

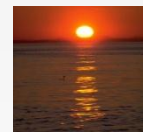


Chapter 2: Relational Model

Database System Concepts, 5th Ed.

©Silberschatz, Korth and Sudarshan

See www.db-book.com for conditions on re-use





Chapter 2: Relational Model

- Structure of Relational Databases
- Fundamental Relational-Algebra-Operations
- Additional Relational-Algebra-Operations
- Extended Relational-Algebra-Operations
- Null Values
- Modification of the Database





Example of a Relation

<i>account_number</i>	<i>branch_name</i>	<i>balance</i>
A-101	Downtown	500
A-102	Perryridge	400
A-201	Brighton	900
A-215	Mianus	700
A-217	Brighton	750
A-222	Redwood	700
A-305	Round Hill	350





Attribute Types

- Each attribute of a relation has a name
- The set of allowed values for each attribute is called the **domain** of the attribute
- Attribute values are (normally) required to be **atomic**; that is, indivisible
 - E.g. the value of an attribute can be an account number, but cannot be a set of account numbers
- Domain is said to be atomic if all its members are atomic
- The special value *null* is a member of every domain
- The null value causes complications in the definition of many operations
 - We shall ignore the effect of null values in our main presentation and consider their effect later





Relation Schema

- Formally, given domains D_1, D_2, \dots, D_n a **relation** r is a subset of

$$D_1 \times D_2 \times \dots \times D_n$$

Thus, a relation is a set of n -tuples (a_1, a_2, \dots, a_n) where each $a_i \in D_i$

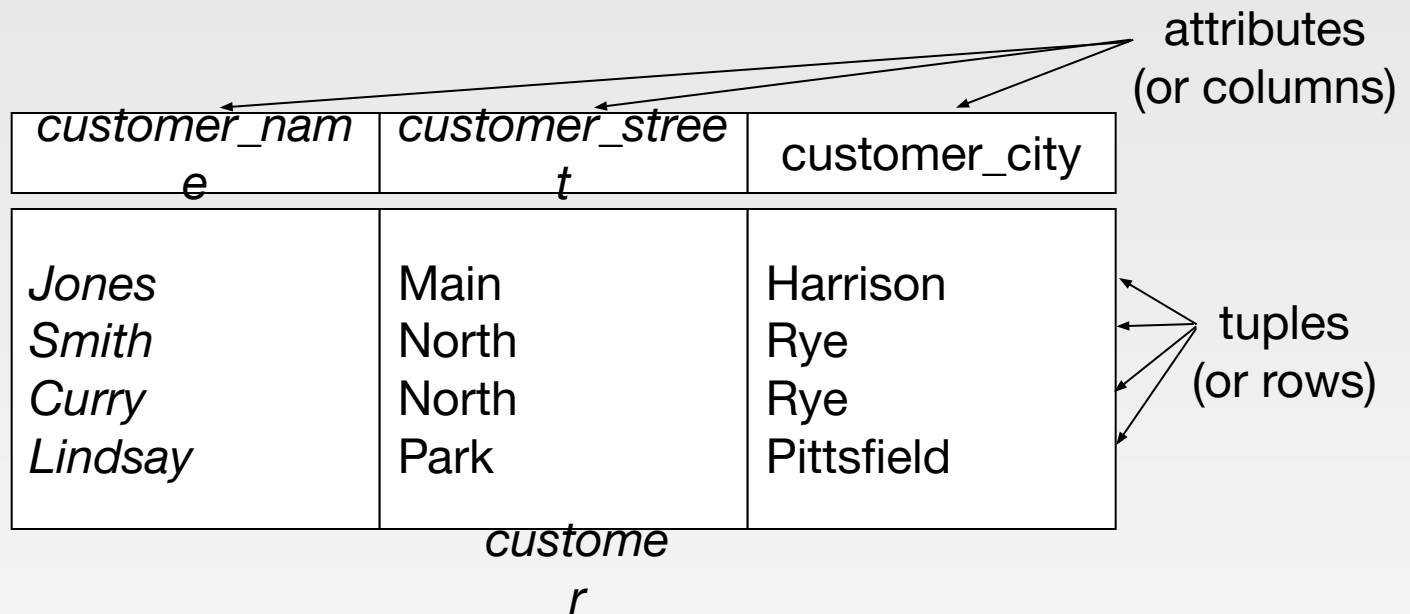
- Schema of a relation consists of
 - attribute definitions
 - name
 - type/domain
 - integrity constraints





