



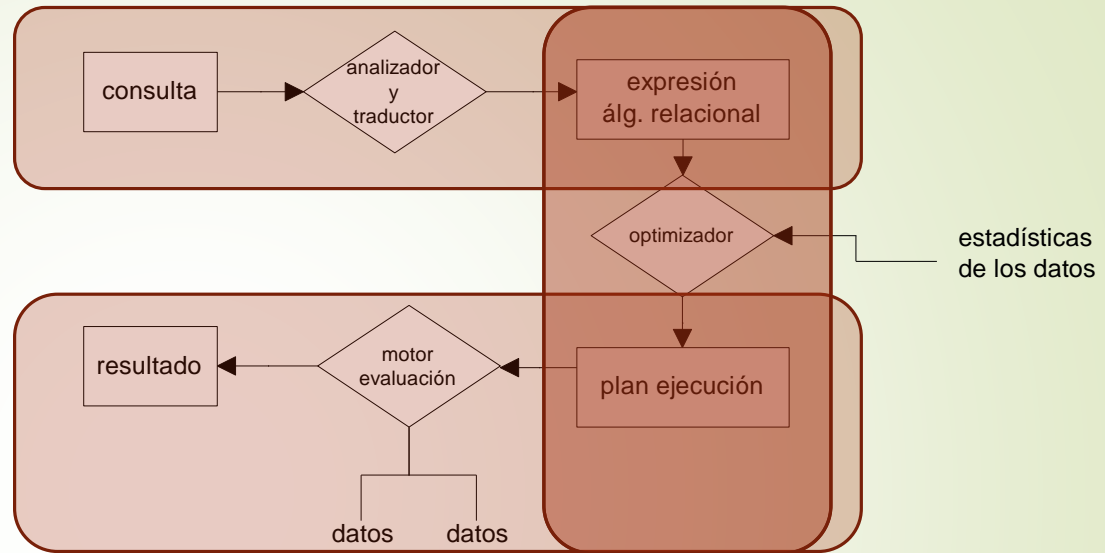
Tema.- Procesamiento y Optimización de consultas

Fundamentos de Bases de Datos (5º edición).

A. Silberschatz. McGraw-Hill [caps. 13-14]

Procesamiento de consultas

➡ Pasos:



Análisis y traducción:

- Comprobación de sintaxis, verificación de nombres
- Construcción del árbol de consulta
- Transformación a una expresión en álgebra relacional extendido

Optimización:

- Indicar cómo evaluar la expresión en álgebra relacional (algoritmo(s) a utilizar, índice(s) a aplicar,...) ⇒ **primitivas de evaluación**
- Determinar el **plan de ejecución/evaluación** (= secuencia de primitivas de evaluación) que minimice el coste de ejecución de la consulta

Evaluación:

- El motor de consultas ejecuta un plan de evaluación

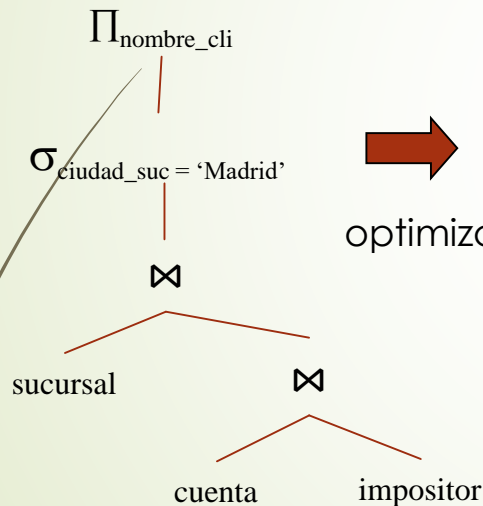
Optimización de consultas

Nombre de clientes con cuenta en una sucursal de Madrid

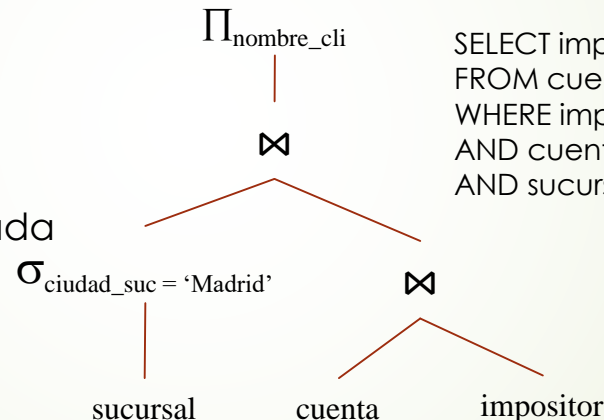
$\Pi_{\text{nombre_cli}} (\sigma_{\text{ciudad_suc} = \text{'Madrid'}} (\text{sucursal} \bowtie (\text{cuenta} \bowtie \text{impositor})))$

↓ optimizada

$\Pi_{\text{nombre_cli}} (\sigma_{\text{ciudad_suc} = \text{'Madrid'}} (\text{sucursal}) \bowtie (\text{cuenta} \bowtie \text{impositor}))$

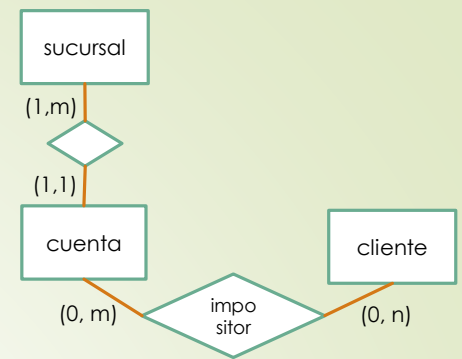


optimizada



```

SELECT impositor.nombre_cli
FROM cuenta, impositor, sucursal
WHERE impositor.num_cta = cuenta.num_cta
AND cuenta.nombre_suc = sucursal.nombre_suc
AND sucursal.ciudad_suc = "Madrid"
  
```



sucursal (nombre_suc, ciudad_suc, activos)
cuenta (num_cta, nombre_suc, saldo)
impositor (nombre_cli, num_cta)
cliente (nombre_cli, edad, ciudad_cli, sueldo)

El optimizador diseña un plan de evaluación de consultas “óptimo” generando planes alternativos

1. **Heurística:** genera expresiones equivalentes mediante **reglas de equivalencia**
2. **Coste:** estima el coste de cada plan (utilizando inf. estadística de tablas, índices, etc.)

Transformación de expresiones relacionales

- Las consultas se pueden expresar de varias maneras, con costes de evaluación diferentes
 - En vez de tomar la expresión relacional original se consideran expresiones alternativas equivalentes.
- 2 expresiones en álgebra relacional son **equivalentes** si generan el mismo conjunto de tuplas
- **Reglas de equivalencia** => 2 expresiones son equivalentes si se puede sustituir la primera expresión por la segunda, o viceversa

Reglas de equivalencia

- **CASCADA DE σ :** las operaciones de selección conjuntivas pueden dividirse en una secuencia de selecciones individuales $\Rightarrow \sigma_{\theta_1 \wedge \theta_2}(E) = \sigma_{\theta_1}(\sigma_{\theta_2}(E))$

ej: $\sigma_{\text{saldo} > 10 \wedge \text{nombre_suc} = \text{"Madrid"}}(\text{cuenta}) = \sigma_{\text{saldo} > 10}(\sigma_{\text{nombre_suc} = \text{"Madrid"}}(\text{cuenta}))$

- **CONMUTATIVIDAD DE σ :** las operaciones de selección son conmutativas $\Rightarrow \sigma_{\theta_1}(\sigma_{\theta_2}(E)) = \sigma_{\theta_2}(\sigma_{\theta_1}(E))$

ej: $\sigma_{\text{saldo} > 10}(\sigma_{\text{nombre_suc} = \text{"Madrid"}}(\text{cuenta})) = \sigma_{\text{nombre_suc} = \text{"Madrid"}}(\sigma_{\text{saldo} > 10}(\text{cuenta}))$

- **CASCADA DE Π :** solo es necesaria la última operación de proyección, las demás pueden omitirse $\Rightarrow \Pi_{L_1}(\Pi_{L_2}(\dots \Pi_{L_i}(E))\dots) = \Pi_{L_1}(E)$

ej: $\Pi_{\text{num_cta}}(\Pi_{\text{num_cta}, \text{nombre_suc}}(\Pi_{\text{num_cta}, \text{nombre_suc}, \text{saldo}}(\text{cuenta}))) = \Pi_{\text{num_cta}}(\text{cuenta})$

Reglas de equivalencia

- Las selecciones pueden combinarse con productos cartesianos y con zeta-joins:

$$\sigma_{\theta}(E_1 \times E_2) = E_1 \bowtie_{\theta} E_2 \quad \sigma_{\theta_1}(E_1 \bowtie_{\theta_2} E_2) = E_1 \bowtie_{\theta_1 \wedge \theta_2} E_2$$

ej: $\sigma_{c.nombre_suc = s.nombre_suc}(\text{cuenta} \times \text{sucursal}) = \text{cuenta} \bowtie_{c.nombre_suc = s.nombre_suc} \text{sucursal}$

ej: $\sigma_{c.saldo > s.activos}(\text{cuenta} \bowtie_{c.nombre_suc = s.nombre_suc} \text{sucursal}) =$
 $= \text{cuenta} \bowtie_{c.nombre_suc = s.nombre_suc \wedge c.saldo > s.activos} \text{sucursal}$

- CONMUTATIVIDAD DE \bowtie :** los zeta-join son conmutativos $\Rightarrow E_1 \bowtie_{\theta} E_2 = E_2 \bowtie_{\theta} E_1$

ej: $\text{cuenta} \bowtie_{c.nombre_suc = s.nombre_suc} \text{sucursal} = \text{sucursal} \bowtie_{c.nombre_suc = s.nombre_suc} \text{cuenta}$

- ASOCIATIVIDAD DE \bowtie :** el join natural es asociativo $\Rightarrow (E_1 \bowtie E_2) \bowtie E_3 = E_1 \bowtie (E_2 \bowtie E_3)$

ej: $(\text{cuenta} \bowtie \text{sucursal}) \bowtie \text{impositor} = \text{cuenta} \bowtie (\text{sucursal} \bowtie \text{impositor})$

Los zeta-join son asociativos en el siguiente sentido:

$$(E_1 \bowtie_{\theta_1} E_2) \bowtie_{\theta_2 \wedge \theta_3} E_3 = E_1 \bowtie_{\theta_1 \wedge \theta_3} (E_2 \bowtie_{\theta_2} E_3) \text{ donde } \theta_2 \text{ implica solo atributos de } E_2 \text{ y } E_3$$

ej: $(\text{cuenta} \bowtie_{c.num_cta=i.num_cta} \text{impositor}) \bowtie_{i.nombre_cli=cl.nombre_cli \wedge c.saldo > cl.sueldo} \text{cliente} =$
 $= \text{cuenta} \bowtie_{c.num_cta=i.num_cta \wedge c.saldo > cl.sueldo} (\text{impositor} \bowtie_{i.nombre_cli=cl.nombre_cli} \text{cliente})$

Reglas de equivalencia

- **Desplazar σ hacia hojas en \bowtie :** La operación selección se distribuye en el zeta-join bajo 2 condiciones:

