# Lezione 15 - Decomposizioni che hanno un join senza perdita

Prof.ssa Maria De Marsico demarsico@di.uniroma1.it



# Cosa significa join senza perdita



Quando uno schema di relazione viene di solito decomposto:

- quando non è in 3NF
- per motivi di <u>efficienza degli accessi</u>
  - più piccola è la taglia delle tuple maggiore è il numero che riusciamo a caricare in memoria nella stessa operazione di lettura
  - se le informazioni della tupla non vengono utilizzate dallo stesso tipo di operazioni nella base di dati meglio decomporre lo schema
  - Esempio: lo schema Studente potrebbe essere decomposto separando le informazioni anagrafiche (CF, Nome, Cognome, DataNascita, LuogoNascita, ecc.) da quelle accademiche (Matricola, CorsoLaurea, AnnoCorso, ecc.)

# Cosa significa join senza perdita



- Abbiamo visto che quando uno schema viene decomposto, non basta che i sottoschemi siano in 3NF
- Rivediamo i due esempi (uno più «astratto», l'altro più «concreto»)



- •Consideriamo lo schema R=ABC con l'insieme di dipendenze funzionali  $F=\A\to B$ ,  $C\to B\A$  (lo schema non è in 3NF per la presenza in  $F^+$  delle dipendenze parziali  $A\to B$  e  $C\to B$ , dato che la chiave è AC. Tale schema può essere decomposto in:
- $R1=AB \operatorname{con} \{A \rightarrow B\}$  e
- $R2=BC \operatorname{con} \{C \rightarrow B\}$ .
- •Tale decomposizione pur preservando tutte le dipendenze in *F*+ non è soddisfacente.



Consideriamo l'istanza legale di R

sono veri i due fatti (a1,b1,c1) e (a2,b1,c2) e **non altri** 

 In base alla decomposizione data, questa istanza si decompone in

•E dovrebbe essere possibile ricostruirla **esattamente** tramite join ... invece ...



• ... e invece se si effettua il join delle due istanze legali risultanti dalla decomposizione si ottiene

R	A	В	С		
	a1	b1	c1		
	a2	b1	c2		
	a1	b1	c2	_	tuple estranee alla realtà di interesse quindi
	a2	b1	c1		perdita di informazione

 Occorre garantire che il join delle istanze risultanti dalla decomposizione non riveli perdita di informazione



Consideriamo ora lo schema

R=(Matricola, Provincia, Comune) con l'insieme di dipendenze funzionali

 $F=\mbox{Matricola}\rightarrow\mbox{Provincia, Comune}\rightarrow\mbox{Provincia}\mbox{$\stackrel{/}{\sim}$}$  (lo schema non è in 3NF per la presenza in  $F^+$  delle dipendenze parziali  $\mbox{Matricola}\rightarrow\mbox{Provincia}$  e  $\mbox{Comune}\rightarrow\mbox{Provincia}$ , dato che la chiave è (Matricola, Comune) (Comune non è determinato da nessun altro attributo!)

Tale schema può essere decomposto in:

R1= (Matricola, Provincia) con  $\langle$  Matricola $\rightarrow$ Provincia $\rangle$  e R2= (Provincia, Comune) con  $\langle$  Comune $\rightarrow$ Provincia $\rangle$ .

Tale schema **pur preservando tutte le dipendenze in** *F*<sup>+</sup> non è soddisfacente.



## Consideriamo l'i stanza **legale** di R

R

Matricola	Provincia	Comune
501	Roma	Tivoli
502	Roma	Mandela

sono veri i due fatti (501,Roma,Tivoli) e (501, Roma,Mandela) e **non altri** 

In base alla decomposizione data, questa istanza si decompone in

**R1** 

Matricola	Provincia
501	Roma
502	Roma

**R2** 

Provincia	Comune
Roma	Tivoli
Roma	Mandela

E dovrebbe essere possibile ricostruirla esattamente tramite join ...

invece ...