Relation Instance

- The current values (*relation instance*) of a relation are specified by a table
- An element t of r is a *tuple*, represented by a *row* in a table
- Order of tuples is irrelevant (tuples may be stored in an arbitrary order)





Database

- A database consists of multiple relations
- Information about an enterprise is broken up into parts, with each relation storing one part of the information
- E.g.
 - account* : information about accounts
 - depositor* : which customer owns which account
 - customer* : information about customers





The *customer* Relation

<i>customer_name</i>	<i>customer_street</i>	<i>customer_city</i>
Adams	Spring	Pittsfield
Brooks	Senator	Brooklyn
Curry	North	Rye
Glenn	Sand Hill	Woodside
Green	Walnut	Stamford
Hayes	Main	Harrison
Johnson	Alma	Palo Alto
Jones	Main	Harrison
Lindsay	Park	Pittsfield
Smith	North	Rye
Turner	Putnam	Stamford
Williams	Nassau	Princeton





The *depositor* Relation

<i>customer_name</i>	<i>account_number</i>
Hayes	A-102
Johnson	A-101
Johnson	A-201
Jones	A-217
Lindsay	A-222
Smith	A-215
Turner	A-305





Why Split Information Across Relations?

- Storing all information as a single relation such as *bank(account_number, balance, customer_name, ..)* results in
 - repetition of information
 - 4 e.g., if two customers own an account (What gets repeated?)
 - the need for null values
 - 4 e.g., to represent a customer without an account
- Normalization theory (Chapter 7) deals with how to design relational schemas





Keys

- Let $K \subseteq R$
- K is a **superkey** of R if values for K are sufficient to identify a unique tuple of each possible relation $r(R)$
 - by “possible r ” we mean a relation r that could exist in the enterprise we are modeling.
 - Example: $\{customer_name, customer_street\}$ and $\{customer_name\}$
are both superkeys of $Customer$, if no two customers can possibly have the same name
- 4 In real life, an attribute such as $customer_id$ would be used instead of $customer_name$ to uniquely identify customers, but we omit it to keep our examples small, and instead assume customer names are unique.





Keys (Cont.)

- K is a **candidate key** if K is minimal
Example: $\{customer_name\}$ is a candidate key for *Customer*, since it is a superkey and no subset of it is a superkey.
- **Primary key:** a candidate key chosen as the principal means of identifying tuples within a relation
 - Should choose an attribute whose value never, or very rarely, changes.
 - E.g. email address is unique, but may change





