

基本的哈希表数据结构

马浩琨

Sep 14, 2025

在日常计算中，我们经常需要快速查找和存储数据，例如在字典中查询单词、在电话簿中寻找联系人，或在搜索引擎中检索信息。这些场景都要求高效的数据访问机制。传统的数据结构如数组或链表，在查找操作时可能需要遍历整个集合，导致时间复杂度为 $O(n)$ ，这在数据量大的情况下效率低下。二叉搜索树虽然将查找效率提升到 $O(\log n)$ ，但仍然存在改进空间。哈希表作为一种经典的数据结构，通过巧妙的「时空权衡」策略，在平均情况下能够实现近乎 $O(1)$ 的查找、插入和删除操作，从而成为许多高性能应用的核心组件。本文将带领读者从理论出发，逐步实现一个功能完整的哈希表，使用 Python 语言进行演示，并深入分析其性能特性。

1 哈希表的核心概念

哈希表是一种基于哈希函数构建的数据结构，它通过将键映射到数组中的特定位置来实现快速数据访问。其核心组件包括哈希函数、底层数组和冲突解决策略。哈希函数的作用是将任意大小的输入键转换为一个固定大小的整数值，即哈希值。理想的哈希函数应具备确定性、高效性和均匀性。确定性确保相同的键总是产生相同的哈希值；高效性要求计算速度快；均匀性则保证哈希值分布均匀，以减少冲突的发生。

在哈希表中，我们使用数组作为底层存储，每个数组元素称为一个「哈希桶」。通过哈希函数计算出的哈希值需要映射到数组索引，通常使用取模运算： $index = \text{hash}(key) \bmod \text{capacity}$ ，其中 `capacity` 表示数组的容量。为了管理哈希表的性能，我们引入负载因子的概念，定义为当前元素数量与数组容量的比值： $\text{load_factor} = \frac{\text{number_of_entries}}{\text{capacity}}$ 。负载因子是触发动态扩容的关键指标，当它超过预设阈值时，哈希表需要进行扩容以维持高效操作。

2 哈希冲突：不可避免的挑战

尽管哈希函数旨在均匀分布键，但由于「鸽巢原理」的存在，冲突是不可避免的。鸽巢原理指出，如果键的数量超过数组容量，那么至少有两个键会映射到同一索引。因此，哈希表必须采用有效的冲突解决策略。主流方法包括链地址法和开放定址法。

链地址法通过将每个哈希桶实现为一个链表来处理冲突。所有映射到同一索引的键值对都存储在该链表中。插入操作时，新元素被追加到链表末尾；查找或删除操作时，需要遍历链表以匹配键。这种方法的优点是实现简单且能有效处理冲突，但缺点是需要额外的指针空间，且缓存性能较差。

开放定址法则将所有元素直接存储在数组中，当冲突发生时，按照某种探测序列寻找下一个空闲槽位。常见探测方法包括线性探测、平方探测和双重哈希。线性探测使用公式 $index = (\text{hash}(key) + i) \bmod \text{capacity}$ 进行顺序查找；平方探测使用 $index = (\text{hash}(key) + i^2) \bmod \text{capacity}$ 以减少聚集现象；双重哈希则引入第二个哈希函数计算步长。开放定址法的优点是数据集中在数组中，缓存友好，但实现复杂，删除操作需使用标记

删除法，且容易产生聚集。

3 动手实现一个哈希表（采用链地址法）

我们将使用 Python 实现一个基于链地址法的哈希表。首先定义 `HashMap` 类，初始化时设置初始容量和负载因子阈值。代码中，我们使用列表来表示哈希桶，每个桶初始化为空列表。

```
1 class HashMap:
2     def __init__(self, capacity=10, load_factor_threshold=0.75):
3         self.capacity = capacity
4         self.load_factor_threshold = load_factor_threshold
5         self.size = 0
6         self.buckets = [[] for _ in range(self.capacity)]
```

这里，`capacity` 默认值为 10，`load_factor_threshold` 设置为 0.75。`buckets` 是一个列表的列表，每个元素代表一个哈希桶，初始为空。`size` 记录当前元素数量，用于计算负载因子。