• ... e invece se si effettua il join delle due istanze legali risultanti dalla decomposizione si ottiene

г	_	
F	~	
	-	

Matricola	Provincia	Comune	
501	Roma	Tivoli	
502	Roma	Mandela	
501	Roma	Mandela	tuple estranee alla realtà di interesse
502	Roma	Tivoli	quindi perdita di informazione



In conclusione, quando si decompone uno schema occorre tenere presente il seguente requisito dello schema decomposto:

 deve permettere di ricostruire mediante join naturale ogni istanza legale dello schema originario (senza aggiunta di informazione estranea)

### **Definizione**



Se si decompone uno schema di relazione R si vuole che la decomposizione  $\{R_1, R_2, ..., R_k\}$  ottenuta sia tale che ogni istanza legale r di R sia ricostruibile mediante join naturale ( $\triangleright \triangleleft$ ) da un istanza legale  $\{r_1, r_2, ..., r_k\}$  dello schema decomposto  $\{R_1, R_2, ..., R_k\}$ . Poiché per ricostruire una tupla t di r è necessario che  $t[R_i] \in r_i$ , i=1,...,k, si deve avere  $r_i=\pi_{Ri}$  (r), i=1,...,k.

**Definizione** Sia R uno schema di relazione. Una decomposizione  $\rho = \{R_1, R_2, ..., R_k\}$  di R ha un join senza perdita se **per ogni** istanza **legale** r di R si ha  $r = \pi_{R1}(r) \triangleright \triangleleft \pi_{R2}(r) \triangleright \triangleleft ... \triangleright \triangleleft \pi_{Rk}(r)$ .

#### **Problema**



- Anche in questo caso partiamo da una decomposizione data, e cerchiamo un modo per verificare che goda della proprietà desiderata.
- Anche in questo caso facciamo alcune premesse.

#### Teorema



**Teorema** Sia R uno schema di relazione e  $\rho = \{R_1, R_2, ..., R_k\}$  una decomposizione di R. Per ogni istanza legale r di R, indicato con  $m_{\rho}(r) = \pi_{R1}(r) \triangleright \triangleleft \pi_{R2}(r) \triangleright \triangleleft ... \triangleright \triangleleft \pi_{Rk}(r)$ , si ha:

$$a)r \subseteq m_{\rho}(r)$$

**b**)
$$\pi_{Ri}(m_{\rho}(r)) = \pi_{Ri}(r)$$

$$\mathbf{c})\mathbf{m}_{\rho}(\mathbf{m}_{\rho}(\mathbf{r})) = \mathbf{m}_{\rho}(\mathbf{r}).$$

Dim.

**Prova di a).** Sia **t una tupla di** *r*. Per ogni *i*, i=1,...,k,  $t[R_i] \in \pi_{R_i}(r)$  e quindi  $t \in m_{R_i}(r)$ .

- una qualsiasi, quindi vale per ogni tupla di r
- Il **sottoinsieme** dei suoi valori su attributi di  $R_i$  sarà sicuramente in una tupla della proiezione  $\pi_{Ri}(r)$
- quindi OGNI sottoinsieme dei suoi valori corrispondente ad un elemento della decomposizione dello schema comparirà in una tupla di un elemento della decomposizione dell'istanza
- poiché anche la proiezione è un **insieme** di tuple al più i duplicati corrisponderanno ad una **unica tupla nella proiezione**
- Quindi ogni tupla potrà essere ricostruita tramite il join naturale ...
- ... anzi il join naturale <u>potrebbe</u> ottenere tuple che <u>non erano</u> in r ricomponendo <u>pezzi di tuple diverse</u>!

## **Digressione**



Consideriamo la solita istanza **legale** di R=ABC con l'insieme di dipendenze funzionali F=ABC (no 3NF – dipendenze parziali)

