

# 索引常见面试题 | 小林coding | Java面试学习

 [xiaolincoding.com/mysql/index/index\\_interview.html](https://xiaolincoding.com/mysql/index/index_interview.html)

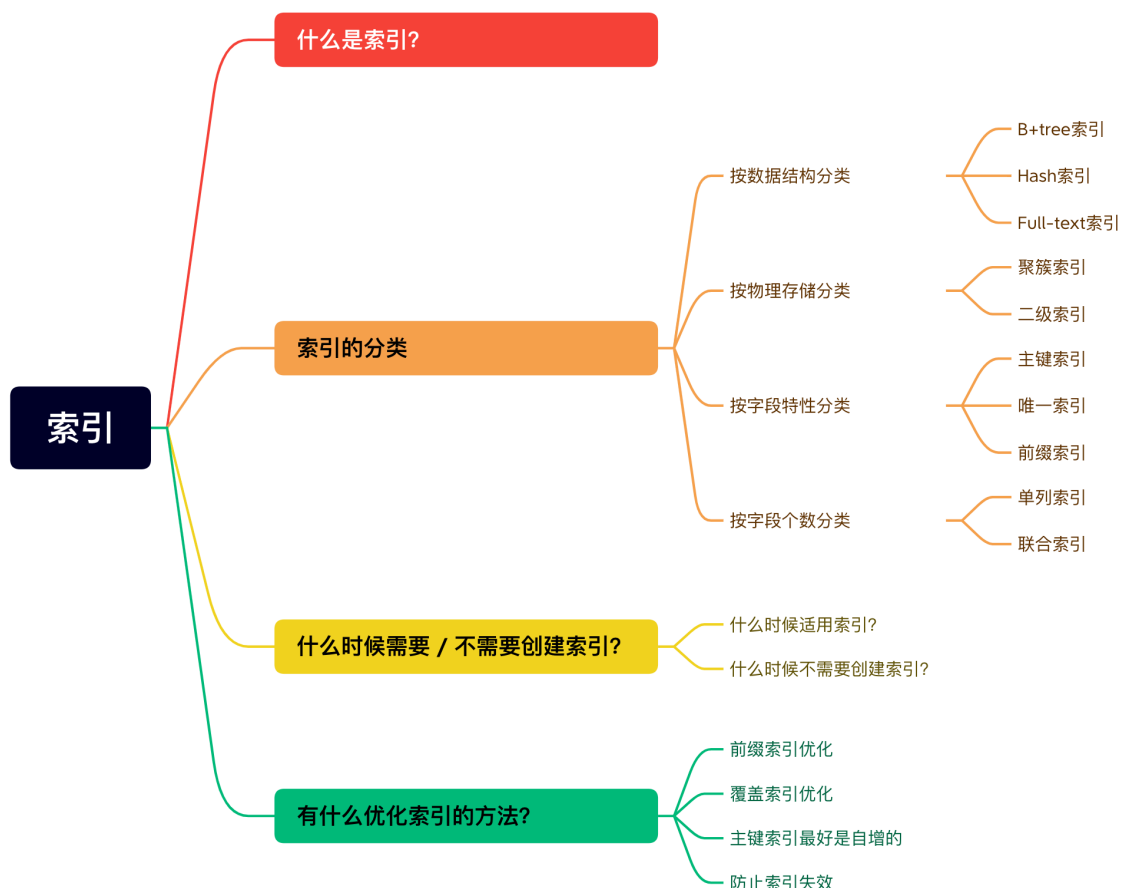
## # 索引常见面试题

大家好，我是小林。

面试中，MySQL 索引相关的问题基本都是一系列问题，都是先从索引的基本原理，再到索引的使用场景，比如：

- 索引底层使用了什么数据结构和算法？
- 为什么 MySQL InnoDB 选择 B+tree 作为索引的数据结构？
- 什么时候适用索引？
- 什么时候不需要创建索引？
- 什么情况下索引会失效？
- 有什么优化索引的方法？
- .....

今天就带大家，夯实 MySQL 索引的知识点。



## # 什么是索引？

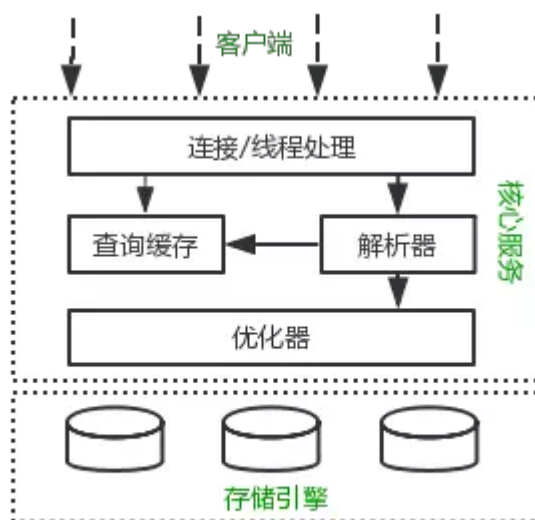
当你想查阅书中某个知识的内容，你会选择一页一页的找呢？还是在书的目录去找呢？

傻瓜都知道时间是宝贵的，当然是选择在书的目录去找，找到后再翻到对应的页。书中的目录，就是充当索引的角色，方便我们快速查找书中的内容，所以索引是以空间换时间的设计思想。

那换到数据库中，索引的定义就是帮助存储引擎快速获取数据的一种数据结构，形象的说就是**索引是数据的目录**。

所谓的存储引擎，说白了就是如何存储数据、如何为存储的数据建立索引和如何更新、查询数据等技术的实现方法。MySQL 存储引擎有 MyISAM、InnoDB、Memory，其中 InnoDB 是在 MySQL 5.5 之后成为默认的存储引擎。

下图是 MySQL 的结构图，索引和数据就是位于存储引擎中：



MySQL逻辑架构，来自：高性能MySQL

## # 索引的分类

你知道索引有哪些吗？大家肯定都能霹雳啪啦地说出聚簇索引、主键索引、二级索引、普通索引、唯一索引、hash索引、B+树索引等等。

然后再问你，你能将这些索引分一下类吗？可能大家就有点模糊了。其实，要对这些索引进行分类，要清楚这些索引的使用和实现方式，然后再针对有相同特点的索引归为一类。

我们可以按照四个角度来分类索引。

- 按「数据结构」分类：**B+tree索引、Hash索引、Full-text索引**。
- 按「物理存储」分类：**聚簇索引（主键索引）、二级索引（辅助索引）**。
- 按「字段特性」分类：**主键索引、唯一索引、普通索引、前缀索引**。
- 按「字段个数」分类：**单列索引、联合索引**。

接下来，按照这些角度来说说各类索引的特点。

## # 按数据结构分类

从数据结构的角度来看，MySQL 常见索引有 B+Tree 索引、HASH 索引、Full-Text 索引。

每一种存储引擎支持的索引类型不一定相同，我在表中总结了 MySQL 常见的存储引擎 InnoDB、MyISAM 和 Memory 分别支持的索引类型。

索引类型	InnoDB 引擎	MyISAM 引擎	Memory 引擎
B+Tree 索引	Yes	Yes	Yes
HASH 索引	No (不支持hash索引，但是在内存结构中有一个自适应hash索引)	No	Yes
Full-Text 索引	Yes (MySQL 5.6 版本后支持)	Yes	No

InnoDB 是在 MySQL 5.5 之后成为默认的 MySQL 存储引擎，B+Tree 索引类型也是 MySQL 存储引擎采用最多的索引类型。

在创建表时，InnoDB 存储引擎会根据不同的场景选择不同的列作为索引：

- 如果有主键，默认会使用主键作为聚簇索引的索引键（key）；
- 如果没有主键，就选择第一个不包含 NULL 值的唯一列作为聚簇索引的索引键（key）；
- 在上面两个都没有的情况下，InnoDB 将自动生成一个隐式自增 id 列作为聚簇索引的索引键（key）；

