

基本的布隆过滤器 (Bloom Filter) 数据结构

王思成

Nov 04, 2025

在当今数据爆炸的时代，处理海量信息的存在性判断成为一个常见挑战。例如，在恶意 URL 检测场景中，我们可能有一个包含十亿个 URL 的黑名单，如何快速判断一个新 URL 是否在其中？如果使用传统的哈希表，内存占用可能高达数十 GB，这在实际应用中往往不可行。数据库查询虽然能节省内存，但查询速度较慢，无法满足实时性要求。布隆过滤器应运而生，它通过巧妙的概率设计，用极小的空间开销和恒定的查询时间，解决了「可能存在」和「肯定不存在」的问题。其代价是允许一定的误判率，但绝不漏判，这使得它在许多场景中成为理想选择。

1 什么是布隆过滤器？

布隆过滤器的核心思想可以概括为一个由随机映射函数和二进制向量组成概率型数据结构。想象一个多指纹采集器，每个元素进来时，会通过多个哈希函数生成多个指纹，并将这些指纹记录在一个比特数组中；判断时，只需检查所有指纹是否都存在。这种设计使得布隆过滤器在空间效率和查询速度上远超一般算法，同时不需要存储元素本身，具有一定保密性。然而，它也存在缺点，例如存在误判率，无法删除元素，且不支持获取所有元素。这些特性决定了它的适用边界，在允许一定误差的场景中表现卓越。

2 深入原理：布隆过滤器如何工作？

布隆过滤器的核心组件包括一个长度为 m 的比特数组和 k 个相互独立且分布均匀的哈希函数。比特数组初始状态所有位都为 0，而哈希函数将输入元素映射到 $[0, m-1]$ 范围内的整数索引。添加元素时，例如添加字符串 geekhour，首先将其输入 k 个哈希函数，得到 k 个哈希值 h_1, h_2, \dots, h_k ，然后将比特数组中对应这些位置的位全部设置为 1。查询元素时，同样计算 k 个哈希值，并检查这些位置是否全部为 1；如果全是 1，则返回「可能存在」；如果有任何一位为 0，则返回「肯定不存在」。这种设计确保了绝不漏判，因为如果一个元素确实被添加过，它设置的位不可能被清零，查询时一定能看到所有位都是 1。误判的发生是由于哈希冲突，导致不同元素的哈希值可能覆盖相同位置。

3 关键参数与误判率分析

布隆过滤器的误判率受三个核心参数影响：预期添加的元素数量 n 、比特数组长度 m 和哈希函数个数 k 。误判率的近似公式为 $P \approx (1 - e^{(-k * n / m)})^k$ ，这个公式直观地展示了参数间的权衡。当 m 增大时，比特数组空间更大，冲突减少，误判率降低；当 n 增大时，元素越多，数组越满，误判率升高；而 k 的个数需要平衡，太少容易冲突，太多则数组快速饱和。实践中，最优哈希函数个数可通过公式 $k = (m / n) * \ln(2)$ 计算。例如，如果

要求误判率为 1%，通常推荐 m 约为 n 的 10 倍，此时最优 k 值约为 7。这种参数选择帮助用户在空间和准确性之间找到平衡点。

4 动手实现：一个简单的布隆过滤器

我们将使用 Python 实现一个简单的布隆过滤器类，命名为 SimpleBloomFilter。这个类包含三个成员变量：`size` 表示比特数组大小 m ，`hash_count` 表示哈希函数个数 k ，以及 `bit_array` 用于模拟比特数组。在 Python 中，我们可以使用内置的 `bytearray` 或第三方库如 `bitarray` 来高效管理位操作，但为了简单起见，这里用整数列表模拟，每个整数代表一个位。

首先，我们需要模拟多个哈希函数。一个高效的方法是使用双重哈希技术，基于一个基础哈希函数生成 k 个不同的哈希值。具体公式为 $\text{hash}_i(x) = (\text{hash1}(x) + i * \text{hash2}(x)) \% \text{size}$ ，其中 `hash1` 和 `hash2` 可以是简单的哈希函数如 FNV-1a 或 MurmurHash 的变体。这种方法确保了哈希值的独立性和均匀分布，同时避免了实现多个独立哈希函数的复杂性。