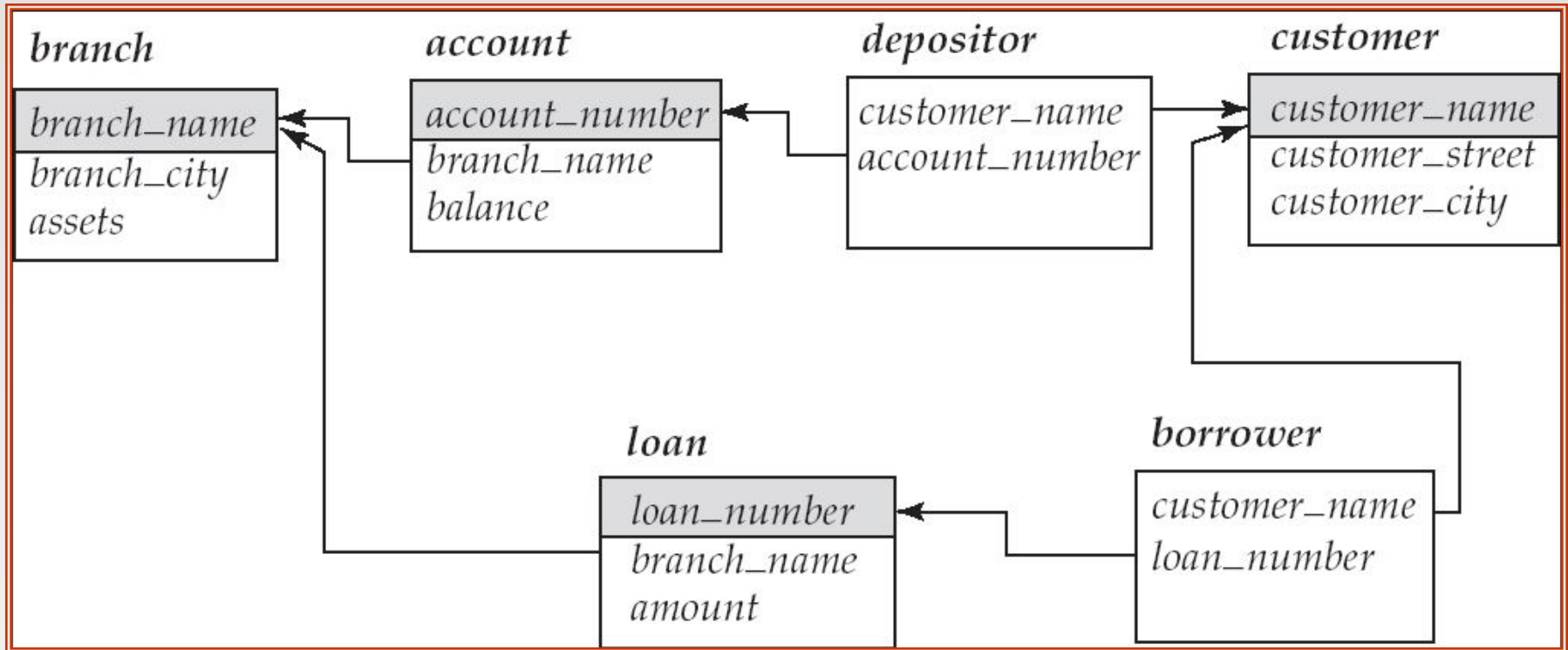
Foreign Keys

- A relation schema may have an attribute that corresponds to the primary key of another relation. The attribute is called a **foreign key**.
 - E.g. *customer_name* and *account_number* attributes of *depositor* are foreign keys to *customer* and *account* respectively.
 - Only values occurring in the primary key attribute of the **referenced relation** may occur in the foreign key attribute of the **referencing relation**.





Schema Diagram





Query Languages

- Language in which user requests information from the database.
- Categories of languages
 - Procedural
 - Non-procedural, or declarative
- “Pure” languages:
 - Relational algebra
 - Tuple relational calculus
 - Domain relational calculus
- Pure languages form underlying basis of query languages that people use.





Relational Algebra

- Procedural language
- Six basic operators
 - select: σ
 - project: π
 - union: \cup
 - set difference: $-$
 - Cartesian product: \times
 - rename: ρ
- The operators take one or two relations as inputs and produce a new relation as a result.





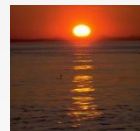
Select Operation – Example

- Relation r

A	B	C	D
a	a	1	7
a	β	5	7
β	β	12	3
β	β	23	10

- $\sigma_{A=B \wedge D > 5}(r)$

A	B	C	D
a	a	1	7
β	β	23	10





Project Operation – Example

- Relation r :

A	B	C
a	10	1
a	20	1
β	30	1
β	40	2

$\Pi_{A,C}(r)$

A	C
a	1
a	1
β	1
β	2

=

A	C
a	1
β	1
β	2





Union Operation – Example

- Relations r, s :

A	B
a	1
a	2
β	1

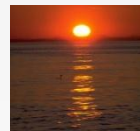
r

A	B
a	2
β	3

s

- $r \cup s$:

A	B
a	1
a	2
β	1
β	3





Set Difference Operation – Example

- Relations r , s :

A	B
a	1
a	2
β	1

r

A	B
a	2
β	3

s

- $r - s$:

A	B
a	1
β	1





Cartesian-Product Operation – Example

- Relations r , s :

A	B
a	1
β	2

r

C	D	E
a	10	a
β	10	a
β	20	b
γ	10	b

s

- $r \times s$:

A	B	C	D	E
a	1	a	10	a
a	1	β	10	a
a	1	β	20	b
a	1	γ	10	b
β	2	a	10	a
β	2	β	10	a
β	2	β	20	b
β	2	γ	10	b





Rename Operation

- Allows us to name, and therefore to refer to, the results of relational-algebra expressions.
- Allows us to refer to a relation by more than one name.
- Example:

$$\rho_X(E)$$

returns the expression E under the name X

- If a relational-algebra expression E has arity n , then

$$\rho_{X(A_1, A_2, \dots, A_n)}(E)$$

returns the result of expression E under the name X , and with the attributes renamed to A_1, A_2, \dots, A_n .