R	
r	

А	В	С
a1	b1	c1
a2	b1	c1
a3	b1	c2

sono veri i fatti (a1,b1,c1) e (a2,b1,c1) e (a3, b1, c2) e **non altri** 

In base ad una delle possibili decomposizioni dello schema, questa

istanza si decompone in

R1	А	В
$\pi_{R1}(r)$	a1	b1
	a2	b1
	а3	b1

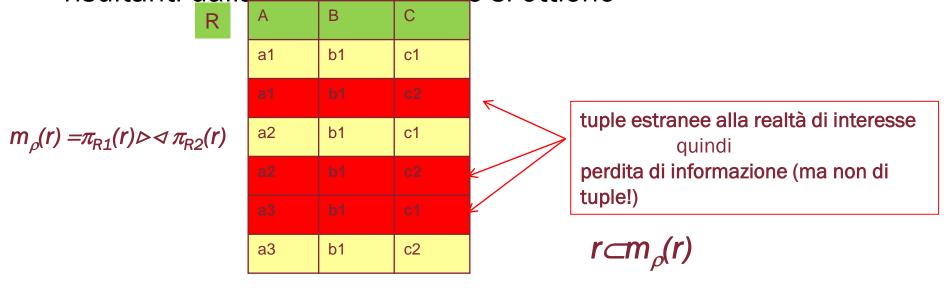
В	С			
b1	c1			
b1	c2			
	B b1			

La prima è ottenuta proiettando l'istanza originale su AB
La seconda è ottenuta proiettando l'istanza originale su BC (notare l'eliminazione del duplicato)
Forse questa eliminazione ci farà perdere tuple originali? NO

Dovrebbe essere possibile ricostruire l'istanza di partenza esattamente tramite join ... invece ...



• ... e invece se si effettua il join delle due istanze legali risultanti dalla decomposizione si ottiene



 Occorre garantire che il join delle istanze risultanti dalla decomposizione non riveli perdita di informazione

#### **Teorema**



## Insiemi di tuple

Prova di b).  $\pi_{Ri}(m_{\rho}(r)) = \pi_{Ri}(r)$ 

Procediamo provando la doppia inclusione

Per a) si ha  $r \subseteq m_{\rho}(r)$  e, quindi,  $\pi_{Ri}(r) \subseteq \pi_{Ri}(m_{\rho}(r))$ .

E' sufficiente, pertanto, mostrare che  $\pi_{Ri}(r) \supseteq \pi_{Ri}(m_{\rho}(r))$ .

Banalmente, per ogni tupla  $t \in m_{\rho}(r)$  e per ogni i, i=1,...,k, deve esistere una tupla  $t' \in r$  tale  $t[R_i] = t'[R_i]$ . Se tale tupla non esistesse, non troveremmo i valori di t su  $R_i$  ( $t[R_i]$ ) nella proiezione  $\pi_{Ri}(r)$  e di conseguenza non li avremmo nel join naturale

# **Digressione**



А	В	С
a1	b1	c1
a2	b1	c1
а3	b1	c2

R1	А	В
$\pi_{R1}(r)$	a1	b1
	a2	b1
	а3	b1

R2	В	С
$\pi_{R2}(r)$	b1	c1
	b1	c2

$m_{\rho}(r)$	$=\pi_{P^1}$	′r)⊳⊲	$\pi_{\rm P2}$	r
$\cdots \rho (\cdot)$	"RI	. /	"K2	• /

А	В	С
a1	b1	c1
a1		
a2	b1	c1
a2		
a3		
а3	b1	c2

	R	A	\	В
		а	11	b1
$\pi_{\!\scriptscriptstyle R1}\!(n)$	$\rho(r)$	а	11	b1
		а	2	b1
		a	2	b1
		· ·	3	b1
		а	3	b1

$\pi_{R2}(r$	$n_{\rho}(r)$

Б	
b1	c1
b1	c2
b1	c1
b1	c2
b1	c1
ALA!	25

duplicati = vengono eliminati dalla proiezione

#### **Teorema**



Prova di c).  $m_{\rho}(m_{\rho}(r)) = m_{\rho}(r)$ .

Per b) si ha  $\pi_{Ri}(m_{\rho}(r)) = \pi_{Ri}(r)$ .

Pertanto , applicando la definizione dell'operatore  $m_{\rho}$  (join delle proiezioni) avremo

$$m_{\rho}(m_{\rho}(r)) = \pi_{R1}(m_{\rho}(r)) \triangleright \triangleleft \pi_{R2}(m_{\rho}(r)) \triangleright \triangleleft \dots \triangleright \triangleleft \pi_{Rk}(m_{\rho}(r)) = \pi_{R1}(r) \triangleright \triangleleft \pi_{R2}(r) \triangleright \triangleleft \dots \triangleright \triangleleft \pi_{Rk}(r) = m_{\rho}(r).$$