其它索引都属于辅助索引（Secondary Index），也被称为二级索引或非聚簇索引。**创建的主键索引和二级索引默认使用的是 B+Tree 索引。**

为了让大家理解 B+Tree 索引的存储和查询的过程，接下来我通过一个简单例子，说明一下 B+Tree 索引在存储数据中的具体实现。

先创建一张商品表，id 为主键，如下：

```
CREATE TABLE `product` (  
  `id` int(11) NOT NULL,  
  `product_no` varchar(20) DEFAULT NULL,  
  `name` varchar(255) DEFAULT NULL,  
  `price` decimal(10, 2) DEFAULT NULL,  
  PRIMARY KEY (`id`) USING BTREE  
) CHARACTER SET = utf8 COLLATE = utf8_general_ci ROW_FORMAT = Dynamic;
```

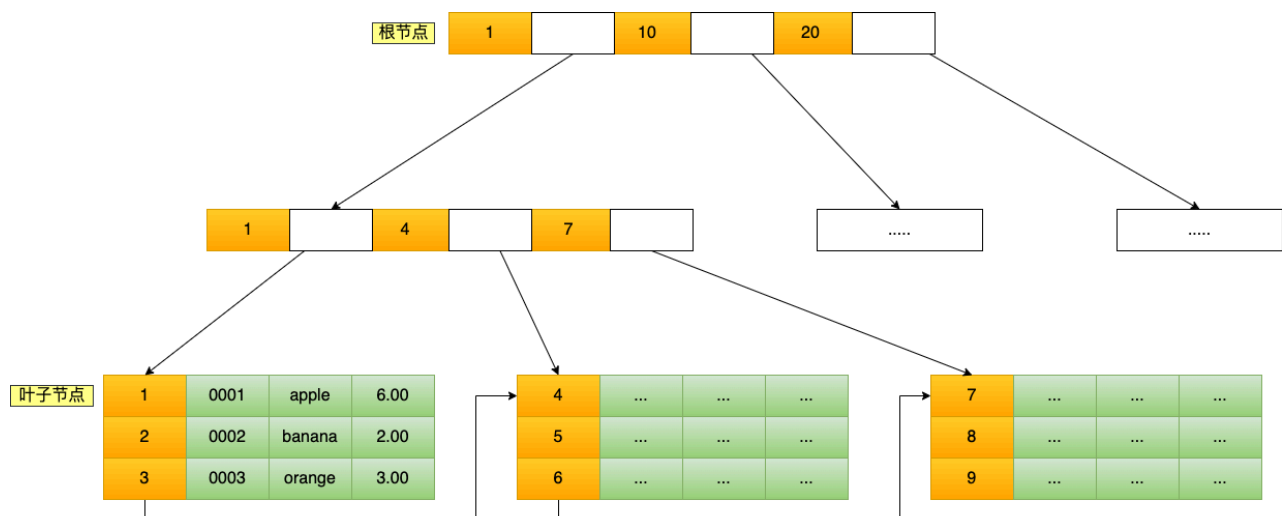
商品表里，有这些行数据：

id	product_no	name	price
1	0001	apple	6.00
2	0002	banana	2.00
3	0003	orange	3.00
4	0004	iphone13	5000.00
5	0005	ipad8	3500.00
6	0006	macbookpro	10000.00
7	0007	ps5	4000.00
8	0008	grape	10.00
9	0009	watermelon	40.00
10	0010	mango	8.00

这些行数据，存储在 B+Tree 索引时是长什么样子的？

B+Tree 是一种多叉树，叶子节点才存放数据，非叶子节点只存放索引，而且每个节点里的数据是**按主键顺序存放**的。每一层父节点的索引值都会出现在下层子节点的索引值中，因此在叶子节点中，包括了所有的索引值信息，并且每一个叶子节点都有两个指针，分别指向下一个叶子节点和上一个叶子节点，形成一个双向链表。

主键索引的 B+Tree 如图所示（图中叶子节点之间我画了单向链表，但是实际上是双向链表，原图我找不到了，修改不了，偷个懒我不重画了，大家脑补成双向链表就行）：



## # 通过主键查询商品数据的过程

比如，我们执行了下面这条查询语句：

```
select * from product where id= 5;
```

这条语句使用了主键索引查询 id 号为 5 的商品。查询过程是这样的，B+Tree 会自顶向下逐层进行查找：

- 将 5 与根节点的索引数据 (1, 10, 20) 比较, 5 在 1 和 10 之间, 所以根据 B+Tree 的搜索逻辑, 找到第二层的索引数据 (1, 4, 7);
- 在第二层的索引数据 (1, 4, 7) 中进行查找, 因为 5 在 4 和 7 之间, 所以找到第三层的索引数据 (4, 5, 6);
- 在叶子节点的索引数据 (4, 5, 6) 中进行查找, 然后我们找到了索引值为 5 的行数据。

数据库的索引和数据都是存储在硬盘的, 我们可以把读取一个节点当作一次磁盘 I/O 操作。那么上面的整个查询过程一共经历了 3 个节点, 也就是进行了 3 次 I/O 操作。

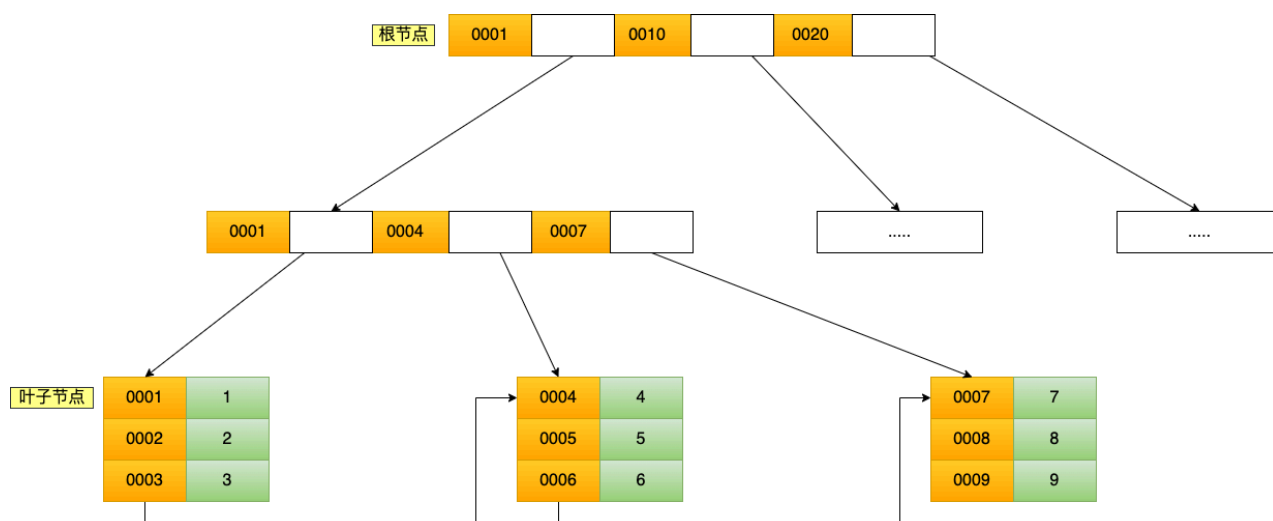
B+Tree 存储千万级的数据只需要 3-4 层高度就可以满足, 这意味着从千万级的表查询目标数据最多需要 3-4 次磁盘 I/O, 所以 **B+Tree 相比于 B 树和二叉树来说, 最大的优势在于查询效率很高, 因为即使在数据量很大的情况, 查询一个数据的磁盘 I/O 依然维持在 3-4 次。**

## # 通过二级索引查询商品数据的过程

主键索引的 B+Tree 和二级索引的 B+Tree 区别如下:

- 主键索引的 B+Tree 的叶子节点存放的是实际数据, 所有完整的用户记录都存放在主键索引的 B+Tree 的叶子节点里;
- 二级索引的 B+Tree 的叶子节点存放的是主键值, 而不是实际数据。

我这里将前面的商品表中的 product\_no (商品编码) 字段设置为二级索引, 那么二级索引的 B+Tree 如下图 (图中叶子节点之间我画了单向链表, 但是实际上是双向链表, 原图我找不到了, 修改不了, 偷个懒我不重画了, 大家脑补成双向链表就行)。