1. Si la selección θ solo implica atributos de una de las expresiones:

$$\sigma_{\theta_1}(E1 \bowtie_{\theta} E2) = \sigma_{\theta_1}(E1) \bowtie_{\theta} E2 \text{ donde } \theta_1 \text{ solo implica atributos de } E1$$

ej: $\sigma_{\text{saldo}>100}(\text{cuenta} \bowtie \text{sucursal}) = \sigma_{\text{saldo}>100}(\text{cuenta}) \bowtie \text{sucursal}$

2. Si la selección θ_1 solo implica atributos de $E1$ y θ_2 solo implica atributos de $E2$:

$$\sigma_{\theta_1 \wedge \theta_2}(E1 \bowtie_{\theta} E2) = \sigma_{\theta_1}(E1) \bowtie_{\theta} \sigma_{\theta_2}(E2)$$

ej: $\sigma_{(\text{saldo}>100) \wedge (\text{activos}=200)}(\text{cuenta} \bowtie \text{sucursal}) = \sigma_{\text{saldo}>100}(\text{cuenta}) \bowtie \sigma_{\text{activos}=200}(\text{sucursal})$

Reglas de equivalencia

- **Desplazar Π hacia las hojas en \bowtie :** La operación proyección se distribuye en el zeta-join bajo 2 condiciones:

1. $\Pi_{L1 \cup L2} (E1 \bowtie_{\theta} E2) = \Pi_{L1} (E1) \bowtie_{\theta} \Pi_{L2} (E2)$ siendo $L1$ atributos de $E1$, $L2$ atributos de $E2$ y donde θ implica solo atributos de $L1 \cup L2$

ej: $\Pi_{\text{nombre_suc, activos}} (\text{sucursal} \bowtie \text{cuenta}) = \Pi_{\text{nombre_suc, activos}} (\text{sucursal}) \bowtie \Pi_{\text{nombre_suc}} (\text{cuenta})$

2. $\Pi_{L1 \cup L2} (E1 \bowtie_{\theta} E2) = \Pi_{L1 \cup L2} (\Pi_{L1 \cup L3} (E1) \bowtie_{\theta} \Pi_{L2 \cup L4} (E2))$

donde $L1$ son atributos de $E1$

$L2$ son atributos de $E2$

$L3$ son atributos de $E1$ que están en θ pero no en $L1$

$L4$ son atributos de $E2$ que están en θ pero no en $L2$

ej: $\Pi_{\text{activos, saldo}} (\text{sucursal} \bowtie \text{cuenta}) =$

$\Pi_{\text{activos, saldo}} (\Pi_{\text{nombre_suc, activos}} (\text{sucursal}) \bowtie \Pi_{\text{nombre_suc, saldo}} (\text{cuenta}))$

Reglas de equivalencia

- Un conjunto de reglas es MÍNIMO si no se puede obtener ninguna regla a partir de la unión de otras
- Los optimizadores utilizan conjuntos mínimos de reglas de equivalencia
- Los optimizadores generan sistemáticamente expresiones equivalentes a la consulta dada
- Una buena ordenación de los joins reduce el tamaño de los resultados

Ej:

$$\sigma_{\text{ciudad_suc} = \text{'Madrid'}} (\text{sucursal}) \bowtie (\text{cuenta} \bowtie \text{impositor})$$

Como es probable que:

- 1) $(\text{cuenta} \bowtie \text{impositor})$ de lugar a una relación muy grande (tantas tuplas como impositores existan)
- 2) el número de cuentas de las sucursales de Madrid es pequeño

sería mejor aplicar la regla de asociatividad del join, dando lugar a la expresión:

$$(\sigma_{\text{ciudad_suc} = \text{'Madrid'}} (\text{sucursal}) \bowtie \text{cuenta}) \bowtie \text{impositor}$$

Tipos de optimización: Heurística

- Mediante la aplicación de reglas de equivalencia, reordena los componentes del árbol de consultas inicial para intentar reducir el coste de la optimización
- Pasos a seguir:

1. Realizar las operaciones de selección tan pronto como sea posible

$$\text{CASCADA DE } \sigma \quad \Rightarrow \sigma_{\theta_1 \wedge \theta_2}(E) = \sigma_{\theta_1}(\sigma_{\theta_2}(E))$$

$$\text{CONMUTATIVIDAD DE } \sigma \quad \Rightarrow \sigma_{\theta_1}(\sigma_{\theta_2}(E)) = \sigma_{\theta_2}(\sigma_{\theta_1}(E))$$

$$\sigma_{\theta_0}(E1 \bowtie_{\theta} E2) = \sigma_{\theta_0}(E1) \bowtie_{\theta} E2 \text{ donde } \theta_0 \text{ implica solo atributos de } E1$$

$$\sigma_{\theta_0}(E1 \times E2) = \sigma_{\theta_0}(E1) \times E2 \text{ donde } \theta_0 \text{ implica solo atributos de } E1$$

$$\sigma_{\theta_1 \wedge \theta_2}(E1 \bowtie_{\theta} E2) = \sigma_{\theta_1}(E1) \bowtie_{\theta} \sigma_{\theta_2}(E2) \text{ donde } \theta_1 \text{ implica solo atributos de } E1 \text{ y } \theta_2 \text{ atributos de } E2$$

$$\sigma_{\theta_1 \wedge \theta_2}(E1 \times E2) = \sigma_{\theta_1}(E1) \times \sigma_{\theta_2}(E2) \text{ donde } \theta_1 \text{ implica solo atributos de } E1 \text{ y } \theta_2 \text{ atributos de } E2$$

Tipos de optimización: Heurística

2. Sustituir el producto cartesiano seguido de σ por \bowtie

$$\sigma_{\theta}(E_1 \times E_2) = E_1 \bowtie_{\theta} E_2$$

$$\sigma_{\theta_1}(E_1 \bowtie_{\theta_2} E_2) = E_1 \bowtie_{\theta_1 \wedge \theta_2} E_2$$

3. Determinar las operaciones σ y \bowtie que producen menos tuplas, y ejecutarlas cuanto antes

$$\text{CONMUTATIVIDAD DE } \sigma \Rightarrow \sigma_{\theta_1}(\sigma_{\theta_2}(E)) = \sigma_{\theta_2}(\sigma_{\theta_1}(E))$$

$$\text{ASOCIATIVIDAD DE } \bowtie \Rightarrow (E_1 \bowtie E_2) \bowtie E_3 = E_1 \bowtie (E_2 \bowtie E_3)$$

$$(E_1 \bowtie_{\theta_1} E_2) \bowtie_{\theta_2 \wedge \theta_3} E_3 = E_1 \bowtie_{\theta_1 \wedge \theta_3} (E_2 \bowtie_{\theta_2} E_3)$$

donde θ_2 implica solo atributos de E_2 y E_3

4. Realizar las proyecciones tan pronto como sea posible

$$\text{CASCADA DE } \Pi \Rightarrow \Pi_{L_1}(\Pi_{L_2}(\dots \Pi_{L_i}(E))\dots) = \Pi_{L_1}(E)$$

$$\Pi_{L_1 \cup L_2}(E_1 \bowtie_{\theta} E_2) = \Pi_{L_1}(E_1) \bowtie_{\theta} \Pi_{L_2}(E_2)$$

$$\Pi_{L_1 \cup L_2}(E_1 \bowtie_{\theta} E_2) = \Pi_{L_1 \cup L_2}(\Pi_{L_1 \cup L_3}(E_1) \bowtie_{\theta} \Pi_{L_2 \cup L_4}(E_2))$$

Ejemplo 1 uso reglas de equivalencia

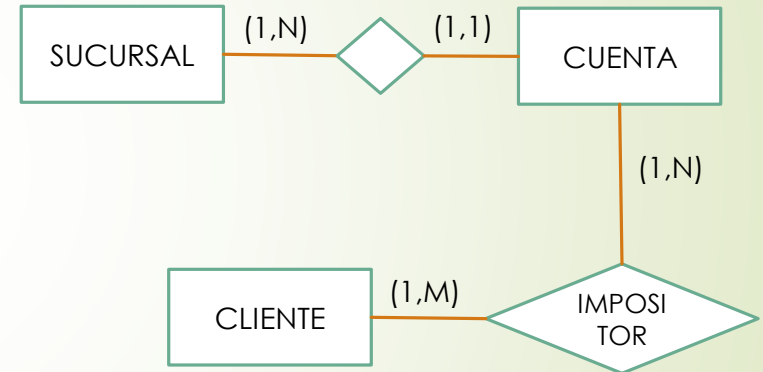
- Representar mediante árboles las expresiones del álgebra relacional de la consulta siguiente y transformarla a una forma más eficiente. Enunciar las reglas de equivalencia utilizadas en cada uno de los pasos del proceso

SUCURSAL (nombre_suc, ciudad_suc, activos)

CUENTA (num_cta, nombre_suc, saldo)

IMPOSITOR (nom_cli, num_cta)

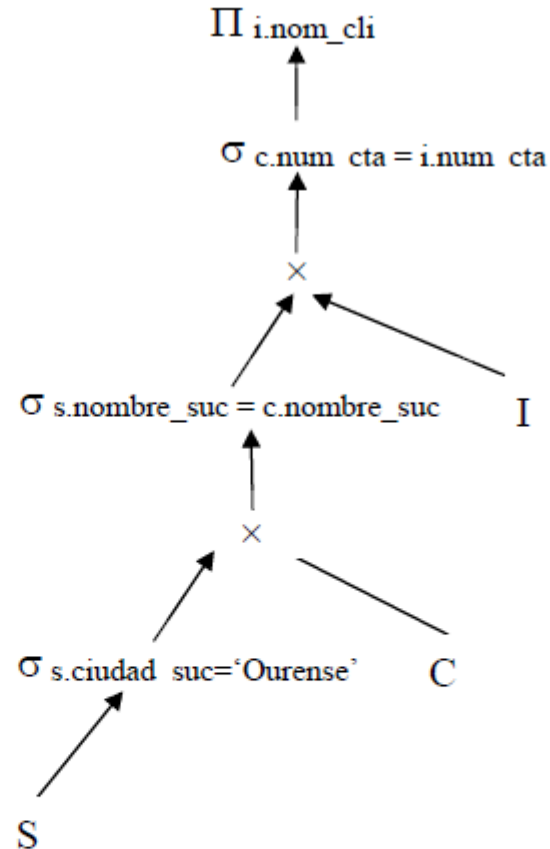
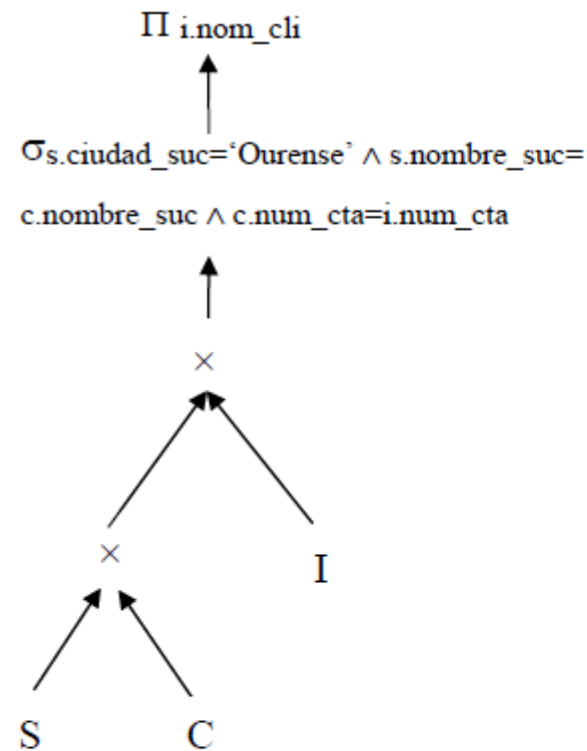
CLIENTE (nom_cli, direccion, activos)



