


```
127.0.0.1:6379> pfadd key "redis"
(integer) 1
127.0.0.1:6379> pfadd key "java" "sql"
(integer) 1
```

```
pfadd key element [element ...]
```

统计不重复的元素

```
127.0.0.1:6379> pfadd key "redis"
(integer) 1
127.0.0.1:6379> pfadd key "sql"
(integer) 1
127.0.0.1:6379> pfadd key "redis"
(integer) 0
127.0.0.1:6379> pfcount key
(integer) 2
```

pfcount **key** [**key** ...]

合并一个或多个 HLL 至新结构

```
127.0.0.1:6379> pfadd k "java" "sql"
(integer) 1
```

```
127.0.0.1:6379> pfadd k2 "redis" "sql"
(integer) 1
127.0.0.1:6379> pfmerge k3 k k2
OK
127.0.0.1:6379> pfcount k3
(integer) 3
```

相关语法：

```
pfmerge destkey sourcekey [sourcekey ...]
```

pfmerge 使用场景

当我们需要合并两个或多个同类页面的访问数据时，我们可以使用 pfmerge 来操作。

代码实战

接下来我们使用 Java 代码来实现 HLL 的三个基础功能，代码如下：

```
import redis.clients.jedis.Jedis;

public class HyperLogLogExample {
    public static void main(String[] args) {
        Jedis jedis = new Jedis("127.0.0.1", 6379);
        // 添加元素
        jedis.pfadd("k", "redis", "sql");
        jedis.pfadd("k", "redis");
        // 统计元素
        long count = jedis.pfcount("k");
        // 打印统计元素
        System.out.println("k: " + count);
        // 合并 HLL
        jedis.pfmerge("k2", "k");
        // 打印新 HLL
        System.out.println("k2: " + jedis.pfcount("k2"));
    }
}
```

以上代码执行结果如下：

```
k: 2
k2: 2
```

HLL 算法原理

HyperLogLog 算法来源于论文 [HyperLogLog the analysis of a near-optimal cardinality estimation algorithm](#)，想要了解 HLL 的原理，先要从伯努利试验说起，伯努利实验说的是抛硬币的事。一次伯努利实验相当于抛硬币，不管抛多少次只要出现一个正面，就称为一次伯努利实验。

我们用 k 来表示每次抛硬币的次数， n 表示第几次抛的硬币，用 k_{\max} 来表示抛硬币的最高次数，最终根据估算发现 n 和 k_{\max} 存在的关系是 $n=2^{(k_{\max})}$ ，但同时我们也发现了另一个问题当试验次数很小的时候，这种估算方法的误差会很大，例如我们进行以下 3 次实验：

- 第 1 次试验：抛 3 次出现正面，此时 $k=3$ ， $n=1$ ；
- 第 2 次试验：抛 2 次出现正面，此时 $k=2$ ， $n=2$ ；
- 第 3 次试验：抛 6 次出现正面，此时 $k=6$ ， $n=3$ 。

对于这三组实验来说， $k_{\max}=6$ ， $n=3$ ，但放入估算公式明显 $3 \neq 2^6$ 。为了解决这个问题 HLL 引入了分桶算法和调和平均数来使这个算法更接近真实情况。

分桶算法是指把原来的数据平均分为 m 份，在每段中求平均数在乘以 m ，以此来消减因偶然性带来的误差，提高预估的准确性，简单来说就是把一份数据分为多份，把一轮计算，分为多轮计算。

而调和平均数指的是使用平均数的优化算法，而非直接使用平均数。

例如小明的月工资是 1000 元，而小王的月工资是 100000 元，如果直接取平均数，那小明的平均工资就变成了 $(1000+100000)/2=50500$ 元，这显然是不准确的，而使用调和平均数算法计算的结果是 $2/(1/1000+1/100000) \approx 1998$ 元，显然此算法更符合实际平均数。

所以综合以上情况，在 Redis 中使用 HLL 插入数据，相当于把存储的值经过 hash 之后，再将 hash 值转换为二进制，存入到不同的桶中，这样就可以用很小的空间存储很多的数据，统计时再去相应的位置进行对比很快就能得出结论，这就是 HLL 算法的基本原理，想要更深入的了解算法及其推理过程，可以看去原版的论文，链接地址在文末。

小结

当需要做大量数据统计时，普通的集合类型已经不能满足我们的需求了，这个时候我们可以借助 Redis 2.8.9 中提供的 HyperLogLog 来统计，它的优点是只需要使用 12k 的空间就能

统计 2^{64} 的数据，但它的缺点是存在 0.81% 的误差，HyperLogLog 提供了三个操作方法 pfadd 添加元素、pfcount 统计元素和 pfmerge 合并元素。

参考文献

- 论文 [HyperLogLog: the analysis of a near-optimal cardinality estimation algorithm](#)

[上一页](#)

[下一页](#)