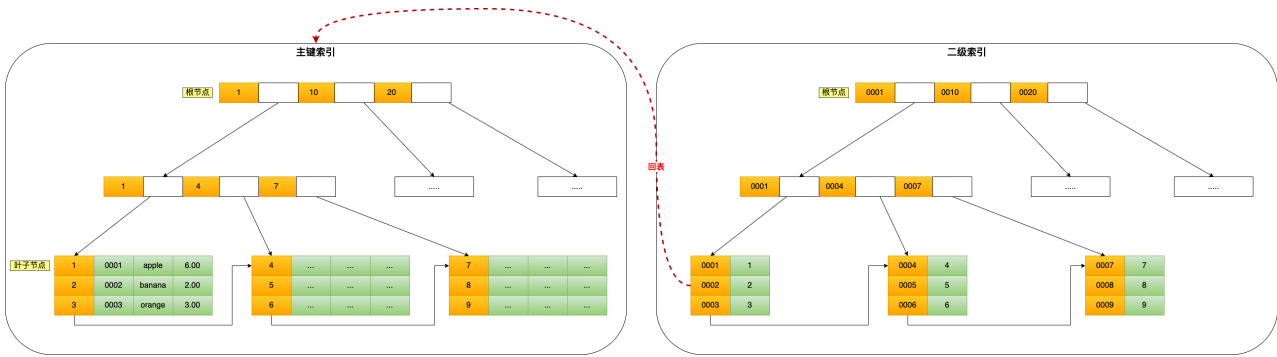
其中非叶子的 key 值是 product\_no (图中橙色部分), 叶子节点存储的数据是主键值 (图中绿色部分)。

如果我用 product\_no 二级索引查询商品, 如下查询语句:

```
select * from product where product_no = '0002';
```

会先检二级索引中的 B+Tree 的索引值 (商品编码, product\_no), 找到对应的叶子节点, 然后获取主键值, 然后再通过主键索引中的 B+Tree 树查询到对应的叶子节点, 然后获取整行数据。这个过程叫「回表」, 也就是说要查两个 B+Tree 才能查到数据。如下图

(图中叶子节点之间我画了单向链表，但是实际上是双向链表，原图我找不到了，修改不了，偷个懒我不重画了，大家脑补成双向链表就行)：



不过，当查询的数据是能在二级索引的 B+Tree 的叶子节点里查询到，这时就不用再查主键索引查，比如下面这条查询语句：

```
select id from product where product_no = '0002';
```

这种在二级索引的 B+Tree 就能查询到结果的过程就叫作「覆盖索引」，也就是只需要查一个 B+Tree 就能找到数据。

# 为什么 MySQL InnoDB 选择 B+tree 作为索引的数据结构？

前面已经讲了 B+Tree 的索引原理，现在就来回答一下 B+Tree 相比于 B 树、二叉树或 Hash 索引结构的优势在哪儿？

之前我也专门写过一篇文章，想详细了解的可以看这篇：[「女朋友问我：为什么 MySQL 喜欢 B+ 树？我笑着画了 20 张图」](#)，这里就简单做个对比。

1、B+Tree vs B Tree

B+Tree 只在叶子节点存储数据，而 B 树 的非叶子节点也要存储数据，所以 B+Tree 的单个节点的数据量更小，在相同的磁盘 I/O 次数下，就能查询更多的节点。

另外，B+Tree 叶子节点采用的是双链表连接，适合 MySQL 中常见的基于范围的顺序查找，而 B 树无法做到这一点。

2、B+Tree vs 二叉树

对于有 N 个叶子节点的 B+Tree，其搜索复杂度为 $O(\log_d N)$ ，其中 d 表示节点允许的最大子节点个数为 d 个。

在实际的应用当中，d 值是大于100的，这样就保证了，即使数据达到千万级别时，B+Tree 的高度依然维持在 3~4 层左右，也就是说一次数据查询操作只需要做 3~4 次的磁盘 I/O 操作就能查询到目标数据。

而二叉树的每个父节点的儿子节点个数只能是 2 个，意味着其搜索复杂度为  $O(\log N)$ ，这已经比 B+Tree 高出不少，因此二叉树检索到目标数据所经历的磁盘 I/O 次数要更多。

3、B+Tree vs Hash

Hash 在做等值查询的时候效率贼快，搜索复杂度为  $O(1)$ 。

但是 Hash 表不适合做范围查询，它更适合做等值的查询，这也是 B+Tree 索引要比 Hash 表索引有着更广泛的适用场景的原因。

## # 按物理存储分类

---

从物理存储的角度来看，索引分为聚簇索引（主键索引）、二级索引（辅助索引）。

这两个区别在前面也提到了：

- 主键索引的 B+Tree 的叶子节点存放的是实际数据，所有完整的用户记录都存放在主键索引的 B+Tree 的叶子节点里；
- 二级索引的 B+Tree 的叶子节点存放的是主键值，而不是实际数据。

所以，在查询时使用了二级索引，如果查询的数据能在二级索引里查询的到，那么就不需要回表，这个过程就是覆盖索引。如果查询的数据不在二级索引里，就会先检索二级索引，找到对应的叶子节点，获取到主键值后，然后再检索主键索引，就能查询到数据了，这个过程就是回表。

## # 按字段特性分类

---

从字段特性的角度来看，索引分为主键索引、唯一索引、普通索引、前缀索引。

### # 主键索引

---

主键索引就是建立在主键字段上的索引，通常在创建表的时候一起创建，一张表最多只有一个主键索引，索引列的值不允许有空值。

在创建表时，创建主键索引的方式如下：

```
CREATE TABLE table_name (  
    ....  
    PRIMARY KEY (index_column_1) USING BTREE  
);
```

### # 唯一索引

---

唯一索引建立在 UNIQUE 字段上的索引，一张表可以有多个唯一索引，索引列的值必须唯一，但是允许有空值。

在创建表时，创建唯一索引的方式如下：

```
CREATE TABLE table_name (  
    ....  
    UNIQUE KEY(index_column_1,index_column_2,...)  
);
```

建表后，如果要创建唯一索引，可以使用下面这条命令：

```
CREATE UNIQUE INDEX index_name  
ON table_name(index_column_1,index_column_2,...);
```



## # 普通索引

---

普通索引就是建立在普通字段上的索引，既不要求字段为主键，也不要求字段为 UNIQUE。

在创建表时，创建普通索引的方式如下：

```
CREATE TABLE table_name (  
    ....  
    INDEX(index_column_1,index_column_2,...)  
);
```

建表后，如果要创建普通索引，可以使用这面这条命令：

```
CREATE INDEX index_name  
ON table_name(index_column_1,index_column_2,...);
```

## # 前缀索引

---

前缀索引是指对字符类型字段的前几个字符建立的索引，而不是在整个字段上建立的索引，前缀索引可以建立在字段类型为 char、varchar、binary、varbinary 的列上。

使用前缀索引的目的是为了减少索引占用的存储空间，提升查询效率。

在创建表时，创建前缀索引的方式如下：

```
CREATE TABLE table_name(  
    column_list,  
    INDEX(column_name(length))  
);
```

建表后，如果要创建前缀索引，可以使用这面这条命令：

```
CREATE INDEX index_name  
ON table_name(column_name(length));
```

