

## 《源码探秘 CPython》41. 集合支持的操作是怎么实现的？

原创 古明地觉 古明地觉的编程教室 2022-03-05 09:00

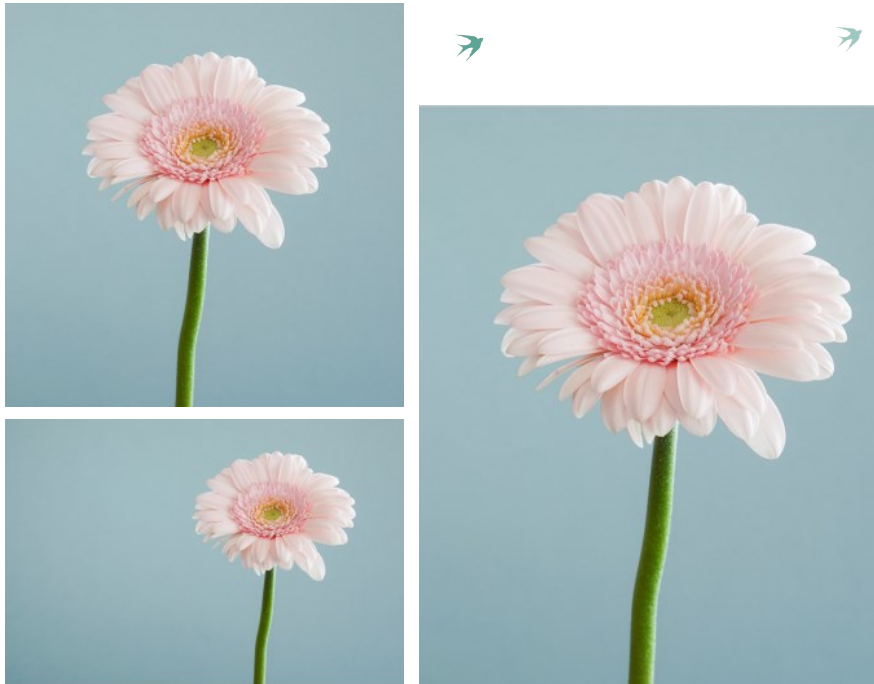


微信扫一扫  
关注该公众号

收录于合集

#CPython

97个 >



楔子

下面我们来分析一下集合的操作是怎么实现的，比如元素的添加、删除，以及集合的扩容等等。

并且集合还支持交集、并集、差集等运算，它们又是如何实现的呢？那么就一起来看一看吧。



添加元素

添加元素，会调用PySet\_Add函数。

```
1 int
2 PySet_Add(PyObject *anyset, PyObject *key)
3 {
4     //参数是两个指针
5     //类型检测
6     if (!PySet_Check(anyset) &&
7         (!PyFrozenSet_Check(anyset) || Py_REFCNT(anyset) != 1)) {
8         PyErr_BadInternalCall();
9         return -1;
10    }
11    //本质上调用了set_add_key
12    return set_add_key((PySetObject *)anyset, key);
13 }
```

在进行了参数检测之后，又调用了set\_add\_key。

```

1 static int
2 set_add_key(PySetObject *so, PyObject *key)
3 {
4     //声明一个变量, 用于保存哈希值
5     Py_hash_t hash;
6
7     //类型检测, 看看是否是ASCII字符串
8     if (!PyUnicode_CheckExact(key) ||
9         (hash = ((PyASCIIObject *) key)->hash) == -1) {
10         //如果不是ASCII字符串
11         //那么计算哈希值
12         hash = PyObject_Hash(key);
13         //如果计算之后的哈希值为-1
14         //在表示该对象不可被哈希, Python层面显然会报错
15         if (hash == -1)
16             return -1;
17     }
18     //底层又调用了set_add_entry, 并把hash也作为参数传了进去
19     return set_add_entry(so, key, hash);
20 }

