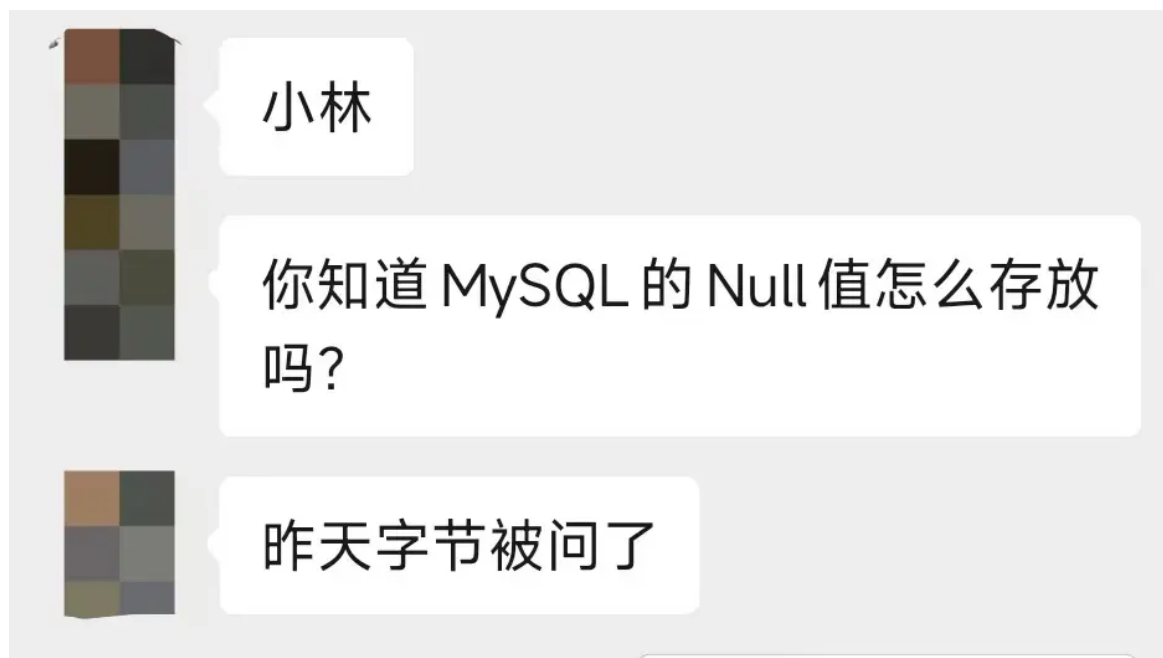


MySQL 一行记录是怎么存储的？

大家好，我是小林。

之前有位读者在面字节的时候，被问到这么个问题：



如果你知道 MySQL 一行记录的存储结构，那么这个问题对你没什么难度。

如果你不知道也没关系，这次我跟大家聊聊 [MySQL 一行记录是怎么存储的？](#)

知道了这个之后，除了能应解锁前面这道面试题，你还会解锁这些面试题：

- MySQL 的 NULL 值会占用空间吗？
- MySQL 怎么知道 varchar(n) 实际占用数据的大小？
- varchar(n) 中 n 最大取值为多少？
- 行溢出后，MySQL 是怎么处理的？

这些问题看似毫不相干，其实都是在围绕「MySQL 一行记录的存储结构」这一个知识点，所以攻破了这个知识点后，这些问题就引刃而解了。

好了，话不多说，发车！



目录



侧边栏



夜间



技术群



资料



支持我



上一篇



下一篇

MySQL 的数据存放在哪个文件？

大家都知道 MySQL 的数据都是保存在磁盘的，那具体是保存在哪个文件呢？

MySQL 存储的行为是由存储引擎实现的，MySQL 支持多种存储引擎，不同的存储引擎保存的文件自然也不同。

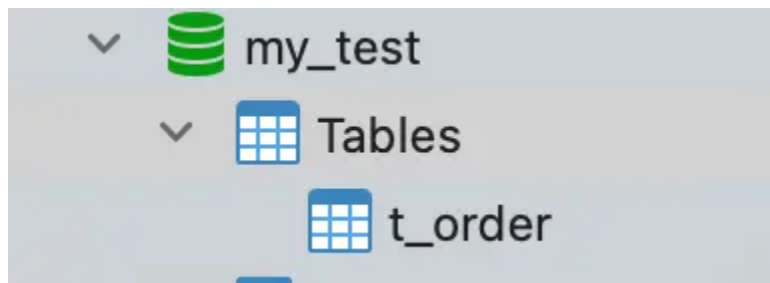
InnoDB 是我们常用的存储引擎，也是 MySQL 默认的存储引擎。所以，本文主要以 InnoDB 存储引擎展开讨论。

先来看看 MySQL 数据库的文件存放在哪个目录？

```
mysql> SHOW VARIABLES LIKE 'datadir';
+-----+-----+
| Variable_name | Value           |
+-----+-----+
| datadir       | /var/lib/mysql/ |
+-----+-----+
1 row in set (0.00 sec)
```

我们每创建一个 database（数据库）都会在 /var/lib/mysql/ 目录里面创建一个以 database 为名的目录，然后保存表结构和表数据的文件都会存放在这个目录里。