Nombre de clientes con cuenta en alguna sucursal de Ourense

```
SELECT      nom_cli
FROM        SUCURSAL s, CUENTA c, IMPOSITOR i
WHERE       s.ciudad_suc = 'Ourense' AND
            s.nombre_suc = c.nombre_suc AND
            c.num_cta = i.num_cta;
```

Ejemplo 1 uso reglas de equivalencia



1. DESPLAZAR σ HACIA HOJAS EN X:

CASCADA DE σ : $\sigma_{\theta_1 \wedge \theta_2}(E) = \sigma_{\theta_1}(\sigma_{\theta_2}(E))$

CONMUTATIVIDAD DE σ : $\sigma_{\theta_1}(\sigma_{\theta_2}(E)) = \sigma_{\theta_2}(\sigma_{\theta_1}(E))$

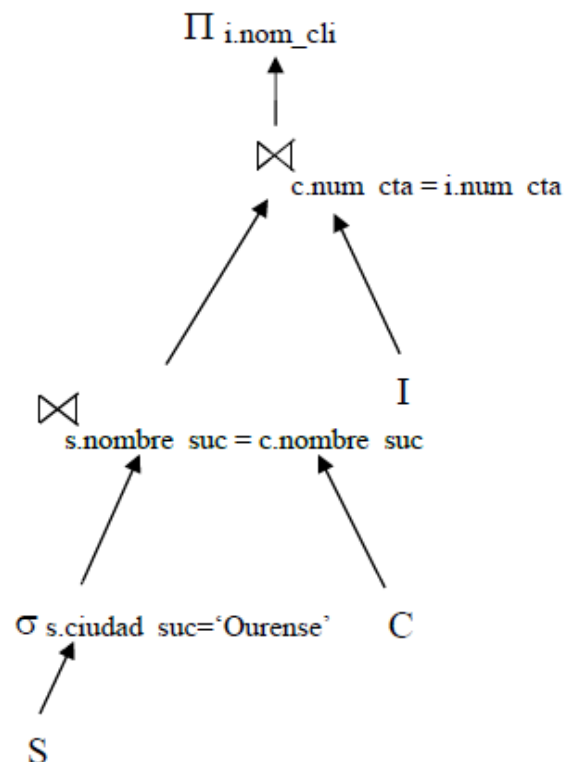
$\sigma_{\theta_0}(E_1 X E_2) = \sigma_{\theta_0}(E_1) X E_2$, θ_0 implica solo atrib. de E1

En este caso, no se utiliza $\sigma_{\theta_1 \wedge \theta_2}(E_1 X E_2) = \sigma_{\theta_1}(E_1) X \sigma_{\theta_2}(E_2)$ donde θ_1 implica solo atributos de E1 y θ_2 atributos de E2

0. $\Pi_{i.nom_cli} (\sigma_{s.ciudad_suc='Ourense' \wedge s.nombre_suc=c.nombre_suc \wedge c.num_cta=i.num_cta} (sucursal X cuenta X impositor))$

1. $\Pi_{i.nom_cli} (\sigma_{c.num_cta=i.num_cta} (\sigma_{s.nombre_suc=c.nombre_suc} (\sigma_{ciudad_suc='Ourense'} (sucursal) X cuenta) X impositor))$

Ejemplo 1 uso reglas de equivalencia



2. Sustituir el producto cartesiano (X) seguido de σ por join:

$$\sigma_{\theta}(E1 \times E2) = E1 \bowtie_{\theta} E2$$

3. Ejecutar antes las selecciones y los joins que producen menos tuplas

$$\text{CONMUTATIVIDAD DE } \sigma \Rightarrow \sigma_{\theta_1}(\sigma_{\theta_2}(E)) = \sigma_{\theta_2}(\sigma_{\theta_1}(E))$$

$$\text{ASOCIATIVIDAD DE } \bowtie \Rightarrow$$

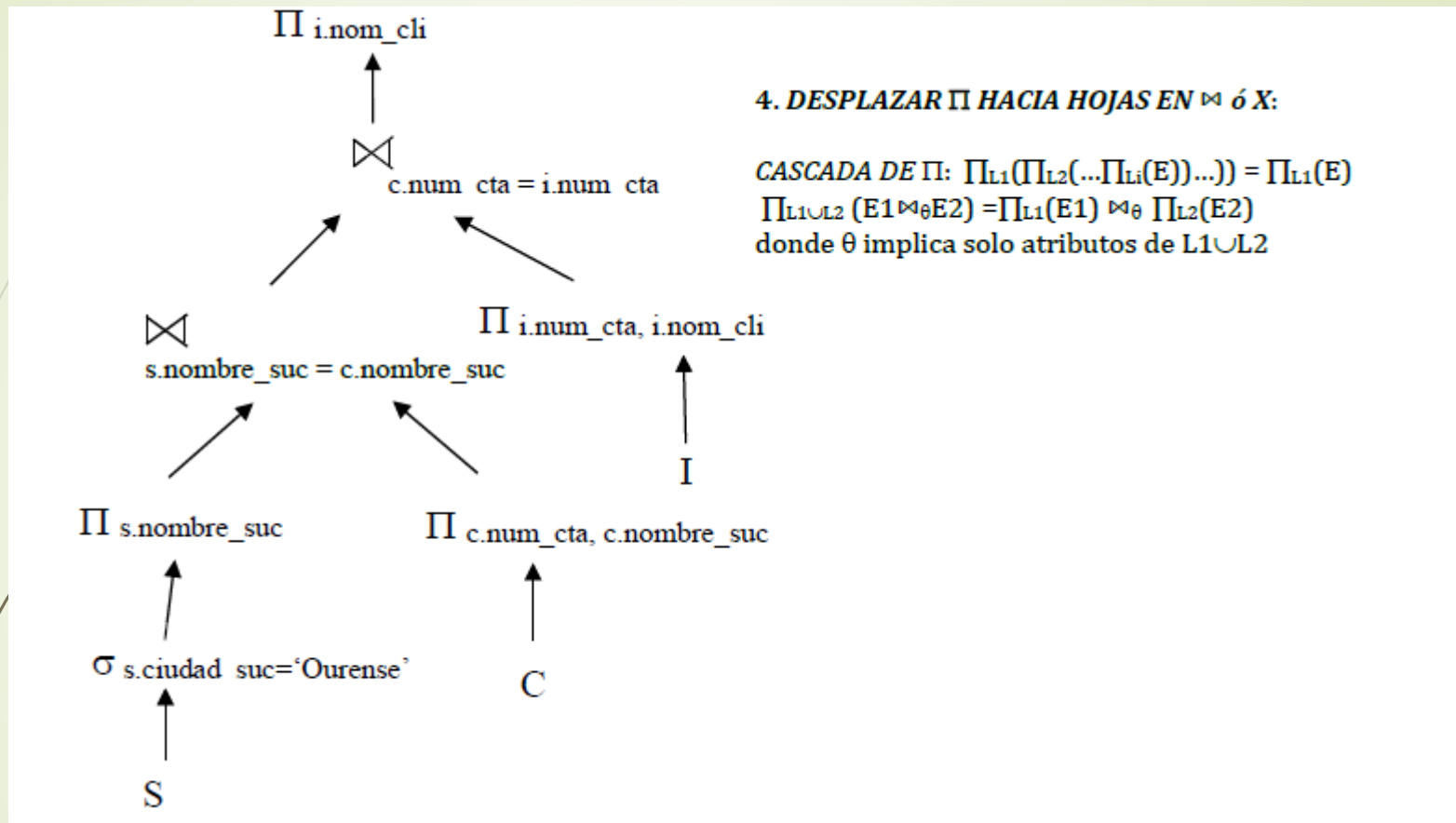
$$(E1 \bowtie E2) \bowtie E3 = E1 \bowtie (E2 \bowtie E3)$$

En este caso, no es necesario aplicar las reglas

2. $\Pi_{nom_cli} ((\sigma_{ciudad_suc='Ourense'} (sucursal) \bowtie cuenta) \bowtie impositor)$

3. $\Pi_{nom_cli} ((\sigma_{ciudad_suc='Ourense'} (sucursal) \bowtie cuenta) \bowtie impositor)$

Ejemplo 1 uso reglas de equivalencia



Expresión final optimizada

$\Pi_{nom_cli} ((\Pi_{nombre_suc} (\sigma_{ciudad_suc='Ourense'} (sucursal))) \bowtie \Pi_{num_cta,nombre_suc} (cuenta)) \bowtie \Pi_{num_cta,nom_cli} (impositor))$

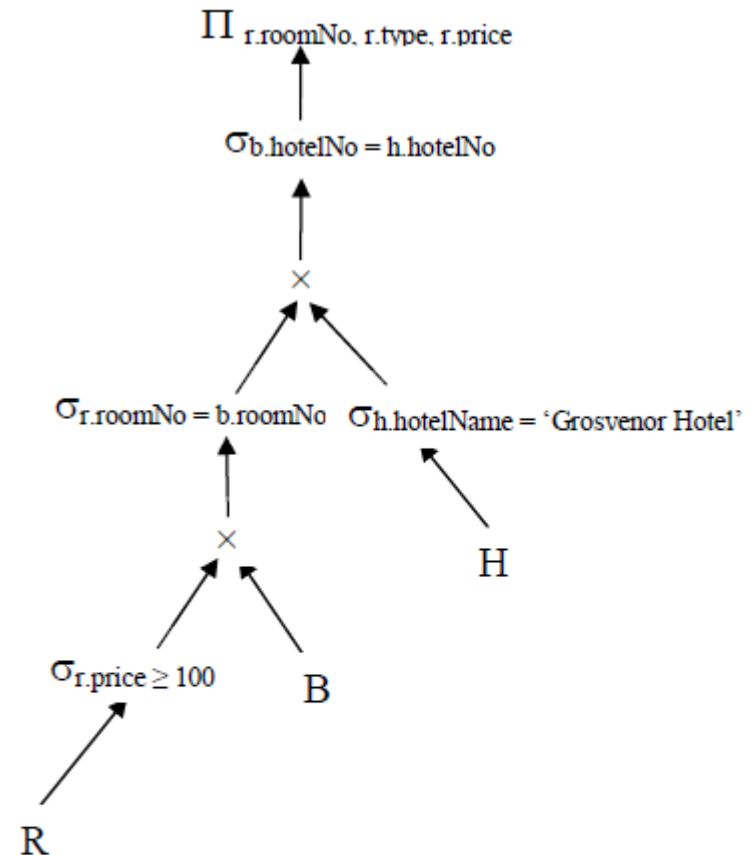
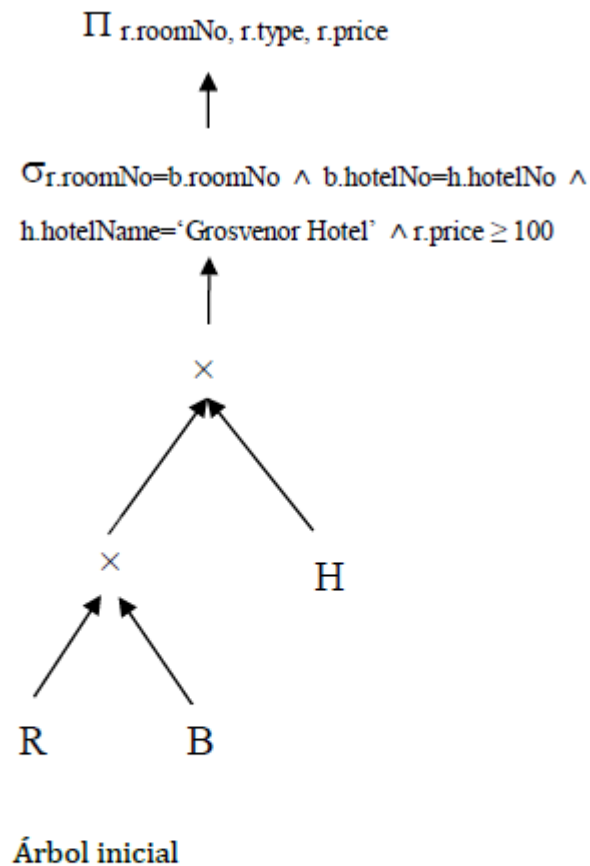
Ejemplo 2 uso reglas de equivalencia