Composition of Operations

- Can build expressions using multiple operations
- Example: $\sigma_{A=C}(r \times s)$

- $r \times s$

<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
<i>a</i>	1	<i>a</i>	10	<i>a</i>
<i>a</i>	1	β	10	<i>a</i>
<i>a</i>	1	β	20	<i>b</i>
<i>a</i>	1	γ	10	<i>b</i>
β	2	<i>a</i>	10	<i>a</i>
β	2	β	10	<i>a</i>
β	2	β	20	<i>b</i>
β	2	γ	10	<i>b</i>

- $\sigma_{A=C}(r \times s)$

<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
<i>a</i>	1	<i>a</i>	10	<i>a</i>
β	2	β	10	<i>a</i>
β	2	β	20	<i>b</i>





Banking Example

branch (branch_name, branch_city, assets)

customer (customer_name, customer_street, customer_city)

account (account_number, branch_name, balance)

loan (loan_number, branch_name, amount)

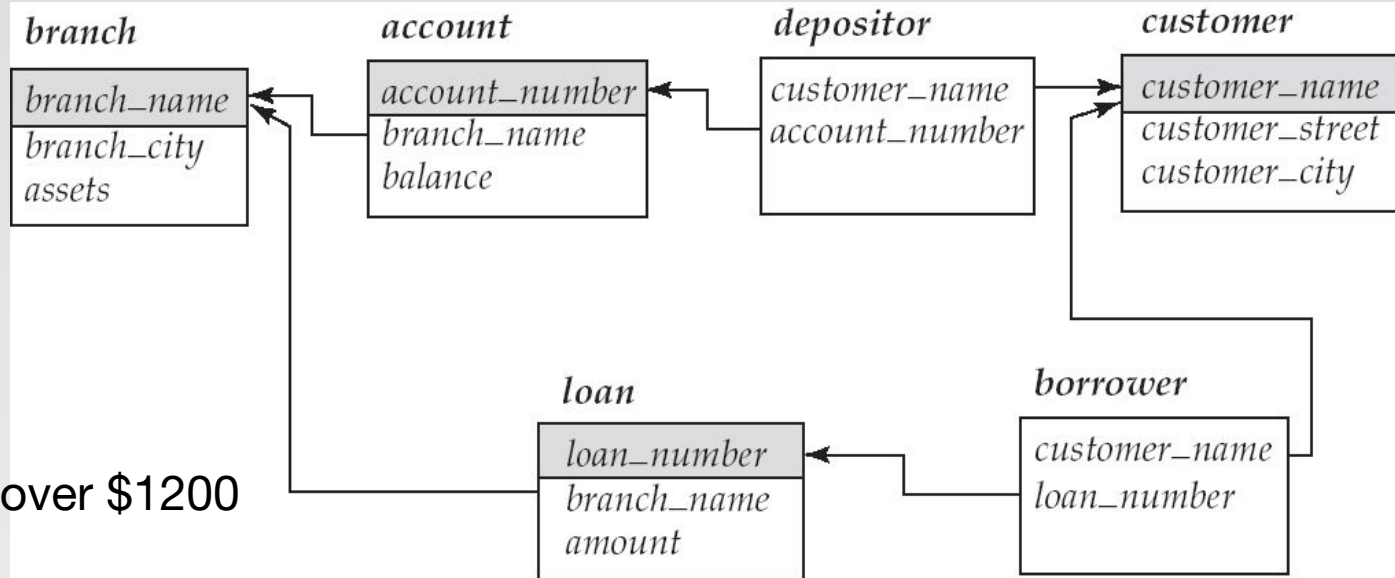
depositor (customer_name, account_number)

borrower (customer_name, loan_number)