```

和字典类似, 这一步也不是添加元素的真正逻辑, 只是计算了哈希值。显然下面的 `set_add_entry` 就是具体的逻辑了。

```

1 static int
2 set_add_entry(PySetObject *so, PyObject *key, Py_hash_t hash)
3 {
4     //...
5     restart:
6         //获取mask
7         mask = so->mask;
8         //hash和mask进行按位与, 得到一个索引
9         i = (size_t)hash & mask;
10        //获取对应的entry指针
11        entry = &so->table[i];
12        if (entry->key == NULL)
13            //如果entry->key == NULL
14            //表示当前位置没有被使用
15            //直接跳到found_unused标签
16            goto found_unused;
17
18        //否则说明该位置已经存储entry
19        freeslot = NULL;
20        perturb = hash; // 将perturb设置为hash
21
22        //接下来就要改变规则, 重新映射了
23        while (1) {
24            //获取已存在entry的hash字段的值
25            //如果和我们当前的哈希值一样的话
26            if (entry->hash == hash) {
27                //获取已存在entry的key
28                PyObject *startkey = entry->key;
29                //entry里面的key不可以为dummy态
30                //因为这相当于删除(伪删除)了, 那么hash应该为-1
31                assert(startkey != dummy);
32                //如果startkey和key相等, 说明指向了同一个对象
33                //那么两者视为相等, 而集合内的元素不允许重复
34                if (startkey == key)
35                    //直接跳转到found_active标签
36                    goto found_active;
37                //如果不是同一个对象, 再比较维护的值是否相等
38                //快分支, 假设两者都是字符串, 然后进行比较
39                if (PyUnicode_CheckExact(startkey)

```

```

40         && PyUnicode_CheckExact(key)
41         && _PyUnicode_EQ(startkey, key))
42         //如果一样, 跳转到found_active标签
43         goto found_active;
44
45     //到这里说明两者不是同一个对象, 也不都是字符串
46     //那么只能走通用的比较逻辑了
47     table = so->table;
48     //增加startkey的引用计数
49     Py_INCREF(startkey);
50     //比较两个对象维护的值是否一致
51     cmp = PyObject_RichCompareBool(startkey, key, Py_EQ);
52     //减少startkey的引用计数
53     Py_DECREF(startkey);
54     //如果cmp大于0, 比较成功
55     if (cmp > 0)
56         //说明两个值是相同的
57         //跳转到found_active标签
58         goto found_active;
59     if (cmp < 0)
60         //小于0 说明比较失败
61         //跳转到comparison_error标签
62         goto comparison_error;
63     //拿到当前的mask
64     mask = so->mask;
65 }
66 //如果不能hash
67 else if (entry->hash == -1)
68     //则设置为freeslot
69     freeslot = entry;
70
71 //如果当前索引值加上9小于等于当前的mask
72 //define LINEAR_PROBES 9
73 if (i + LINEAR_PROBES <= mask) {
74     //循环9次, 这里逻辑我们一会单独说
75     for (j = 0 ; j < LINEAR_PROBES ; j++) {
76         // .....
77     }
78 }
79
80 //程序走到这里说明索引冲突了
81 //改变规则, 重新计算索引值
82 perturb >>= PERTURB_SHIFT;
83 //我们看到计算规则和字典是一样的
84 i = (i * 5 + 1 + perturb) & mask;
85 //获取新索引对应的entry
86 entry = &so->table[i];
87 //如果对应的key为NULL, 说明重新计算索引之后找到了可以存储的地方
88 if (entry->key == NULL)
89     //跳转到found_unused_or_dummy
90     goto found_unused_or_dummy;
91 //否则说明比较倒霉, 改变规则重新映射之后, 索引依旧冲突
92 //那么继续循环, 比较key是否一致等等
93 }
94
95 found_unused_or_dummy:
96     //如果这个freeslot为NULL, 说明是可用的
97     if (freeslot == NULL)
98         //跳转
99         goto found_unused;
100 //否则, 说明为dummy态
101 //那么我们依旧可以使用, 正好废物利用
102 //将used数量加一
103 so->used++;

```

```

104 //设置key和hash值
105 freeslot->key = key;
106 freeslot->hash = hash;
107 return 0;
108
109 //发现未使用的
110 found_unused:
111 //将fill和used个数+1
112 so->fill++;
113 so->used++;
114 //设置key和hash值
115 entry->key = key;
116 entry->hash = hash;
117 //检查active态+dummy的entry个数是否小于mask的3/5
118 if ((size_t)so->fill*5 < mask*3)
119 //是的话, 表示无需扩容
120 return 0;
121 //否则要进行扩容
122 //如果active态的entry大于50000, 那么两倍扩容, 否则四倍扩容
123 return set_table_resize(so, so->used>50000 ? so->used*2 : so->used*4
124 );
125
126 //如果是found_active, 表示key重复了
127 //直接减少一个引用计数即可
128 found_active:
129 Py_DECREF(key);
130 return 0;
131
132 //比较失败, 同样减少引用计数, 返回-1
133 comparison_error:
134 Py_DECREF(key);
135 return -1;
}

