

# Basi di Dati 1

Marco Casu



# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduzione</b>                           | <b>3</b>  |
| <b>2</b> | <b>Il Modello Relazionale</b>                 | <b>5</b>  |
| 2.1      | Notazione con Indice . . . . .                | 5         |
| 2.2      | Rappresentazione come Funzioni . . . . .      | 6         |
| 2.3      | Integrità dei Dati . . . . .                  | 7         |
| 2.4      | Le Chiavi . . . . .                           | 8         |
| 2.4.1    | La Chiave Esterna . . . . .                   | 9         |
| 2.5      | Le Dipendenze Funzionali . . . . .            | 9         |
| <b>3</b> | <b>Algebra relazionale</b>                    | <b>10</b> |
| 3.0.1    | Proiezione . . . . .                          | 10        |
| 3.0.2    | Selezione . . . . .                           | 10        |
| 3.0.3    | Unione . . . . .                              | 11        |
| 3.0.4    | Differenza . . . . .                          | 11        |
| 3.0.5    | Intersezione . . . . .                        | 12        |
| 3.0.6    | Prodotto Cartesiano . . . . .                 | 12        |
| 3.0.7    | Join . . . . .                                | 13        |
| 3.0.8    | $\Theta$ -Join . . . . .                      | 13        |
| 3.0.9    | Condizioni Negative . . . . .                 | 13        |
| 3.1      | Quantificazione Universale . . . . .          | 14        |
| 3.2      | Esempi di Esercizi . . . . .                  | 14        |
| 3.2.1    | Esempio 1 . . . . .                           | 14        |
| 3.2.2    | Esempio 2 . . . . .                           | 15        |
| <b>4</b> | <b>Design di un Database</b>                  | <b>16</b> |
| 4.1      | Definizioni Formali . . . . .                 | 17        |
| 4.2      | Gli Assiomi di Armstrong . . . . .            | 18        |
| 4.2.1    | Chiusura di un Insieme di Attributi . . . . . | 19        |
| 4.3      | Teorema Fondamentale : $F^+ = F^A$ . . . . .  | 19        |
| 4.3.1    | Dimostrazione . . . . .                       | 19        |
| 4.4      | La Terza Forma Normale . . . . .              | 21        |
| 4.4.1    | Dipendenze Parziali e Transitive . . . . .    | 21        |
| 4.4.2    | Decomposizione in più Schemi . . . . .        | 22        |
| 4.4.3    | Forma Normale Boyce-Codd . . . . .            | 22        |
| 4.5      | Calcolo per la Chiusura di $X$ . . . . .      | 23        |
| 4.5.1    | L'Algoritmo . . . . .                         | 23        |
| 4.5.2    | La Ricerca delle Chiavi . . . . .             | 25        |

# 1 Introduzione

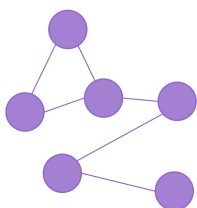
L'informazione memorizzata nei sistemi elettronici può essere di due tipi, **strutturata** e **non strutturata**. In questo corso ci occuperemo dell'informazione strutturata, ossia composta da oggetti matematici ben definiti. Un **sistema informativo** connette e contiene le informazioni, alla quale si può accedere da diversi componenti. Si prenda come esempio di sistema informativo l'archivio fisico dei documenti (ossia i dati) di una azienda, ad esso, possono accedere i vari reparti, come la sezione commercio o le risorse umane, ognuno ha a disposizione l'accesso ad un determinato sotto-insieme di dati (permessi). Prima di adoperare i sistemi informativi, ogni reparto aveva il suo personale archivio, ciò faceva sì che molti dei dati fossero duplicati e presenti per più reparti, causando un'elevata **ridondanza** di dati. Inoltre, due dati da 2 archivi diversi potrebbero dipendere tra loro. È quindi importante mantenere l'informazione **centralizzata**. I dati vanno organizzati, gestiti, e regolamentati da permessi di accesso a secondo dell'utente che vuole accedervi.

Un sistema informativo è composto dai seguenti componenti :

- Database (DB)
- Database Management System (DBMS) - Ossia il software per il mantenimento dei dati
- Application Software
- Computer Hardware - La memoria nella quale è contenuto

È importante mantenere una visione astratta del modello di sistema informativo, che sia indipendente dall'hardware in uso. Fissiamo un modello che utilizzi tipi di dati che non possano variare nel tempo, strutturati in maniera completa, è necessario pensare ad un metodo formale per organizzare i dati. Vogliamo farlo in un ambiente condiviso, preoccupandoci quindi di quali e quanti dati vengono condivisi con i singoli utenti, che avranno permessi diversi e potranno accedere a dati diversi (ad esempio, uno studente può accedere ai suoi esami registrati, ma non a quelli degli altri). Inoltre bisogna anche amministrare i metodi con cui si accede ai dati, preoccupandoci della concorrenza (accedere allo stesso dato nello stesso momento). I dati vengono organizzati in maniera omogenea, esiste un tipo di dato per descrivere un'entità (ad esempio uno studente), e tutte le volte che si vuole immagazzinare nel sistema uno studente, bisogna ricorrere allo stesso tipo di dato. L'informazione viene rappresentata da un aggregamento di più dati *grezzi*, ossia, rispettivamente la stringa "Maurizio Ernesti" e l'intero "3761523746", insieme rappresentano l'informazione di un professore (nome, cognome e numero di telefono).