1 – Representar gráficamente mediante árboles las expresiones del álgebra relacional de la consulta siguiente, y transformarla a una forma más eficiente. Enunciar las reglas de equivalencia utilizadas en cada uno de los pasos del proceso:

```
SELECT r.roomNo, r.type, r.price
FROM Room r, Booking b, Hotel h
WHERE r.roomNo = b.roomNo AND b.hotelNo = h.hotelNo AND
      h.hotelName = "Grosvenor Hotel" AND r.price ≥ 100;
```



Ejemplo 2 uso reglas de equivalencia



1. DESPLAZAR σ HACIA HOJAS EN X:

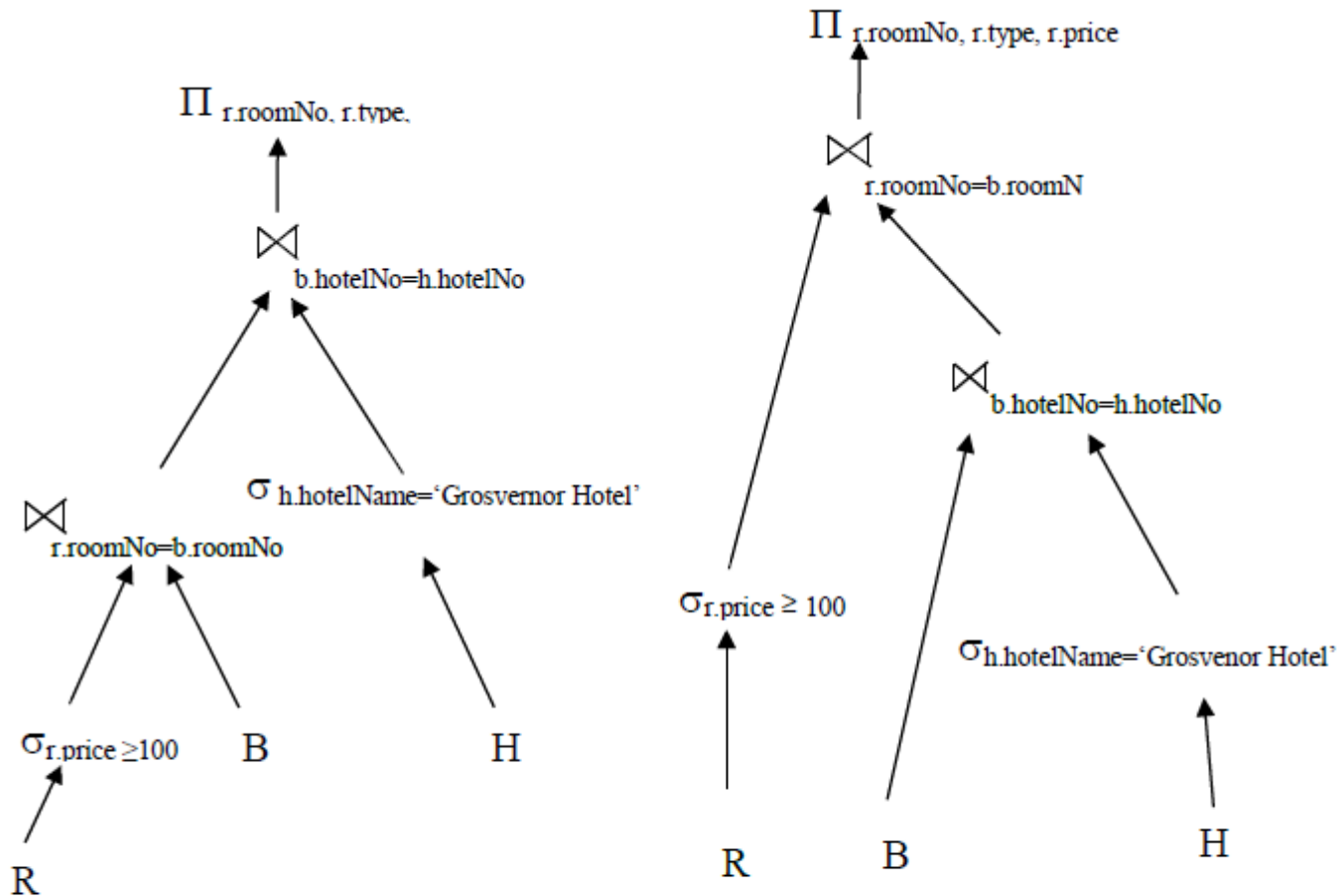
CASCADA DE σ : $\sigma_{\theta_1 \wedge \theta_2}(E) = \sigma_{\theta_1}(\sigma_{\theta_2}(E))$

CONMUTATIVIDAD DE σ : $\sigma_{\theta_1}(\sigma_{\theta_2}(E)) = \sigma_{\theta_2}(\sigma_{\theta_1}(E))$

$\sigma_{\theta_0}(E_1 \times E_2) = \sigma_{\theta_0}(E_1) \times E_2$, θ_0 implica solo atrib. de E_1

$\sigma_{\theta_1 \wedge \theta_2}(E_1 \times E_2) = \sigma_{\theta_1}(E_1) \times \sigma_{\theta_2}(E_2)$ donde θ_1 implica solo atributos de E_1 y θ_2 atributos de E_2

Ejemplo 2 uso reglas de equivalencia



2. Sustituir el producto cartesiano (X) seguido de σ por join:

$$\sigma_{\theta}(E1 \times E2) = E1 \bowtie_{\theta} E2$$

3. Ejecutar antes las selecciones y los joins que producen menos tuplas

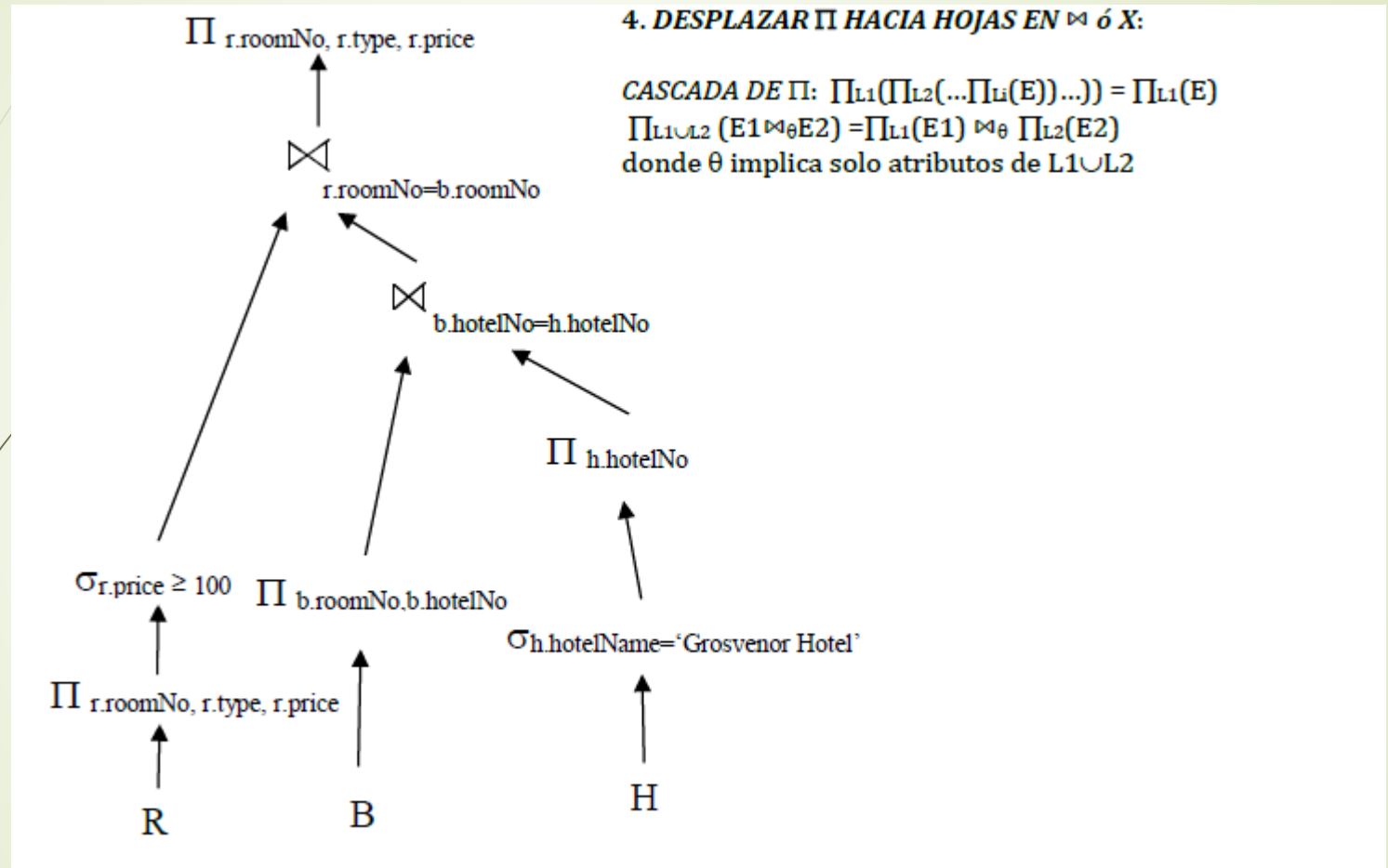
(en este caso, el n° hoteles es inferior al n° habitaciones)

CONMUTATIVIDAD DE $\sigma \Rightarrow \sigma_{\theta_1}(\sigma_{\theta_2}(E)) = \sigma_{\theta_2}(\sigma_{\theta_1}(E))$

