

# 数据库系统原理笔记

陈鸿峥

2019.11\*

## 目录

<b>1 数据库系统概述</b>	<b>1</b>
1.1 高层概述 . . . . .	1
1.2 数据库语言 . . . . .	2
1.3 关系型数据库 . . . . .	2
<b>2 SQL简介</b>	<b>3</b>
2.1 数据定义语言 . . . . .	3
2.2 数据查询语言 . . . . .	3
<b>3 关系代数</b>	<b>7</b>
<b>4 关系型数据库设计</b>	<b>7</b>

本课程采用书目Avi Silberschatz, Henry F. Korth, S. Sudarshan, *Database System Concepts (6th ed)*<sup>1</sup>。

## 1 数据库系统概述

### 1.1 高层概述

早期的数据库直接建立在文件系统上，但这会导致：

- 数据冗余与不一致
- 访问数据非常麻烦
- 完整性问题：难以添加限制（如年龄为非负整数）
- 更新的原子性

---

\*Build 20191120

<sup>1</sup><http://www.db-book.com>

- 多用户的并发访问
- 安全性问题：权限

数据抽象包括：

- 物理层：存储块、物理组织
- 逻辑层：通过类型定义进行描述，同时记录类型之间的相互关系
- 视图层：屏蔽数据类型细节的一组应用程序，同时提供了访问权限

查询过程：解释编译+求值(evaluation)

## 1.2 数据库语言

- 数据操纵语言(Data Manipulation Language, DML)
- 数据定义语言(Data Definition Language, DDL)：DDL的输出放在数据字典中，数据字典包含元数据(metadata)，需要满足一致性约束
  - 域约束：如整型、字符型等
  - 参照完整性(referential integrity)：一个关系中给定属性集上的取值也在另一关系的某一属性集的取值中出现
  - 断言(assertion)：域约束和参照完整性只是断言的特殊形式，如“每一学期每一个系必须至少开设5门课程”
  - 权限(authorization)

## 1.3 关系型数据库

每一个表(table)都是一个关系(relation)，纵向为属性(attributes/columns)，横向为元组(tuples/rows)。一组特定的行称为关系实例(instance)。

ID	name	dept_name
12345	Chen	CS
10001	Bob	Biology
10101	Alice	Physics

注意关系都是无序的，元组可以以任意顺序存储。

数据库模式(schema)是数据库的逻辑设计，如instructor(ID,name,dept\_name,salary)；数据库实例(instance)则是给定时刻数据库中数据的一个快照，与具体数据相关。

- 超码(superkey)：某些属性的集合使得在一个关系中可以唯一地标识一个元组，如ID属性
- 候选码(candidate key)：最小超码，即其任意真子集都不能成为超码
- 主码(primary key)：用在一个关系中区分不同元组的候选码
- 外码(foreign key)：外部参照关系

## 2 SQL简介

SQL即结构化查询语言(Structured Query Language, SQL)

### 2.1 数据定义语言

数据类型

- `char(n)`: 固定长度字符串
- `varchar(n)`: 可变长度字符串, 用得最多
- `int`
- `smallint`
- `numeric(p,d)`: 定点数, 共 $p$ 位 (包括符号位), 其中 $d$ 位在小数点右侧
- `real, double precision`
- `float(n)`: 精度至少为 $n$ 位的浮点数

```
create table instructor (  
    ID char(5),  
    name varchar(20),  
    dept_name varchar(20),  
    salary numeric(8,2),  
    primary key (ID), -- integrity constraints  
    foreign key (dept_name) references department);
```

### 2.2 数据查询语言

基本的形式如下所示

```
select A1, A2, ..., An  
from r1, r2, ..., rm  
where P
```

其中 $A_i$ 为属性(attribute)、 $R_i$ 为一个关系(relation)/一个表、 $P$ 是谓词(predicate), 返回满足关系的数组。

- 自然连接(natural join): 将两个关系的各个元组进行连接, 并且只考虑那些在两个关系模式中都出现的属性上取值相同的元组对 (相当于取交)
- 连接(join): 只考虑 $A$ 这个属性进行连接, `join ... using (A)`

