```
着字典大小的增加而变慢。
 田 面向对象高级编程
                                  dict就是第二种实现方式,给定一个名字,比如 'Michael',dict在内部就可以直接计算出 Michael 对应的存放成绩的"页码",也就是 95 这个数字存放的内
 田 错误、调试和测试
                                  存地址,直接取出来,所以速度非常快。
 田 IO编程
                                  你可以猜到,这种key-value存储方式,在放进去的时候,必须根据key算出value的存放位置,这样,取的时候才能根据key直接拿到value。
 田 进程和线程
                                  把数据放入dict的方法,除了初始化时指定外,还可以通过key放入:
   正则表达式
 田 常用内建模块
                                   >>> d['Adam'] = 67
 田 常用第三方模块
                                   >>> d['Adam']
                                   67
   virtualenv
 田 图形界面
                                  由于一个key只能对应一个value, 所以, 多次对一个key放入value, 后面的值会把前面的值冲掉:
 田 网络编程
                                   >>> d['Jack'] = 90
 田 电子邮件
                                   >>> d['Jack']
 田 访问数据库
                                   >>> d['Jack'] = 88
 田 Web开发
                                   >>> d['Jack']
                                   88
 田 异步IO
 田 实战
                                  如果key不存在, dict就会报错:
   FAQ
                                   >>> d['Thomas']
   期末总结
                                   Traceback (most recent call last):
                                     File "<stdin>", line 1, in <module>
                                   KeyError: 'Thomas'
关于作者
                                  要避免key不存在的错误,有两种办法,一是通过 in 判断key是否存在:
        廖雪峰 🗸 北京 朝阳区
                                   >>> <mark>'Thomas' in d</mark>
                                   False
                                  二是通过dict提供的 get() 方法,如果key不存在,可以返回 None ,或者自己指定的value:
                                   >>> d.get('Thomas')
                                   >>> d.get('Thomas', -1)
                                   -1
                                   75
                                  法(Hash)。
                                  要保证hash的正确性,作为key的对象就不能变。在Python中,字符串、整数等都是不可变的,因此,可以放心地作为key。而list是可变的,就不能作为
                                  key:
                                   >>> key = [1, 2, 3]
                                   >>> d[key] = 'a list'
                                   Traceback (most recent call last):
                                    File "<stdin>", line 1, in <module>
                                   TypeError: unhashable type: 'list'
                                  set
                                  set和dict类似,也是一组key的集合,但不存储value。由于key不能重复,所以,在set中,没有重复的key。
```

```
注意:返回None的时候Python的交互环境不显示结果。
要删除一个key,用 pop(key)方法,对应的value也会从dict中删除:
>>> d.pop('Bob')
 {'Michael': 95, 'Tracy': 85}
请务必注意,dict内部存放的顺序和key放入的顺序是没有关系的。
和list比较, dict有以下几个特点:
1. 查找和插入的速度极快,不会随着key的增加而变慢;
2. 需要占用大量的内存,内存浪费多。
而list相反:
1. 查找和插入的时间随着元素的增加而增加;
2. 占用空间小,浪费内存很少。
所以,dict是用空间来换取时间的一种方法。
dict可以用在需要高速查找的很多地方,在Python代码中几乎无处不在,正确使用dict非常重要,需要牢记的第一条就是dict的key必须是不可变对象。
这是因为dict根据key来计算value的存储位置,如果每次计算相同的key得出的结果不同,那dict内部就完全混乱了。这个通过key计算位置的算法称为哈希算
```

```
>>> s = set([1, 2, 3])
 >>> S
 \{1, 2, 3\}
注意,传入的参数[1, 2, 3]是一个list,而显示的 {1, 2, 3} 只是告诉你这个set内部有1, 2, 3这3个元素,显示的顺序也不表示set是有序的。。
重复元素在set中自动被过滤:
 >>> s = set([1, 1, 2, 2, 3, 3])
 >>> S
 \{1, 2, 3\}
```

```
通过 remove(key) 方法可以删除元素:
 >>> s.remove(4)
>>> S
\{1, 2, 3\}
set可以看成数学意义上的无序和无重复元素的集合,因此,两个set可以做数学意义上的交集、并集等操作:
>>> s1 = set([1, 2, 3])
```

```
set和dict的唯一区别仅在于没有存储对应的value,但是,set的原理和dict一样,所以,同样不可以放入可变对象,因为无法判断两个可变对象是否相等,也
就无法保证set内部"不会有重复元素"。试试把list放入set,看看是否会报错。
```

```
>>> a.sort()
>>> a
['a', 'b', 'c']
```

>>> a = ['c', 'b', 'a']

>>> s.add(4)

 $\{1, 2, 3, 4\}$ >>> s.add(4)

 $\{1, 2, 3, 4\}$

>>> s2 = set([2, 3, 4])

>>> s1 & s2

>>> s1 | s2 $\{1, 2, 3, 4\}$

再议不可变对象

{2, 3}

>>> S

>>> S

而对于不可变对象,比如str,对str进行操作呢:

上面我们讲了,str是不变对象,而list是可变对象。

对于可变对象,比如list,对list进行操作,list内部的内容是会变化的,比如:

要创建一个set,需要提供一个list作为输入集合:

通过 add(key) 方法可以添加元素到set中,可以重复添加,但不会有效果:

```
>>> a = 'abc'
>>> a.replace('a', 'A')
'Abc'
>>> a
'abc'
```

>>> a 'abc'

虽然字符串有个 replace() 方法,也确实变出了 'Abc',但变量 a 最后仍是 'abc',应该怎么理解呢?

```
象的内容才是 'abc':
```

'abc'

'Abc'

我们先把代码改成下面这样:

>>> b = a.replace('a', 'A')

>>> a = 'abc'

>>> b 'Abc'

'abc'

所以,对于不变对象来说,调用对象自身的任意方法,也不会改变该对象自身的内容。相反,这些方法会创建新的对象并返回,这样,就保证了不可变对象

当我们调用 a. replace('a', 'A') 时,实际上调用方法 replace 是作用在字符串对象 'abc'上的,而这个方法虽然名字叫 replace,但却没有改变字符

串'abc'的内容。相反,replace 方法创建了一个新字符串'Abc'并返回,如果我们用变量b指向该新字符串,就容易理解了,变量a仍指向原有的字符

要始终牢记的是,「a」是变量,而「abc」才是字符串对象!有些时候,我们经常说,对象「a」的内容是「abc」,但其实是指,「a」本身是一个变量,它指向的对

```
小结
使用key-value存储结构的dict在Python中非常有用,选择不可变对象作为key很重要,最常用的key是字符串。
tuple虽然是不变对象,但试试把 (1, 2, 3) 和 (1, [2, 3]) 放入dict或set中,并解释结果。
```

本身永远是不可变的。

参考源码

the_dict.py

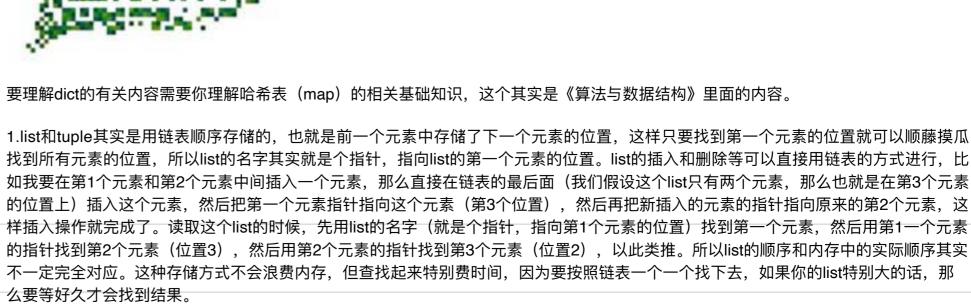
the_set.py

串 'abc', 但变量 b 却指向新字符串 'Abc' 了:

读后有收获可以支付宝请作者喝咖啡,读后有疑问请加微信群讨论:



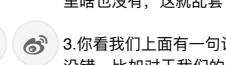
Comments



的指针找到第2个元素(位置3),然后用第2个元素的指针找到第3个元素(位置2),以此类推。所以list的顺序和内存中的实际顺序其实 不一定完全对应。这种存储方式不会浪费内存,但查找起来特别费时间,因为要按照链表一个一个找下去,如果你的list特别大的话,那

2.dict则为了快速查找使用了一种特别的方法,哈希表。哈希表采用哈希函数从key计算得到一个数字(哈希函数有个特点:对于不同的 key,有很大的概率得到的哈希值也不同),然后直接把value存储到这个数字所对应的地址上,比如key='ABC',value=10,经过哈希函

数得到key对应的哈希值为123,那么就申请一个有1000个地址(从0到999)的内存,然后把10存放在地址为123的地方。类似的,对于 Make a comment key='BCD', value=20, 得到key的哈希值为234, 那么就把20存放在地址为234的地方。对于这样的表查找起来是非常方便的。只要给出 key,计算得到哈希值,然后直接到对应的地址去找value就可以了。无论有几个元素,都可以直接找到value,无需遍历整个表。不过虽 Sign in to make a comment 然dict查找速度快,但内存浪费严重,你看我们只存储了两个元素,都要申请一个长度为1000的内存。 3.现在你知道为啥key要用不可变对象了吧?因为不可变对象是常量,每次的哈希值算出来都是固定的,这样就不会出错。比如 key='ABC', value=10, 存储地址为123, 假设我突发奇想, 把key改成'BCD', 那么当查找'BCD'的value的时候就会去234的地址找, 但那 里啥也没有,这就乱套了。



友情链接: 中华诗幕後 - 阿里云 - SICP - 4clojure

3.你看我们上面有一句话:对于不同的key,有很大的概率得到的哈希值也不同。那么有很小的概率不同的key可以得到相同的哈希值了? 没错,比如对于我们的例子来说,哈希值只有3位,那么只要元素个数超过1000,就一定会有至少两个key的哈希值相同(鸽笼原理) 这种情况叫"冲突",设计哈希表的时候要采取办法减少冲突,实在冲突了也要想办法补救。不过这是编译器的事情,况且对于初学者的我 们来说碰到的冲突的概率基本等于零,就不用操心了。

```
本网站运行在阿里云上并使用阿里云CDN加速。
```

廖雪峰的官方网站©2019

Powered by iTranswarp