```

代码很多, 我们还删除了一部分, 整个流程总结一下就是:

- 传入hash值, 计算出索引值, 通过索引值找到对应的entry;
- 如果entry->key=NULL, 那么将hash和key存到对应的entry;
- 如果entry->key != NULL, 那么就比较两个key是否相同;
- 如果相同, 则不插入, 直接减少引用计数。因为不是字典, 不存在更新一说;
- 如果不相同, 那么从该索引往后遍历9个entry, 如果存在key为NULL的entry, 那么设置进去;
- 如果以上条件都不满足, 则改变策略重新计算索引值, 直到找到一个满足key为NULL的entry;
- 判断容量问题, 如果active态+dummy态的entry个数不小于3/5\*mask, 那么扩容, 扩容的规则是active态的entry个数是否大于50000, 是的话就二倍扩容, 否则4倍扩容;

最后是if (i + LINEAR\_PROBES <= mask), 这一部分代码我们省略了, 那它是做什么的呢? 首先哈希值相同但是key不同时, 按照学习字典的思路, 肯定是映射一个新的索引。

但是问题来了, 这样是不能有效地利用CPU缓存的, L1 Cache加载数据会一次性加载64字节, 称为一个cache line。如果两个位置间隔比较远, 因为映射出来的索引是随机的, 对应的entry可能不在cache中, 从而导致CPU下一次需要重新读取。

所以Python中引入了**LINEAR\_PROBES**, 从当前的entry开始, 查找前面的9个entry。如果还找不到可用位置, 然后才重新计算, 从而提高cache的稳定性。

所以集合和字典在解决哈希冲突的时候采取的策略是一样的, 只不过集合多考虑了CPU的cache。

删除元素会调用`set_remove`函数，但是删除的核心逻辑位于`set_discard_entry`函数中。

```

1 static int
2 set_discard_entry(PySetObject *so, PyObject *key, Py_hash_t hash)
3 { //传入集合、key、以及计算的哈希值
4
5     setentry *entry;
6     PyObject *old_key;
7     //通过传入的key和hash找到该entry
8     //并且entry->key要和传入的key是一样的
9     entry = set_lookkey(so, key, hash);
10    //如果entry为NULL, 说明不存在此key
11    //直接返回-1
12    if (entry == NULL)
13        return -1;
14    //如果entry不为NULL, 但是对应的key为NULL
15    //返回DISCARD_NOTFOUND
16    if (entry->key == NULL)
17        return DISCARD_NOTFOUND;
18    //获取要删除的key
19    old_key = entry->key;
20    //并将entry设置为dummy
21    entry->key = dummy;
22    //hash值设置为-1
23    entry->hash = -1;
24    //减少使用数量
25    so->used--;
26    //减少引用计数
27    Py_DECREF(old_key);
28    //返回DISCARD_FOUND
29    return DISCARD_FOUND;
30 }

```

如果找到了指定的key，在`set_remove`函数里面会返回None，否则报出KeyError。

可以看到集合添加、删除元素和字典是有些相似的，毕竟底层都是使用了哈希表嘛。

## 集合的扩容

当集合的容量不够时，会自动扩容，具体的逻辑位于`set_table_resize`函数中。

```

1 static int
2 set_table_resize(PySetObject *so, Py_ssize_t minused)
3 { //显然参数是:PySetObject *指针以及容量大小
4
5     //三个setentry *指针
6     setentry *oldtable, *newtable, *entry;
7     //oldmask
8     Py_ssize_t oldmask = so->mask;
9     //newmask
10    size_t newmask;
11
12    //是否为其申请过内存
13    int is_oldtable_malloced;
14    //将PySet_MINSIZE个entry直接copy过来