字符串运算: 模式匹配

- `upper(s)`、`lower(s)`
- `%`: 匹配任意字符串
- `_`: 匹配任意一个字符

- like: 以这些字符串进行匹配

```
select distinct dept_name
from instructor

select * -- all attributes
from instructor

select '437' as F00

select ID, name, salary/12 as monthly_salary

select name
from instructor
where dept_name = 'Comp. Sci.' and salary > 80000

select *
from instructor, teaches -- Cartesian product

select distinct T.name
from instructor as T, instructor as S -- rename, Cartesian product
where T.salary > S.salary and S.dept_name = 'Comp. Sci.'
-- salary between 90000 and 100000

/**
percent ( % ). The % character matches any substring.
underscore ( _ ). The _ character matches any character.
**/

select name
from instructor
where name like '%dar%' matches any string containing "dar" as a substring

select distinct name
from instructor
order by name

select *
from instructor
order by salary desc, name asc;

(select course_id from section where sem = 'Fall' and year = 2017)
union -- intersect, except
(select course_id from section where sem = 'Spring' and year = 2018)
```

```

select name
from instructor
where salary is null

-- arg, min, max, sum, count
select dept_name, avg (salary) as avg_salary
from instructor
where dept_name= 'Comp. Sci.'
group by dept_name;

select count (distinct ID)
from teaches
where semester = 'Spring' and year = 2018;

-- group constraints
select dept_name, avg(salary) as avg_salary
from instructor
group by dept_name
having avg(salary) > 42000;

```

三值逻辑，添加了Unknown

AND	OR	NOT
$T \wedge U = U$	$T \vee U = T$	$\neg U = U$
$F \wedge U = F$	$F \vee U = U$	
$U \wedge U = U$	$U \vee U = U$	

集合成员资格通过in来判断，同时有some和all的关系测试

```

select name
from instructor
where salary > some (select salary
                     from instructor
                     where dept_name = 'Biology')

select S.ID, S.name
from student as S
where not exists ((select course_id
                  from course
                  where dept_name = 'Biology')
except
(select T.course_id
 from takes as T
 where S.ID = T.ID))

```

```

select T.course_id
from course as T
where 1 >= (select count(R.course_id)
            from section as R
            where T.course_id = R.course_id and R.year = 2009);
-- not unique

with max_budget(value) as
    (select max(budget)
     from department)
select budget
from department, max_budget
where department.budget = max_budget.value;

select dept_name,
       (select count(*)
        from instructor
        where department.dept_name = instructor.dept_name)
       as num_instructors
from department; -- scalar subquery

```

数据库的修改操作如下

```

delete from instructor
where dept_name = 'Finance';

insert into course(course_id, title)
values ('CS-437', 'Database Systems');

update instructor
set salary = salary * 1.05
where salary < 70000;

update instructor
set salary = case
    when salary <= 1000 then salary * 1.05
    when salary <= 2000 then salary * 2.05
    else salary * 1.03
end

```

### 3 关系代数

选择	$\sigma$	挑选出符合一定性质的元组 $\sigma_{\text{Sub}=\text{"Phy"} \wedge \text{Age} > 30}(\text{teachers})$
投影	$\Pi$	只选出对应属性 $\Pi_{\text{ID, name, salary}}(\text{teachers})$
笛卡尔积	$\times$	将两个关系整合（简单并置，需要进一步筛选）
合并	$r \bowtie_{\theta} s = \sigma_{\theta}(r \times s)$	
并集	$\cup$	数目应相同，属性可兼容
交集	$\cap$	
差集	$-$	
赋值	$\leftarrow$	
重命名	$\rho_x(E)$	给 $E$ 的返回值赋名为 $x$

### 4 关系型数据库设计

数据库设计的范式(normal form)