ASOCIATIVIDAD DE $\bowtie \Rightarrow$

$$(E1 \bowtie E2) \bowtie E3 = E1 \bowtie (E2 \bowtie E3)$$

Ejemplo 2 uso reglas de equivalencia



Expresión final optimizada

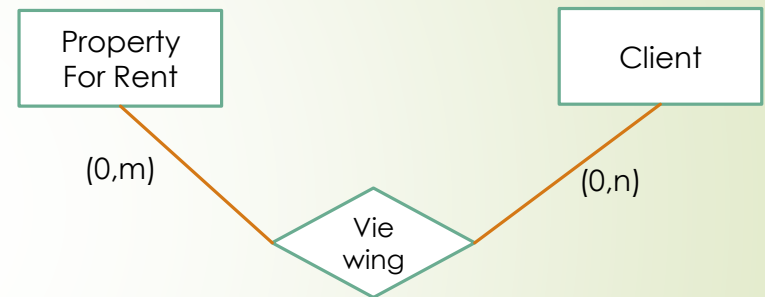
$\Pi_{\text{roomNo, type, price}} (\sigma_{\text{price} \geq 100} (\Pi_{\text{roomNo, type, price}} (\text{Room}))) \bowtie (\Pi_{\text{roomNo, hotelNo}} (\text{Booking}) \bowtie \Pi_{\text{hotelNo}} (\sigma_{\text{hotelName} = \text{'Grosvenor Hotel'}} (\text{Hotel})))$

Ejemplo 3 uso reglas de equivalencia

PropertyForRent (propertyNo, street, city, postcode, type, rooms, rent, ownerNo)

Client (clientNo, fName, lName, address, telNo, prefType, maxRent)

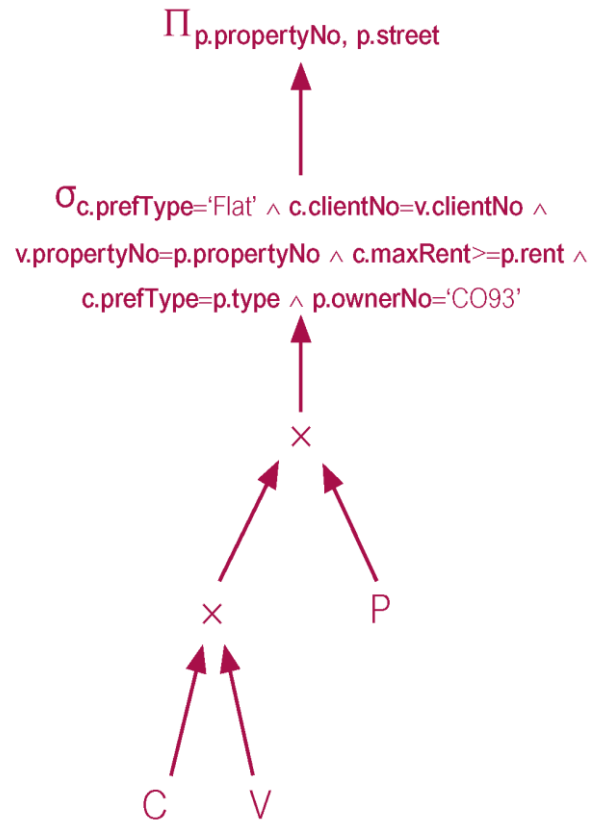
Viewing (clientNo, propertyNo, comment, viewDate)



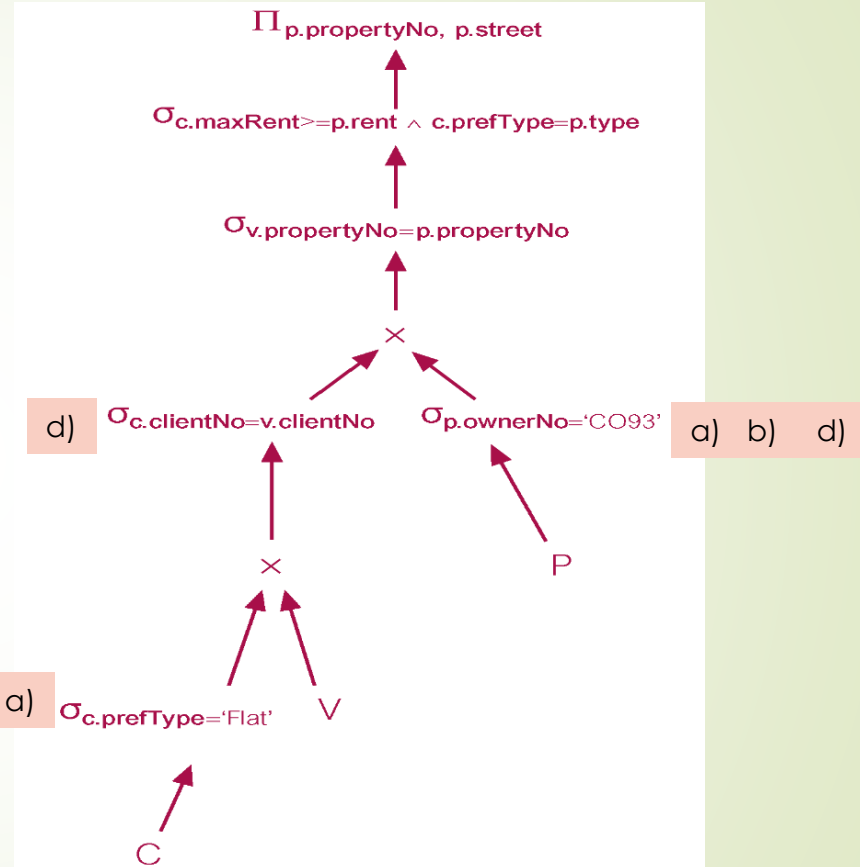
Para los clientes que buscan pisos, localizar los inmuebles que satisfacen sus requisitos y pertenecen al propietario CO93

```
SELECT      p.propertyNo, p.Street
FROM        Client c, Viewing v, PropertyForRent p
WHERE       c.prefType = 'Flat' AND
            c.clientNo = v.clientNo AND
            v.propertyNo = p.propertyNo AND
            c.maxRent >= p.rent AND
            c.prefType = p.type AND
            p.ownerNo = 'CO93';
```


Ejemplo 3 uso reglas de equivalencia



Árbol inicial



1. Desplazar σ hacia las hojas del árbol

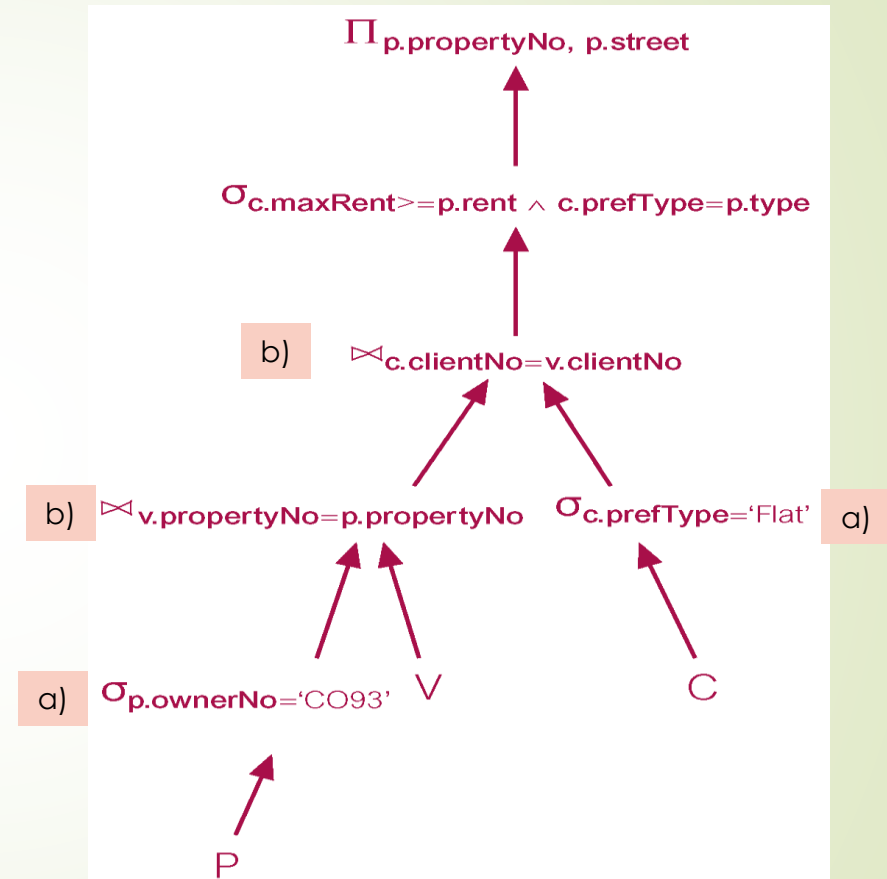
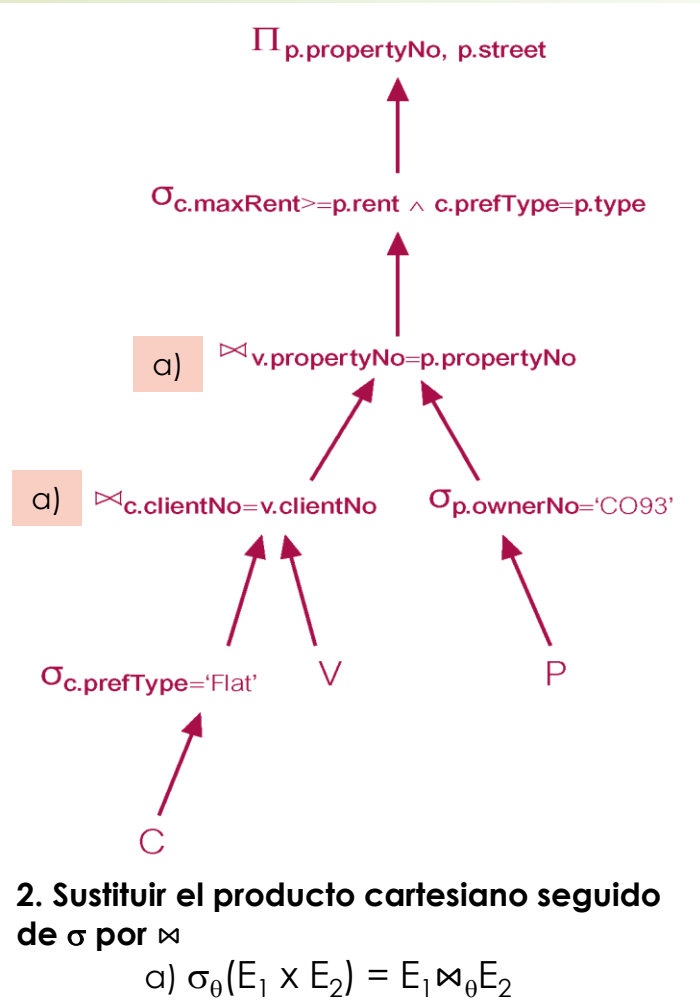
a) CASCADA DE σ $\Rightarrow \sigma_{\theta_1 \wedge \theta_2}(E) = \sigma_{\theta_1}(\sigma_{\theta_2}(E))$

b) CONMUTATIVIDAD DE σ $\Rightarrow \sigma_{\theta_1}(\sigma_{\theta_2}(E)) = \sigma_{\theta_2}(\sigma_{\theta_1}(E))$

