

# Part 11 数据库的设计规范

为什么要有数据库的设计规范？解决数据冗余，更新异常，插入异常，删除异常等问题。

## 1. 范式

### 1.1 键和相关属性的概念

- **超键**：能唯一标识元组的属性集叫做超键。
- **候选键**：如果超键不包括多余的属性，那么这个超键就是候选键。·**主键**：用户可以从候选键中选择一个作为主键。
- **外键**：如果数据表R1中的某属性集不是R1的主键，而是另一个数据表R2的主键，那么这个属性集就是数据表R1的外键。
- **主属性**：包含在任一候选键中的属性称为主属性。
- **非主属性**：与主属性相对，指的是不包含在任何一个候选键中的属性。

通常，我们也将候选键称之为“**码**”，把主键也称为“**主码**”。因为键可能是由多个属性组成的，针对单个属性，我们还可以用主属性和非主属性来进行区分。

### 1.2 第一范式

第一范式主要是确保数据表中每个字段的值必须具有**原子性**，也就是说数据表中每个字段的值为**不可再次拆分**的最小数据单元。

### 1.3 第二范式(2 NF)

第二范式要求，在满足第一范式的基础上，还要**满足数据表里的每一条数据记录，都是可唯一标识的。而且所有非主键字段，都必须完全依赖主键，不能只依赖主键的一部分**。如果知道主键的所有属性的值，就可以检索到任何元组(行)的任何属性的任何值。(要求中的主键，其实可以拓展替换为候选键)。

对于非主属性来说，并非完全依赖候选键。这样会产生怎样的问题呢？

- 1. 数据冗余**：如果一个球员可以参加  $m$  场比赛，那么球员的姓名和年龄就重复了  $m - 1$  次。一个比赛也可能会有  $n$  个球员参加，比赛的时间和地点就重复了  $n - 1$  次。
- 2. 插入异常**：如果我们想要添加一场新的比赛，但是这时还没有确定参加的球员都有谁，那么就没法插入。
- 3. 删除异常**：如果我要删除某个球员编号，如果没有单独保存比赛表的话，就会同时把比赛信息删除掉。
- 4. 更新异常**：如果我们调整了某个比赛的时间，那么数据表中所有这个比赛的时间都需要进行调整，否则就会出现一场比赛时间不同的情况。

1NF 告诉我们字段属性需要是原子性的，而 2NF 告诉我们一张表就是一个独立的对象，一张表只表达一个意思。

小结：2NF 要求实体的属性完全依赖主关键字。如果存在不完全依赖，那么这个属性和主关键字的这一部分应该分离出来形成一个新的实体，新实体与元实体之间是一对多的关系。

## 1.4 第三范式(3NF)

第三范式是在第二范式的基础上，确保数据表中的每一个非主键字段都和主键字段直接相关，也就是说，**要求数据表中的所有非主键字段不能依赖于其他非主键字段**。(即，不能存在非主属性A依赖于非主属性B，非主属性B依赖于主键C的情况，即存在“A→B→C”的决定关系)通俗地讲，该规则的意思是所有**非主键属性**之间不能有依赖关系，必须**相互独立**。

这里的主键可以拓展为候选键。

符合3NF后的数据模型通俗地讲，2NF和3NF通常以这句话概括：“每个非键属性依赖于键，依赖于整个键，并且除了键别无他物”。

## 1.5 小结

- 第一范式(1NF)，确保每列保持**原子性**
  - 数据库的每一列都是不可分割的原子数据项，不可再分的最小数据单元，而不能是集合、数组、记录等非原子数据项。
- 第二范式(2NF)，确保每列都和主键**完全依赖**
  - 尤其在复合主键的情况下，非主键部分不应该依赖于部分主键。
- 第三范式(3NF) 确保每列都和主键列**直接相关**，而不是间接相关

**范式的优点:** 数据的标准化有助于消除数据库中的**数据冗余**，第三范式(3NF) 通常被认为在性能、扩展性和数据完整性方面达到了最好的平衡。

**范式的缺点:** 范式的使用，可能**降低查询的效率**。因为范式等级越高，设计出来的数据表就越多、越精细，数据的冗余度就越低，进行数据查询的时候就可能需要**关联多张表**，这不但代价昂贵，也可能使一些**索引策略无效**。

范式只是提出了设计的标准，实际上设计数据表时，未必一定要符合这些标准。开发中，我们会出现为了性能和读取效率违反范式化的原则，通过**增加少量的冗余或重复的数据**来提高数据库的**读性能**，减少关联查询，join表的次数，实现空间换取时间的目的。因此在实际的设计过程中要理论结合实际，灵活运用。