比如，我这里有一个名为 my_test 的 database，该 database 里有一张名为 t_order 数据库表。



然后，我们进入 /var/lib/mysql/my_test 目录，看看里面有什么文件？



目录



侧边栏



夜间



技术群



资料



支持我



上一篇



下一篇

t_order.frm
t_order.ibd

可以看到，共有三个文件，这三个文件分别代表着：

- db.opt，用来存储当前数据库的默认字符集和字符校验规则。
- t_order.frm，t_order 的**表结构**会保存在这个文件。在 MySQL 中建立一张表都会生成一个.frm 文件，该文件是用来保存每个表的元数据信息的，主要包含表结构定义。
- t_order.ibd，t_order 的**表数据**会保存在这个文件。表数据既可以存在共享表空间文件（文件名：ibdata1）里，也可以存放在独占表空间文件（文件名：表名字.idb）。这个行为是由参数 innodb_file_per_table 控制的，若设置了参数 innodb_file_per_table 为 1，则会将存储的数据、索引等信息单独存储在一个独占表空间，从 MySQL 5.6.6 版本开始，它的默认值就是 1 了，因此从这个版本之后，MySQL 中每一张表的数据都存放在一个独立的 .ibd 文件。

好了，现在我们知道了一张数据库表的数据是保存在「表名字.idb」的文件里的，这个文件也称为独占表空间文件。

表空间文件的结构是怎么样的？

表空间由段 (segment)、区 (extent)、页 (page)、行 (row) 组成，InnoDB存储引擎的逻辑存储结构大致如下图：



目录



侧边栏



夜间



技术群



资料



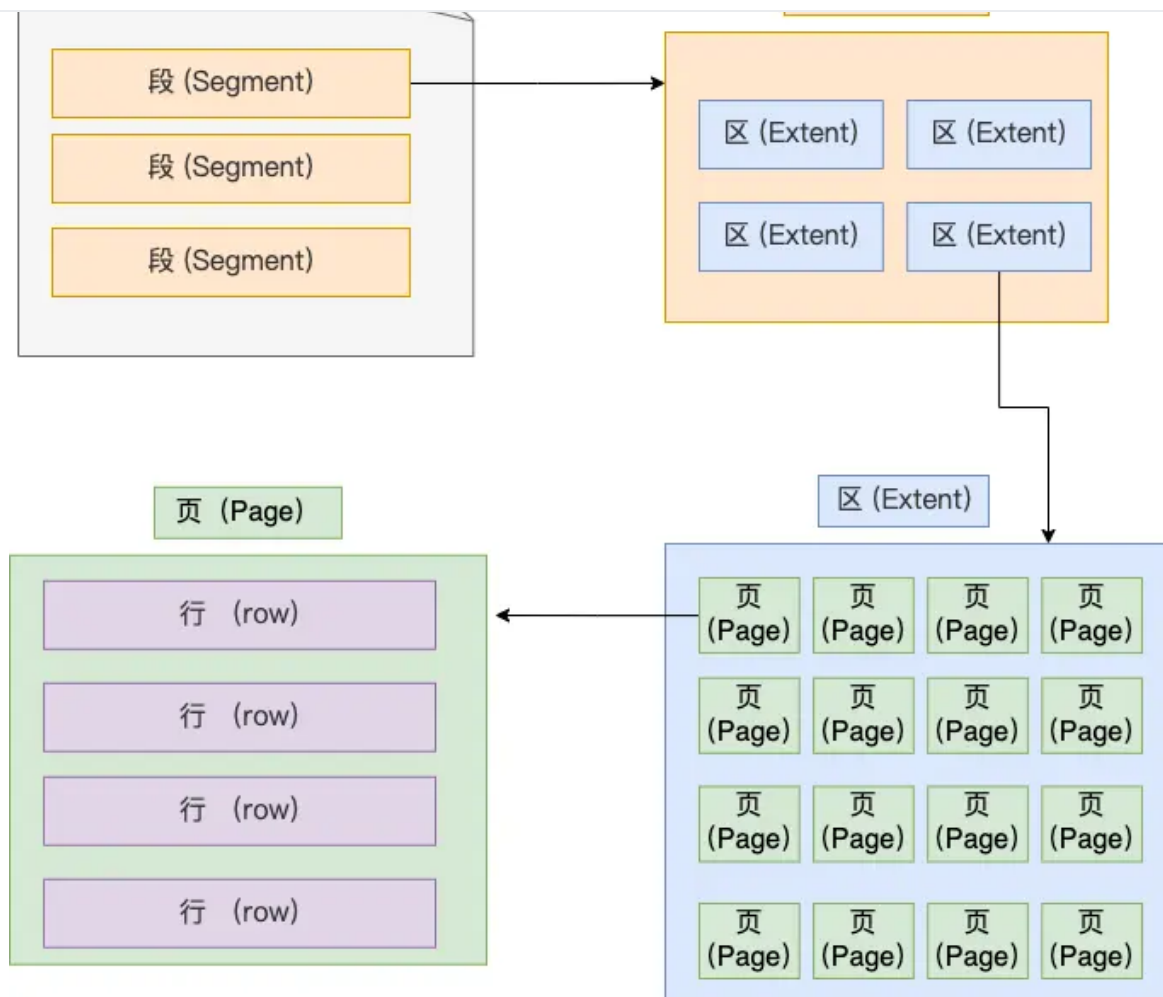
支持我



上一篇



下一篇



目录



侧边栏



夜间



技术群



资料



支持我



上一篇



下一篇

下面我们从下往上一个个看看。

1、行 (row)

数据库表中的记录都是按行 (row) 进行存放的，每行记录根据不同的行格式，有不同的存储结构。

后面我们详细介绍 InnoDB 存储引擎的行格式，也是本文重点介绍的内容。

2、页 (page)