# **Digressione**



$m_{o}(r)$	$=\pi_{R1}$	(r)⊳⊲	$\pi_{R2}(r)$

А	В	С
a1	b1	c1
a1	b1	c2
a2	b1	c1
a2	b1	c2
а3	b1	c1
а3	b1	c2

А	В	С
a1	b1	c1
a1	b1	c2
a2	b1	c1
a2	b1	c2
а3	b1	c1
а3	b1	c2

R A B a1 b1 a2 b1 a3 b1

R

 $\pi_{R2}(m_{\rho}(r)$ 

В	С
b1	c1
b1	c2

 $m_{\rho}(m_{\rho}(r)) = \pi_{R1}(m_{\rho}(r)) \triangleright \triangleleft \pi_{R2}(m_{\rho}(r))$ 

#### Che si fa?



- Abbiamo uno schema di relazione R, un insieme di dipendenze funzionali F e una decomposizione  $\rho$ .
- Come facciamo a <u>verificare</u> che la decomposizione data ha un join senza perdita?
- L'algoritmo che mostriamo permette di effettuare la verifica in tempo polinomiale.
- F ci serve perché ... ci serviremo di nuovo di una particolare istanza legale di R

## Algoritmo di verifica



## Algoritmo – verifica che una decomposizione abbia del join senza perdita

**Input** uno schema di relazione R, un insieme F di dipendenze funzionali su R, una decomposizione  $\rho = \{R_1, R_2, ..., R_k\}$  di R;

**Output** decide se  $\rho$  ha un join senza perdita;

#### begin

Costruisci una tabella *r* nel modo seguente:

*r* ha /R/ colonne e  $/\rho/$  righe

all'incrocio dell'i-esima riga e della j-esima colonna metti

il simbolo  $a_i$  se l'attributo  $A_i \in R_i$ 

il simbolo  $b_{ii}$  altrimenti

repeat

for every  $X \rightarrow Y \in F$ 

l'attributo  $A_j$  fa parte del sottoschema  $R_i$ 

Indice i = elemento della decomposizione = riga

una colonna per ogni attributo di R e una riga per

ogni elemento della decomposizione (sottoschema)

Indice j = attributo = colonna

**do if** ci sono due tuple  $t_1$  e  $t_2$  in r tali che  $t_1[X] = t_2[X]$  e  $t_1[Y] \neq t_2[Y]$  **then for every attribute**  $A_j$  in Y **do if**  $t_1[A_j] = a_j$  **then**  $t_2[A_j] := t_1[A_j]$ **else**  $t_1[A_i] := t_2[A_i]$  gestiamo correttamente anche il caso in cui  $t_2[A_i]=a_i$ 

until r ha una riga con tutte 'a' or r non è cambiato; if r ha una riga con tutte 'a' then  $\rho$  ha un join senza perdita else  $\rho$  non ha un join senza perdita

anche in questo caso
l'algoritmo termina <u>sempre!</u>
occorre poi verificare se in r
c'è la tupla che cerchiamo

### **Osservazioni**



- Possiamo considerare gli a<sub>j</sub> come valori particolari appartenenti al dominio dell'attributo A<sub>i</sub>
- Possiamo considerare i b<sub>ij</sub> come valori particolari appartenenti al dominio dell'attributo A<sub>i</sub>
- Possiamo considerare tutti i valori a<sub>j</sub> uguali tra di loro
- Il valore  $\mathbf{b}_{ij}$  è diverso da  $\mathbf{a}_{j}$  e da un altro valore  $\mathbf{b}_{kj}$  anche se appartengono tutti allo stesso dominio (quello dell'attributo  $\mathbf{A}_{i}$ )
- Come conseguenza, la nostra r iniziale è una particolare istanza dello schema R

#### Osservazioni



- Per i nostri scopi, <u>trasformiamo</u> questa istanza iniziale in una <u>istanza legale</u>
- Per fare questo, <u>dobbiamo fare in modo</u> che **tutte** le dipendenze in F siano **soddisfatte**
- Allora, ogni volta che troviamo due tuple uguali sugli attributi a sinistra, facciamo in modo che siano uguali anche su quelli a destra (attributo per attributo!)
- Nel fare questo, diamo la precedenza alle a (non diventano MAI b)