接下来实现 `set` 方法用于插入或更新键值对。首先计算键的哈希值并映射到索引，然后遍历对应桶的链表。如果找到相同键，则更新值；否则追加新键值对。完成后检查负载因子，若超过阈值则调用 `_resize` 方法扩容。

```
1     def set(self, key, value):
2         index = hash(key) % self.capacity
3         bucket = self.buckets[index]
4         for i, (k, v) in enumerate(bucket):
5             if k == key:
6                 bucket[i] = (key, value)
7                 return
8         bucket.append((key, value))
9         self.size += 1
10        if self.load_factor > self.load_factor_threshold:
11            self._resize()
```

`hash(key)` 是 Python 内置函数，返回键的哈希值。取模运算确保索引在数组范围内。遍历桶中元素时，使用枚举来访问索引和键值对。如果键已存在，更新值并返回；否则添加新元素并增加 `size`。负载因子通过属性计算，我们稍后定义。

`get` 方法用于查找键对应的值。计算索引后，遍历桶链表，找到匹配键则返回值，否则返回 `None`。

```
1     def get(self, key):
2         index = hash(key) % self.capacity
3         bucket = self.buckets[index]
4         for k, v in bucket:
5             if k == key:
6                 return v
7         return None
```

`delete` 方法类似，遍历桶链表找到并删除键值对。

```

1 def delete(self, key):
2     index = hash(key) % self.capacity
3     bucket = self.buckets[index]
4     for i, (k, v) in enumerate(bucket):
5         if k == key:
6             del bucket[i]
7             self.size -= 1
8             return

```

动态扩容通过 `_resize` 方法实现。当负载因子超过阈值时，创建一个新数组，容量通常翻倍，然后重新哈希所有元素到新数组中。

```

1 def _resize(self):
2     new_capacity = self.capacity * 2
3     new_buckets = [[] for _ in range(new_capacity)]
4     for bucket in self.buckets:
5         for key, value in bucket:
6             index = hash(key) % new_capacity
7             new_buckets[index].append((key, value))
8     self.buckets = new_buckets
9     self.capacity = new_capacity

```

这里，新容量为原容量的两倍。遍历所有旧桶中的键值对，重新计算哈希索引并插入新桶。最后更新哈希表的 `buckets` 和 `capacity`。

此外，我们可以添加辅助方法如 `load_factor` 属性和 `len` 方法。

```

1 @property
2 def load_factor(self):
3     return self.size / self.capacity
4
5 def __len__(self):
6     return self.size

```

`load_factor` 使用属性装饰器实现计算属性，返回当前负载因子。`len` 方法返回元素数量，便于使用 `len()` 函数。

4 性能分析与进阶话题

哈希表的性能高度依赖于哈希函数和冲突解决策略。在最佳情况下，当键均匀分布时，操作时间复杂度为 $O(1)$ ；最坏情况下，所有键冲突，性能退化为 $O(n)$ 。平均情况下，基于负载因子，时间复杂度通常为 $O(1)$ 。负载因子的选择直接影响性能；例如，阈值 0.75 在空间和时间效率之间提供了平衡。

设计一个好的哈希函数是关键。虽然加密哈希函数如 MD5 或 SHA-1 提供均匀分布，但计算开销大，不适合哈希表。实用哈希函数如 FNV-1 或 DJB2 针对字符串优化，能高效产生均匀哈希值。例如，DJB2 哈希函数使用迭代计算： $hash = ((hash \ll 5) + hash) + char$ ，其中 $char$ 是字符串的字符代码。

冲突解决策略的选择也影响性能。链地址法简单可靠，但额外空间开销；开放定址法缓存友好，但容易聚集。在实际应用中，如 Python 的 dict 或 Java 的 HashMap，都采用优化策略结合动态扩容以确保高效性。

哈希表以其高效的查找、插入和删除操作，成为计算机科学中不可或缺的数据结构。优势在于平均 $O(1)$ 的时间复杂度，但劣势包括不支持顺序遍历和性能对哈希函数的依赖。经典应用场景包括数据库索引、缓存系统如 Redis、字典实现以及对象属性存储。通过亲手实现哈希表，读者可以更深入理解其内部机制，并尝试优化如哈希函数选择或冲突策略，以提升实践能力。哈希表不仅是理论上的经典，更是工程中的利器，值得每一位开发者掌握。