- 第一范式：原子性，即所有元素都是不可分的单元

•

**定义 1** (函数依赖(functional dependency)). 设 $R$ 为关系模式,  $\alpha \subset R, \beta \subset R$ , 函数依赖 $\alpha \rightarrow \beta$ 在 $R$ 上满足, 当且仅当对于任意合法的关系 $r(R)$ , 任何两个关于 $r$ 的数对 $t_1$ 和 $t_2$ , 如果满足属性 $\alpha$ , 那么它们也满足属性 $\beta$ , 即

$$t_1[\alpha] = t_2[\alpha] \implies t_1[\beta] = t_2[\beta]$$

实际上就是函数单射的概念。

考虑下例 $r(A, B)$ 的关系

A	B
1	4
1	5
3	7

关系 $A \rightarrow B$ 不成立, 但关系 $B \rightarrow A$ 成立。

**定义 2** (多值依赖(multivalued dependency)). 设 $R$ 为关系模式,  $\alpha \subset R, \beta \subset R$ , 多值依赖在 $R$ 上满足 $\alpha \twoheadrightarrow \beta$

$\beta$ , 当且仅当对于所有数对 $t_1$ 和 $t_2$ , 使得 $t_1[\alpha] = t_2[\alpha]$ , 都存在数对 $t_3$ 和 $t_4$ 使得

$$\begin{aligned} t_1[\alpha] &= t_2[\alpha] = t_3[\alpha] = t_4[\alpha] \\ t_3[\beta] &= t_1[\beta] \\ t_3[R - \beta] &= t_2[R - \beta] \\ t_4[\beta] &= t_2[\beta] \\ t_4[R - \beta] &= t_1[R - \beta] \end{aligned}$$

简而言之, 设 $(x, y, z)$ 有对应着属性 $\alpha, \beta, R - \alpha - \beta$ , 则对于任意 $R$ 中的元素 $(a, b, c)$ 和 $(a, d, e)$ , 都有 $(a, b, e)$ 和 $(a, d, c)$ 在 $I$

## 5 事务

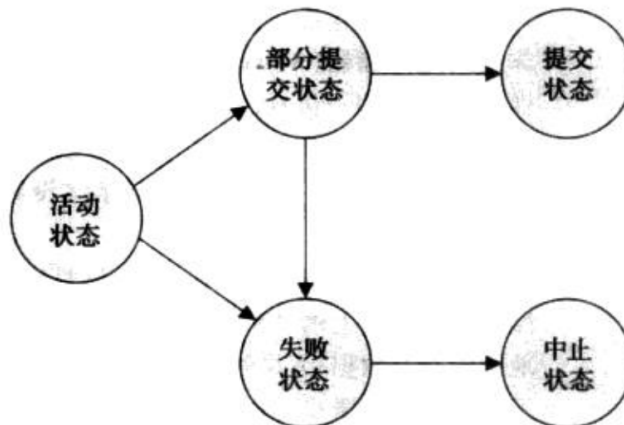
**定义 3** (事务(transaction)). 构成单一逻辑工作单元的操作集合称为事务。即使有故障, 数据库系统也必须保证数据库的正确执行——要么执行整个事务, 要么属于该事务的操作一个也不执行。

事务具有以下的基本性质:

- 原子性: 要么执行完, 要么不执行, 不能执行到一半。
- 一致性: 除了基本的数据完整性约束, 还有更多的一致性约束。
- 隔离性: 每个事务都察觉不到系统中有其他事务在并发执行, 一定是完成一个再进行下一个。
- 持久性: 一个事务成功完成对数据库的改变是永久的, 即使出现系统故障。

事务的基本状态:

- 活动的(active): 初始状态
- 部分提交的(partially committed): 最后一条语句执行后
- 失败的(failed): 执行出错
- 中止的(aborted): 事务回滚且数据库已恢复到事务开始执行前
- 提交的(committed): 成功完成





**定义 4 (冲突).** 当 $I$ 和 $J$ 是不同事务在相同数据项上的操作, 并且其中至少有一个是write指令时, 则称 $I$ 与 $J$ 是冲突的。

**定义 5 (冲突等价(conflict equivalent)).** 如果调度 $S$ 可以经过一系列非冲突指令交换转换为 $S'$ , 则称 $S$ 和 $S'$ 是冲突等价的。若 $S'$ 为串行调度, 则 $S$ 为冲突可串行化的(*conflict serializable*)。

可以通过构造一个优先图(precedence graph)来判断是否冲突可串行化。如果无环 (用环检测算法), 则冲突可串行化, 可以通过拓扑排序得到串行化顺序。