## # 按字段个数分类

---

从字段个数的角度来看，索引分为单列索引、联合索引（复合索引）。

- 建立在单列上的索引称为单列索引，比如主键索引；
- 建立在多列上的索引称为联合索引；

## # 联合索引

---

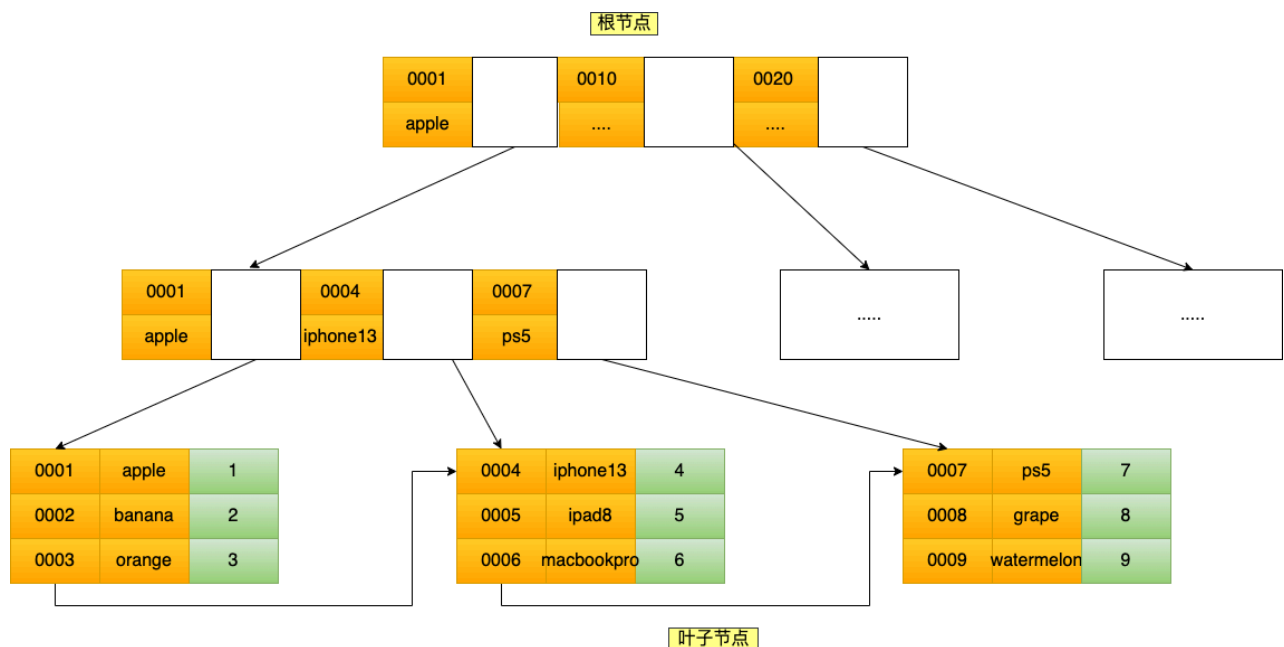
通过将多个字段组合成一个索引，该索引就被称为联合索引。

比如，将商品表中的 product\_no 和 name 字段组合成联合索引 (product\_no, name)，创建联合索引的方式如下：

```
CREATE INDEX index_product_no_name ON product(product_no, name);
```



联合索引(**product\_no**, **name**) 的 B+Tree 示意图如下 (图中叶子节点之间我画了单向链表, 但是实际上是双向链表, 原图我找不到了, 修改不了, 偷个懒我不重画了, 大家脑补成双向链表就行)。



可以看到, 联合索引的非叶子节点用两个字段的值作为 B+Tree 的 key 值。当在联合索引查询数据时, 先按 **product\_no** 字段比较, 在 **product\_no** 相同的情况下再按 **name** 字段比较。

也就是说, 联合索引查询的 B+Tree 是先按 **product\_no** 进行排序, 然后再 **product\_no** 相同的情况再按 **name** 字段排序。

因此, 使用联合索引时, 存在**最左匹配原则**, 也就是按照最左优先的方式进行索引的匹配。在使用联合索引进行查询的时候, 如果不遵循「最左匹配原则」, 联合索引会失效, 这样就无法利用到索引快速查询的特性了。

比如, 如果创建了一个 (**a**, **b**, **c**) 联合索引, 如果查询条件是以下这几种, 就可以匹配上联合索引:

- where a=1 ;
- where a=1 and b=2 and c=3 ;
- where a=1 and b=2 ;

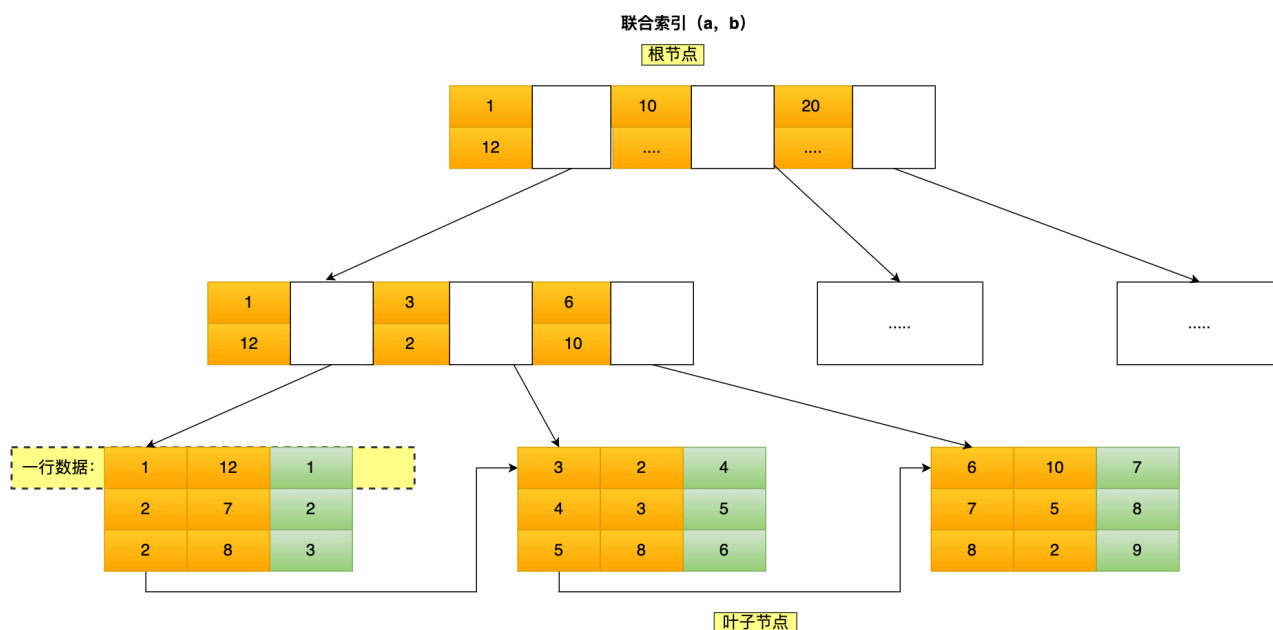
需要注意的是, 因为有查询优化器, 所以 **a** 字段在 where 子句的顺序并不重要。

但是, 如果查询条件是以下这几种, 因为不符合最左匹配原则, 所以就无法匹配上联合索引, 联合索引就会失效:

- where b=2 ;
- where c=3 ;
- where b=2 and c=3 ;

上面这些查询条件之所以会失效，是因为(a, b, c) 联合索引，是先按 a 排序，在 a 相同的情况再按 b 排序，在 b 相同的情况再按 c 排序。所以，b 和 c 是全局无序，局部相对有序的，这样在没有遵循最左匹配原则的情况下，是无法利用到索引的。

我这里举联合索引 (a, b) 的例子，该联合索引的 B+ Tree 如下（图中叶子节点之间我画了单向链表，但是实际上是双向链表，原图我找不到了，修改不了，偷个懒我不重画了，大家脑补成双向链表就行）。



可以看到，a 是全局有序的 (1, 2, 2, 3, 4, 5, 6, 7, 8)，而 b 是全局是无序的 (12, 7, 8, 2, 3, 8, 10, 5, 2)。因此，直接执行where b = 2这种查询条件没有办法利用联合索引的，利用索引的前提是索引里的 key 是有序的。

只有在 a 相同的情况才，b 才是有序的，比如 a 等于 2 的时候，b 的值为 (7, 8)，这时就是有序的，这个有序状态是局部的，因此，执行where a = 2 and b = 7是 a 和 b 字段能用到联合索引的，也就是联合索引生效了。

## # 联合索引范围查询

联合索引有一些特殊情况，并不是查询过程使用了联合索引查询，就代表联合索引中的所有字段都用到了联合索引进行索引查询，也就是可能存在部分字段用到联合索引的 B+Tree，部分字段没有用到联合索引的 B+Tree 的情况。

这种特殊情况就发生在范围查询。联合索引的最左匹配原则会一直向右匹配直到遇到「范围查询」就会停止匹配。也就是范围查询的字段可以用到联合索引，但是在范围查询字段的后面的字段无法用到联合索引。

范围查询有很多种，那到底是哪些范围查询会导致联合索引的最左匹配原则会停止匹配呢？

接下来，举例几个范围查例子。

Q1: select \* from t\_table where a > 1 and b = 2，联合索引 (a, b) 哪一个字段用到了联合索引的 B+Tree？