记录是按照行来存储的，但是数据库的读取并不以「行」为单位，否则一次读取（也就是一次 I/O 操作）只能处理一行数据，效率会非常低。

因此，**InnoDB 的数据是按「页」为单位来读写的**，也就是说，当需要读一条记录的时候，并不是将这个行记录从磁盘读出来，而是以页为单位，将其整体读入内存。

的，一次最少从磁盘中读取 16K 的内容到内存中，一次最少把内存中的 16K 内容刷新到磁盘中。

页的类型有很多，常见的有数据页、undo 日志页、溢出页等等。数据表中的行记录是用「数据页」来管理的，数据页的结构这里我就不讲细说了，之前文章有说过，感兴趣的可以去看这篇文章：[换一个角度看 B+ 树](#)

总之知道表中的记录存储在「数据页」里面就行。

3、区 (extent)

我们知道 InnoDB 存储引擎是用 B+ 树来组织数据的。

B+ 树中每一层都是通过双向链表连接起来的，如果是以页为单位来分配存储空间，那么链表中相邻的两个页之间的物理位置并不是连续的，可能离得非常远，那么磁盘查询时就会有大量的随机 I/O，随机 I/O 是非常慢的。

解决这个问题也很简单，就是让链表中相邻的页的物理位置也相邻，这样就可以使用顺序 I/O 了，那么在范围查询（扫描叶子节点）的时候性能就会很高。

那具体怎么解决呢？

在表中数据量大的时候，为某个索引分配空间的时候就不再按照页为单位分配了，而是按照区 (extent) 为单位分配。每个区的大小为 1MB，对于 16KB 的页来说，连续的 64 个页会被划为一个区，这样就使得链表中相邻的页的物理位置也相邻，就能使用顺序 I/O 了。

4、段 (segment)

表空间是由各个段 (segment) 组成的，段是由多个区 (extent) 组成的。段一般分为数据段、索引段和回滚段等。

- 索引段：存放 B + 树的非叶子节点的区的集合；
- 数据段：存放 B + 树的叶子节点的区的集合；
- 回滚段：存放的是回滚数据的区的集合，之前讲[事务隔离](#)的时候就介绍到了 MVCC 利用了回滚段实现了多版本查询数据。

好了，终于说完表空间的结构了。接下来，就具体讲一下 InnoDB 的行格式了。

之所以要绕一大圈才讲行记录的格式，主要是想让大家知道行记录是存储在哪个文件，以及行记录在这个表空间文件中的哪个区域，有一个从上往下切入的视角，这样理解起来不会觉



目录



侧边栏



夜间



技术群



资料



支持我



上一篇



下一篇

InnoDB 行格式有哪些？

行格式 (row_format) ，就是一条记录的存储结构。

InnoDB 提供了 4 种行格式，分别是 Redundant、Compact、Dynamic和 Compressed 行格式。

- Redundant 是很古老的行格式了，MySQL 5.0 版本之前用的行格式，现在基本没人用了。
- 由于 Redundant 不是一种紧凑的行格式，所以 MySQL 5.0 之后引入了 Compact 行记录存储方式，Compact 是一种紧凑的行格式，设计的初衷就是为了让一个数据页中可以存放更多的行记录，从 MySQL 5.1 版本之后，行格式默认设置成 Compact。
- Dynamic 和 Compressed 两个都是紧凑的行格式，它们的行格式都和 Compact 差不多，因为都是基于 Compact 改进一点东西。从 MySQL5.7 版本之后，默认使用 Dynamic 行格式。

Redundant 行格式我这里就不讲了，因为现在基本没人用了，这次重点介绍 Compact 行格式，因为 Dynamic 和 Compressed 这两个行格式跟 Compact 非常像。

所以，弄懂了 Compact 行格式，之后你们在去了解其他行格式，很快也能看懂。

COMPACT 行格式长什么样？

先跟 Compact 行格式混个脸熟，它长这样：



可以看到，一条完整的记录分为「记录的额外信息」和「记录的真实数据」两个部分。

接下来，分别详细说下。

记录的额外信息

记录的额外信息包含 3 个部分：变长字段长度列表、NULL 值列表、记录头信息。



目录



侧边栏



夜间



技术群



资料



支持我



上一篇



下一篇

1. 变长字段长度列表

varchar(n) 和 char(n) 的区别是什么，相信大家都非常清楚，char 是定长的，varchar 是变长的，变长字段实际存储的数据的长度（大小）不固定的。

所以，在存储数据的时候要把这些数据占用的字节数也存起来，存到「变长字段长度列表」里面，读取数据的时候才能根据这个「变长字段长度列表」去读取对应长度的数据。其他 TEXT、BLOB 等变长字段也是这么实现的。

为了展示「变长字段长度列表」具体是怎么保存变长字段占用的字节数，我们先创建这样一张表，字符集是 ascii（所以每一个字符占用的 1 字节），行格式是 Compact，t_user 表中 name 和 phone 字段都是变长字段：

```
CREATE TABLE `t_user` (  
  `id` int(11) NOT NULL,  
  `name` VARCHAR(20) NOT NULL,  
  `phone` VARCHAR(20) DEFAULT NULL,  
  `age` int(11) DEFAULT NULL,  
  PRIMARY KEY (`id`) USING BTREE  
) ENGINE = InnoDB DEFAULT CHARACTER SET = ascii ROW_FORMAT = COMPACT;
```

现在 t_user 表里有这三条记录：

id	name	phone	age
1	a	123	18
2	bb	1234	(NULL)
3	ccc	(NULL)	(NULL)

