

# 为什么GROUP BY比DISTINCT快3倍？90%的程序员都踩过这个坑！

原创 杨乐乐哉 乐哥聊编程 2025年03月12日 21:00 安徽



本站300+篇技术文档已打包  
本站300+篇技术文档已打包  
本站300+篇技术文档已打包



乐哥聊编程

专注编程分享，聊程序人生

444篇原创内容

公众号

## 引言：去重操作，你真的用对了吗？

“为什么同样的去重需求，同事的代码运行如飞，你的却慢如蜗牛？”

一位开发者在处理百万级数据时，发现 `SELECT DISTINCT` 耗时高达10秒，而改用 `GROUP BY` 后仅需3秒——**性能差距超过3倍！**

**90%的程序员误以为两者功能完全相同**，但背后的机制天差地别。本文将彻底揭秘它们的核心区别，并教你如何避免性能陷阱！

## 一、功能解析：GROUP BY和DISTINCT的职责边界

### 1. DISTINCT：简单粗暴的去重工具

- **核心功能**：对查询结果的所有列进行全局去重。
- **典型场景**：

```
-- 查询所有不重复的城市
SELECT DISTINCT city FROM users;

-- 查询不重复的城市和性别组合
SELECT DISTINCT city, gender FROM users;
```

- **本质**：将结果集加载到临时表，排序后去除重复行。

## 2. GROUP BY：分组聚合的多面手

- **核心功能**：按指定字段分组，通常结合聚合函数（如 COUNT，SUM）使用。
- **典型场景**：

```
-- 统计每个城市的用户数
SELECT city, COUNT(*) FROM users GROUP BY city;

-- 计算每个部门的平均薪资
SELECT department, AVG(salary) FROM employees GROUP BY department;
```

- **隐藏技能**：当不搭配聚合函数时，可实现类似 DISTINCT 的去重效果。

## 二、本质区别：为何GROUP BY能吊打DISTINCT？

### 1. 执行计划对比：临时表与索引的生死博弈

- **DISTINCT的执行流程**（无索引时）：

```
-- EXPLAIN结果: Using temporary; Using filesort
EXPLAIN SELECT DISTINCT city FROM users;
```

1. **全表扫描**：读取所有数据到内存或磁盘临时表。
2. **排序去重**：对临时表按所有字段排序，遍历去除重复行。
3. **资源消耗**：若数据量超过 tmp\_table\_size，临时表写入磁盘，触发I/O风暴。

- **GROUP BY的执行流程**（有索引时）：

```
-- 添加索引
ALTER TABLE users ADD INDEX idx_city (city);

-- EXPLAIN结果: Using index (无临时表，无排序)
EXPLAIN SELECT city FROM users GROUP BY city;
```

1. **索引扫描**：直接遍历 city 字段的B+树索引（有序）。
2. **跳跃去重**：利用索引有序性，跳过重复值，无需全量排序。
3. **零回表**：若仅查询分组字段，直接从索引取数据，无需访问原表。

2. 性能差距的三大核心原因

- 原因1：索引利用率
  - GROUP BY 可利用索引跳过排序和临时表，而 DISTINCT 即便有索引也可能触发全表扫描。
- 原因2：内存与磁盘的较量
  - DISTINCT 的临时表默认使用磁盘存储（MyISAM引擎），GROUP BY 优先使用内存（MEMORY引擎）。
- 原因3：隐式排序的代价

```
SELECT city FROM users GROUP BY city ORDER BY NULL; -- 性能再提升20%
```

- DISTINCT 必须全排序，而 GROUP BY 可通过 ORDER BY NULL 禁用排序：

3. 实测数据：千万级表性能对比

| 查询语句                                 | 无索引耗时 | 有索引耗时 |
|--------------------------------------|-------|-------|
| SELECT DISTINCT city FROM users      | 12.3s | 9.8s  |
| SELECT city FROM users GROUP BY city | 11.7s | 2.1s  |

- 结论：
  - 无索引时：两者性能接近，但 GROUP BY 略优（省去显式排序）。
  - 有索引时：GROUP BY 性能提升超过3倍！

三、避坑指南：这些场景GROUP BY反而更危险！

1. 分组字段无索引时

- 问题：GROUP BY 同样需要全表扫描+临时表，性能与 DISTINCT 无异。
- 优化方案：

```
-- 为分组字段添加索引
ALTER TABLE orders ADD INDEX idx_product (product_id);
```

2. 查询包含非分组字段

- 错误示例：

```
-- 查询结果不确定 (MySQL可能随机返回user_id)
SELECT city, user_id FROM users GROUP BY city;
```

- 正确写法：

```
SELECT city, MAX(user_id) FROM users GROUP BY city; -- 明确聚合规则
```

### 3. 分组的字段含复杂计算

- 错误示例：

```
-- 索引失效，触发全表扫描
SELECT YEAR(create_time) FROM orders GROUP BY YEAR(create_time);
```

- 优化方案：

```
-- 新增冗余字段并添加索引
ALTER TABLE orders ADD COLUMN create_year INT, ADD INDEX idx_year (create_year);
UPDATE orders SET create_year = YEAR(create_time);
```

## 四、终极实战：亿级数据去重优化方案

### 1. 索引设计黄金法则

- 覆盖索引：确保查询字段全部包含在索引中。

```
-- 为city和gender创建联合索引
ALTER TABLE users ADD INDEX idx_city_gender (city, gender);
-- 以下查询直接走索引
SELECT city, gender FROM users GROUP BY city, gender;
```

### 2. 参数调优：强制临时表驻留内存

```
-- 调整会话级内存参数（单位：字节）
SET tmp_table_size = 256 * 1024 * 1024; -- 256MB
SET max_heap_table_size = 256 * 1024 * 1024;
```

### 3. 分布式方案：Elasticsearch预聚合

- 适用场景：实时统计海量数据的唯一值。
- 操作步骤：

```
GET /users/_search
{
  "size": 0,
  "aggs": {
```

```
"unique_cities": {  
  "terms": { "field": "city.keyword", "size": 1000 }  
}  
}  
}
```

1. 将去重字段（如 city）同步到Elasticsearch。
2. 使用 terms aggregation 实现毫秒级去重：

## 五、灵魂拷问：你真的需要去重吗？

- **场景1**：仅统计不重复数量 → `SELECT COUNT(DISTINCT city)`。
- **场景2**：需要列出所有不重复值 → 优先用 `GROUP BY` +索引。
- **场景3**：去重后需复杂计算 → 使用窗口函数（如 `ROW_NUMBER()`）。

## 结尾互动

你在使用GROUP BY和DISTINCT时踩过哪些坑？欢迎留言分享！



乐哥聊编程

专注编程分享，聊程序人生

444篇原创内容

公众号

面试 55    MySQL 精讲系列 9

面试 · 目录

上一篇

高并发下MySQL扛不住了？6招让大表性能飙升300%！

下一篇

高并发崩溃的元凶竟是它？揭秘Java中volatile的致命陷阱！