由于联合索引（二级索引）是先按照 a 字段的值排序的，所以符合  $a > 1$  条件的二级索引记录肯定是相邻，于是在进行索引扫描的时候，可以定位到符合  $a > 1$  条件的第一条记录，然后沿着记录所在的链表向后扫描，直到某条记录不符合  $a > 1$  条件位置。所以 a 字段可以在联合索引的 B+Tree 中进行索引查询。

但是在符合  $a > 1$  条件的二级索引记录的范围里，b 字段的值是无序的。比如前面图的联合索引的 B+ Tree 里，下面这三条记录的 a 字段的值都符合  $a > 1$  查询条件，而 b 字段的值是无序的：

- a 字段值为 5 的记录，该记录的 b 字段值为 8；
- a 字段值为 6 的记录，该记录的 b 字段值为 10；
- a 字段值为 7 的记录，该记录的 b 字段值为 5；

因此，我们不能根据查询条件  $b = 2$  来进一步减少需要扫描的记录数量（b 字段无法利用联合索引进行索引查询的意思）。

所以在执行 Q1 这条查询语句的时候，对应的扫描区间是  $(2, +\infty)$ ，形成该扫描区间的边界条件是  $a > 1$ ，与  $b = 2$  无关。

因此，Q1 这条查询语句只有 a 字段用到了联合索引进行索引查询，而 b 字段并没有使用到联合索引。

我们也可以在执行计划中的 key\_len 知道这一点，在使用联合索引进行查询的时候，通过 key\_len 我们可以知道优化器具体使用了多少个字段的搜索条件来形成扫描区间的边界条件。

举例个例子，a 和 b 都是 int 类型且不为 NULL 的字段，那么 Q1 这条查询语句执行计划如下，可以看到 key\_len 为 4 字节（如果字段允许为 NULL，就在字段类型占用的字节数上加 1，也就是 5 字节），说明只有 a 字段用到了联合索引进行索引查询，而且可以看到，即使 b 字段没用到联合索引，key 为 idx\_a\_b，说明 Q1 查询语句使用了 idx\_a\_b 联合索引。

1

EXPLAIN select \* from t\_table where a > 1 and b = 2

只有 a 字段用到了联合索引

MessageResult 1ProfileStatus

id	select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
▶	1	SIMPLE	t_table	(NULL)	range	idx_a_b	idx_a_b	4	(NULL)	8	11.11 Using where; Using index

通过 Q1 查询语句我们可以知道，a 字段使用了 > 进行范围查询，联合索引的最左匹配原则在遇到 a 字段的范围查询（>）后就停止匹配了，因此 b 字段并没有使用到联合索引。

Q2: `select * from t_table where a >= 1 and b = 2`，联合索引 (a, b) 哪一个字段用到了联合索引的 B+Tree？

Q2 和 Q1 的查询语句很像，唯一的区别就是 a 字段的查询条件「大于等于」。

由于联合索引（二级索引）是先按照 a 字段的值排序的，所以符合  $\geq 1$  条件的二级索引记录肯定是相邻，于是在进行索引扫描的时候，可以定位到符合  $\geq 1$  条件的第一条记录，然后沿着记录所在的链表向后扫描，直到某条记录不符合  $a \geq 1$  条件位置。所以 a 字段可以

在联合索引的 B+Tree 中进行索引查询。

虽然在符合  $a \geq 1$  条件的二级索引记录的范围里，b 字段的值是「无序」的，但是对于符合  $a = 1$  的二级索引记录的范围里，b 字段的值是「有序」的（因为对于联合索引，是先按照 a 字段的值排序，然后在 a 字段的值相同的情况下，再按照 b 字段的值进行排序）。

于是，在确定需要扫描的二级索引的范围时，当二级索引记录的 a 字段值为 1 时，可以通过  $b = 2$  条件减少需要扫描的二级索引记录范围（b 字段可以利用联合索引进行索引查询的意思）。也就是说，从符合  $a = 1$  and  $b = 2$  条件的第一条记录开始扫描，而不需要从第一个 a 字段值为 1 的记录开始扫描。

所以，Q2 这条查询语句 a 和 b 字段都用到了联合索引进行索引查询。

我们也可以在执行计划中的 key\_len 知道这一点。执行计划如下，可以看到 key\_len 为 8 字节，说明优化器使用了 2 个字段的查询条件来形成扫描区间的边界条件，也就是 a 和 b 字段都用到了联合索引进行索引查询。

1 EXPLAIN select * from t_table where a >= 1 and b = 2											
a 和 b 字段都用到了联合索引											
Message Result 1 Profile Status											
id	select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
1	SIMPLE	t_table	(NULL)	index	idx_a_b	idx_a_b	8	(NULL)	9	11.11	Using where; Using index

通过 Q2 查询语句我们可以知道，虽然 a 字段使用了  $\geq$  进行范围查询，但是联合索引的最左匹配原则并没有在遇到 a 字段的范围查询 ( $\geq$ ) 后就停止匹配了，b 字段还是可以用到了联合索引的。

Q3: SELECT \* FROM t\_table WHERE a BETWEEN 2 AND 8 AND b = 2，联合索引 (a, b) 哪一个字段用到了联合索引的 B+Tree？

Q3 查询条件中 a BETWEEN 2 AND 8 的意思是查询 a 字段的值在 2 和 8 之间的记录。不同的数据库对 BETWEEN ... AND 处理方式是有差异的。在 MySQL 中，BETWEEN 包含了 value1 和 value2 边界值，类似于  $\geq$  and  $\leq$ 。而有的数据库则不包含 value1 和 value2 边界值（类似于  $>$  and  $<$ ）。

这里我们只讨论 MySQL。由于 MySQL 的 BETWEEN 包含 value1 和 value2 边界值，所以类似于 Q2 查询语句，因此 Q3 这条查询语句 a 和 b 字段都用到了联合索引进行索引查询。

我们也可以在执行计划中的 key\_len 知道这一点。执行计划如下，可以看到 key\_len 为 8 字节，说明优化器使用了 2 个字段的查询条件来形成扫描区间的边界条件，也就是 a 和 b 字段都用到了联合索引进行索引查询。

1 EXPLAIN select * from t_table where a BETWEEN 2 and 8 and b = 2;											
a 字段和 b 字段都用到了联合索引											
Message Result 1 Profile Status											
id	select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
1	SIMPLE	t_table	(NULL)	range	idx_a_b	idx_a_b	8	(NULL)	8	11.11	Using where; Using index

通过 Q3 查询语句我们可以知道，虽然 a 字段使用了 BETWEEN 进行范围查询，但是联合索引的最左匹配原则并没有在遇到 a 字段的范围查询（BETWEEN）后就停止匹配了，b 字段还是可以用到了联合索引的。

Q4: `SELECT * FROM t_user WHERE name like 'j%' and age = 22`，联合索引（name, age）哪一个字段用到了联合索引的 B+Tree？

由于联合索引（二级索引）是先按照 name 字段的值排序的，所以前缀为 'j' 的 name 字段的二级索引记录都是相邻的，于是在进行索引扫描的时候，可以定位到符合前缀为 'j' 的 name 字段的第一条记录，然后沿着记录所在的链表向后扫描，直到某条记录的 name 前缀不为 'j' 为止。

所以 a 字段可以在联合索引的 B+Tree 中进行索引查询，形成的扫描区间是 ['j', 'k')。注意，j 是闭区间。如下图：

条件 like 'j%'