Example Queries



- Find all loans of over \$1200

$$\sigma_{amount > 1200} (loan)$$

- Find the loan number for each loan of an amount greater than \$1200

$$\Pi_{loan_number} (\sigma_{amount > 1200} (loan))$$

- Find the names of all customers who have a loan, an account, or both, from the bank

$$\Pi_{customer_name} (borrower) \cup \Pi_{customer_name} (depositor)$$



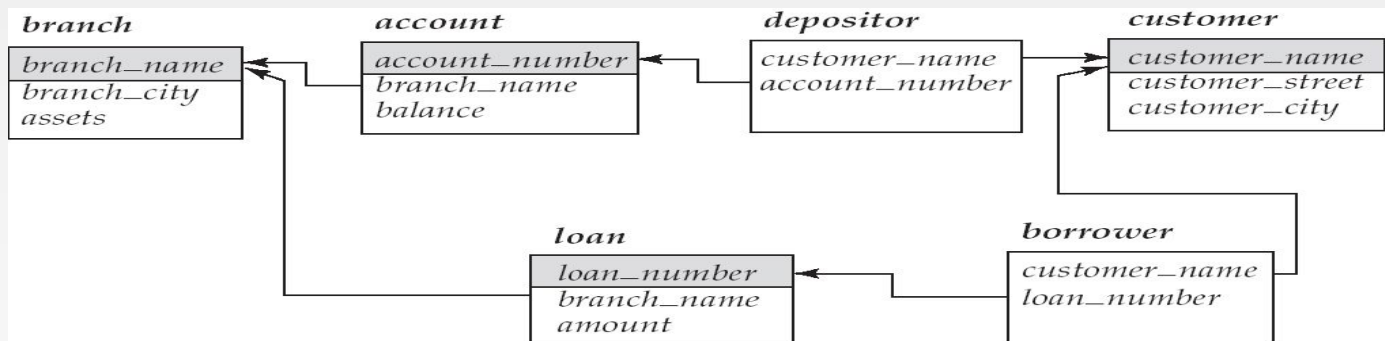


Example Queries

- Find the names of all customers who have a loan at the Perryridge branch.

$$\Pi_{customer_name} (\sigma_{branch_name="Perryridge"} (\sigma_{borrower.loan_number = loan.loan_number} (borrower \times loan)))$$

- Find the names of all customers who have a loan at the Perryridge branch but do not have an account at any branch of the bank.

$$\Pi_{customer_name} (\sigma_{branch_name = "Perryridge"}$$
$$(\sigma_{borrower.loan_number = loan.loan_number} (borrower \times loan))) -$$
$$\Pi_{customer_name} (depositor)$$


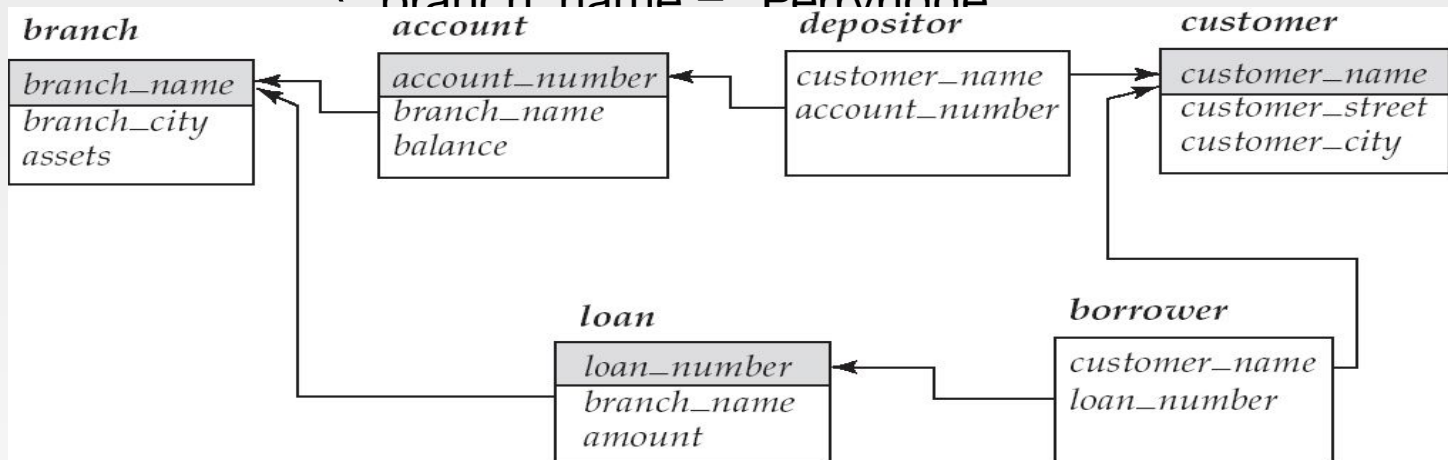


Example Queries

- Find the names of all customers who have a loan at the Perryridge branch.