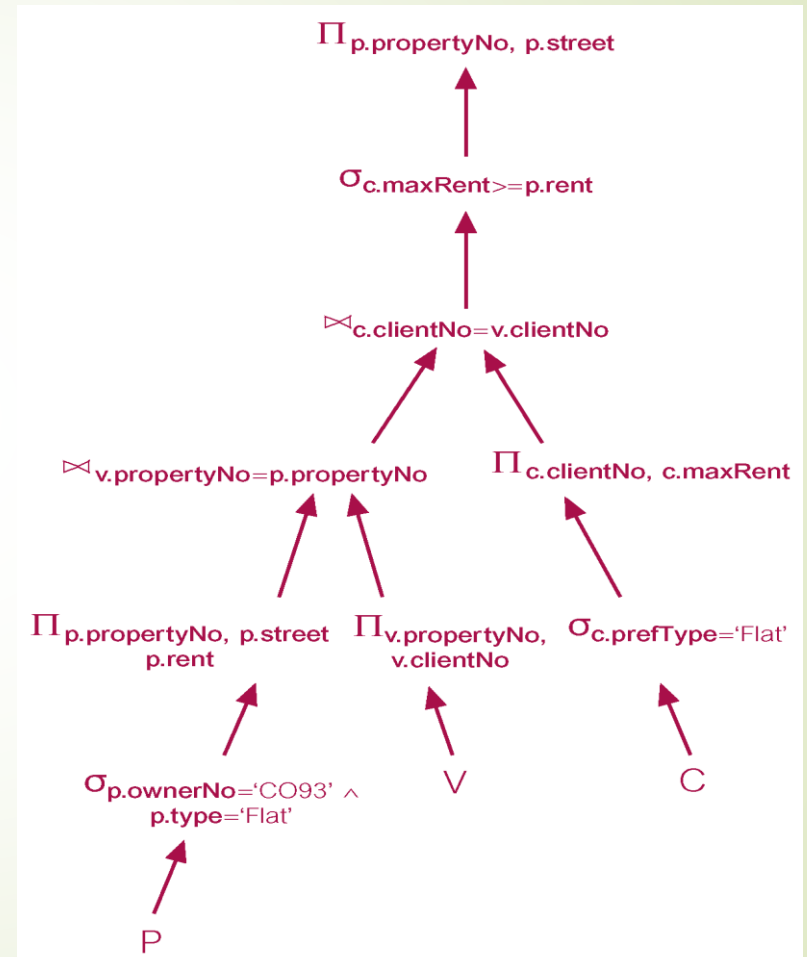
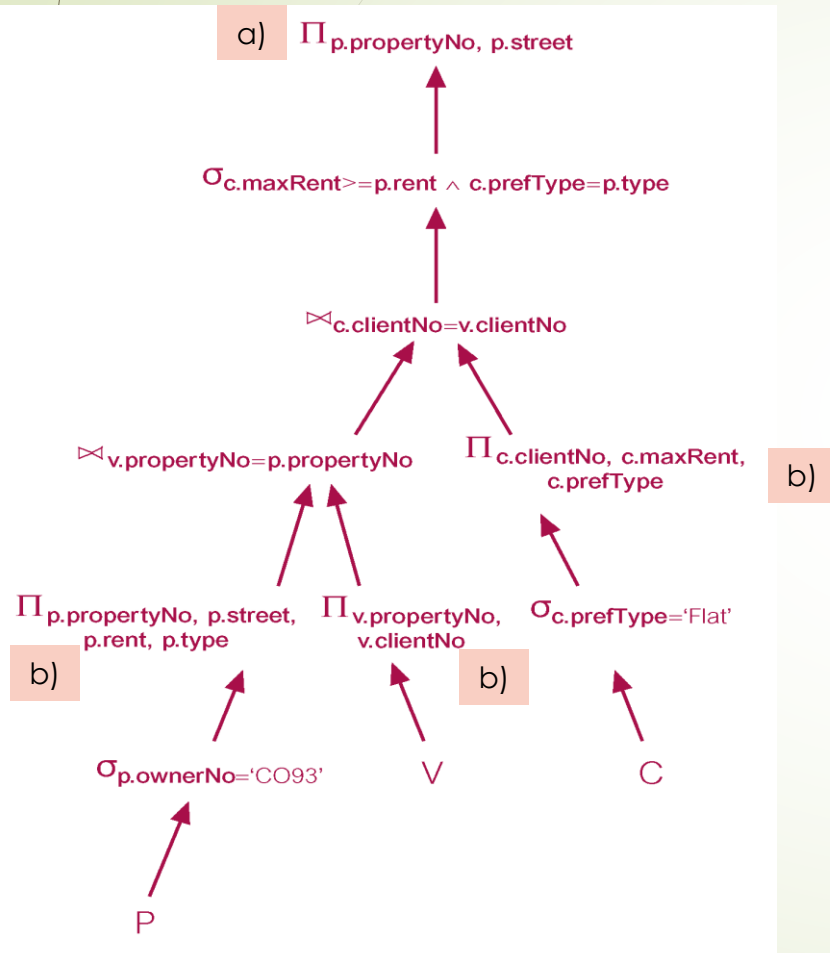
c) $\sigma_{\theta_0}(E_1 \times E_2) = \sigma_{\theta_0}(E_1) \times E_2$ donde θ_0 implica solo atributos de E_1

d) $\sigma_{\theta_1 \wedge \theta_2}(E_1 \times E_2) = \sigma_{\theta_1}(E_1) \times \sigma_{\theta_2}(E_2)$ donde θ_1 implica solo atributos de E_1 y θ_2 de E_2

Ejemplo 3 uso reglas de equivalencia



Ejemplo 3 uso reglas de equivalencia



4. Desplazar Π hacia las hojas del árbol

a) $CASCADE\ DE\ \Pi \Rightarrow \Pi_{L1}(\Pi_{L2}(\dots \Pi_{Li}(E))\dots) = \Pi_{L1}(E)$

b) $\Pi_{L1 \cup L2}(E1 \bowtie_{\theta} E2) = \Pi_{L1 \cup L2}(\Pi_{L1 \cup L3}(E1) \bowtie_{\theta} \Pi_{L2 \cup L4}(E2))$

5. Sustitución de $c.prefType = p.type$ por $(p.type = 'Flat')$

Tipos de optimización: basada en coste

- Se genera una gama de planes de evaluación empleando reglas de equivalencia
- Se realiza una estimación estadística de los resultados de las expresiones
- Ver apartado 14.3 Silberstchatz, A.; Korth, H.; Sudarshan, S. Fundamentos de Bases de Datos. Mc Graw-Hill (5ª edición)
- Se escoge el plan de evaluación de coste mínimo

Estimación estadística de tamaños

➤ Información del catálogo:

➤ n_r (n° registros de r)

➤ $V(A, r)$ (n° valores distintos para el (conjunto de) atributo(s)
 A en la relación r)

Estimación estadística de tamaños

➤ Tamaño del filtro de atributos (proyección)

$$\Pi_A(r) \Rightarrow V(A, r) \text{ porque elimina duplicados}$$

Ej: $\Pi_{\text{edad}}(\text{empleado})$ donde el rango de edad es [20, 60)

$\Rightarrow V(\text{edad}, \text{empleado}) = 40$ valores diferentes

➤ Tamaño de la selección de tuplas (igualdad)

$$\sigma_{A=v}(r) \Rightarrow n_r \frac{\text{casos favorables}}{\text{casos posibles}} = n_r \frac{1}{V(A,r)}$$

Ej: $\sigma_{\text{ciudad_suc} = \text{"Madrid"}}(\text{sucursal})$

donde $n_{\text{sucursal}} = 1000$ tuplas y $V(\text{ciudad_suc}, \text{sucursal}) = 5$ equiprobables

Probabilidad de que 1 tupla satisfaga la condición: $\frac{\text{casos favorables}}{\text{casos posibles}} = \frac{1}{5}$

Nº de tuplas que satisfacen la condición: $1000 \frac{1}{5} = \mathbf{200 \text{ tuplas}}$

Estimación estadística de tamaños

► Tamaño de la selección de tuplas (rango)

$$\sigma_{A < v}(r) \Rightarrow n_r \left(\frac{v - \min(A, r)}{\max(A, r) - \min(A, r)} \right) \text{ si se conoce } v \text{ y hay distribución uniforme de valores}$$

Ej: $\sigma_{\text{edad} < 30}$ (empleado) edad $[20, 60)$, $n_{\text{empleado}} = 2000$

Probabilidad de que 1 tupla satisfaga la condición: $\frac{\text{casos favorables}}{\text{casos posibles}} = \frac{30-20}{60-20} = \frac{1}{4}$

Nº de tuplas que satisfacen la condición: $2000 \cdot \frac{1}{4} = \mathbf{500 \text{ tuplas}}$

$$\sigma_{A < v}(r) \Rightarrow \frac{n_r}{2} \text{ si NO se conoce } v$$

Ej: $\sigma_{\text{edad} < 'X'}$ (empleado) edad $[20, 60)$, $n_{\text{empleado}} = 2000$

Probabilidad de que 1 tupla satisfaga la condición: $\frac{1}{2}$

Nº de tuplas que satisfacen la condición: $2000 \cdot \frac{1}{2} = \mathbf{1000 \text{ tuplas}}$

Estimación estadística de tamaños

- Tamaño de la selección de tuplas (rango)

$$\sigma_{A \geq v}(r) \Rightarrow n_r \left(\frac{\max(A,r) - v}{\max(A,r) - \min(A,r)} \right)$$

Ej: $\sigma_{\text{edad} \geq 30}(\text{empleado})$ edad $[20, 60]$, $n_{\text{empleado}} = 2000$

Probabilidad de que 1 tupla satisfaga la condición: $\frac{\text{casos favorables}}{\text{casos posibles}} = \frac{60-30}{60-20} = \frac{3}{4}$

Nº de tuplas que satisfacen la condición: $2000 \cdot \frac{3}{4} = \mathbf{1500 \text{ tuplas}}$

Estimación estadística de tamaños

► Tamaño de la selección de tuplas (conjunción de condiciones)

$$\sigma_{\theta_1 \wedge \theta_2 \wedge \dots \wedge \theta_m}(r)$$

Probabilidad $P_{\theta t}$ de que 1 tupla t satisfaga la condición θ_1 : $\frac{\text{casos favorables}}{\text{casos posibles}}$

Si las condiciones son independientes, la probabilidad de que 1 tupla cumpla todas las condiciones es el producto de las probabilidades de las condiciones: $P_{\theta_1} * P_{\theta_2} * \dots * P_{\theta_m}$

Nº tuplas que satisface la selección completa: $n_r(P_{\theta_1} * P_{\theta_2} * \dots * P_{\theta_m})$

Ej: $\sigma_{(\text{edad} < 30) \wedge (\text{salario} = 120)}(\text{empleado})$ donde edad $[20, 60)$, $n_{\text{empleado}} = 2000$, salario $[50, 150)$

Probabilidad de que 1 tupla satisfaga la condición (edad < 30): $\frac{\text{casos favorables}}{\text{casos posibles}} = \frac{1}{4}$

Probabilidad de que 1 tupla satisfaga la condición (salario = 120): $\frac{\text{casos favorables}}{\text{casos posibles}} = \frac{1}{150-50} = \frac{1}{100}$