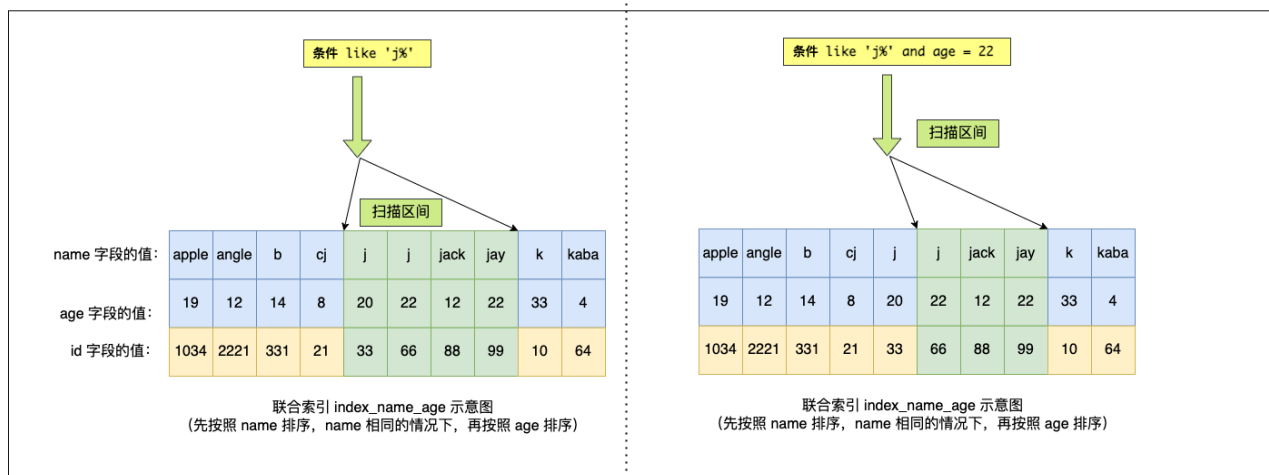
扫描区间

name 字段的值:	apple	angle	b	cj	j	j	jack	jay	k	kaba
age 字段的值:	19	12	14	8	20	22	12	22	33	4
id 字段的值:	1034	2221	331	21	33	66	88	99	10	64

联合索引 index\_name\_age 示意图  
(先按照 name 排序，name 相同的情况下，再按照 age 排序)

虽然在符合前缀为 'j' 的 name 字段的二级索引记录的范围里，age 字段的值是「无序」的，但是对于符合 name = j 的二级索引记录的范围里，age 字段的值是「有序」的（因为对于联合索引，是先按照 name 字段的值排序，然后在 name 字段的值相同的情况下，再按照 age 字段的值进行排序）。

于是，在确定需要扫描的二级索引的范围时，当二级索引记录的 name 字段值为 'j' 时，可以通过 age = 22 条件减少需要扫描的二级索引记录范围（age 字段可以利用联合索引进行索引查询的意思）。也就是说，从符合 name = 'j' and age = 22 条件的第一条记录时开始扫描，而不需要从第一个 name 为 j 的记录开始扫描。如下图的右边：



所以，Q4 这条查询语句 a 和 b 字段都用到了联合索引进行索引查询。

我们也可以在执行计划中的 key\_len 知道这一点。本次例子中：

- name 字段的类型是 varchar(30) 且不为 NULL，数据库表使用了 utf8mb4 字符集，一个字符集为 utf8mb4 的字符是 4 个字节，因此 name 字段的实际数据最多占用的存储空间长度是 120 字节（30 x 4），然后因为 name 是变长类型的字段，需要再加 2 字节（用于存储该字段实际数据的长度值），也就是 name 的 key\_len 为 122。
- age 字段的类型是 int 且不为 NULL，key\_len 为 4。

## TIP

可能有的同学对于「因为 name 是变长类型的字段，需要再加 2 字节」这句话有疑问。之前这篇[文章](#)说「如果变长字段允许存储的最大字节数小于等于 255 字节，就会用 1 字节表示变长字段的长度」，而这里为什么是 2 字节？

key\_len 的显示比较特殊，行格式是由 innodb 存储引擎实现的，而执行计划是在 server 层生成的，所以它不会去问 innodb 存储引擎可变字段的长度占用多少字节，而是不管三七二十一都使用 2 字节表示可变字段的长度。

毕竟 key\_len 的目的只是为了告诉你索引查询中用了哪些索引字段，而不是为了准确告诉这个字段占用多少字节空间。

Q4 查询语句的执行计划如下，可以看到 key\_len 为 126 字节，name 的 key\_len 为 122，age 的 key\_len 为 4，说明优化器使用了 2 个字段的查询条件来形成扫描区间的边界条件，也就是 name 和 age 字段都用到了联合索引进行索引查询。

```
1 EXPLAIN select * from t_user where name like "j%" and age = 22;
```

Message Result 1 Profile Status											
id	select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
1	SIMPLE	t_user	(NULL)	range	index_name_age	index_name_age	126	(NULL)	1	10.00	Using index condition



通过 Q4 查询语句我们可以知道，虽然 name 字段使用了 like 前缀匹配进行范围查询，但是联合索引的最左匹配原则并没有在遇到 name 字段的范围查询（like 'j%') 后就停止匹配了，age 字段还是可以用到了联合索引的。

综上所述，联合索引的最左匹配原则，在遇到范围查询（如 >、<）的时候，就会停止匹配，也就是范围查询的字段可以用到联合索引，但是在范围查询字段的后面的字段无法用到联合索引。注意，对于 >=、<=、BETWEEN、like 前缀匹配的范围查询，并不会停止匹配，前面我也用了四个例子说明了。

## # 索引下推

现在我们知道，对于联合索引 (a, b)，在执行 `select * from table where a > 1 and b = 2` 语句的时候，只有 a 字段能用到索引，那在联合索引的 B+Tree 找到第一个满足条件的主键值（ID 为 2）后，还需要判断其他条件是否满足（看 b 是否等于 2），那是在联合索引里判断？还是回主键索引去判断呢？

- 在 MySQL 5.6 之前，只能从 ID2（主键值）开始一个个回表，到「主键索引」上找出数据行，再对比 b 字段值。
- 而 MySQL 5.6 引入的索引下推优化（index condition pushdown），可以在联合索引遍历过程中，对联合索引中包含的字段先做判断，直接过滤掉不满足条件的记录，减少回表次数。

当你的查询语句的执行计划里，出现了 Extra 为 `Using index condition`，那么说明使用了索引下推的优化。

## # 索引区分度

另外，建立联合索引时的字段顺序，对索引效率也有很大影响。越靠前的字段被用于索引过滤的概率越高，实际开发工作中建立联合索引时，要把区分度大的字段排在前面，这样区分度大的字段越有可能被更多的 SQL 使用到。

区分度就是某个字段 column 不同值的个数「除以」表的总行数，计算公式如下：

$$\text{区分度} = \frac{\text{distinct}(\text{column})}{\text{count}(*)}$$

比如，性别的区分度就很小，不适合建立索引或不适合排在联合索引列的靠前的位置，而 UUID 这类字段就比较适合做索引或排在联合索引列的靠前的位置。



因为如果索引的区分度很小，假设字段的值分布均匀，那么无论搜索哪个值都可能得到一半的数据。在这些情况下，还不如不要索引，因为 MySQL 还有一个查询优化器，查询优化器发现某个值出现在表的数据行中的百分比（惯用的百分比界线是"30%"）很高的时候，它一般会忽略索引，进行全表扫描。

## # 联合索引进行排序

---

这里出一个题目，针对下面这条 SQL，你怎么通过索引来提高查询效率呢？

```
select * from order where status = 1 order by create_time asc
```

有的同学会认为，单独给 status 建立一个索引就可以了。

但是更好的方式给 status 和 create\_time 列建立一个联合索引，因为这样可以避免 MySQL 数据库发生文件排序。

因为在查询时，如果只用到 status 的索引，但是这条语句还要对 create\_time 排序，这时就要用文件排序 filesort，也就是在 SQL 执行计划中，Extra 列会出现 Using filesort。

所以，要利用索引的有序性，在 status 和 create\_time 列建立联合索引，这样根据 status 筛选后的数据就是按照 create\_time 排好序的，避免在文件排序，提高了查询效率。

## # 什么时候需要 / 不需要创建索引？

---

索引最大的好处是提高查询速度，但是索引也是有缺点的，比如：

- 需要占用物理空间，数量越大，占用空间越大；
- 创建索引和维护索引要耗费时间，这种时间随着数据量的增加而增大；
- 会降低表的增删改的效率，因为每次增删改索引，B+ 树为了维护索引有序性，都需要进行动态维护。

所以，索引不是万能钥匙，它也是根据场景来使用的。

## # 什么时候适用索引？

---

- 字段有唯一性限制的，比如商品编码；
- 经常用于 WHERE 查询条件的字段，这样能够提高整个表的查询速度，如果查询条件不是一个字段，可以建立联合索引。
- 经常用于 GROUP BY 和 ORDER BY 的字段，这样在查询的时候就不需要再去做一次排序了，因为我们都已经知道了建立索引之后在 B+Tree 中的记录都是排序好的。

