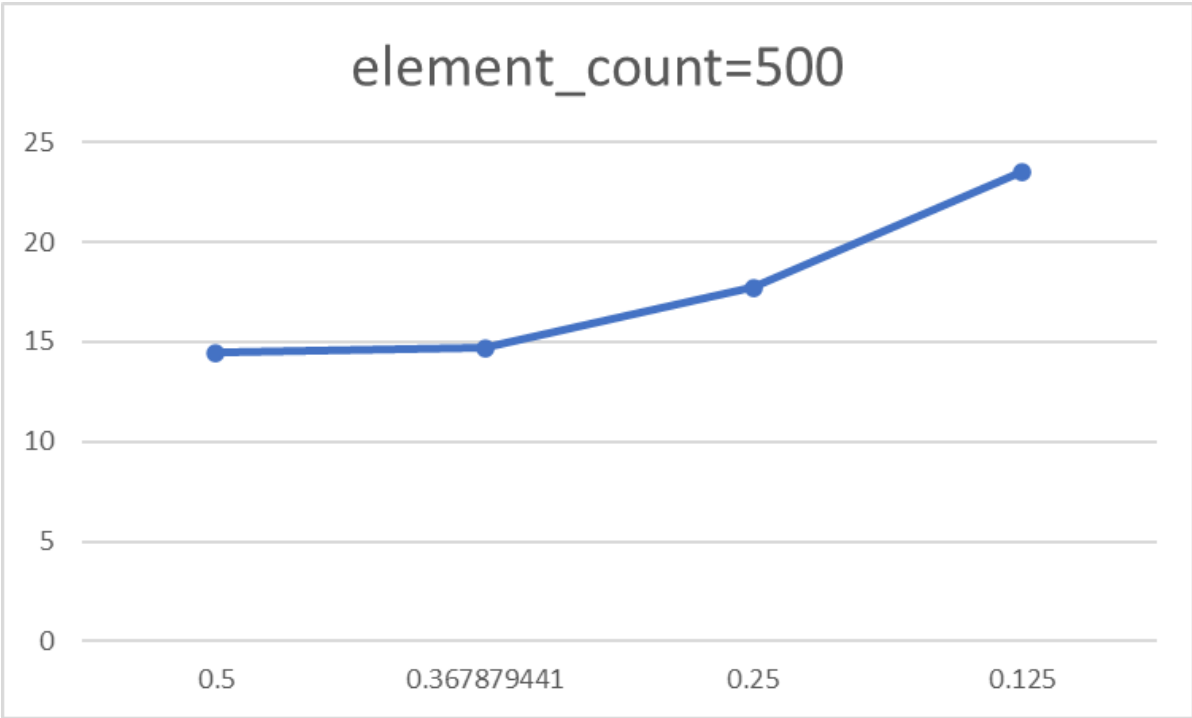
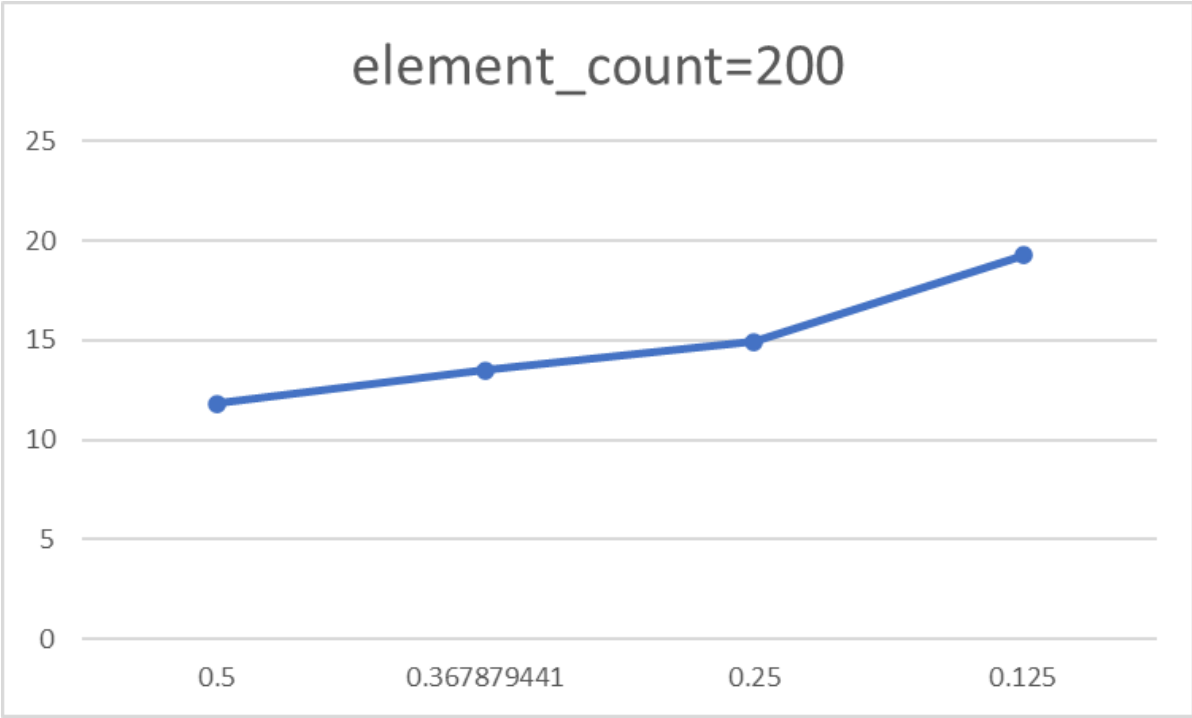
Section 2 —— 探究增长率对跳表性能的影响