范式本身没有优劣之分，只有适用场景不同。没有完美的设计，只有合适的设计，我们在数据表的设计中，还需要根据需求将范式和反范式混合使用。

## 2. 反范式化

### 2.1 概述

规范化 vs 性能

- 为满足某种商业目标，数据库性能比规范化数据库更重要。
- 在数据规范化的同时，要综合考虑数据库的性能。
- 通过在给定的表中添加额外的字段，以大量减少需要从中搜索信息所需的时间。
- 通过在给定的表中插入计算列，以方便查询。

## 2.2 举例

想要查询课程 ID 为 10001 的前 1000 条评论。

如果评论表中没有存储学生的姓名，只存储了学生的id，则显示评论时，需要将评论表和学生表连接，效率较低。

如果评论表中冗余地存储了学生姓名，则无需连接即可展示结果，效率变高。

## 2.3 反范式的新问题

- 存储空间变大了（**数据冗余**）
- 一个表中字段做了修改，另一个表中冗余的字段也需要做同步修改，否则数据不一致（**更新问题**）
- 若采用**存储过程**来支持数据的更新、删除等额外操作，如果更新频繁，会非常消耗系统资源
- 在**数据量小**的情况下，反范式不能体现性能的优势，可能还会让数据库的设计更加复杂

## 2.4 反范式的适用场景

当冗余信息有价值或者能**大幅度提高查询效率**的时候，我们才会采取反范式的优化。

- 增加冗余字段的建议
- 历史快照、历史数据的需要
  - 在现实生活中，我们经常需要一些冗余信息，比如**订单中的收货人信息**，包括姓名、电话和地址等。每次发生的 订单收货信息 都属于 历史快照，需要进行保存，但用户可以随时修改自己的信息，这时保存这些冗余信息是非常有必要的。
  - 反范式优化也常用在 **数据仓库** 的设计中，因为数据仓库通常 存储历史数据，对增删改的实时性要求不强，对历史数据的分析需求强。这时适当允许数据的冗余度，更方便进行数据分析。

## 3. BCNF(巴斯范式)

主属性（仓库名）对于候选键（管理员，物品名）是部分依赖的关系，这样就有可能导致异常情况。因此引入BCNF，它在 3NF 的基础上消除了主属性对候选键的部分依赖或者传递依赖关系。

## 4. ER模型

### 4.1 ER 模型包括那些要素？

ER 模型中有三个要素，分别是**实体**、**属性**和**关系**。

- **实体**，可以看做是数据对象，往往对应于现实生活中的真实存在的个体。在 ER 模型中，用 **矩形** 来表示。实体分为两类，分别是 **强实体** 和 **弱实体**。强实体是指不依赖于其他实体的实体；弱实体是指对另一个实体有很强的依赖关系的实体。
- **属性**，则是指实体的特性。比如超市的地址、联系电话、员工数等。在 ER 模型中用 **椭圆形** 来表示。
- **关系**，则是指实体之间的联系。比如超市把商品卖给顾客，就是一种超市与顾客之间的联系。在 ER 模型中用 **菱形** 来表示。

注意：实体和属性不容易区分。这里提供一个原则：我们要从系统整体的角度出发去看，**可以独立存在的是实体，不可再分的是属性**。也就是说，属性不能包含其他属性。

## 4.2 关系的类型

在 ER 模型的 3 个要素中，关系又可以分为 3 种类型，分别是一对一、一对多、多对多。

- **一对一**：指实体之间的关系是一一对应的。
- **一对多**：指一边的实体通过关系，可以对应多个另外一边的实体。相反，另外一边的实体通过这个关系，则只能对应唯一的一边的实体。
- **多对多**：指关系两边的实体都可以通过关系对应多个对方的实体。

## 5. 数据表的设计原则

数据表设计的一般原则："三少一多"