接下来，我们看看看看这三条记录的行格式中的「变长字段长度列表」是怎样存储的。

先来看第一条记录：

- name 列的值为 a，长度是 1 字节，十六进制 0x01
- phone 列的值为 123，长度是 3 字节，十六进制 0x03
- age 列和 id 列不是变长字段，所以这里不用管。



目录



侧边栏



夜间



技术群



资料



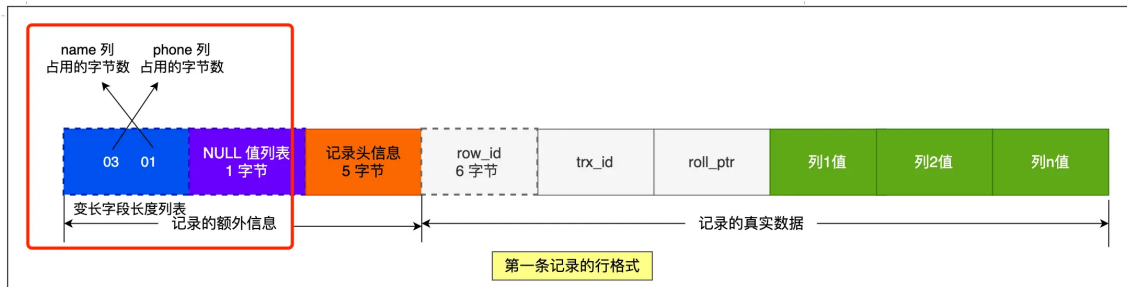
支持我



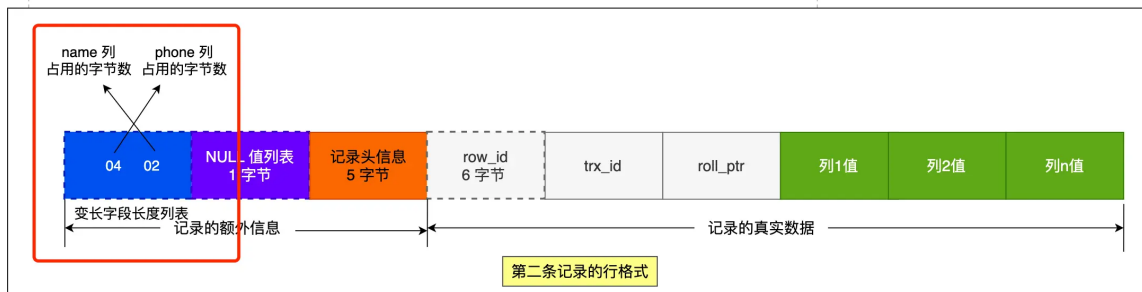
上一篇



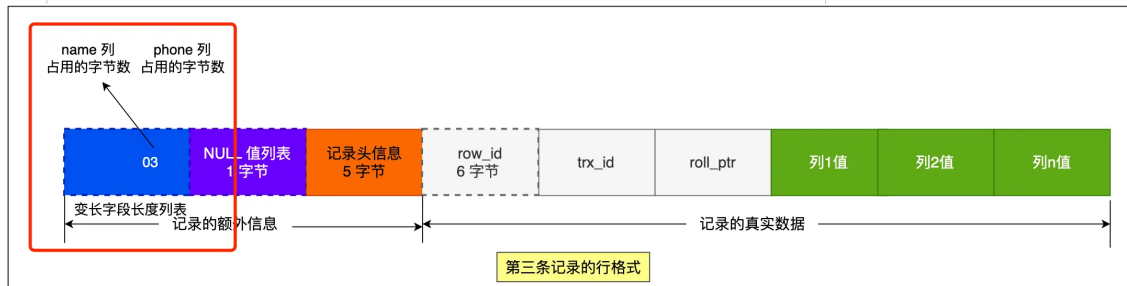
下一篇



同样的道理，我们也可以得出**第二条记录**的行格式中，「变长字段长度列表」里的内容是「04 02」，如下图：



第三条记录中 phone 列的值是 NULL，**NULL 是不会存放在行格式中记录的真实数据部分里的**，所以「变长字段长度列表」里不需要保存值为 NULL 的变长字段的长度。



为什么「变长字段长度列表」的信息要按照逆序存放？

这个设计是有想法的，主要是因为「记录头信息」中指向下一个记录的指针，指向的是下一条记录的「记录头信息」和「真实数据」之间的位置，这样的好处是向左读就是记录头信息，向右读就是真实数据，比较方便。

「变长字段长度列表」中的信息之所以要逆序存放，是因为这样可以**使得位置靠前的记录的真实数据和数据对应的字段长度信息可以同时在一个 CPU Cache Line 中，这样就可以提高 CPU Cache 的命中率。**



目录



侧边栏



夜间



技术群



资料



支持我



上一篇



下一篇

每个数据库表的行格式都有「变长字段字节数列表」吗？

其实变长字段字节数列表不是必须的。

当数据表没有变长字段的时候，比如全部都是 int 类型的字段，这时候表里的行格式就不会有「变长字段长度列表」了，因为没必要，不如去掉以节省空间。

所以「变长字段长度列表」只出现在数据表有变长字段的时候。

2. NULL 值列表

表中的某些列可能会存储 NULL 值，如果把这些 NULL 值都放到记录的真实数据中会比较浪费空间，所以 Compact 行格式把这些值为 NULL 的列存储到 NULL 值列表中。

如果存在允许 NULL 值的列，则每个列对应一个二进制位 (bit)，二进制位按照列的顺序逆序排列。

- 二进制位的值为 1 时，代表该列的值为 NULL。
- 二进制位的值为 0 时，代表该列的值不为 NULL。

另外，NULL 值列表必须用整数个字节的位表示（1 字节 8 位），如果使用的二进制位个数不足整数个字节，则在字节的高位补 0。

还是以 t_user 表的这三条记录作为例子：

id	name	phone	age	
1	a	123	18	
2	bb	1234	(NULL)	
3	ccc	(NULL)	(NULL)	