```

```

15 setentry small_copy[PySet_MINSIZE];
16 //minused必须大于等于0
17 assert(minused >= 0);
18 //newsize不断扩大二倍, 直到大于minused
19 //所以我们刚才说的大于50000, 二倍扩容, 否则四倍扩容
20 //实际上是最终的newsize是比二倍或者四倍扩容的结果要大的
21 size_t newsize = PySet_MINSIZE;
22 while (newsize <= (size_t)minused) {
23     //newsize最大顶多也就是PY_SSIZE_T_MAX+1
24     //但一般不可能有这么多元素
25     newsize <<= 1;
26 }
27 //为新的table申请空间
28 oldtable = so->table;
29 assert(oldtable != NULL);
30 is_oldtable_malloced = oldtable != so->smalltable;
31
32 //如果newsize和PySet_MINSIZE(这里的8)相等
33 if (newsize == PySet_MINSIZE) {
34
35     //拿到smalltable, 就是默认初始化8个entry数组的那个成员
36     newtable = so->smalltable;
37     //如果oldtable和newtable一样
38     if (newtable == oldtable) {
39         //并且没有dummy态的entry
40         if (so->fill == so->used) {
41             //那么无需做任何事情
42             return 0;
43         }
44         //否则的话, dummy的个数一定大于0
45         assert(so->fill > so->used);
46         //扔掉dummy态, 只把oldtable中active态的entry拷贝过来
47         memcpy(small_copy, oldtable, sizeof(small_copy));
48         //将small_copy重新设置为oldtable
49         oldtable = small_copy;
50     }
51 }
52 else {
53     //否则的话, 肯定大于8, 申请newsize个setentry所需要的空间
54     newtable = PyMem_NEW(setentry, newsize);
55     //如果newtable为NULL, 那么申请内存失败, 返回-1
56     if (newtable == NULL) {
57         PyErr_NoMemory();
58         return -1;
59     }
60 }
61 //newtable肯定不等于oldtable
62 assert(newtable != oldtable);
63 //创建一个能容纳newsize个entry的空set
64 memset(newtable, 0, sizeof(setentry) * newsize);
65 //将mask设置为newsize-1
66 //将table设置为newtable
67 so->mask = newsize - 1;
68 so->table = newtable;
69 //获取newmask
70 newmask = (size_t)so->mask;
71 //遍历旧table的setentry数组
72 //将setentry的key和hash全部设置到新的table里面
73 //如果fill==used, 说明没有dummy态的entry
74 if (so->fill == so->used) {
75     for (entry = oldtable; entry <= oldtable + oldmask; entry++) {
76         if (entry->key != NULL) {
77             //设置元素的逻辑在此函数中
78             set_insert_clean(newtable, newmask, entry->key, entry->h

```

```

79 ash);
80     }
81 }
82 } else {
83     //逻辑和上面一样, 但是存在dummy态的entry
84     //判断时需要多一个条件:entry->key != dummy
85     //由于会丢弃dummy态的entry, 因此扩容后fill和used相等
86     //所以这里将used赋值给fill
87     so->fill = so->used;
88     //另外估计有人觉得这里的代码有点啰嗦
89     //代码是类似的, 没必要分成两个分支
90     //其实这是Python为了性能考虑的
91     //如果fill==used, 说明不存在dummy态的entry
92     //那么遍历时就无需加上entry->key != dummy这个条件了
93     for (entry = oldtable; entry <= oldtable + oldmask; entry++) {
94         if (entry->key != NULL && entry->key != dummy) {
95             set_insert_clean(newtable, newmask, entry->key, entry->h
96 ash);
97         }
98     }
99 }
100
101 //如果已经为旧的table申请了内存, 那么要将其归还给系统堆
102 if (is_oldtable_mallocated)
103     PyMem_DEL(oldtable);
104
105 return 0;
106 }