Il modello è il modo in cui decidiamo di organizzare e collegare i dati, esistono modelli **logici**, indipendenti dalla struttura fisica, come il **modello relazionale**, ed esistono modelli **concettuali**, ossia rappresentazioni ancora più astratte indipendenti dal modello logico, i primi modelli sono stati introdotti negli anni 60.



un esempio è il modello *Mesh*, rappresentato con un grafo, dove i nodi sono i dati (i record), e gli archi le loro relazioni.

In questo modello i collegamenti sono esplicitati fisicamente (possiamo immaginare con dei puntatori), i modelli relazionali, diversamente, hanno relazioni rappresentate implicitamente dai valori stessi che contengono.

| Studenti  |         |        |            | Esami    |      |       |
|-----------|---------|--------|------------|----------|------|-------|
| Matricola | Cognome | Nome   | Compleanno | Studente | Voto | Corso |
| 276545    | Gialli  | Lucia  | 25/11/1980 | 348191   | 25   | 01    |
| 176515    | Rossi   | Mario  | 13/09/1982 | 176515   | 18   | 02    |
| 348191    | Verdi   | Andrea | 04/07/1981 | 176515   | 27   | 02    |

La relazione tra gli esami e gli studenti (ogni studente ha  $n$  esami registrati) è data implicitamente dalla presenza del campo *Studente* nella tabella *Esami*, che equivale al campo *Matricola* della tabella *Studenti*, il modello relazionale è basato su oggetti, classi ed attributi. Ad esempio, se volessi sapere che voto ha preso *Mario Rossi* al corso numero 02, mi basterebbe consultare la tabella *Esami*, e controllare il campo che ha la matricola equivalente a quella di *Mario Rossi*.

Lo schema logico descrive la presenza di tutte le entità con i loro rispettivi attributi, esistono poi gli schemi esterni, ossia sotto-insiemi degli schemi, destinati a determinati tipi di utenti che ne hanno l'accesso. Anche lo schema logico completo può essere schema esterno, il super-amministratore di un sistema informativo, come schema esterno, avrà l'intero schema logico, avendo accesso a tutte le entità.

Come già detto, lo schema logico rappresenta la struttura degli oggetti/entità, ed è invariata nel tempo, ma la sua istanza, ossia gli effettivi campi delle tabelle, possono variare, ed anche rapidamente. Per gestire gli schemi si utilizzano dei veri e propri linguaggi.

- **Data Definition Language (DDL)** - Per la definizione degli schemi logici ed altre operazioni generali.
- **Data Manipulation Language (DML)** - Per interrogare lo schema logico, leggerne i valori ed eventualmente modificarli.

È largamente utilizzato il linguaggio **SQL (Structured Query Language)**, che funge sia da DDL che da DML.

Un base di dati deve essere :

- Manipolabile
- Modificabile
- Centralizzata
- Il minimo ridondante
- Sicura

I dati molto spesso devono soddisfare certi vincoli (ad esempio, ogni studente ha una sola residenza) , tali vincoli son chiamati **dipendenze funzionali**, e possono riguardare anche il dominio di certi attributi (ad esempio, il voto verbalizzato di un esame deve essere maggiore o uguale a 18 e minore o uguale a 30). I dati devono essere protetti da accessi non autorizzati, è necessaria la dichiarazione di regole di accesso.

Definiamo adesso uno specifico tipo di operazione sulle basi di dati, ossia le **transizioni**, che non sono altr che sequenze ordinate di operazioni che vanno obbligatoriamente eseguite

insieme in sequenza, con il divieto assoluto che ne venga eseguita solo una parte. Facciamo un esempio, si dia il caso che su una base di dati bancaria, si vogliano trasferire 1000 euro dal conto *C1* al conto *C2*, le seguenti operazioni definite informalmente sono :

- Cercare il conto *C1*
- Sottrarne al bilancio 1000
- Cercare il conto *C2*
- Aggiungerne al bilancio 1000

Se per errore si eseguono solo i primi 2 passi, ci si ritroveranno 1000 euro persi, sottratti al primo conto, ma non addizionati al secondo. Una transizione va quindi completamente eseguita, se non dovesse essere così, l'intera sequenza di operazioni va abortita. Inoltre, è importante notare che data la concorrenzialità, in una base di dati potrebbe accadere di accedere allo stesso dato nello stesso momento, ciò potrebbe causare errori, è quindi importante evitare di lavorare contemporaneamente su uno stesso specifico campo.

## 2 Il Modello Relazionale

Il modello relazionale è basato sulla relazione intesa in senso matematico, una relazione non è altro che un insieme di tuple, tutte della stessa lunghezza, con elementi appartenenti a diversi domini. Il dominio è l'insieme dei possibili valori che gli elementi delle tuple possono assumere. Se prendiamo una lista di  $k$  domini, il prodotto cartesiano di tutti i  $k$  domini è l'insieme di tuple di lunghezza  $k$ .

$$D_1 \times D_2 \times D_3 \dots \times D_k = \{(v_1, v_2, v_3, \dots, v_k) | v_1 \in D_1, v_2 \in D_2, v_3 \in D_3, \dots, v_k \in D_k\} \quad (1)$$

Per esempio, per  $k = 2$  consideriamo i seguenti domini  $D_1 = \{White, Black\}$  e  $D_2 = \{0, 1, 2\}$ , si ha che :

$$D_1 \times D_2 = \{(White, 0), (White, 1), (White, 2), (Black, 0), (Black, 1), (Black, 2)\} \quad (2)$$

È una relazione di grado 2, perchè ogni tupla ha 2 coordinate.

**Teorema 1** *Il grado di una relazione è equivalente al numero di elementi di ogni tupla appartenente a tale relazione.*

In tale prodotto cartesiano è possibile costruire  $2^6$  possibili relazioni, ossia sotto-insiemi del prodotto cartesiano.