接下来，我们看看看看这三条记录的行格式中的 NULL 值列表是怎样存储的。

先来看**第一条记录**，第一条记录所有列都有值，不存在 NULL 值，所以用二进制来表示是酱紫的：



目录



侧边栏



夜间



技术群



资料



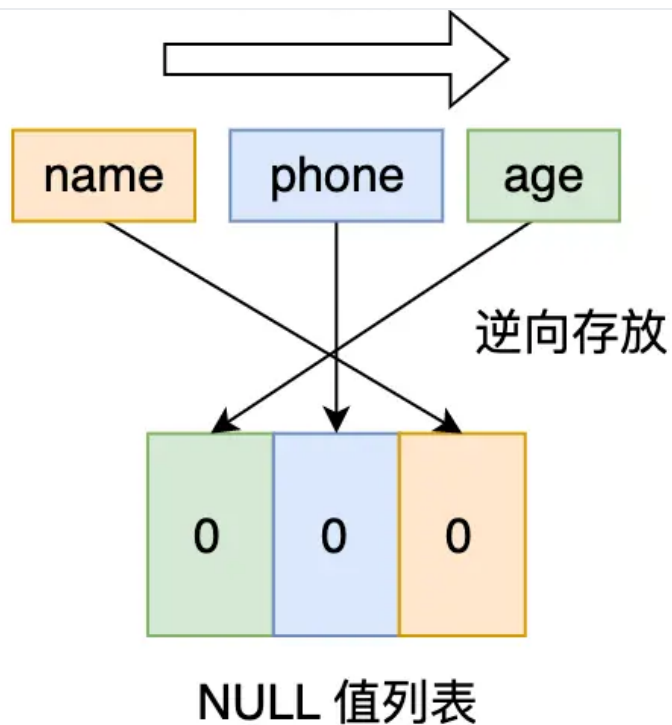
支持我



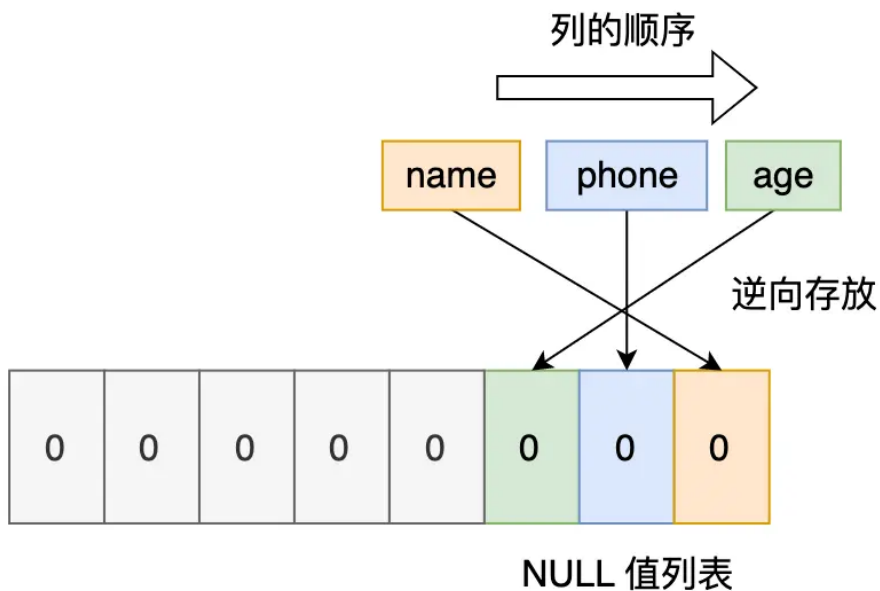
上一篇



下一篇



但是 InnoDB 是用整数字节的二进制位来表示NULL值列表的，现在不足 8 位，所以要在高位补 0，最终用二进制来表示是酱紫的：



所以，对于第一条数据，NULL 值列表用十六进制表示是 0x00。

接下来看**第二条记录**，第二条记录 age 列是 NULL 值，所以，对于第二条数据，NULL值列表



目录



侧边栏



夜间



技术群



资料



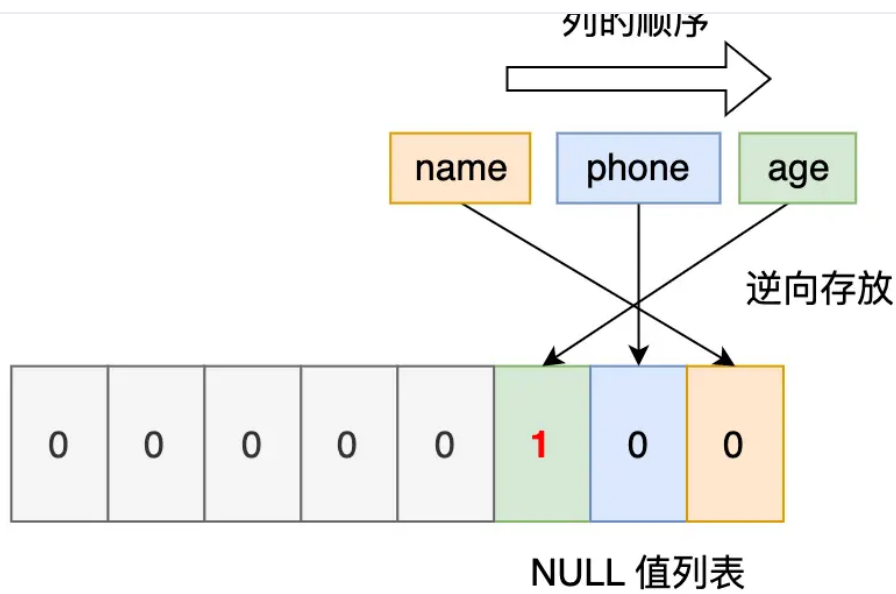
支持我