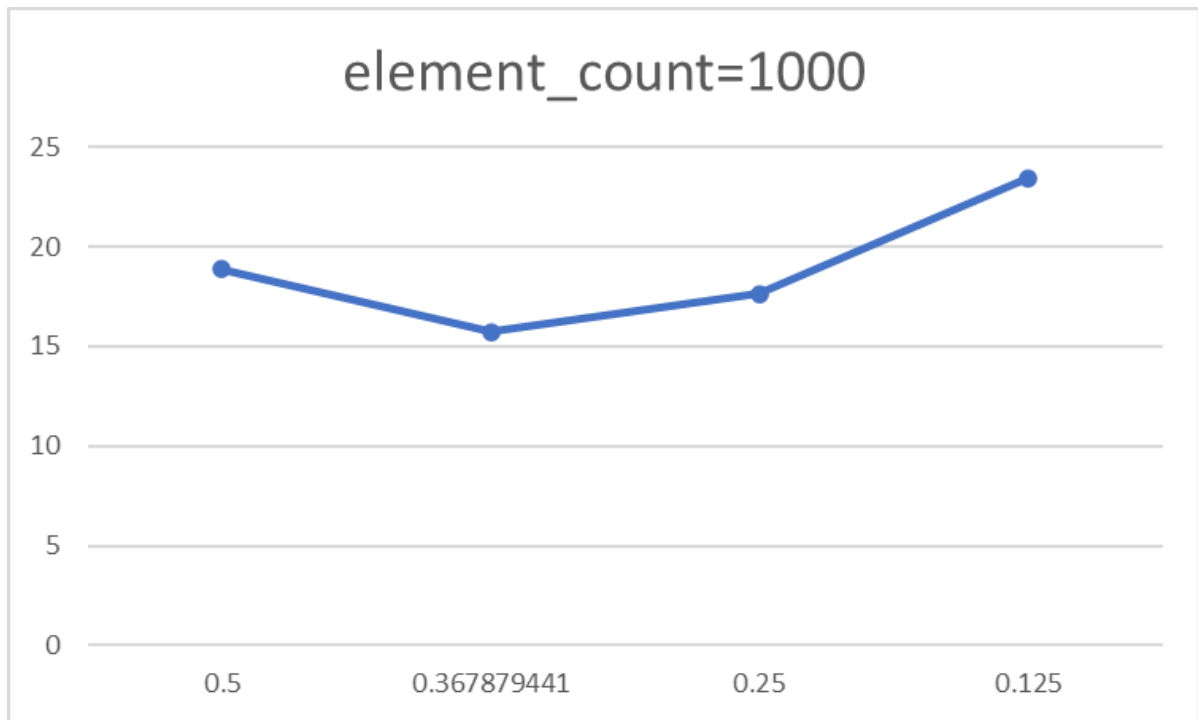
选择5个 `element_count` 即跳表元素个数(50, 100, 200, 500, 1000), 4个 `p` 值即概率 ($1/2$, $1/e$, $1/4$, $1/8$), 分别测试搜索10000次的 `average query distance` 即平均搜索长度, 结果如下

element_count	p	average query distance
50	0.5	9.93635028
50	0.367879441	10.68415525
50	0.25	11.67262122
50	0.125	15.76735273
100	0.5	10.70715588
100	0.367879441	12.26985
100	0.25	13.64880171
100	0.125	18.87082646
200	0.5	11.80721189
200	0.367879441	13.47543519
200	0.25	14.92413028
200	0.125	19.2740173
500	0.5	14.45684688
500	0.367879441	14.69501132
500	0.25	17.71731112
500	0.125	23.54042721
1000	0.5	18.87121512

element_count	p	average query distance
1000	0.367879441	15.72349454
1000	0.25	17.64183862
1000	0.125	23.45992068







分析

由上述折线图可知：

- 首先，在 p 一定的前提下，随着 `element_count` 的增大，`average query distance` 逐渐增加，因为整个链表变长了，搜索的长度会增加
- 然后，在 `element_count` 一定的前提下，随着 p 的减小，`average query distance` 总体上呈现出先减后增的趋势，并且会在 $1/e$ 处取得极值：
 - $L(n) = \log_{1/p} n$ ，代表 `MAX_LEVEL`，搜索的时间复杂度 $C(i) = i/p$ ， $C(p) = \frac{\ln(n)}{-p \ln(p)}$ 。
对 $H(x) = x \ln(x)$ 求导： $H'(x) = \ln(x) + 1$ ，在 $(0, \frac{1}{e}]$ 单调递减， $[\frac{1}{e}, +\infty)$ 单调递增，因此会在 $1/e$ 处取得极值
- 在 `element_count` 较小时（例如 50, 100, 200, 500 时），可能会出现偶然状况，比如查找的节点都大量的集中在表头部分，此时结论并不严格成立，但是随 `element_count` 增大（例如 1000 时），结论会越来越明显