- **数据表的个数越少越好**
- **数据表中的字段个数越少越好**
- **数据表中联合主键的字段个数越少越好**
- **使用主键和外键越多越好**

注意：这个原则并不是绝对的，有时候我们需要牺牲数据的冗余度来换取数据处理的效率。

## 6. 数据库对象编写建议

### 6.1 关于库

1. 【强制】库的名称必须控制在32个字符以内，只能使用英文字母、数字和下划线，建议以英文字母开头。
2. 【强制】库名中英文**一律小写**，不同单词采用**下划线**分割。须见名知意。
3. 【强制】库的名称格式：业务系统名称\_子系统名。
4. 【强制】库名禁止使用关键字（如type,order等）。
5. 【强制】创建数据库时必须**显式指定字符集**，并且字符集只能是utf8或者utf8mb4。创建数据库SQL举例：CREATE DATABASE crm\_fund **DEFAULT CHARACTER SET 'utf8'**;
6. 【建议】对于程序连接数据库账号，遵循**权限最小原则**。使用数据库账号只能在一个DB下使用，不准跨库。程序使用的账号**原则上不准有drop权限**。
7. 【建议】临时库以**tmp\_**为前缀，并以日期为后缀；备份库以**bak\_**为前缀，并以日期为后缀。

### 6.2 关于表、列

1. 【强制】表和列的名称必须控制在32个字符以内，表名只能使用英文字母、数字和下划线，建议以**英文字母开头**。
2. 【强制】**表名、列名一律小写**，不同单词采用下划线分割。须见名知意。

3. 【强制】表名要求有模块名强相关，同一模块的表名尽量使用 统一前缀。比如：crm\_fund\_item
4. 【强制】创建表时必须 显式指定字符集 为utf8或utf8mb4。
5. 【强制】表名、列名禁止使用关键字（如type,order等）。
6. 【强制】创建表时必须 显式指定表存储引擎 类型。如无特殊需求，一律为InnoDB。
7. 【强制】建表必须有comment。
8. 【强制】字段命名应尽可能使用表达实际含义的英文单词或 缩写。如：公司 ID，不要使用 corporation\_id, 而用corp\_id 即可。
9. 【强制】布尔值类型的字段命名为 is\_描述。如member表上表示是否为enabled的会员的字段命名为 is\_enabled。
10. 【强制】禁止在数据库中存储图片、文件等大的二进制数据。通常文件很大，短时间内造成数据量快速增长，数据库进行数据库读取时，通常会进行大量的随机IO操作，文件很大时，IO操作很耗时。通常存储于文件服务器，数据库只存储文件地址信息。
11. 【建议】建表时关于主键：表必须有主键
  1. 强制要求主键为id，类型为int或bigint，且为auto\_increment建议使用unsigned无符号型。
  2. 标识表里每一行主体的字段不要设为主键，建议设为其他字段如user\_id，order\_id等，并建立unique key索引。因为如果设为主键且主键值为随机插入，则会导致innodb内部页分裂和大量随机I/O，性能下降。
12. 【建议】核心表（如用户表）必须有行数据的 创建时间字段（create\_time）和 最后更新时间字段（update\_time），便于查问题。
13. 【建议】表中所有字段尽量都是 NOT NULL 属性，业务可以根据需要定义 DEFAULT值。因为使用 NULL值会存在每一行都会占用额外存储空间、数据迁移容易出错、聚合函数计算结果偏差等问题。
14. 【建议】所有存储相同数据的 列名和列类型必须一致（一般作为关联列，如果查询时关联列类型不一致会自动进行数据类型隐式转换，会造成列上的索引失效，导致查询效率降低）。
15. 【建议】中间表（或临时表）用于保留中间结果集，名称以 tmp\_ 开头。备份表用于备份或抓取源表快照，名称以 bak\_ 开头。中间表和备份表定期清理。
16. 【示范】一个较为规范的建表语句：

```
1 CREATE TABLE user_info (  
2     `id` int unsigned NOT NULL AUTO_INCREMENT COMMENT '自增主键',  
3     `user_id` bigint(11) NOT NULL COMMENT '用户id',  
4     `username` varchar(45) NOT NULL COMMENT '真实姓名',  
5     `email` varchar(30) NOT NULL COMMENT '用户邮箱',  
6     `nickname` varchar(45) NOT NULL COMMENT '昵称',  
7     `birthday` date NOT NULL COMMENT '生日',  
8     `sex` tinyint(4) DEFAULT '0' COMMENT '性别',  
9     `short_introduce` varchar(150) DEFAULT NULL COMMENT '一句话介绍自己，最多50  
    个汉字',  
10    `user_resume` varchar(300) NOT NULL COMMENT '用户提交的简历存放地址',  
11    `user_register_ip` int NOT NULL COMMENT '用户注册时的源ip',  
12    `create_time` timestamp NOT NULL DEFAULT CURRENT_TIMESTAMP COMMENT '创建  
    时间',  
13    `update_time` timestamp NOT NULL DEFAULT CURRENT_TIMESTAMP ON UPDATE  
    CURRENT_TIMESTAMP COMMENT '修改时间',  
14    `user_review_status` tinyint NOT NULL COMMENT '用户资料审核状态，1为通过，2为  
    审核中，3为未 通过，4为还未提交审核',  
15    PRIMARY KEY (`id`),  
16    UNIQUE KEY `uniq_user_id` (`user_id`),
```

```

17     KEY `idx_username`(`username`),
18     KEY `idx_create_time_status`(`create_time`,`user_review_status`)
19 ) ENGINE=InnoDB DEFAULT CHARSET=utf8 COMMENT='网站用户基本信息'