**Teorema 2** *Una relazione è un qualsiasi sottoinsieme del prodotto cartesiano.*

### 2.1 Notazione con Indice

Sia  $r$  una relazione di grado  $k$ , data  $t$  una tupla appartenente alla relazione  $r$ , ed  $i$  un intero da 1 a  $k$ , con  $t[i]$  si intende l'elemento alla coordinata  $i$ -esima della tupla  $t$ .

| Tabella |   |
|---------|---|
| Black   | 0 |
| White   | 0 |

Se prendiamo  $t$  come la prima riga della tabella, si avrà che  $t[1] = \text{"Black"}$  e  $t[2] = 0$ . Ricordando la notazione tabellare, in basi di dati le intestazioni delle colonne ed il loro dominio è denominato **attributo** (ad Esempio  $Name : string$ ). In una tabella, due attributi distinti non possono avere lo stesso nome.

## 2.2 Rappresentazione come Funzioni

Sia  $R$  un oggetto definito come insieme di attributi, una tupla su  $R$ , ossia un istanza di tale oggetto, può essere visto come una funzione definita su  $R$  che associa ad ogni attributo  $A$ , un valore presente nel dominio di  $A$ . Considerando ciò, presa  $t$  una tupla di  $R$ , ed  $A$  uno dei suoi attributi, indichiamo con  $t(A)$ , il valore (ossia l'istanza) di quell'attributo preso dalla funzione  $t$  sulla variabile  $A$ .

Ad esempio, la relazione  $R$  ha la tupla  $t_1 = (Paolo, Rossi, 2, 26.5)$ , considerando  $t_1$  come una funzione, essa associa ad ogni attributo  $A$ , un elemento del suo dominio :

$$f : (Nome, Cognome, Esami, Media) \rightarrow string \cup string \cup int \cup real \quad (3)$$

Da qui si ha  $t_1(Nome) = Paolo$   $t_1(Cognome) = Rossi$   $t_1(Media) = 26.5$

Lo schema logico, è il dominio di tale funzione, l'istanza sono le tuple, le istanze sono insiemi di tuple, ossia una relazione.

**Teorema 3** *L'istanza è un sotto-insieme di tuple*

Ogni riga della tabella è una tupla distinta, ed ogni colonna corrisponde al dominio. Possiamo quindi rappresentare un oggetto della base di dati con la seguente notazione :

$$R(A_1, A_2, A_3, \dots, A_k) \quad (4)$$

Qui  $R$  è una relazione dello schema, quindi uno schema di basi di dati non è altro che un insieme di relazioni  $(R_1, R_2, R_3, \dots, R_k)$  (invarianti nel tempo) per le quali, ognuna di esse possiede un istanza (variante nel tempo). Da qui in poi, utilizzeremo  $R$  per denominare le relazioni, ed  $r$  per le loro istanze.

Riprendendo la notazione con indice vista in precedenza 2.1, si può utilizzare piuttosto che un indice intero, l'intestazione dell'attributo per il quale si voglia leggere l'istanza.

| Luoghi |           |             |
|--------|-----------|-------------|
| Città  | Regione   | Popolazione |
| Roma   | Lazio     | 3000000     |
| Milano | Lombardia | 1500000     |
| Genova | Liguria   | 150000      |

Data  $t_1$  la prima riga dell'istanza, si ha  $t_1[Città] = \text{"Roma"}$ . È possibile anche farlo con sottoinsiemi di attributi, ossia  $t_1[Regione, Popolazione] = (\text{"Lazio"}, 3000000)$ . Non è importante l'ordine degli attributi come argomenti, come "risultato" riceviamo una sotto-tupla della tupla  $t_1$ , detta *restrizione*. Abbiamo visto come le istanze delle tabelle non sono altro che insiemi di tuple, esistono tante possibili istanze quanto la cardinalità del prodotto delle cardinalità dei domini.

Può succedere in certi casi, che nell'istanza di una relazione, sia presente una riga in cui un attributo è **sconosciuto**, per rappresentare tale campo nelle basi di dati, si utilizzi il valore polimorfico<sup>1</sup> *NULL*, utilizzato per riempire gli spazi vuoti, ad esempio in una tabella contenente

---

<sup>1</sup>Appartente a tutti i domini

i dati degli utenti iscritti ad un sito, è possibile che alcuni utenti abbiano omesso il numero di telefono, per loro il campo avrà valore *NULL*, ossia sconosciuto. Si ricordi che *NULL* è diverso da 0.

| Luoghi    |         |             |
|-----------|---------|-------------|
| Matricola | Nome    | Cellulare   |
| 1039      | Luca    | 3475746371  |
| 4316      | Giorgio | <i>NULL</i> |
| 1499      | Sandro  | 3857482845  |

## 2.3 Integrità dei Dati

La presenza di un valore *NULL* può causare alcuni errori, vedremo come sono presenti alcuni attributi, le quali istanze devono per forza essere dichiarate e non sconosciute. Esistono però diversi tipi di errori, come dei campi identificativi duplicati o valori fuori dominio.

| Studenti  |         |       |
|-----------|---------|-------|
| Matricola | Nome    | Media |
| 1039      | Luca    | 28    |
| 4316      | Giorgio | 33    |
| 1039      | Sandro  | 22    |

- **Errore 1** - Nella prima e nella terza riga sono presenti due studenti con la stessa matricola. Il campo matricola identifica ogni singolo e distinto studente, e non può essere duplicato.
- **Errore 2** - Uno studente ha come media dei voti 33, è impossibile dato che i voti sono compresi tra 18 e 30, è quindi un valore fuori dominio.