## # 什么时候不需要创建索引？

---

- WHERE 条件，GROUP BY，ORDER BY 里用不到的字段，索引的价值是快速定位，如果起不到定位的字段通常是不需要创建索引的，因为索引是会占用物理空间的。
- 字段中存在大量重复数据，不需要创建索引，比如性别字段，只有男女，如果数据库表中，男女的记录分布均匀，那么无论搜索哪个值都可能得到一半的数据。在这些情况下，还不如不要索引，因为 MySQL 还有一个查询优化器，查询优化器发现某个值出现在表的数据行中的百分比很高的时候，它一般会忽略索引，进行全表扫描。

- 表数据太少的时候，不需要创建索引；
- 经常更新的字段不用创建索引，比如不要对电商项目的用户余额建立索引，因为索引字段频繁修改，由于要维护 B+Tree 的有序性，那么就需要频繁的重建索引，这个过程是会影响数据库性能的。

## # 有什么优化索引的方法？

---

这里说一下几种常见优化索引的方法：

- 前缀索引优化；
- 覆盖索引优化；
- 主键索引最好是自增的；
- 防止索引失效；

### # 前缀索引优化

---

前缀索引顾名思义就是使用某个字段中字符串的前几个字符建立索引，那我们为什么需要使用前缀来建立索引呢？

使用前缀索引是为了减小索引字段大小，可以增加一个索引页中存储的索引值，有效提高索引的查询速度。在一些大字符串的字段作为索引时，使用前缀索引可以帮助我们减小索引项的大小。

不过，前缀索引有一定的局限性，例如：

- order by 就无法使用前缀索引；
- 无法把前缀索引用作覆盖索引；

### # 覆盖索引优化

---

覆盖索引是指 SQL 中 query 的所有字段，在索引 B+Tree 的叶子节点上都能找得到的那些索引，从二级索引中查询得到记录，而不需要通过聚簇索引查询获得，可以避免回表的操作。

假设我们只需要查询商品的名称、价格，有什么方式可以避免回表呢？

我们可以建立一个联合索引，即「商品ID、名称、价格」作为一个联合索引。如果索引中存在这些数据，查询将不会再次检索主键索引，从而避免回表。

所以，使用覆盖索引的好处就是，不需要查询出包含整行记录的所有信息，也就减少了大量的 I/O 操作。

### # 主键索引最好是自增的

---

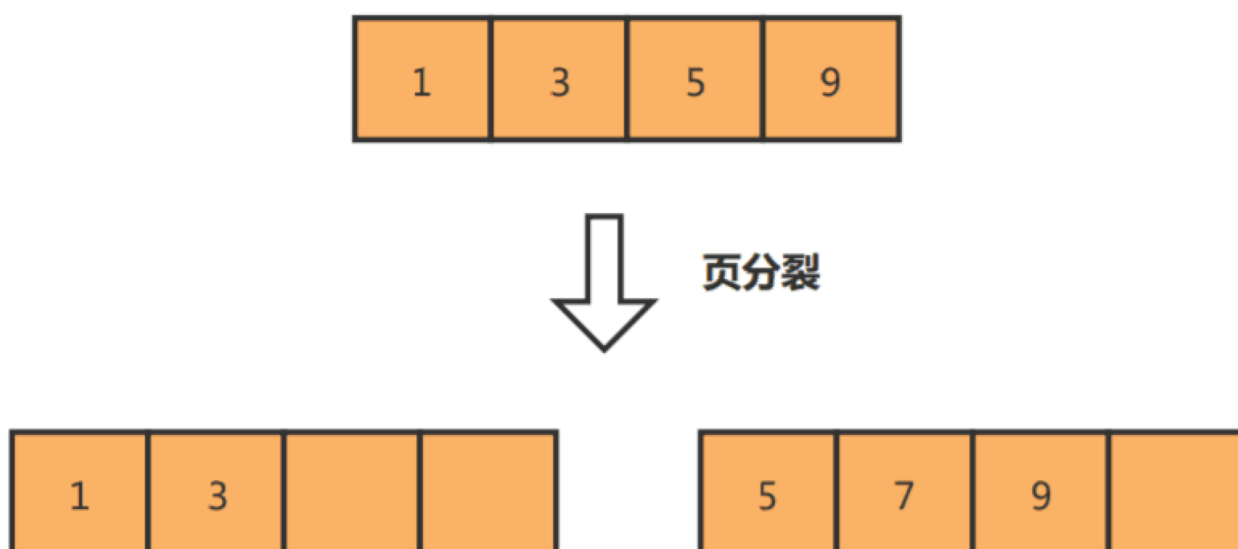
我们在建表的时候，都会默认将主键索引设置为自增的，具体为什么要这样做呢？又有什么好处？

InnoDB 创建主键索引默认为聚簇索引，数据被存放在 B+Tree 的叶子节点上。也就是说，同一个叶子节点内的各个数据是按主键顺序存放的，因此，每当有一条新的数据插入时，数据库会根据主键将其插入到对应的叶子节点中。

**如果我们使用自增主键**，那么每次插入的新数据就会按顺序添加到当前索引节点的位置，不需要移动已有的数据，当页面写满，就会自动开辟一个新页面。因为每次**插入一条新记录，都是追加操作，不需要重新移动数据**，因此这种插入数据的方法效率非常高。

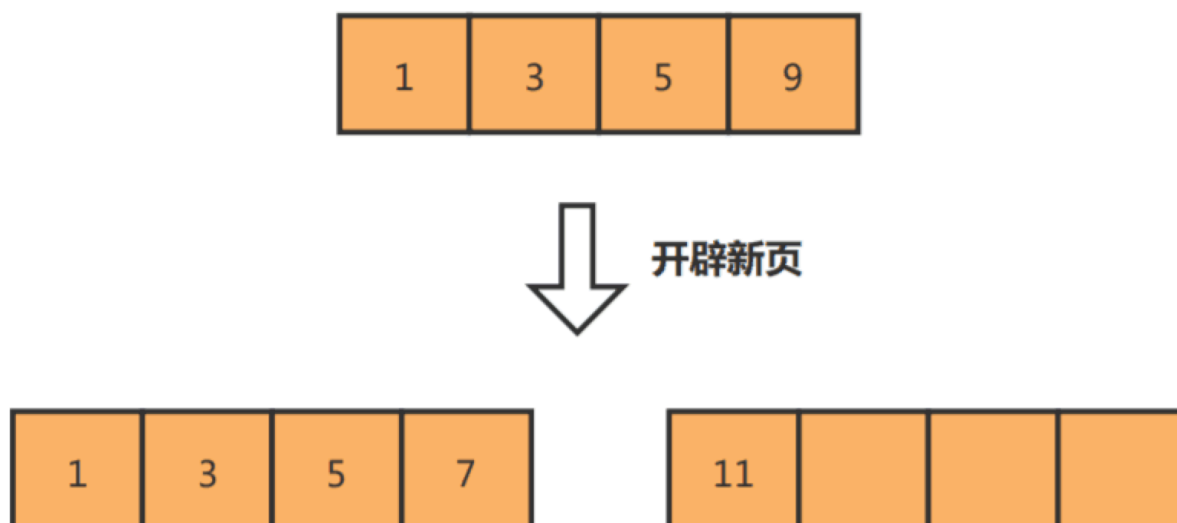
**如果我们使用非自增主键**，由于每次插入主键的索引值都是随机的，因此每次插入新的数据时，就可能会插入到现有数据页中间的某个位置，这将不得不移动其它数据来满足新数据的插入，甚至需要从一个页面复制数据到另外一个页面，我们通常将这种情况称为**页分裂**。页分裂还有可能会造成大量的内存碎片，导致索引结构不紧凑，从而影响查询效率。