```

整个逻辑还是不难理解的, 该函数内部负责申请内存, 初始化成员。但是设置元素的核心逻辑位于set\_insert\_clean中, 我们看一下。

```

1 static void
2 set_insert_clean(setentry *table, size_t mask, PyObject *key, Py_hash_t
3 hash)
4 {
5     setentry *entry;
6     //perturb初始值为hash
7     size_t perturb = hash;
8     //计算索引
9     size_t i = (size_t)hash & mask;
10    size_t j;
11
12    while (1) {
13        //获取当前entry
14        entry = &table[i];
15        if (entry->key == NULL)
16            //如果为空则跳转found_null设置key与hash
17            goto found_null;
18        if (i + LINEAR_PROBES <= mask) {
19            //否则还是老规矩, 遍历之后的9个entry
20            for (j = 0; j < LINEAR_PROBES; j++) {
21                entry++;
22                //找到空的entry, 那么跳转到found_null设置key与hash
23                if (entry->key == NULL)
24                    goto found_null;
25            }
26        }
27        // 没有找到, 那么改变规则, 重新计算索引
28        perturb >>= PERTURB_SHIFT;
29        i = (i * 5 + 1 + perturb) & mask;
30    }
31 found_null:
32     //设置key与hash
33     entry->key = key;

```

```
34     entry->hash = hash;
    }
}
```

以上就是集合的扩容，我们又看到了字典的影子。



## 集合的交集运算

我们在使用集合的时候，可以取两个集合的交集、并集、差集、对称差集等等。这里介绍一下交集，其余的可以自己参考源码研究一下，源码位于setobject.c中。

```
1  static PyObject *
2  set_intersection(PySetObject *so, PyObject *other)
3  {
4      //result, 集合运算之后会产生新的集合
5      PySetObject *result;
6      PyObject *key, *it, *tmp;
7      Py_hash_t hash;
8      int rv;
9
10     //如果两个对象相同
11     if ((PyObject *)so == other)
12         //直接返回其中一个的拷贝即可
13         return set_copy(so);
14
15     //这行代码表示创建一个空的PySetObject *
16     result = (PySetObject *)make_new_set_basetype(Py_TYPE(so), NULL);
17     //如果result == NULL, 说明创建失败
18     if (result == NULL)
19         return NULL;
20
21     //检测other是不是PySetObject *
22     if (PyAnySet_Check(other)) {
23         //初始索引为0
24         Py_ssize_t pos = 0;
25         //setentry *
26         setentry *entry;
27
28         //如果other元素的个数大于so
29         if (PySet_GET_SIZE(other) > PySet_GET_SIZE(so)) {
30             //就把so和other进行交换
31             tmp = (PyObject *)so;
32             so = (PySetObject *)other;
33             other = tmp;
34         }
35
36         //从少的那一方的开始遍历
37         while (set_next((PySetObject *)other, &pos, &entry)) {
38             //拿到key和hash
39             key = entry->key;
40             hash = entry->hash;
41             //传入other的key和hash, 在so中去找
42             rv = set_contains_entry(so, key, hash);
43             if (rv < 0) {
44                 //如果rv<0, 说明不存在
45                 Py_DECREF(result);
46                 return NULL;
47             }
48             if (rv) {
```



```

49         //存在的话设置进result里面
50         if (set_add_entry(result, key, hash)) {
51             Py_DECREF(result);
52             return NULL;
53         }
54     }
55 }
56 //直接返回
57 return (PyObject *)result;
58 }
59 //...
60 }

```

逻辑比我们想象中的要单纯，假设有两个集合S1和S2，遍历元素少的集合，然后判断元素在另一个集合中是否存在。如果存在，则添加进要返回的集合中，否则遍历下一个。



以上就是集合相关的内容，它的效率也是非常高的，能够以O(1)的复杂度去查找某个元素。最关键的是，它用起来也特别的方便。

此外Python里面还有一个frozenset，也就是不可变的集合。但frozenset对象和set对象都是同一个结构体，只有PySetObject，没有PyFrozenSetObject。

我们在看PySetObject的时候，发现该对象里面也有一个hash成员，如果是不可变集合，那么hash值是不为-1的，因为它不可以添加、删除元素，是不可变对象。由于比较相似，因此frozenset就不再说了，可以自己对着源码简单看一下，源码还是setobject.c。

收录于合集 [#CPython 97](#)

[← 上一篇](#)

《源码探秘 CPython》42. 迭代器的实现原理

[下一篇 →](#)

《源码探秘 CPython》40. 集合是怎么实现的？

喜欢此内容的人还喜欢

python 7天进阶之路-对象和json转换  
缪斯之子



[系列]微服务·深入理解 gRPC - Part2  
走向架构师的每一天



Abaqus python脚本开发 第三章 各类指令的方法对象变量 (3)  
山石结构