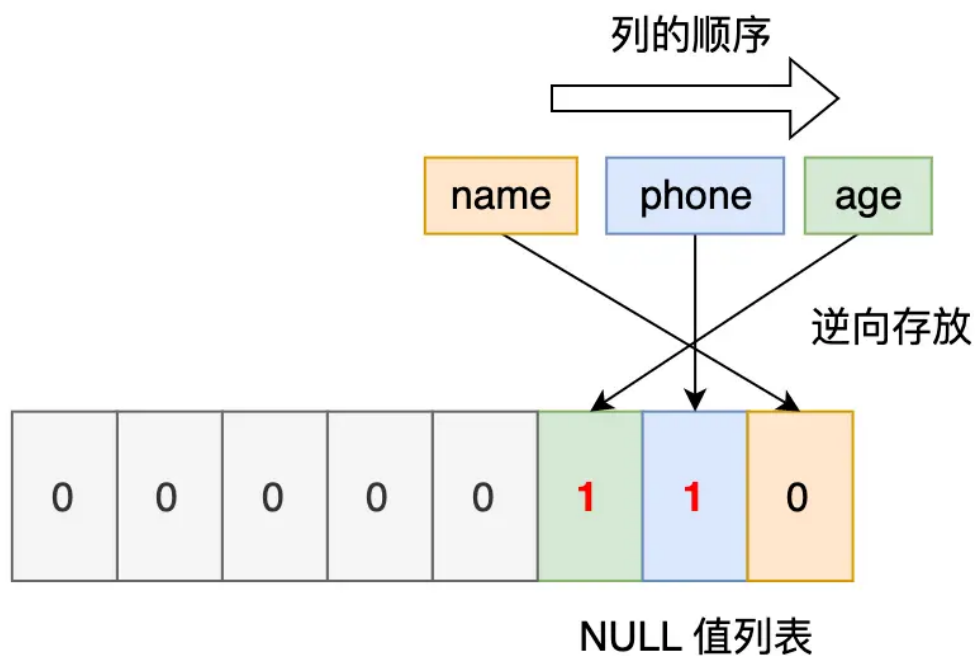
上一篇



下一篇



最后**第三条记录**，第三条记录 phone 列和 age 列是 NULL 值，所以，对于第三条数据，NULL 值列表用十六进制表示是 0x06。



我们把三条记录的 NULL 值列表都填充完毕后，它们的行格式是这样的：



目录



侧边栏



夜间



技术群



资料



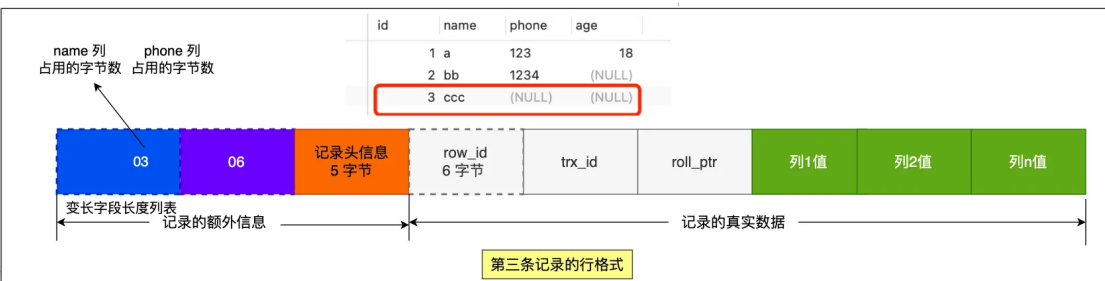
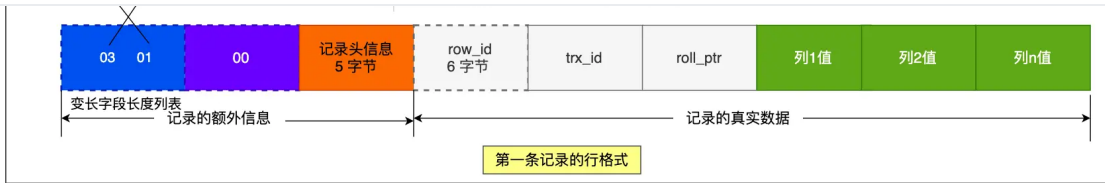
支持我



上一篇



下一篇



每个数据库表的行格式都有「NULL 值列表」吗？

NULL 值列表也不是必须的。

当数据表的字段都定义成 NOT NULL 的时候，这时候表里的行格式就不会有 NULL 值列表了。所以在设计数据库表的时候，通常都是建议将字段设置为 NOT NULL，这样可以节省 1 字节的空间（NULL 值列表占用 1 字节空间）。