```

17. 【建议】创建表时，可以使用可视化工具。这样可以确保表、字段相关的约定都能设置上。实际上，我们通常很少自己写 DDL 语句，可以使用一些可视化工具来创建和操作数据库和数据表。可视化工具除了方便，还能直接帮我们将数据库的结构定义转化成 SQL 语言，方便数据库和数据表结构的导出和导入。

## 6.3 关于索引

1. 【强制】InnoDB表必须主键为id int/bigint auto\_increment，且主键值 禁止被更新。
2. 【强制】InnoDB和MyISAM存储引擎表，索引类型必须为 BTREE。
3. 【建议】主键的名称以 pk\_ 开头，唯一键以 uni\_ 或 uk\_ 开头，普通索引以 idx\_ 开头，一律使用小写格式，以字段的名称或缩写作为后缀。
4. 【建议】多单词组成的columnname，取前几个单词首字母，加末单词组成column\_name。如：sample 表 member\_id 上的索引：idx\_sample\_mid。
5. 【建议】单个表上的索引个数 不能超过6个。
6. 【建议】在建立索引时，多考虑建立 联合索引，并把区分度最高的字段放在最前面。
7. 【建议】在多表 JOIN 的SQL里，保证被驱动表的连接列上有索引，这样JOIN 执行效率最高。
8. 【建议】建表或加索引时，保证表里互相不存在 冗余索引。比如：如果表里已经存在key(a,b)，则 key(a)为冗余索引，需要删除。

## 6.4 SQL编写

1. 【强制】程序端SELECT语句必须指定具体字段名称，禁止写成\*。
2. 【建议】程序端insert语句指定具体字段名称，不要写成INSERT INTO t1 VALUES(...)。
3. 【建议】除静态表或小表（100行以内），DML语句必须有WHERE条件，且使用索引查找。
4. 【建议】INSERT INTO...VALUES(XX),(XX),(XX).. 这里XX的值不要超过5000个。值过多虽然上线很快，但会引起主从同步延迟。
5. 【建议】SELECT语句不要使用UNION，推荐使用UNION ALL，并且UNION子句个数限制在5个以内。
6. 【建议】线上环境，多表 JOIN 不要超过5个表。
7. 【建议】减少使用ORDER BY，和业务沟通能不排序就不排序，或将排序放到程序端去做。ORDER BY、GROUP BY、DISTINCT 这些语句较为耗费CPU，数据库的CPU资源是极其宝贵的。
8. 【建议】包含了ORDER BY、GROUP BY、DISTINCT 这些查询的语句，WHERE 条件过滤出来的结果集请保持在1000行以内，否则SQL会很慢。
9. 【建议】对单表的多次alter操作必须合并为一次。对于超过100W行的大表进行alter table，必须经过DBA审核，并在业务低峰期执行，多个alter需整合在一起。因为alter table会产生 表锁，期间阻塞对于该表的所有写入，对于业务可能会产生极大影响。
10. 【建议】批量操作数据时，需要控制事务处理间隔时间，进行必要的sleep。
11. 【建议】事务里包含SQL不超过5个。因为过长的事务会导致锁数据较久，MySQL内部缓存、连接消耗过多等问题。

12. 【建议】事务里更新语句尽量基于主键或UNIQUE KEY，如UPDATE... WHERE id=XX;否则会产生间隙锁，内部扩大锁定范围，导致系统性能下降，产生死锁。