Tali errori vengono definiti problemi di **integrità dei dati**, per mantenere tale integrità è necessario che le istanze delle relazioni soddisfino delle determinate proprietà dette **vincoli**. Vedremo che esistono :

- Vincoli di chiave
- Vincoli di dominio
- Vincoli funzionali
- Vincoli di esistenza

| Impiegati        |         |         |                |            |              |
|------------------|---------|---------|----------------|------------|--------------|
| Codice Impiegato | Nome    | Cognome | Ruolo          | Assunzione | Dipartimento |
| 01               | Luca    | Rossi   | Analista       | 1785       | 01           |
| 02               | Mario   | Verdi   | Amministratore | 1980       | 02           |
| 02               | Giorgio | Neri    | Ricercatore    | 1985       | 05           |

| Dipartimenti |                 |
|--------------|-----------------|
| Numero       | Nome            |
| 01           | Managment       |
| 02           | Amministrazione |

Vediamo come nello schema logico appena mostrato ci sono diversi vincoli da definire che non sono rispettati. Ad esempio, va definito il vincolo di dominio per cui il valore

"Assunzione"  $\geq 1980$ , ossia tale valore deve essere strettamente maggiore di una certa data (possibile data di nascita dell'azienda). Si noti che non si sta rispettando un vincolo di chiave, dato che il "Codice Impiegato" presente nella seconda e nella terza riga lo stesso valore, essendo esso l'attributo identificativo, non deve essere duplicato ("Codice Impiegato" *UNIQUE*). Un altro errore meno evidente, è che alla terza riga della tabella "Impiegati", è presente un campo dipartimento con codice 05, tale campo dovrebbe collegare quella riga con il suo rispettivo dipartimento presente su un'altra tabella, ma notiamo che nella tabella "Dipartimenti", non è presente alcuna riga con "Numero" identificativo 05. Un altro vincolo noto è il vincolo di esistenza, che impone ad un certo attributo di non accettare valore *NULL*.

Un'altra importante distinzione da fare tra vincoli è di suddividerli in :

- **Intra-relazionali** - Vincoli definiti e da soddisfare all'interno della singola tabella (Ad esempio, un vincolo di dominio per il quale un valore deve essere sufficientemente grande).
- **Inter-relazionali** - Vincoli soddisfatti dai collegamenti di più tabelle (Un chiaro esempio è il sopra-citato errore sulla tabella del dipartimento).

## 2.4 Le Chiavi

In un istanza  $r$  di una relazione  $R$ , per ogni tupla è necessario che vi sia un attributo (o un insieme di attributi)  $X$  che la identifichi e distingua dalle altre tuple (che quindi non sia mai duplicato). Tale attributo/insieme di attributi è detto **chiave**, un chiaro esempio può essere il campo "Matricola" all'interno di un'ipotetica tabella studenti. Vediamo una definizione più formale.

### Teorema 4

**Punto 1** - Per ogni istanza della relazione  $R$ , non esistono due tuple  $t_1, t_2$  che hanno gli stessi valori per tutti i singoli attributi, preso un insieme di attributi  $X$ , vale sempre  $t_1[X] \neq t_2[X]$ .

$$\text{sia } X \in R(A_1, A_2, \dots, A_k) | \forall t_1, t_2 \in r \text{ se } t_1[X] = t_2[X] \implies t_1 = t_2 \quad (5)$$

**Punto 2** - Inoltre, non esistono sotto-insiemi di  $X$  che soddisfino la condizione sopra-citata.

$$\forall X' \subseteq X, \exists t_1, t_2 \in r \text{ tale che } t_1[X'] = t_2[X'] \wedge t_1 \neq t_2 \quad (6)$$

Approfondendo il punto 2, se esiste una relazione  $R$  con chiave  $X = (A_1, A_2, A_3)$ , è chiaro che, nella sua istanza  $r$ , non esisteranno due righe con gli stessi valori assegnati  $X$ , quindi, preso il sotto-insieme  $X' = (A_1, A_2)$ , se esistono due tuple  $t_1, t_2$  tali che  $t_1[X'] = t_2[X']$ , per forza di cose esse non saranno la stessa tupla in quanto, per il punto 1, differiranno per l'attributo  $A_3$ . È importante per una relazione che abbia una chiave significativa basata sull'informazione che rappresenta (Una relazione che rappresenta degli studenti, non può avere come chiave il campo del nome, dato che potrebbero esserci 2 studenti con lo stesso nome, bensì si predilige la matricola), si sceglie quindi una **chiave primaria**, ovviamente con vincolo di esistenza.



### 2.4.1 La Chiave Esterna

Può esistere inoltre un attributto nelle relazioni detto **chiave esterna** o **foreign key**, essa identifica all'interno di una relazione, un attributo associato ad un'altra relazione (un'altra tabella), identificandone la chiave primaria.

| Studenti  |         |        |            | Esami    |      |       |
|-----------|---------|--------|------------|----------|------|-------|
| Matricola | Cognome | Nome   | Compleanno | Studente | Voto | Corso |
| 276545    | Gialli  | Lucia  | 25/11/1980 | 348191   | 25   | 01    |
| 176515    | Rossi   | Mario  | 13/09/1982 | 176515   | 18   | 02    |
| 348191    | Verdi   | Andrea | 04/07/1981 | 176515   | 27   | 02    |