### **Osservazioni**



- Se due tuple sono uguali sulla parte sinistra ma sono diverse in un attributo della parte a destra della dipendenza, e una delle due ha una a come valore di quell'attributo, facciamo diventare a anche il valore dell'attributo nell'altra tupla
- Se due tuple sono uguali sulla parte sinistra ma sono diverse in un attributo della parte a destra della dipendenza, e nessuna delle due ha una a come valore di quell'attributo, facciamo diventare uguali le due b (ad esempio, se abbiamo b<sub>ij</sub> e b<sub>kj</sub> facciamo diventare entrambi i valori b<sub>ij</sub> oppure b<sub>kj</sub>)
- In pratica due valori saranno uguali se sono entrambi a oppure se hanno una b con lo stesso pedice
- L'algoritmo di ferma quando TUTTE le coppie di tuple soddisferanno la condizione che <u>se</u> gli attributi delle parti sinistre delle dipendenze sono uguali, <u>lo sono anche</u> gli attributi delle parti destre
- Ma questo significa che alla fine r è stata trasformata in un'istanza legale di R

#### **Teorema**



- Teorema Sia R uno schema di relazione, F un insieme di dipendenze funzionali su R e ρ={R₁, R₂, ..., Rκ} una decomposizione di R. L'Algoritmo di verifica decide correttamente se ρ ha un join senza perdita.
- Dim. Occorre dimostrare che:
- $\rho$  ha un join senza perdita ( $m_{\rho}(r) = r$  per ogni r legale)
  - se e solo se
- quando l'algoritmo termina la tabella r ha una tupla con tutte 'a'.

## **Teorema: dimostrazione**



#### Parte solo se.

•Supponiamo per assurdo che  $\rho$  abbia un join senza perdita  $(m_{\rho}(r)=r)$  e che quando l'algoritmo termina la tabella r non abbia una tupla con tutte 'a'. La tabella r può essere interpretata come un'istanza legale di R (basta sostituire ai simboli 'a' e 'b' valori presi dai domini dei corrispondenti attributi in modo tale che ad uno stesso simbolo venga sostituito lo stesso valore) in quanto l'algoritmo termina quando non ci sono più violazioni delle dipendenze in F. Poiché nessun simbolo 'a' che compare nella tabella costruita inizialmente viene mai modificato dall'algoritmo, per ogni  $i, i=1,...,k, \pi_{Ri}(r)$  contiene (fin dall'inizio!) una tupla con tutte 'a' (quella ottenuta proiettando l'istanza r sugli attributi di R<sub>i</sub> e precisamente nella riga corrispondente al sottoschema  $R_i$ ) pertanto  $m_o(r)$  contiene sicuramente una tupla con tutte 'a' e, quindi,  $m_{\rho}(r)\neq r$  (contraddizione).

### **Teorema: dimostrazione**



#### Parte se

- •Per uno sketch della prova della parte "se" consultare il testo J. D. Ullman, "Principles of database and knowledge-base systems", vol. I, Computer Science Press, 1988.
- •Una scansione delle pagine rilevanti verrà fornita insieme alle dispense del corso.

## Corollario



- Corollario Sia R uno schema di relazione, F un insieme di dipendenze funzionali su R e ρ={R₁, R₂} una decomposizione di R.
  - •Se  $R1 \cap R2 \rightarrow R1 R2$  o  $R1 \cap R2 \rightarrow R2 R1$  allora  $\rho$  ha un join senza perdita.

Osservazione banale:  $R_1$ =  $(R_1 - R_2) \cup (R_1 \cap R_2)$  e  $R_2$ =  $(R_2 - R_1) \cup (R_1 \cap R_2)$ 

	$R_1 - R_2$	$R_1 \cap R_2$	R <sub>2</sub> -R <sub>1</sub>
R <sub>1</sub>	<u>a</u>	<u>a</u>	<u>b</u> <sub>2-1</sub>
$R_2$	<u>b</u> <sub>1-2</sub>	<u>a</u>	<u>a</u>

Possiamo riordinare gli attributi in modo da avere un quadro più immediato della distribuzione dei valori nell'istanza di R

 $\underline{a}$ ,  $\underline{b}_{2-1}$ , e  $\underline{b}_{1-2}$  sono vettori di valori del tipo  $a_j$  o  $b_{ij}$ In base a  $R1 \cap R2 \rightarrow R1 - R2$  o  $R1 \cap R2 \rightarrow R2 - R1$  applicando l'algoritmo una delle due tuple (rispettivamente la seconda o la prima) diventerà di tutte a

Non è necessario che  $R1 \cap R2 \rightarrow R1 - R2$  o  $R1 \cap R2 \rightarrow R2 - R1$  siano in F, è sufficiente che siano in  $F^+$  ...