3. 记录头信息

记录头信息中包含的内容很多，我就不一一列举了，这里说几个比较重要的：

- delete_mask：标识此条数据是否被删除。从这里可以知道，我们执行 delete 删除记录的时候，并不会真正的删除记录，只是将这个记录的 delete_mask 标记为 1。
- next_record：下一条记录的位置。从这里可以知道，记录与记录之间是通过链表组织的。在前面我也提到了，指向的是下一条记录的「记录头信息」和「真实数据」之间的位置，这样的好处是向左读就是记录头信息，向右读就是真实数据，比较方便。
- record_type：表示当前记录的类型，0表示普通记录，1表示B+树非叶子节点记录，2表示



目录



侧边栏



夜间



技术群



资料



支持我



上一篇



下一篇

记录的真实数据

记录真实数据部分除了我们定义的字段，还有三个隐藏字段，分别为：row_id、trx_id、roll_pointer，我们来看下这三个字段是什么。



- row_id

如果我们建表的时候指定了主键或者唯一约束列，那么就没有 row_id 隐藏字段了。如果既没有指定主键，又没有唯一约束，那么 InnoDB 就会为记录添加 row_id 隐藏字段。row_id 不是必需的，占用 6 个字节。

- trx_id

事务id，表示这个数据是由哪个事务生成的。trx_id 是必需的，占用 6 个字节。

- roll_pointer

这条记录上一个版本的指针。roll_pointer 是必需的，占用 7 个字节。

如果你熟悉 MVCC 机制，你应该就清楚 trx_id 和 roll_pointer 的作用了，如果你还不知道 MVCC 机制，可以看完[这篇文章](#)，一定要掌握，面试也很经常问 MVCC 是怎么实现的。

varchar(n) 中 n 最大取值为多少？

varchar(n) 字段类型的 n 代表的是最多存储的字符数量，那 n 最大能设置多少？

这个问题要考虑两个因素：

- 行格式中「变长字段长度列表」最大能表示多少字节？知道了这个才能知道，一行数据最大能存储多少字节的数据。
- 数据库表的字符集，确定了这个，才能知道 1 个字符占用多少字节。

行格式中「变长字段长度列表」有时候是占用 1 字节，有时候是占用 2 字节：



目录



侧边栏



夜间



技术群



资料



支持我



上一篇



下一篇

- 如果变长字段允许存储的最大字节数大于 255 字节，「变长字段长度列表」就占用 2 个字节；

可以看到，「[变长字段长度列表](#)」占用的字节数最大不会超过 2 字节。

2 个字节的最大值是 65535（十进制），从这里可以推测一行记录最大能存储 65535 字节的数据，实际上真的是这样吗？

我这里以 ascii 字符集作为例子，这意味着 1 个字符占用 1 字节。那么 varchar(65535) 就意味着最多可存储 65535 个 ascii 字符，刚好满足一行记录最大能存储 65535 字节的数据。

我们定义一个 varchar(65535) 类型的字段，字符集为 ascii 的数据库表。

```
CREATE TABLE test (  
  `name` VARCHAR(65535) NULL  
) ENGINE = InnoDB DEFAULT CHARACTER SET = ascii ROW_FORMAT = COMPACT;
```

看能不能成功创建一张表：

```
ERROR 1118 (42000): Row size too large. The maximum row size for the used table type, not counting BLOBs, is 65535. This includes storage overhead, check the manual. You have to change some columns to TEXT or BLOBs
```

可以看到，创建失败了。

从报错信息就可以知道一行数据的最大字节数是 65535（不包含 TEXT、BLOBs 这种大对象类型），其中包含了 storage overhead。

问题来了，这个 storage overhead 是什么呢？其实就是变长字段长度列表和 NULL 值列表，也就是说[一行数据的最大字节数 65535，其实是包含「变长字段长度列表」和「NULL 值列表」所占用的字节数的。](#)

我们存储字段类型为 varchar(n) 的数据时，其实分成了三个部分来存储：

- 真实数据
- 真实数据占用的字节数
- NULL 标识，如果不允许为 NULL，这部分不需要

前面我创建表的时候，字段是允许为 NULL 的，所以会占用 1 字节来存储 NULL 标识，字段是变长字段且变长字段允许存储的最大字节数大于 255 字节，所以会占用 2 字节存储真实数



目录



侧边栏



夜间



技术群



资料



支持我



上一篇



下一篇

```
mysql> CREATE TABLE test ( `name` VARCHAR(65533) NULL ) ENGINE = InnoDB DEFAULT CHARACTER SET =  
ascii ROW_FORMAT = COMPACT;  
ERROR 1118 (42000): Row size too large. The maximum row size for the used table type, not countin  
g BLOBs, is 65535. This includes storage overhead, check the manual. You have to change some colu  
mns to TEXT or BLOBs
```

可以看到，还是不行，接下来看看 varchar(65532) 是否可行？

```
mysql> CREATE TABLE test ( `name` VARCHAR(65532) NULL ) ENGINE = InnoDB DEFAULT CHARACTER SET =  
ascii ROW_FORMAT = COMPACT;  
Query OK, 0 rows affected (0.01 sec)
```

可以看到，创建成功了。

当然，我上面这个例子是针对字符集为 ascii 情况，如果采用的是 UTF-8，varchar(n) 最多能存储的数据计算方式就不一样了：

- 在 UTF-8 字符集下，一个字符串最多需要三个字节，varchar(n) 的 n 最大取值就是 $65532/3 = 21844$ 。

上面所说的只是针对于一个字段的计算方式。

如果有多个字段的话，要保证所有字段的长度 + 变长字段字节数列表所占用的字节数 + NULL值列表所占用的字节数 <= 65535。

行溢出后，MySQL 是怎么处理的？

MySQL 中磁盘和内存交互的基本单位是页，一个页的大小一般是 16KB，也就是 16384 字节，而一个 varchar(n) 类型的列最多可以存储 65532 字节，一些大对象如 TEXT、BLOB 可能存储更多的数据，这时一个页可能就存不了一条记录。这个时候就会**发生行溢出，多的数据就会存到另外的「溢出页」中。**