- $\Pi_{\text{customer_name}} (\sigma_{\text{branch_name} = \text{"Perryridge"}} (\sigma_{\text{borrower.loan_number} = \text{loan.loan_number}} (\text{borrower} \times \text{loan})))$

- $\Pi_{\text{customer_name}} (\sigma_{\text{loan.loan_number} = \text{borrower.loan_number}} (\sigma_{\text{branch_name} = \text{"Perryridge"}} (\text{loan}) \times \text{borrower}))$





Additional Operations

- Additional Operations
 - Set intersection
 - Natural join
 - Aggregation
 - Outer Join
 - Division
- **All above, other than aggregation, can be expressed using basic operations we have seen earlier**





Set-Intersection Operation – Example

- Relation r, s :

A	B
α	1
α	2
β	1

r

A	B
α	2
β	3

s

- $r \cap s$

A	B
α	2





Natural Join Operation – Example

- Relations r , s :

A	B	C	D
a	1	a	a
β	2	γ	a
γ	4	β	b
a	1	γ	a
δ	2	β	b

r

B	D	E
1	a	a
3	a	β
1	a	γ
2	b	δ
3	b	ϵ

s

- $r \bowtie s$

A	B	C	D	E
a	1	a	a	a
a	1	a	a	γ
a	1	γ	a	a
a	1	γ	a	γ
δ	2	β	b	δ





Natural-Join Operation

- Notation: $r \bowtie s$
- Let r and s be relations on schemas R and S respectively. Then, $r \bowtie s$ is a relation on schema $R \cup S$ obtained as follows:
 - Consider each pair of tuples t_r from r and t_s from s .
 - If t_r and t_s have the same value on each of the attributes in $R \cap S$, add a tuple t to the result, where

4 t has the same value as t_r on r

4 t has the same value as t_s on s

- Example:

$R = (A, B, C, D)$

$S = (E, B, D)$

- Result schema = (A, B, C, D, E)

- $r \bowtie s$ is defined as:

$$\pi_{r.A, r.B, r.C, r.D, s.E} (\sigma_{r.B = s.B \wedge r.D = s.D} (r \times s))$$





Bank Example Queries

- Find the largest account balance
 - Strategy:
 - 4 Find those balances that are *not* the largest
 - Rename *account* relation as *d* so that we can compare each account balance with all others
 - 4 Use set difference to find those account balances that were *not* found in the earlier step.
 - The query is:

$$\Pi_{balance}(account) - \Pi_{account.balance}(\sigma_{account.balance < d.balance}(account \times \rho_d(account)))$$

<i>account</i>
<i>account_number</i>
<i>branch_name</i>
<i>balance</i>





Aggregate Functions and Operations

- **Aggregation function** takes a collection of values and returns a single value as a result.

avg: average value

min: minimum value

max: maximum value

sum: sum of values

count: number of values

- **Aggregate operation** in relational algebra

$$G_1, G_2, \dots, G_n \quad \mathcal{G} \quad F_1(A_1), F_2(A_2), \dots, F_n(A_n) (E)$$

E is any relational-algebra expression

- G_1, G_2, \dots, G_n is a list of attributes on which to group (can be empty)
- Each F_i is an aggregate function
- Each A_i is an attribute name





Aggregate Operation – Example

- Relation r :

A	B	C
a	a	7
a	β	7
β	β	3
β	β	10

- $g_{\text{sum}(c)}(r)$

sum(c)
27

- Question: Which aggregate operations cannot be expressed using basic relational operations?





Aggregate Operation – Example

- Relation *account* grouped by *branch-name*:

<i>branch_name</i>	<i>account_number</i>	<i>balance</i>
Perryridge	A-102	400
Perryridge	A-201	900
Brighton	A-217	750
Brighton	A-215	750
Redwood	A-222	700

branch_name g **sum**(*balance*) (*account*)

<i>branch_name</i>	sum (<i>balance</i>)
Perryridge	1300
Brighton	1500
Redwood	700





Aggregate Functions (Cont.)