Nella tabella "Esami", l'attributo "Studente" è una chiave esterna che identifica e collega tale tabella con la relazione "Studenti", tramite la sua chiave primaria "Matricola". Dato che ad ogni esame è associato uno studente che l'ha sostenuto. Il vincolo di integrità inter-relazionale precedentemente citato impone che per ogni valore presente su un attributo di una chiave esterna, esista il suo corrispondente campo con tale valore nella relazione alla quale fa riferimento. Si ricordi che tale vincolo non è violato dalla presenza di un valore *NULL*.

| Multe  |            |             |         | Ufficiali |           |         |
|--------|------------|-------------|---------|-----------|-----------|---------|
| Codice | Data       | Ufficiale   | Targa   | Codice    | Nome      | Cognome |
| 4312   | 01/12/1988 | 001         | AA123AA | 001       | Giancarlo | Pozzi   |
| 1351   | 12/04/1989 | <i>NULL</i> | KA194AR | 002       | Sara      | Tua     |
| 9572   | 10/10/1990 | 002         | ND193MF | 003       | Nicola    | Canti   |

Può quindi esistere un'istanza dove la chiave esterna ha valore *NULL*, conseguentemente non avrà nessun riferimento nella relazione alla quale è collegata.

## 2.5 Le Dipendenze Funzionali

Come si possono definire facilmente i vincoli di integrità dei dati<sup>2.3</sup> visti nei paragrafi precedenti? Nelle basi di dati si utilizzano le note **dipendenze funzionali**, ossia un insieme di attributi che dipende da un altro insieme di attributi all'interno dello stesso schema. Tale definizione può sembrare poco esplicativa, ma è di vitale importanza in questo paragrafo che lo studente abbia ben saldo in mente il concetto di dipendenza funzionale, in quanto centrale nel corso di *Basi di Dati*. Formalmente, una dipendenza funzionale stabilisce un collegamento semantico tra due distinti insiemi di attributi  $X$  e  $Y$  appartenenti allo stesso schema. Si scrive:

$$X \rightarrow Y \quad (7)$$

e si legge: "X **determina** Y", stabilendone dei vincoli di integrità.

**Teorema 5** Sia  $r$  un'istanza della relazione  $R$ , la dipendenza funzionale  $X \rightarrow Y$  è soddisfatta se:

- Sia  $X$  che  $Y$  sono due sotto-insieme distinti di  $R$ .
- Le tuple di  $r$  che sono identiche per  $X$ , sono anche identiche per  $Y$ .

$$\forall t_1, t_2 \in R \text{ se } t_1[X] = t_2[X] \implies t_1[Y] = t_2[Y] \quad (8)$$

### 3 Algebra relazionale

Tramite l'**algebra relazionale** possiamo fare interrogazioni alla nostra base di dati per ottenere informazioni su una porzione di essa, le *query* sono scritte in linguaggio SQL, e vengono tradotte nel linguaggio formale e procedurale dell'algebra relazionale. L'algebra relazionale fornisce degli operatori che lavorano sulle istanze delle nostre relazioni, vi sono 4 tipi di operatori :

- Operatori di rimozione sulle singole relazioni
- Operatori insiemistici
- Operatori che combinano tuple da relazioni diverse
- Operatore di rinomina

Vediamo nel dettaglio tutti gli operatori che abbiamo a disposizione :

#### 3.0.1 Proiezione

L'operatore di proiezione, indicato con  $\pi$ , esegue un *taglio verticale* su una tabella, selezionando un sotto-insieme degli attributi, creando una tabella come quella iniziale, ma con esclusivamente le istanze degli attributi selezionati.

$$\pi_{A_1, A_2, \dots, A_k}(R)$$

*Esempio :*

|       |                 |           |      |   |      |          |
|-------|-----------------|-----------|------|---|------|----------|
|       | <b>Customer</b> |           |      |   |      |          |
|       | Code            | Name      | Town |   |      |          |
| Sia : | 001             | Giancarlo | Roma | ho che $\pi_{Name}(Customer) =$ <table><tr><td>Roma</td></tr><tr><td>Cagliari</td></tr></table> | Roma | Cagliari |
|       | Roma            |           |      |   |      |          |
|       | Cagliari        |           |      |   |      |          |
| 002   | Sara            | Cagliari  |      |   |      |          |
| 003   | Nicola          | Roma      |      |   |      |          |

#### 3.0.2 Selezione

La selezione, indicata con  $\sigma$ , esegue un *taglio orizzontale* sulla tabella, ossia, seleziona tutte le righe che soddisfano un determinato vincolo  $C$ , tale vincolo è un'espressione booleana della forma  $A \Theta B$ , dove  $\Theta \in \{<, >, =, \leq, ge\}$ . La condizione ovviamente per esser valida, ha bisogno che  $A$  e  $B$  abbiano lo stesso dominio (non posso comparare un intero con una stringa).

$$\sigma_C(R)$$

*Esempio :*

| Sia : | Customer |           |          | ho che : |
|-------|----------|-----------|----------|----------|
|       | Code     | Name      | Town     |          |
|       | 001      | Giancarlo | Roma     |          |
|       | 002      | Sara      | Cagliari |          |
|       | 003      | Nicola    | Roma     |          |

|                                    |      |           |      |
|------------------------------------|------|-----------|------|
| $\sigma_{Town="Roma"}(Customer) =$ | Code | Name      | Town |
|                                    | 001  | Giancarlo | Roma |
|                                    | 003  | Nicola    | Roma |

Ovviamente posso utilizzare gli operatori logici per eseguire richieste multiple :

$$\sigma_{Town="Roma" \wedge Code=001}(Customer) =$$

| Code | Name      | Town |
|------|-----------|------|
| 001  | Giancarlo | Roma |