如果一个数据页存不了一条记录，InnoDB 存储引擎会自动将溢出的数据存放到「溢出页」中。在一般情况下，InnoDB 的数据都是存放在「数据页」中。但是当发生行溢出时，溢出的数据会存放到「溢出页」中。

当发生行溢出时，在记录的真实数据处只会保存该列的一部分数据，而把剩余的数据放在「溢出页」中，然后真实数据处用 20 字节存储指向溢出页的地址，从而可以找到剩余数据所在的页。大致如下图所示。



目录



侧边栏



夜间



技术群



资料



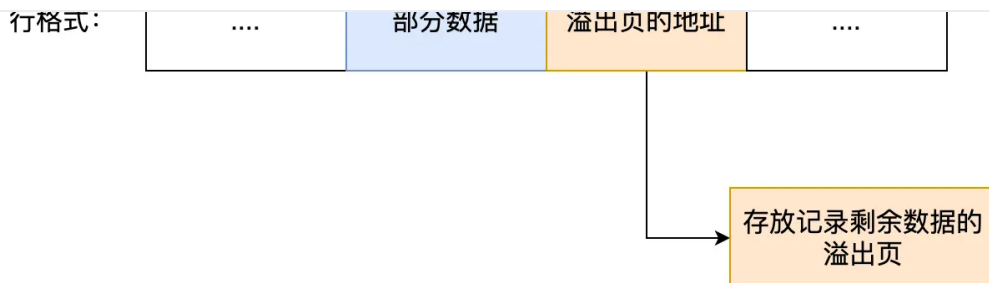
支持我



上一篇



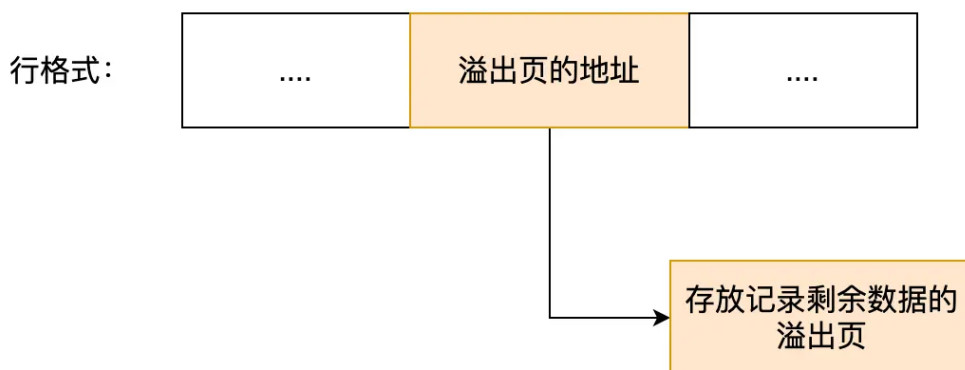
下一篇



上面这个是 Compact 行格式在发生行溢出后的处理。

Compressed 和 Dynamic 这两个行格式和 Compact 非常类似，主要的区别在于处理行溢出数据时有些区别。

这两种格式采用完全的行溢出方式，记录的真实数据处不会存储该列的一部分数据，只存储 20 个字节的指针来指向溢出页。而实际的数据都存储在溢出页中，看起来就像下面这样：



总结

MySQL 的 NULL 值是怎么存放的？

MySQL 的 Compact 行格式中会用「NULL值列表」来标记值为 NULL 的列，NULL 值并不会存储在行格式中的真实数据部分。

NULL值列表会占用 1 字节空间，当表中所有字段都定义成 NOT NULL，行格式中就不会有 NULL值列表，这样可节省 1 字节的空间。



目录



侧边栏



夜间



技术群



资料



支持我



上一篇



下一篇

MySQL 的 Compact 行格式中会用「变长字段长度列表」存储变长字段实际占用的数据大小。

varchar(n) 中 n 最大取值为多少？

一行记录最大能存储 65535 字节的数据，但是这个包含「变长字段字节数列表所占用的字节数」和「NULL值列表所占用的字节数」。

如果一张表只有一个 varchar(n) 字段，且允许为 NULL，字符集为 ascii。varchar(n) 中 n 最大取值为 65532。

计算公式：65535 - 变长字段字节数列表所占用的字节数 - NULL值列表所占用的字节数 = 65535 - 2 - 1 = 65532

行溢出后，MySQL 是怎么处理的？

如果一个数据页存不了一条记录，InnoDB 存储引擎会自动将溢出的数据存放到「溢出页」中。

Compact 行格式针对行溢出的处理是这样的：当发生行溢出时，在记录的真实数据处只会保存该列的一部分数据，而把剩余的数据放在「溢出页」中，然后真实数据处用 20 字节存储指向溢出页的地址，从而可以找到剩余数据所在的页。

Compressed 和 Dynamic 这两种格式采用完全的行溢出方式，记录的真实数据处不会存储该列的一部分数据，只存储 20 个字节的指针来指向溢出页。而实际的数据都存储在溢出页中。

参考资料：

- 《MySQL 是怎样运行的》
- 《MySQL技术内幕 InnoDB存储引擎》

最新的图解文章都在公众号首发，别忘记关注哦！！如果你想加入百人技术交流群，扫码下方二维码回复「加群」。



目录



侧边栏



夜间



技术群



资料



支持我



上一篇



下一篇