

# CS2102

## Database Systems

AY2022/23 Semester 2

Notes by Jonathan Tay

Last updated on January 24, 2023

---

## Contents

<b>I</b>	<b>Relational Model</b>	<b>1</b>
<b>1</b>	<b>Relational Algebra</b>	<b>1</b>
1.1	Equivalence and Compatibility . . . . .	1
1.2	Basic Operators . . . . .	1
1.3	Set Operators . . . . .	2
1.4	Join Operators . . . . .	2

---

# Part I

## Relational Model

Data in **relational databases** are stored in **relations** (tables). Column headers are **attributes**, and rows are **tuples**.

The **degree** is the number of columns, and the **cardinality** is the number of rows.

The **domain** of an attribute  $A_i$ , denoted as  $\text{dom}(A_i)$ , is the set of all possible *atomic* values for  $A_i$ . NULL is an additional special value for unknown or invalid values.

keys

A **superkey** is a subset of attributes that uniquely identifies a tuple. A **key** is a *minimal* superkey.

The **candidate keys** is the set of all keys for a relation, of which one is selected as a **primary key**.

Primary key values must be non-NULL.

foreign keys

A **foreign key** is a subset of attributes of a *referencing relation* that refers to the primary key of a *referenced relation*:

(referencing attributes)  $\rightsquigarrow$  (referenced attributes)

Because the names of the attributes are not necessarily unique, each attribute is prefixed with the name of the relation, like so:

(<relation name> · <attribute name>, ...)  $\rightsquigarrow$  ...

Foreign keys must appear as a primary key in the referenced table, NULL, or a tuple containing NULL.

The key constraints above are not intrinsic properties of a relation; rather, they are specified by the database designer to avoid problematic but otherwise valid data.

### 1 Relational Algebra

Relation are always the **operands** in relational algebra, on which **operators** are applied.

Because relations are closed under *any combination of operators*, the result of an operation is *always* a relation, and no other output is possible (**closure property**).

3-valued logic (false, true, and NULL)

Any operation involving NULL will result in NULL. Hence,  $\equiv$  and  $\neq$  are needed to compare (in)equality of NULL values directly:

NULL = NULL produces NULL, but  $\text{NULL} \equiv \text{NULL}$  produces true.

$a \wedge b$		$a$		
		true	false	NULL
$b$	true	true	false	NULL
	false	false	false	false
	NULL	NULL	false	NULL

$a \vee b$		$a$		
		true	false	NULL
$b$	true	true	true	true
	false	true	false	NULL
	NULL	true	NULL	NULL

#### 1.1 Equivalence and Compatibility

Two *expressions* are **equivalent** if either *both* produce an error, or *both* produce the same result.

Errors occur if attributes are missing (e.g. by projection or renaming), or are incompatible — there is no implicit type conversion.

The order of types in a tuple matters — (int, text) is not equivalent to (text, int).

Two *relations* are **union-compatible** if they have the same number of attributes with the same domain (type).

#### 1.2 Basic Operators

Note that logical conjunction ( $\wedge$ ) has greater precedence than logical disjunction ( $\vee$ ).

selection —  $\sigma_{[c]}(R)$

Filters the rows of relation  $R$ , returning the set of tuples that satisfy the condition  $c$ .

Conditions which evaluate to NULL are excluded from the result.

The condition  $c$  must only specify attributes of  $R$ .

projection —  $\pi_{[l]}(R)$

Maps the relation  $R$ , returning the set of tuples with the attributes in the ordered list  $l$ , in the order specified by  $l$ .

Equivalent to a column filter, but the rows which were previously unique tuples may no longer be unique, and are therefore de-duplicated.

The elements of  $l$  must not be operations, must be attributes of  $R$ , and must be unique.

**renaming** —  $\rho_{[\mathcal{R}]}(R)$ 

Renames the attributes of relation  $R$  to the attributes in the unordered list  $\mathcal{R}$ .

Elements of  $\mathcal{R}$  must be in the following form:  
 $\langle \text{new name} \rangle \leftarrow \langle \text{old name} \rangle$ .

New attribute names must be unique, and existing attributes must only be renamed at most once per operation (i.e.  $\rho_{[A \leftarrow B, B \leftarrow C]}$  is invalid).

**1.3 Set Operators**

The typical set union ( $\cup$ ), set intersection ( $\cap$ ), and set difference ( $-$ ) operators are omitted for brevity.

**cross product** —  $R \times S$ 

For every tuple in  $R$ , concatenate it with every tuple in  $S$  to form a new relation, such that the cardinality (number of rows) of the result is  $|R| \times |S|$ .

The set of attributes in  $R$  and  $S$  must be disjoint, such that the degree (number of columns) of the result is  $\text{deg}(R) + \text{deg}(S)$ .

Cross products are associative —  $R \times (S \times T) = (R \times S) \times T$ .

**1.4 Join Operators**

**Joins** are a composite operator, composing cross product, selection, and projection on a relation.

This concatenates two tables and removes unwanted/redundant rows and columns from the result of a cross product.

**theta-join ( $\theta$ -join)** —  $R \bowtie_{[\theta]} S$ 

Cross  $R$  and  $S$ , then filter (by selection) the rows using the condition  $\theta$ .

**equi-join** —  $R \bowtie_{=} S$ 

A special case of theta-join, where only equality ( $=$ ) conditions are allowed (c.f.  $\theta$ -join which allows  $\equiv$ ,  $\leq$ ,  $<>$ , etc.).

This may be more performant versus a theta-join as hashing can be used internally.

**natural inner join** —  $R \bowtie S$ 

First, find the set of attributes that are common to both  $R$  and  $S$ .

Then, cross  $R$  and  $S$ , and retain (by selection) the rows for which the **common attributes** are equal —

this also eliminates rows with any NULL value.

If there are no common attributes, then this is simply the cross product (by vacuous truth).

Finally, de-duplicate (by projection) the columns by their attribute names.

Theta-joins, equi-joins, and natural inner joins are collectively known as **inner joins**.

If we wish to perform a join but still retain *all* rows from the *left* table, *right* table, or even *both* tables and simply pad missing values with NULL, we can use **outer joins**.

**full outer join** —  $L \bowtie_{[\theta]}^{\text{full}} R$ 

attributes of $L$		attributes of $R$	
values from $L$	...	values from $R$	...
$\vdots$	$\ddots$	$\vdots$	$\ddots$
values from $L$	...	NULL	...
$\vdots$	$\ddots$	$\vdots$	$\ddots$
NULL	...	values from $R$	...
$\vdots$	$\ddots$	$\vdots$	$\ddots$

Inner joins only retain rows (in **red**) from which both  $L$  and  $R$  have a value for which the condition  $\theta$  evaluates to **true** (during the selection).

However, rows (in **green** and **blue**) which do not satisfy  $\theta$  may also be desirable, and can be retained by outer joins.

The **full outer join** retains all of the rows above.

**left outer join** —  $L \bowtie_{[\theta]}^{\text{left}} R$ 

The left outer join retains the **red** and **green** rows.

**right outer join** —  $L \bowtie_{[\theta]}^{\text{right}} R$ 

The right outer join retains the **red** and **blue** rows.

**natural outer joins**

The natural keyword can be prefixed to the left, right, or full outer joins, e.g. "natural left outer join".

The equality condition is implicitly defined over the set of attributes that are common to both  $L$  and  $R$ .