

However if $a \rightarrow b$

then we can do:

$$|\gamma^a R| = |\gamma^{a,b} R| = |\pi_{a,b}(R \bowtie \gamma^a R)|$$

i.e. all queries return the same number of tuples

And:

$$\gamma^{a,b} R = \pi_{a,b}(R \bowtie \gamma^a R)$$

But more frequently you will need:

Assume $R(\underline{a}, b, c)$, $S(\underline{a}, \underline{d})$

$$\gamma_{\text{count}(d)}^{a,b}(R \bowtie S) =$$

$$\pi_{a,b,\text{count}(d)}(R \bowtie \gamma_{\text{count}(d)}^a S)$$

But only because $a \rightarrow b$!!

Ex: Find id and name of student and the number of courses she/he is registered in.

Aggregation.

Frequently it is necessary to summarize a set of tuples into only one.

Ex:

- How many tuples satisfy this condition?
- What is the average of this attribute?

γ group-by operator

In its simplest form $\gamma_{\langle \text{seq. aggr. exps} \rangle} R$ computes a sequence of aggregation expressions on a relation R .

Aggregation functions.

Given a set of tuples or attributes, compute a single value.

$\text{count}(x)$ Count number of tuples in set.

$\text{count}(att)$ Count number of tuples with attribute not NULL

$\text{sum}(att)$ Sums the value of attr.

$$\text{avg}(att) = \frac{\text{sum}(att)}{\text{count}(att)}$$

$\text{max}(att)$, $\text{min}(att)$.

Example

$R(a, b, c)$

a	b	c
7	a	1 ← NULL
2	x	-1
5	y	5

$\gamma_{\text{count}(*), \text{sum}(a), \text{count}(c)} R$

"count(*)"	"sum(a)"	"count(c)"
3	12	2

$\gamma_{\frac{\text{sum}(a)}{\text{count}(a)}} \rightarrow$ a avg R rename attribute.

a avg
4

Grouping

Sometimes we need to make summaries of different subsets of tuples.

Ex: How many courses is each student taking?

- What is the average price of each part?

Be careful:

$\gamma_b^a R, \pi_b \gamma_b^a R, \pi_b \gamma^a R$
are all illegal

Remember: the schema of γ does not contain attributes of R not listed in the grouping attributes $\gamma_{\langle \text{list att} \rangle}$

my SQL allows this:

$\gamma_b^a R$

Value of b is non deterministic.
Chosen at random from one tuple in grouping subset.

We don't like NON DETERMINISM
Unless you know what you're doing.

Instead use:

$\pi_{a,b} [R \bowtie \gamma^a R]$

We can combine operations:

$\pi_{count(c)}$ $\sigma_{count(c) > 1}$ $\gamma_{\overline{a}}$ $\sigma_{b > 3}$ R
 SELECT count(c) FROM
 (SELECT a, count(c)
 FROM R
 WHERE b > 3
 GROUP BY a) AS X
 WHERE count(c) > 1;

Any subquery requires a name
 ↑ selection on result of aggregation

$\Pi \Sigma_q$ of $\gamma \Sigma_p$ is so common that SQL has syntactic sugar for it:

```
SELECT count(c)
FROM R
WHERE b > 3
GROUP BY a
HAVING count(c) > 1.
```

Ex: Find the student id of students who are taking 3 or more courses.

$$\gamma_{\langle \text{att list} \rangle_R}$$

Creates one tuple for each different value of the list of attributes.

Ex. $R(a, b, c)$

a	b	c
3	9	1
2	5	4
3	9	5
2	1	8

$$\gamma^{a,b} \in \mathbb{R}$$

a	b
3	9
2	5
2	1

 $\gamma^a R$
$$\begin{array}{r} 9 \\ \hline 3 \\ 2 \end{array}$$
 $\gamma^c R$
$$\begin{array}{r} c \\ \hline 1 \\ 4 \\ 5 \\ 8 \end{array}$$

Warning : This is my notation.

In fact, our textbook does not even include γ in its RA chapter.

SQL

$\gamma_{\text{count}(*), \text{count}(a)} R$

Remember, it returns only one tuple

SELECT count(*), count(a)
FROM R;

This is not a

$\Pi_{\text{count}(*), \text{count}(a)}$ But it can be interpreted as

$\Pi_{\text{count}(*), \text{count}(a)} \gamma_{\text{count}(*), \text{count}(a)} R$

Redundant in this case.

$\gamma^{a,b} R$

SELECT a, b FROM R
GROUP BY a, b

Yes, redundant but necessary

REMOVES DUPLICATES !!

Equivalent to:

SELECT DISTINCT a, b FROM R

$\Rightarrow \Pi_{a,b} R = \gamma^{a,b} R$

only in RA (relations are sets) 4

Combining both:

$\gamma_{\langle \text{list of attr} \rangle \langle \text{list of agg expr} \rangle} R$

Computes the expressions on each subset of different values of attributes.

Ex:

	a	b	c
$R(a,b,c)$	3	9	1
	2	5	4
	3	9	5
	2	1	8

$\gamma^a_{\text{avg}(c), \text{count}(*)} R$

a	"avg(c)"	"count(*)"
3	5	2
2	6	2

SELECT count(c), count(*)
FROM R
GROUP BY a