Vale la proprietà :  $\sigma_C(\sigma_{C'}(R)) = \sigma_{C \wedge C'}(R)$

### 3.0.3 Unione

L'operazione di unione fra due relazioni, denominata con il simbolo  $\cup$ , come nel senso insiemistico, crea una nuova istanza della relazione, contenente le tuple da entrambe le relazioni. Non si possono unire due relazioni qualsiasi, è necessario che esse siano **union-compatibili**, ossia che abbiano lo stesso numero di attributi, e che gli attributi corrispondenti abbiano lo stesso dominio.

|         |                   |           |              |        |             |
|---------|-------------------|-----------|--------------|--------|-------------|
| Siano : | <b>Insegnanti</b> |           | <b>Admin</b> |        | Si ha che : |
|         | Code              | Name      | Code         | Name   |             |
|         | 001               | Giancarlo | 001          | Luca   |             |
|         | 002               | Sara      | 004          | Andrea |             |
|         | 003               | Nicola    | 005          | Carlo  |             |

$Insegnanti \cup Admin =$

| Code | Name      |
|------|-----------|
| 001  | Giancarlo |
| 002  | Sara      |
| 003  | Nicola    |
| 004  | Andrea    |
| 005  | Carlo     |

Si noti come la riga 

|     |      |
|-----|------|
| 001 | Luca |
|-----|------|

 non appare nell'unione, dato che condivide la stessa chiave con la riga 

|     |           |
|-----|-----------|
| 001 | Giancarlo |
|-----|-----------|

, quindi, quando si uniscono due relazioni, bisogna fare attenzione che due righe da esse non condividano la stessa chiave, altrimenti si avrà una perdita di informazioni. Si noti che, se una delle due tabelle avesse presentato un attributo di troppo, esso sarebbe potuto essere stato rimosso con una proiezione, rendendole comunque union-compatibili.

### 3.0.4 Differenza

Come per l'unione, è necessario che le due relazioni sulla quale voglio applicare la differenza siano union-compatibili. Si indica con il simbolo  $-$ , ed il risultato fra due relazioni, conterrà le tuple del primo operando, che non stanno anche nel secondo operando.

|         |                   |           |              |      |             |
|---------|-------------------|-----------|--------------|------|-------------|
| Siano : | <b>Insegnanti</b> |           | <b>Admin</b> |      | Si ha che : |
|         | Code              | Name      | Code         | Name |             |
|         | 001               | Giancarlo | 002          | Sara |             |
|         | 002               | Sara      |              |      |             |

$$Insegnanti - Admin =$$

| Code | Name      |
|------|-----------|
| 001  | Giancarlo |

La differenza non è commutativa! Infatti, in questo caso  $Admin - Insegnanti = \emptyset$

### 3.0.5 Intersezione

Anch'essa richiede che i due operandi siano union-compatibili. Si indica con  $\cap$ , ed il risultato conterrà esclusivamente le tuple che fanno parte sia della prima che della seconda relazione. (si noti che  $A \cap B = A - (A - B)$ ).

| Siano : | Insegnanti |           | Si ha che : | Admin |        |
|---------|------------|-----------|-------------|-------|--------|
|         | Code       | Name      |             | Code  | Name   |
|         | 001        | Giancarlo |             | 002   | Sara   |
|         | 002        | Sara      |             | 004   | Andrea |

| $Insegnanti \cap Admin =$ |  | Code | Name |
|---------------------------|--|------|------|
|                           |  | 002  | Sara |

Vediamo adesso gli operatori della terza categoria, che creano **relazioni multiple** combinando tabelle eterogenee.

### 3.0.6 Prodotto Cartesiano

Il prodotto cartesiano è indicato con il simbolo  $\times$ , e rappresenta tutte le possibili combinazioni di righe fra gli elementi del primo operando con gli elementi del secondo, tale operazione risulta parecchio dispendiosa, e potrebbe associare elementi che non hanno nessuna correlazione sensata, ad esempio :

| Siano : | Customer |           | Si ha che : | Order |          |
|---------|----------|-----------|-------------|-------|----------|
|         | Code     | Name      |             | CO    | Customer |
|         | 001      | Giancarlo |             | AX00  | 001      |
|         | 002      | Sara      |             | AX01  | 002      |

| $Customer \times Order =$ |  | Code | Name      | CO   | Customer |
|---------------------------|--|------|-----------|------|----------|
|                           |  | 001  | Giancarlo | AX00 | 001      |
|                           |  | 001  | Giancarlo | AX01 | 002      |
|                           |  | 002  | Sara      | AX00 | 001      |
|                           |  | 002  | Sara      | AX01 | 002      |

Logicamente però, stiamo associando a dei clienti, degli ordini che non gli appartengono (il campo *Customer* di *Order* è la foreign key che identifica il campo *Code* della tabella *Customer*). Quindi, se volessimo una tabella con ogni cliente e gli ordini ad esso associati, combineremo anche una selezione del tipo :

| $\sigma_{Customer.Code=Order.Customer}(Customer \times Order) =$ |  | Code | Name      | CO   | Customer |
|--|--|------|-----------|------|----------|
|  |  | 001  | Giancarlo | AX00 | 001      |
|  |  | 002  | Sara      | AX01 | 002      |

È chiaro che, per tabelle con numerosi elementi, il prodotto cartesiano da computare risulta parecchio dispendioso, dato che, se  $R_1$  ha  $n$  elementi, ed  $R_2$  ha  $m$  elementi, il prodotto cartesiano  $R_1 \times R_2$  avrà  $n \cdot m$  elementi, e quando si opera su basi di dati con milioni di righe, tale operazione deve essere evitata quando possibile, per questo esiste un'operazione simile al prodotto cartesiano, che seleziona automaticamente le righe da combinare secondo parametri ben precisi.