Probabilidad de que 1 tupla satisfaga (edad < 30) \wedge (salario = 120): $\frac{1}{4} \cdot \frac{1}{100} = \frac{1}{400}$

Nº de tuplas que satisface (edad < 30) \wedge (salario = 120): $2000 \cdot \frac{1}{400} = \mathbf{5 \text{ tuplas}}$

Estimación estadística de tamaños

► Tamaño de la selección de tuplas (disyunción de condiciones)

$$\sigma_{\theta_1 \vee \theta_2 \vee \dots \vee \theta_m} (r)$$

Probabilidad $P_{\theta t}$ de que 1 tupla t satisfaga la condición θ_1 : $\frac{\text{casos favorables}}{\text{casos posibles}}$

Si las condiciones son independientes, la probabilidad de que 1 tupla satisfaga la disyunción es:

(1 menos la probabilidad de que no satisfaga ninguna) $1 - (1 - P_{\theta_1}) * (1 - P_{\theta_2}) * \dots * (1 - P_{\theta_m})$

Nº tuplas que satisface la selección completa: $n_r [1 - (1 - P_{\theta_1}) * (1 - P_{\theta_2}) * \dots * (1 - P_{\theta_m})]$

Ej: $\sigma_{(\text{edad} < 30) \vee (\text{salario} = 120)} (\text{empleado})$ donde edad $[20, 60)$, $n_{\text{empleado}} = 2000$, salario $[50, 150)$

Probabilidad de que 1 tupla satisfaga la condición (edad < 30): $\frac{\text{casos favorables}}{\text{casos posibles}} = \frac{1}{4}$

Probabilidad de que 1 tupla satisfaga la condición (salario = 120): $\frac{\text{casos favorables}}{\text{casos posibles}} = \frac{1}{150-50} = \frac{1}{100}$

Probabilidad de que 1 tupla satisfaga (edad < 30) o (salario = 120)

*es (1 menos la probabilidad de que no satisfaga **ninguna**) : $1 - \left(1 - \frac{1}{4}\right) * \left(1 - \frac{1}{100}\right) = \frac{103}{400}$*

Nº de tuplas que satisface (edad < 30) o (salario = 120): $2000 \frac{103}{400} = 515 \text{ tuplas}$

Estimación estadística de tamaños

- Join de relaciones $r \bowtie s$, siendo **R** el esquema de **r** y **S** el esquema de **s**

- a) Si $R \cap S = \emptyset \Rightarrow$ **no tienen atributos en común** $\Rightarrow n_r * n_s$
- b) Si $R \cap S$ es **clave de R** \Rightarrow cada tupla de **s** se combina como máximo con 1 tupla de **r** \Rightarrow el nº de tuplas del join NO es mayor que las tuplas de **s**

No Titulados/as

c_emp	...
e2	...
e3	...
e6	...

Titulados/as

c_emp	...
e1	...
e4	...
e5	...

Proyectos

n_proy	jefe
p1	e2
p2	e1
p3	e4
p4	e1

Titulados/as \bowtie Proyectos

c_emp	...	n_proy	jefe
e1	...	p2	e1
e4	...	p3	e4
e1	...	p4	e1

Estimación estadística de tamaños

- Join de relaciones $r \bowtie s$, siendo **R** el esquema de **r** y **S** el esquema de **s**

- a) Si $R \cap S = \emptyset \Rightarrow$ **no tienen atributos en común** $\Rightarrow n_r * n_s$
- b) Si $R \cap S$ es clave de **R** \Rightarrow cada tupla de **s** se combina como máximo con 1 tupla de **r** \Rightarrow el nº de tuplas del join NO es mayor que las tuplas de **s**

Si, además, los atributos comunes son **clave foránea de S**, entonces el nº de tuplas del join es EXACTAMENTE el nº de tuplas de **s**

Empleados/as

c_emp	...
e1	...
e2	...
e3	...
e4	...

Proyectos

n_proy	jefe
p1	e2
p2	e1
p3	e1

Empleados/as \bowtie Proyectos

c_emp	...	n_proy	jefe
e2	...	p1	e2
e1	...	p2	e1
e1	...	p3	e1

Estimación estadística de tamaños

- Join de relaciones $r \bowtie s$, siendo **R** el esquema de **r** y **S** el esquema de **s**

c) Si $R \cap S$ **NO es clave de R ni de S**, suponiendo que todos los valores son equiprobables, entonces:

Cli	
n_cli	edad
c1	10
c2	10
c3	20
c4	20
c5	25
c6	25
c7	30
c8	30

$$V(\text{edad}, \text{clientes}) = 4$$
$$n_{\text{clientes}} = 8$$

Prov	
n_prov	edad
p1	15
p2	15
p3	15
p4	20
p5	20
p6	20

$$V(\text{edad}, \text{proveedores}) = 2$$
$$n_{\text{proveedores}} = 6$$

$$1 \text{ tupla de Cli produce como máximo: } \frac{n_{\text{Prov}}}{V(\text{edad}, \text{Prov})} = \frac{6}{2} = 3 \text{ tuplas}$$

$$\text{Todas las tuplas de Cli producen como máximo } n_{\text{Cli}} \frac{n_{\text{Prov}}}{V(\text{edad}, \text{Prov})} = 8 \times 3 = 24 \text{ tuplas}$$

$$1 \text{ tupla de Prov produce como máximo: } \frac{n_{\text{Cli}}}{V(\text{edad}, \text{Cli})} = \frac{8}{4} = 2 \text{ tuplas}$$

$$\text{Todas las tuplas de Prov producen como máximo } n_{\text{Prov}} \frac{n_{\text{Cli}}}{V(\text{edad}, \text{Cli})} = 6 \times 2 = 12 \text{ tuplas}$$

$$n_{\text{Cli} \bowtie \text{Prov}} = \min(24, 12) = 12 \text{ tuplas}$$

Estimación estadística de tamaños

- Join de relaciones $r \bowtie s$, siendo **R** el esquema de **r** y **S** el esquema de **s**

c) Si $R \cap S$ NO es clave de R ni de S, suponiendo que todos los valores son equiprobables, entonces:

- a) cada tupla de R produce $\frac{n_s}{V(A,s)}$ tuplas en el join, donde A son los atributos comunes
=> todas las tuplas de R producen $n_r \frac{n_s}{V(A,s)}$ tuplas
- b) cada tupla de S produce $\frac{n_r}{V(A,r)}$ tuplas en el join, donde A son los atributos comunes
=> todas las tuplas de S producen $n_s \frac{n_r}{V(A,r)}$ tuplas
- c) si $V(A,r) \neq V(A,s)$, el nº de tuplas del join es **min** $(\frac{n_r n_s}{V(A,r)}, \frac{n_r n_s}{V(A,s)})$

Ejemplo 1 cálculo de estadísticas

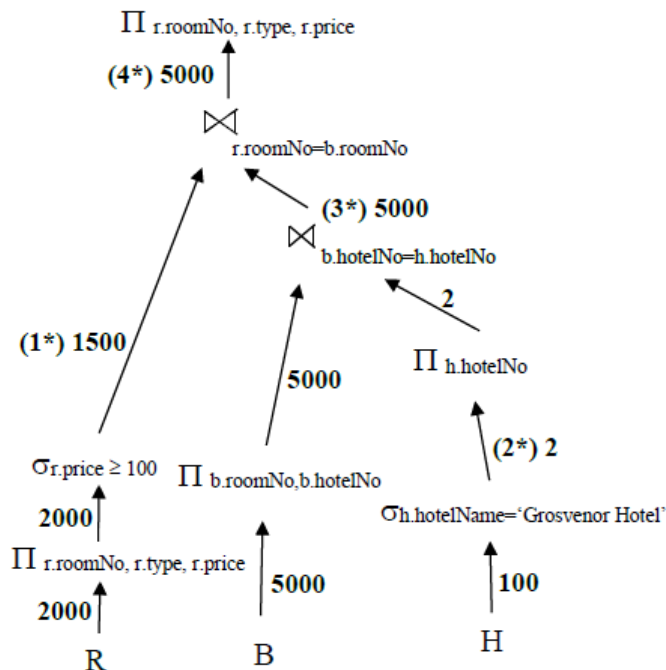
```
SELECT r.roomNo, r.type, r.price
FROM Room r, Booking b, Hotel h
WHERE r.roomNo = b.roomNo AND b.hotelNo = h.hotelNo
AND h.hotelName = "Grosvenor Hotel" AND r.price ≥ 100;
```

Estimar el tamaño de las operaciones que aparecen en el árbol final, teniendo en cuenta lo siguiente:

nRoom = 2000,
V(hotelName, Hotel) = 50,

nHotel = 100,
min(r.price, Room) = 50,

nBooking = 5000,
max(r.price, Room) = 250



(1*) $\sigma_{r.price \geq 100}$, dado que se conoce a priori el valor de comparación se calcula como:

$$2000 (250-100)/(250-50) = 1500 \text{ tuplas}$$

(2*) $\sigma_{h.hotelName = 'Grosvenor Hotel'}$ si se supone una distribución uniforme de los valores del atributo h.hotelName, dado que $V(hotelName, Hotel) = 50$, la probabilidad de que una tupla cumpla la selección es $\frac{\text{casos favorables}}{\text{casos posibles}} = \frac{1}{50}$

Por tanto, el nº de tuplas que satisfacen la condición: $100 \times 1/50 = 2$ tuplas

(3*) $B \bowtie H$, donde $B \cap H = hotelNo$, que es clave de H y foránea de B, entonces el número de tuplas no es mayor que las de B = 5000.

(4*) $(B \bowtie H) \bowtie R$, donde $B \bowtie H \cap R = roomNo$, que es clave de R y foránea de B. Entonces el número de tuplas no es mayor que las de $B \bowtie H = 5000$.