举个例子，假设某个数据页中的数据是1、3、5、9，且数据页满了，现在准备插入一个数据7，则需要把数据页分割为两个数据页：



出现页分裂时，需要将一个页的记录移动到另外一个页，性能会受到影响，同时页空间的利用率下降，造成存储空间的浪费。

而如果记录是顺序插入的，例如插入数据11，则只需开辟新的数据页，也就不会发生页分裂：



因此，在使用 InnoDB 存储引擎时，如果没有特别的业务需求，建议使用自增字段作为主键。

另外，主键字段的长度不要太大，因为主键字段长度越小，意味着二级索引的叶子节点越小（二级索引的叶子节点存放的数据是主键值），这样二级索引占用的空间也就越小。

## # 索引最好设置为 NOT NULL

为了更好的利用索引，索引列要设置为 NOT NULL 约束。有两个原因：

- 第一原因：索引列存在 NULL 就会导致优化器在做索引选择的时候更加复杂，更加难以优化，因为可为 NULL 的列会使索引、索引统计和值比较都更复杂，比如进行索引统计时，count 会省略值为 NULL 的行。
- 第二个原因：NULL 值是一个没意义的值，但是它会占用物理空间，所以会带来的存储空间的问题，因为 InnoDB 存储记录的时候，如果表中存在允许为 NULL 的字段，那么行格式中至少会用 1 字节空间存储 NULL 值列表，如下图的紫色部分：



## # 防止索引失效

用上了索引并不意味着查询的时候会使用到索引，所以我们心里要清楚有哪些情况会导致索引失效，从而避免写出索引失效的查询语句，否则这样的查询效率是很低的。

我之前写过索引失效的文章，想详细了解的可以去看这篇文章：[谁还没碰过索引失效呢？](#)

这里简单说一下，发生索引失效的情况：

- 当我们使用左或者左右模糊匹配的时候，也就是 `like %xx` 或者 `like %xx%` 这两种方式都会造成索引失效；
- 当我们在查询条件中对索引列做了计算、函数、类型转换操作，这些情况下都会造成索引失效；
- 联合索引要能正确使用需要遵循最左匹配原则，也就是按照最左优先的方式进行索引的匹配，否则就会导致索引失效。
- 在 `WHERE` 子句中，如果在 `OR` 前的条件列是索引列，而在 `OR` 后的条件列不是索引列，那么索引会失效。

我上面说的是常见的索引失效场景，实际过程中，可能会出现其他的索引失效场景，这时我们就需要查看执行计划，通过执行计划显示的数据判断查询语句是否使用了索引。

如下图，就是一个没有使用索引，并且是一个全表扫描的查询语句。

```
1 explain select * from t_user where id + 1 = 10;
```

Message Result 1 Profile Status											
id	select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
1	SIMPLE	t_user	(NULL)	ALL	(NULL)	(NULL)	(NULL)	(NULL)	9	100.00	Using where

对于执行计划，参数有：

- `possible_keys` 字段表示可能用到的索引；
- `key` 字段表示实际用的索引，如果这一项为 `NULL`，说明没有使用索引；
- `key_len` 表示索引的长度；
- `rows` 表示扫描的数据行数。
- `type` 表示数据扫描类型，我们需要重点看这个。

`type` 字段就是描述了找到所需数据时使用的扫描方式是什么，常见扫描类型的执行效率从低到高的顺序为：

- `All`（全表扫描）；
- `index`（全索引扫描）；
- `range`（索引范围扫描）；
- `ref`（非唯一索引扫描）；
- `eq_ref`（唯一索引扫描）；
- `const`（结果只有一条的主键或唯一索引扫描）。

在这些情况里，`all` 是最坏的情况，因为采用了全表扫描的方式。`index` 和 `all` 差不多，只不过 `index` 对索引表进行全扫描，这样做的好处是不再需要对数据进行排序，但是开销依然很大。所以，要尽量避免全表扫描和全索引扫描。

`range` 表示采用了索引范围扫描，一般在 `where` 子句中使用 `<`、`>`、`in`、`between` 等关键词，只检索给定范围的行，属于范围查找。从这一级别开始，索引的作用会越来越明显，因此我们需要尽量让 **SQL 查询可以使用到 `range` 这一级别及以上的 `type` 访问方式。**

ref 类型表示采用了非唯一索引，或者是唯一索引的非唯一性前缀，返回数据返回可能是多条。因为虽然使用了索引，但该索引列的值并不唯一，有重复。这样即使使用索引快速查找到了第一条数据，仍然不能停止，要进行目标值附近的小范围扫描。但它的好处是它并不需要扫全表，因为索引是有序的，即便有重复值，也是在一个非常小的范围内扫描。

eq\_ref 类型是使用主键或唯一索引时产生的访问方式，通常使用在多表联查中。比如，对两张表进行联查，关联条件是两张表的 user\_id 相等，且 user\_id 是唯一索引，那么使用 EXPLAIN 进行执行计划查看的时候，type 就会显示 eq\_ref。

const 类型表示使用了主键或者唯一索引与常量值进行比较，比如 select name from product where id=1。

需要说明的是 const 类型和 eq\_ref 都使用了主键或唯一索引，不过这两个类型有所区别，**const 是与常量进行比较，查询效率会更快，而 eq\_ref 通常用于多表联查中。**

除了关注 type，我们也要关注 extra 显示的结果。

这里说几个重要的参考指标：

- Using filesort：当查询语句中包含 group by 操作，而且无法利用索引完成排序操作的时候，这时不得不选择相应的排序算法进行，甚至可能会通过文件排序，效率是很低的，所以要避免这种问题的出现。
- Using temporary：使了用临时表保存中间结果，MySQL 在对查询结果排序时使用临时表，常见于排序 order by 和分组查询 group by。效率低，要避免这种问题的出现。
- Using index：所需数据只需在索引即可全部获得，不须要再到表中取数据，也就是使用了覆盖索引，避免了回表操作，效率不错。

## # 总结

---

这次主要介绍了索引的原理、分类和使用。我把重点总结在了下面这个表格



<p>为什么 MySQL InnoDB 选择 B+tree 作为索引的数据结构?</p>	<p>B+Tree vs B Tree:</p> <ul style="list-style-type: none"> <li>• 存储相同数据量级别的情况下, B+Tree 树高比 B Tree 低, 磁盘 I/O 次数更少。</li> <li>• B+Tree 叶子节点用双向链表串起来, 适合范围查询, B Tree 无法做到这点</li> </ul> <p>B+Tree vs 二叉树:</p> <ul style="list-style-type: none"> <li>• 随着数据量的增加, 二叉树的树高会越来越高, 磁盘 I/O 次数也会更多, B+Tree 在千万级别的数据量下, 高度依然维持在 3~4 层左右, 也就是说一次数据查询操作只需要做 3~4 次的磁盘 I/O 操作就能查询到目标数据。</li> </ul> <p>B+Tree vs Hash:</p> <ul style="list-style-type: none"> <li>• 虽然 Hash 的等值查询效率很高, 但是无法做范围查询</li> </ul>
<p>什么时候适用索引?</p>	<ul style="list-style-type: none"> <li>• 字段有唯一性限制的, 比如商品编码;</li> <li>• 经常用于 WHERE 查询条件的字段;</li> <li>• 经常用于 GROUP BY 和 ORDER BY 的字段;</li> </ul>
<p>什么时候不需要创建索引?</p>	<ul style="list-style-type: none"> <li>• WHERE 条件, GROUP BY, ORDER BY 里用不到的字段;</li> <li>• 字段中存在大量重复数据, 不需要创建索引;</li> <li>• 表数据太少的时候, 不需要创建索引;</li> <li>• 经常更新的字段不用创建索引;</li> </ul>
<p>什么时候索引会失效?</p>	<ul style="list-style-type: none"> <li>• 当我们使用左或者左右模糊匹配的时候, 也就是 `like %xx` 或者 `like %xx%` 这两种方式都会造成索引失效;</li> <li>• 当我们在查询条件中对索引列做了计算、函数、类型转换操作, 会导致索引失效</li> <li>• 联合索引要能正确使用需要遵循最左匹配原则, 也就是按照最左优先的方式进行索引的匹配, 否则就会导致索引失效。</li> <li>• 在 WHERE 子句中, 如果在 OR 前的条件列是索引列, 而在 OR 后的条件列不是索引列, 那么索引会失效。</li> <li>• 为了更好的利用索引, 索引列要设置为 NOT NULL 约束。</li> </ul>
<p>有什么优化索引的方法?</p>	<ul style="list-style-type: none"> <li>• 前缀索引优化;</li> <li>• 覆盖索引优化;</li> <li>• 主键索引最好是自增的;</li> <li>• 防止索引失效;</li> </ul>