### 3.0.7 Join

L'operatore join, indicato dal simbolo  $\bowtie$ , seleziona le tuple del prodotto cartesiano che soddisfano una precisa condizione  $C$ , ossia che gli attributi con lo stesso nome, debbano avere anche lo stesso valore, è quindi importante che si denominino in maniera precisa gli attributi di relazioni che si vogliono combinare con il join.

$$R_1 \bowtie R_2 = \sigma_C(R_1 \times R_2) \quad (9)$$

$$C = R_1.A_1 = R_2.A_1 \wedge R_1.A_2 = R_2.A_2 \wedge \dots \wedge R_1.A_k = R_2.A_k \quad (10)$$

dove  $A_1, A_2, \dots, A_k$  sono gli attributi condivisi con lo stesso nome fra  $R_1$  e  $R_2$ . La tabella risultante, non presenterà due volte gli attributi da  $R_1$  e  $R_2$  con lo stesso nome, ma li unirà in un attributo solo.

Se dovesse capitare che le due relazioni non hanno alcun attributo in comune, il risultato sarà il prodotto cartesiano.

Siano :

| Customer |           | Order |      |        |
|----------|-----------|-------|------|--------|
| Code     | Name      | CO    | Code | Item   |
| 001      | Giancarlo | AX00  | 001  | Glue   |
| 002      | Sara      | AX01  | 002  | Bricks |
| 003      | Lucia     | AX02  | 001  | Shoes  |

Si ha che :

$Customer \bowtie Order =$

| Customer.Code | Name      | CO   | Item   |
|---------------|-----------|------|--------|
| 001           | Giancarlo | AX00 | Glue   |
| 001           | Giancarlo | AX02 | Shoes  |
| 002           | Sara      | AX01 | Bricks |

### 3.0.8 $\Theta$ -Join

Il  $\Theta$ -Join (Che si pronuncia "theta join"), seleziona le tuple risultanti dal prodotto cartesiano, che soddisfino la condizione  $A\Theta B$  dove  $\Theta \in \{<, >, =, \leq, ge\}$  ed  $A, B$  sono attributi rispettivamente della prima, e seconda relazione (il dominio di  $A$  deve essere lo stesso del dominio di  $B$ ).

$$R_1 \bowtie_{A\Theta B} R_2 = \sigma_{A\Theta B}(R_1 \times R_2) \quad (11)$$

### 3.0.9 Condizioni Negative

È possibile utilizzare il simbolo  $\neg C$  per intendere che si vogliono selezionare tutte le tuple che non soddisfino la condizione  $C$ .

Sia :

| Customer |           |          |
|----------|-----------|----------|
| Code     | Name      | Town     |
| 001      | Giancarlo | Roma     |
| 002      | Sara      | Cagliari |
| 003      | Nicola    | Roma     |
| 004      | Gianluca  | Venezia  |

ho che :

$\sigma_{\neg(Town="Roma")}(Customer) =$

|     |          |          |
|-----|----------|----------|
| 002 | Sara     | Cagliari |
| 004 | Gianluca | Venezia  |

### 3.1 Quantificazione Universale

Per ora tutti gli operatori che abbiamo visto, implicano una quantificazione esistenziale, quando seleziono da una relazione le righe che soddisfano una condizione  $C$ , seleziona le righe anche se vi è un altro elemento che non soddisfa tale condizione.

Se seleziono tutti i clienti che hanno ordini con prezzo superiore a 100 euro, selezionerà tutti i clienti che hanno **almeno** un ordine superiore a 100 euro, non solo quelli che hanno **esclusivamente** ordini superiori a 100 euro.

Vogliamo scrivere query che implicino la quantificazione universale e non esistenziale, ossia, che tale condizione valga **"per ogni"**, e ciò è equivalente a **"non c'è ne sono tali che"**:

- I clienti che hanno solo ordini superiori a 100 euro = I clienti che non hanno ordini inferiori o uguali a 100 euro

Siano :

| Customer |           | Order |      |       |
|----------|-----------|-------|------|-------|
| Code     | Name      | CO    | Code | Price |
| 001      | Giancarlo | AX00  | 001  | 90    |
| 002      | Sara      | AX01  | 002  | 120   |
| 003      | Lucia     | AX02  | 001  | 200   |

Vogliamo elencare nome e codice dei clienti che hanno fatto ordini esclusivamente con prezzo superiore a 100 euro, quindi, *Giancarlo* non sarà incluso, dato che ha fatto sì un ordine da 200 euro, ma anche uno da 90 euro. La nostra query  $Q$  sarà :

$$Q = Customer - \pi_{Code, Name}(\sigma_{Price \leq 100}(Customer \bowtie Order)) \quad (12)$$

Quindi, prendere tutti coloro che hanno fatto solo ordini superiori a 100 euro, equivale a prendere il totale, e sottrarne coloro che hanno fatto almeno un ordine inferiore o uguale a 100 euro.

$$Q =$$

| Code | Name |
|------|------|
| 002  | Sara |

### 3.2 Esempi di Esercizi

#### 3.2.1 Esempio 1

Ci sono situazioni in cui dobbiamo combinare una relazione con se stessa per ottenere paia di tuple della stessa tabella, vediamo un esempio, si osservino le seguenti relazioni :

| Impiegati |     |         |        |
|-----------|-----|---------|--------|
| Name      | Cod | Salario | CodSup |
| Rossi     | C1  | 100     | C3     |
| Pirlo     | C2  | 200     | C3     |
| Bianchi   | C3  | 500     | NULL   |
| Verdi     | C4  | 200     | C2     |
| Neri      | C5  | 150     | C1     |
| Tosi      | C6  | 100     | C1     |