Ejemplo 2 cálculo de estadísticas

Estimar el tamaño de las operaciones que aparecen en el árbol, teniendo en cuenta lo siguiente:

$$n_{\text{Client}} = 1000,$$

$$V(\text{prefType}, \text{Client}) = 5,$$

$$\min(\text{maxRent}, \text{Client}) = 100,$$

$$\min(\text{rent}, \text{PropertyForRent}) = 350,$$

$$V(\text{rent}, \text{PropertyForRent}) = 20$$

$$V(\text{clientNo}, \text{Viewing}) = 250,$$

$$V(\text{type}, \text{PropertyForRent}) = 5$$

$$n_{\text{Viewing}} = 3000,$$

$$V(\text{maxRent}, \text{Client}) = 50$$

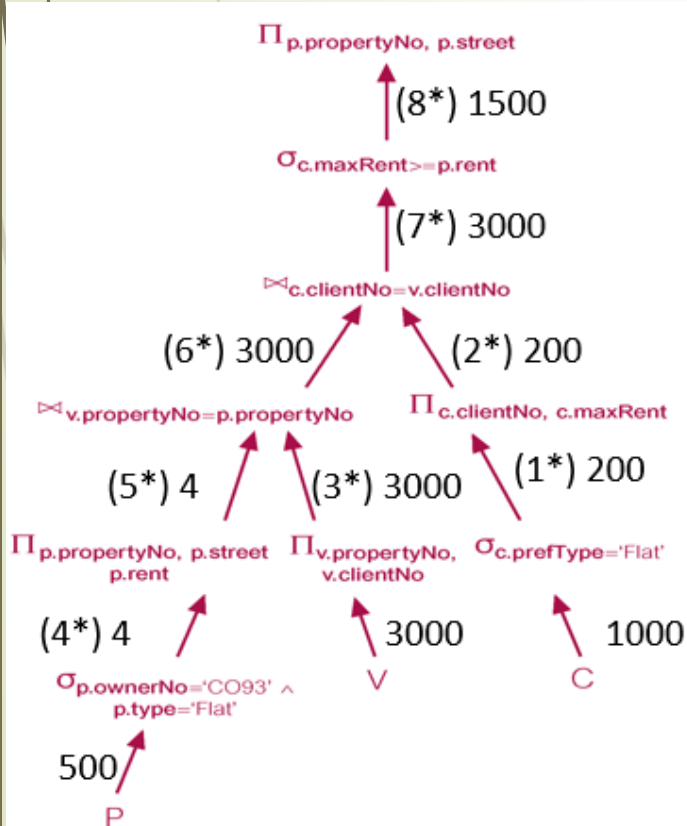
$$\max(\text{maxRent}, \text{Client}) = 3000,$$

$$\max(\text{rent}, \text{PropertyForRent}) = 1000,$$

$$V(\text{ownerNo}, \text{PropertyForRent}) = 30,$$

$$V(\text{propertyNo}, \text{Viewing}) = 400$$

$$n_{\text{PropertyForRent}} = 500,$$



(1*) $\sigma_{c.\text{prefType} = \text{'Flat'}}$ si se supone una distribución uniforme de los valores del atributo $c.\text{prefType}$, la probabilidad de que 1 tupla cumpla la selección es $\frac{\text{casos favorables}}{\text{casos posibles}} = \frac{1}{5}$

Por tanto, el nº de tuplas que satisfacen la condición: $1000 \times \frac{1}{5} = 200 \text{ tuplas}$

(2*) $\pi_{c.\text{clientNo}, c.\text{maxRent}}$ su tamaño estimado es $V(c.\text{clientNo}, c.\text{maxRent}, C)$. Como $c.\text{clientNo}$ es clave en C, no hay posibilidad de la existencia de duplicados, por lo que el número de tuplas se mantiene.

(3*) $\pi_{v.\text{propertyNo}, v.\text{clientNo}}$ su tamaño estimado es $V(v.\text{propertyNo}, v.\text{clientNo}, V)$. Como $v.\text{propertyNo}$ y $v.\text{clientNo}$ forman la clave en V, no hay posibilidad de la existencia de duplicados, por lo que el número de tuplas se mantiene.

Ejemplo 2 cálculo de estadísticas

Estimar el tamaño de las operaciones que aparecen en el árbol final, teniendo en cuenta lo siguiente:

$$n_{\text{Client}} = 1000,$$

$$V(\text{preftype}, \text{Client}) = 5,$$

$$\min(\text{maxRent}, \text{Client}) = 100,$$

$$\min(\text{rent}, \text{PropertyForRent}) = 350,$$

$$V(\text{rent}, \text{PropertyForRent}) = 20$$

$$V(\text{clientNo}, \text{Viewing}) = 250,$$

$$V(\text{type}, \text{PropertyForRent}) = 5$$

$$n_{\text{Viewing}} = 3000,$$

$$V(\text{maxRent}, \text{Client}) = 50$$

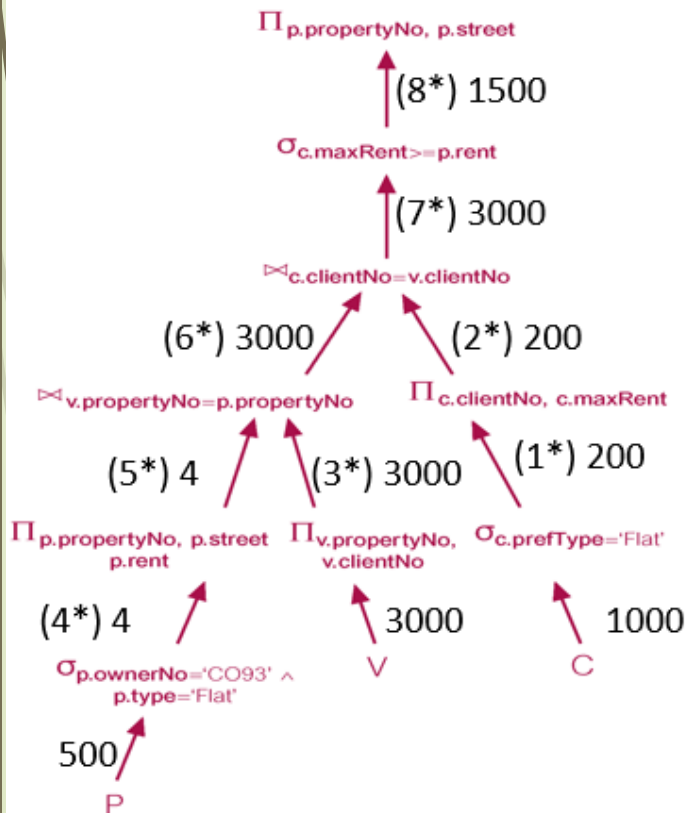
$$\max(\text{maxRent}, \text{Client}) = 3000,$$

$$\max(\text{rent}, \text{PropertyForRent}) = 1000,$$

$$V(\text{ownerNo}, \text{PropertyForRent}) = 30,$$

$$V(\text{propertyNo}, \text{Viewing}) = 400$$

$$n_{\text{PropertyForRent}} = 500,$$



(4*) $\sigma_{p.\text{ownerNo} = 'C093' \wedge p.\text{type} = 'Flat'}$ Debe estimarse el tamaño de cada parte de la conjunción

$\sigma_{p.\text{ownerNo} = 'C093'}$ como se conoce el valor de $p.\text{ownerNo}$, y se supone una distribución uniforme de los valores del atributo, la probabilidad de que 1 tupla cumpla la selección es $\frac{\text{casos favorables}}{\text{casos posibles}} = \frac{1}{30}$

$\sigma_{p.\text{type} = 'Flat'}$ como se conoce el valor de $p.\text{type}$, y se supone una distribución uniforme de los valores del atributo, la probabilidad de que 1 tupla cumpla la selección es $\frac{\text{casos favorables}}{\text{casos posibles}} = \frac{1}{5}$

Dado que las condiciones son independientes, la probabilidad de que 1 tupla satisfaga

$$(\sigma_{p.\text{ownerNo} = 'C093'}) \wedge (\sigma_{p.\text{type} = 'Flat'}) \text{ es } \frac{1}{30} \times \frac{1}{5} = \frac{1}{150}$$

De este modo, el nº de tuplas que satisface ambas condiciones es $500 \times \frac{1}{150} \simeq 4 \text{ tuplas}$

(5*) $\Pi_{p.\text{propertyNo}, p.\text{street}, p.\text{rent}}$ su tamaño estimado es $V(p.\text{propertyNo}, p.\text{street}, p.\text{rent}, P)$. Como $p.\text{propertyNo}$ es clave en P, no hay posibilidad de la existencia de duplicados, por lo que el número de tuplas se mantiene.

Ejemplo 2 cálculo de estadísticas

Estimar el tamaño de las operaciones que aparecen en el árbol final, teniendo en cuenta lo siguiente:

$$n_{\text{Client}} = 1000,$$

$$V(\text{preftype}, \text{Client}) = 5,$$

$$\min(\text{maxRent}, \text{Client}) = 100,$$

$$\min(\text{rent}, \text{PropertyForRent}) = 350,$$

$$V(\text{rent}, \text{PropertyForRent}) = 20$$

$$V(\text{clientNo}, \text{Viewing}) = 250,$$

$$V(\text{type}, \text{PropertyForRent}) = 5$$

$$n_{\text{Viewing}} = 3000,$$

$$V(\text{maxRent}, \text{Client}) = 50$$

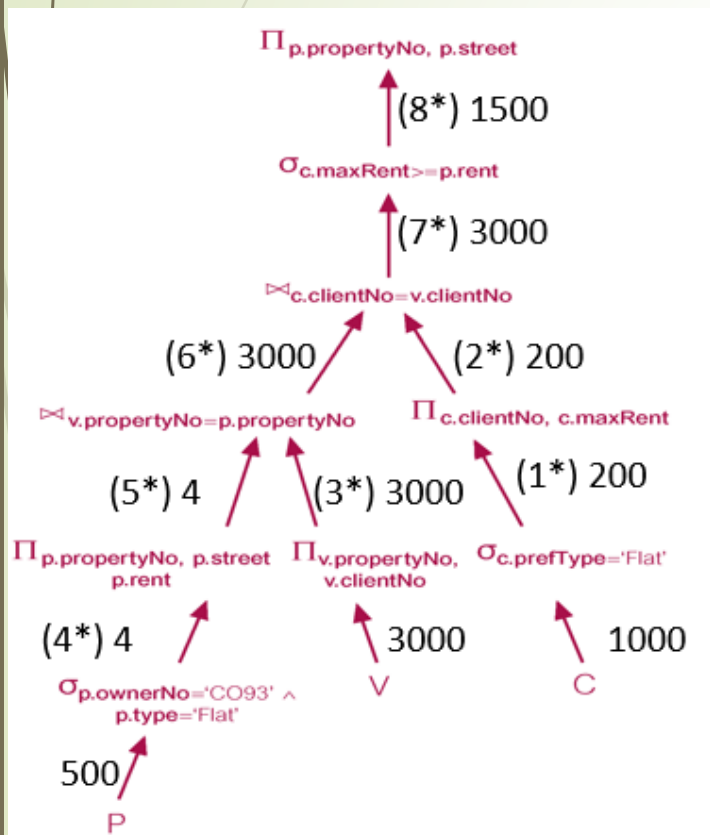
$$\max(\text{maxRent}, \text{Client}) = 3000,$$

$$\max(\text{rent}, \text{PropertyForRent}) = 1000,$$

$$V(\text{ownerNo}, \text{PropertyForRent}) = 30,$$

$$V(\text{propertyNo}, \text{Viewing}) = 400$$

$$n_{\text{PropertyForRent}} = 500,$$



(6*) $\bowtie_{v.\text{propertyNo} = p.\text{propertyNo}}$, como propertyNo es clave de PropertyForRent, cada tupla de V se combina como máximo con una tupla de P. Por lo tanto, $n_{V \bowtie P} \leq n_V$. Como propertyNo es foránea en P, entonces $n_{V \bowtie P} = n_V = 3000$ tuplas

(7*) $\bowtie_{c.\text{clientNo} = v.\text{clientNo}}$, como c.clientNo es clave de Client, cada tupla de V se combina como máximo con una tupla de C. Por lo tanto, $n_{C \bowtie V} \leq n_V$. Como clientNo es foránea en V, entonces $n_{C \bowtie V} = n_V = 3000$ tuplas

(8*) $\sigma_{c.\text{maxRent} \geq p.\text{rent}}$ dado que no se conoce a priori el valor de p.rent la probabilidad de que 1 tupla cumpla la selección es $\frac{\text{casos favorables}}{\text{casos posibles}} = \frac{1}{2}$

Por tanto, el nº de tuplas que satisfacen la condición: $3000 \times \frac{1}{2} = 1500$ tuplas