初始化方法 `__init__(self, size, hash_count)` 负责设置比特数组大小和哈希函数个数，并将比特数组初始化为全 0。在代码中，我们用一个长度为 `size` 的列表表示，每个元素初始为 0。

添加方法 `add(self, item)` 首先将输入项转换为字符串或字节，以确保一致性。然后，计算 k 个位位置，通过循环调用哈希函数生成索引，并将比特数组中对应位置设为 1。在实现中，我们使用位操作来高效设置位，例如通过位或运算。

查询方法 `contains(self, item)` 类似地计算 k 个位位置，并检查这些位置是否全为 1；如果是，返回 `True`，表示可能存在；否则返回 `False`，表示肯定不存在。

以下是一个简单的 Python 实现示例：

```
1 class SimpleBloomFilter:
2     def __init__(self, size, hash_count):
3         self.size = size
4         self.hash_count = hash_count
5         self.bit_array = [0] * size # 初始化比特数组为全 0
6
7     def _hashes(self, item):
8         # 模拟双重哈希生成 k 个哈希值
9         item = str(item).encode('utf-8')
10        h1 = hash(item) # 使用 Python 内置哈希作为基础，实际应用中应使用更均匀的哈希函数
11        h2 = h1 ^ 0xFFFFFFFF # 简单变换生成第二个哈希值
12        hashes = []
13        for i in range(self.hash_count):
14            hashes.append((h1 + i * h2) % self.size)
15        return hashes
16
17    def add(self, item):
18        for pos in self._hashes(item):
```

```
19     self.bit_array[pos] = 1 # 设置对应位为 1  
20  
21     def contains(self, item):  
22         for pos in self._hashes(item):  
23             if self.bit_array[pos] == 0:  
24                 return False # 如果任何一位为 0, 肯定不存在  
25         return True # 所有位为 1, 可能存在
```

在这段代码中，初始化方法创建了一个指定大小的列表作为比特数组。哈希生成方法 `_hashes` 将输入项转换为字节后，使用 Python 内置哈希函数生成基础值，并通过线性组合生成 k 个索引；实际应用中，建议使用更均匀的哈希函数如 FNV-1a 以减少冲突。添加方法遍历哈希值，将对应位置设为 1；查询方法检查所有位置，如果任何一位为 0 则立即返回 `False`。这种实现虽然简单，但清晰地展示了布隆过滤器的核心逻辑。

为了测试这个实现，我们可以创建一个布隆过滤器实例，设置 $m=1000$ 和 $k=7$ ，然后添加一组单词如 `[hello, world, bloom, filter]`。测试 `contains(hello)` 应返回 `True`，而 `contains(foo)` 应返回 `False`。随着添加元素增多，误判可能发生，例如在添加大量元素后，查询一个未添加项可能返回 `True`，这演示了布隆过滤器的概率特性。

5 进阶话题与优化

布隆过滤器有多种变体，例如计数布隆过滤器，它将比特位扩展为计数器，支持元素的删除操作。在计数布隆过滤器中，添加元素时计数器加一，删除时减一，但这增加了空间开销，并可能引发计数溢出问题。应用场景方面，布隆过滤器广泛用于缓存穿透防护，通过在查询数据库前判断键是否存在，避免无效查询；在爬虫 URL 去重中，快速判断 URL 是否已抓取；在垃圾邮件过滤中，检查发件人是否在黑名单内；以及在数据库如 LevelDB 中优化查询，减少不必要的磁盘访问。这些应用凸显了布隆过滤器在高效处理大规模数据时的价值。

回顾布隆过滤器的核心，它通过空间换时间和概率交换确定性的哲学，解决了传统方法在高并发和大数据场景下的瓶颈。使用时需权衡其优缺点：适用于对空间敏感、允许误判的快速查询场景；而不适用于要求百分百准确、需要删除操作或存储元素本身的场景。鼓励读者在项目中探索布隆过滤器的应用，例如在分布式系统中使用变体优化性能。进一步学习可关注高效哈希函数或分布式布隆过滤器的实现，以深化对这一数据结构的理解。