- Result of aggregation does not have a name
 - Can use rename operation to give it a name
 - For convenience, we permit renaming as part of aggregate operation

branch_name **g** **sum**(*balance*) **as** *sum_balance* (*account*)





Outer Join

- An extension of the join operation that avoids loss of information.
 - Computes the join and then adds tuples from one relation that does not match tuples in the other relation to the result of the join.
 - Uses *null* values:
 - *null* signifies that the value is unknown or does not exist
 - All comparisons involving *null* are (roughly speaking) **false** by definition.
- 4 We shall study precise meaning of comparisons with nulls later





Outer Join – Example

- Relation *loan*

<i>loan_number</i>	<i>branch_name</i>	<i>amount</i>
L-170	Downtown	3000
L-230	Redwood	4000
L-260	Perryridge	1700

- Relation *borrower*

<i>customer_name</i>	<i>loan_number</i>
Jones	L-170
Smith	L-230
Hayes	L-155





Outer Join – Example

- Join

loan ⋈ *borrower*

<i>loan_number</i>	<i>branch_name</i>	<i>amount</i>	<i>customer_name</i>
L-170	Downtown	3000	Jones
L-230	Redwood	4000	Smith

- Left Outer Join

loan ⋈_L *borrower*

<i>loan_number</i>	<i>branch_name</i>	<i>amount</i>	<i>customer_name</i>
L-170	Downtown	3000	Jones
L-230	Redwood	4000	Smith
L-260	Perryridge	1700	<i>null</i>





Outer Join – Example

- Right Outer Join

$loan \bowtie_{\square} borrower$

<i>loan_number</i>	<i>branch_name</i>	<i>amount</i>	<i>customer_name</i>
L-170	Downtown	3000	Jones
L-230	Redwood	4000	Smith
L-155	<i>null</i>	<i>null</i>	Hayes

- Full Outer Join

$loan \bowtie_{\square\square} borrower$

<i>loan_number</i>	<i>branch_name</i>	<i>amount</i>	<i>customer_name</i>
L-170	Downtown	3000	Jones
L-230	Redwood	4000	Smith
L-260	Perryridge	1700	<i>null</i>
L-155	<i>null</i>	<i>null</i>	Hayes

- **Question:** can outer joins be expressed using basic relational algebra operations





Null Values

- It is possible for tuples to have a null value, denoted by *null*, for some of their attributes
- *null* signifies an unknown value or that a value does not exist.
- The result of any arithmetic expression involving *null* is *null*.
- Aggregate functions simply ignore null values (as in SQL)
- For duplicate elimination and grouping, null is treated like any other value, and two nulls are assumed to be the same (as in SQL)





Null Values

- Comparisons with null values return the special truth value: *unknown*
 - If *false* was used instead of *unknown*, then $\text{not } (A < 5)$ would not be equivalent to $A \geq 5$
- Three-valued logic using the truth value *unknown*:
 - OR: $(\text{unknown or true}) = \text{true}$,
 $(\text{unknown or false}) = \text{unknown}$
 $(\text{unknown or unknown}) = \text{unknown}$
 - AND: $(\text{true and unknown}) = \text{unknown}$,
 $(\text{false and unknown}) = \text{false}$,
 $(\text{unknown and unknown}) = \text{unknown}$
 - NOT: $(\text{not unknown}) = \text{unknown}$
 - In SQL “*P* is **unknown**” evaluates to true if predicate *P* evaluates to *unknown*
- Result of select predicate is treated as *false* if it evaluates to *unknown*





Bank Example Queries

- Find the names of all customers who have a loan and an account at bank.

$$\Pi_{customer_name} (borrower) \cap \Pi_{customer_name} (depositor)$$

- Find the name of all customers who have a loan at the bank and the loan amount

$$\Pi_{customer_name, loan_number, amount} (borrower \quad loan)$$





Bank Example Queries

- Find all customers who have an account from at least the “Downtown” and the Uptown” branches.

- Query 1

$$\begin{aligned} & \Pi_{customer_name} (\sigma_{branch_name = \text{“Downtown”}} (depositor \bowtie account)) \cap \\ & \Pi_{customer_name} (\sigma_{branch_name = \text{“Uptown”}} (depositor \bowtie account)) \end{aligned}$$

- Query 2

$$\begin{aligned} & \Pi_{customer_name, branch_name} (depositor \bowtie account) \\ & \div \rho_{temp(branch_name)} (\{(\text{“Downtown”}), (\text{“Uptown”})\}) \end{aligned}$$

Note that Query 2 uses a constant relation.