Dove  $CodSup$  identifica il  $Cod$  della riga che ogni impiegato ha come supervisore. Ad esempio, Rossi ha come supervisore Bianchi, e Neri ha come supervisore Rossi. Vogliamo trovare tutti gli impiegati che hanno lo stipendio superiore o uguale al loro supervisore. Per procedere possiamo combinare tramite il prodotto cartesiano la tabella con se stessa, prima però rinominando ogni attributo aggiungendoci una  $C$  davanti, per identificare una tabella che è copia dell'altra:

$$ImpiegatiC = \rho_{Name,Cod,Salario,CodSup \rightarrow CName,CCod,CSalario,CCodSup}(Impiegati) \quad (13)$$

Adesso procediamo col combinare tali tabelle, però selezionando esclusivamente le tuple che hanno il  $CodSup$  ed il  $CCod$  identici, in modo che avremo una lista dei dipendenti con a destra il loro superiore.

$$\sigma_{CodSup=CCod}(Impiegati \times ImpiegatiC) \quad (14)$$

| Name  | Cod | Salario | CodSup | CName   | CCod | CSalario | CCodSup |
|-------|-----|---------|--------|---------|------|----------|---------|
| Rossi | C1  | 100     | C3     | Bianchi | C3   | 500      | NULL    |
| Pirlo | C2  | 200     | C3     | Bianchi | C3   | 500      | NULL    |
| Verdi | C4  | 200     | C2     | Pirlo   | C2   | 200      | C3      |
| Neri  | C5  | 150     | C1     | Rossi   | C1   | 100      | C3      |
| Tosi  | C6  | 100     | C1     | Rossi   | C1   | 100      | C3      |

Fatto ciò adesso, ci basta selezionare quelli che hanno il salario superiore o uguale a quello del proprio supervisore, ossia  $Salario \geq CSalario$  :

$$\sigma_{(CodSup=CCod) \wedge (Salario \geq CSalario)}(Impiegati \times ImpiegatiC) \quad (15)$$

| Name  | Cod | Salario | CodSup | CName | CCod | CSalario | CCodSup |
|-------|-----|---------|--------|-------|------|----------|---------|
| Verdi | C4  | 200     | C2     | Pirlo | C2   | 200      | C3      |
| Neri  | C5  | 150     | C1     | Rossi | C1   | 100      | C3      |
| Tosi  | C6  | 100     | C1     | Rossi | C1   | 100      | C3      |

### 3.2.2 Esempio 2

Vediamo adesso un esempio di un altro tipo, si consideri sempre la stessa relazione di prima :

| Imp     |     |         |        |
|---------|-----|---------|--------|
| Name    | Cod | Salario | CodSup |
| Rossi   | C1  | 100     | C3     |
| Pirlo   | C2  | 200     | C3     |
| Bianchi | C3  | 500     | NULL   |
| Verdi   | C4  | 200     | C2     |
| Neri    | C5  | 150     | C1     |
| Tosi    | C6  | 100     | C1     |

Adesso vogliamo trovare l'impiegato che il salario più alto, ma come possiamo fare? Procederemo con il comparare ogni impiegato con tutti gli altri, selezionando esclusivamente quelli che hanno lo stipendio inferiore all'impiegato con la quale sono stati comparati (selezioniamo esclusivamente il codice impiegato):

$$Imp2 = Imp \quad (16)$$

$$\pi_{Imp.cod}(Imp \bowtie_{Imp.Salario < Imp2.Salario} Imp2) \quad (17)$$

Ottengo la tabella :

| Impiegati.Cod |
|---------------|
| C1            |
| C2            |
| C4            |
| C5            |
| C6            |

Che rappresenta tutti gli impiegati che hanno qualcuno con lo stipendio superiore al loro. Quindi, chi non è presente in questa tabella, sarà l'impiegato con lo stipendio più alto. Prendiamo allora tutti gli impiegati e sottraiamo ad essi la nostra tabella.

$$\pi_{Cod}(Imp) - \pi_{Imp.cod}(Imp \bowtie_{Imp.Salario < Imp2.Salario} Imp2) =$$

$$= \begin{array}{|c|} \hline C1 \\ \hline C2 \\ \hline C3 \\ \hline C4 \\ \hline C5 \\ \hline C6 \\ \hline \end{array} - \begin{array}{|c|} \hline C1 \\ \hline C2 \\ \hline C4 \\ \hline C5 \\ \hline C6 \\ \hline \end{array} = \begin{array}{|c|} \hline C3 \\ \hline \end{array}$$

Quindi l'impiegato di codice *C3*, ossia 

|         |    |     |      |
|---------|----|-----|------|
| Bianchi | C3 | 500 | NULL |
|---------|----|-----|------|

 è colui con lo stipendio più alto.

## 4 Design di un Database

L'obiettivo è capire come creare uno schema in maniera corretta, immaginiamo di dover creare una base di dati per memorizzare le informazioni relative agli studenti di un corso di laurea triennale, ed i relativi esami sostenuti.

| Matricola | SurN    | Name  | BirthD | City  | Prov | ExCode | ExName    | Doc   | Date | Grade |
|-----------|---------|-------|--------|-------|------|--------|-----------|-------|------|-------|
| 01        | Rossi   | Mario | ...    | Roma  | Roma | 10     | Physics   | Pippo | ...  | 28    |
| 02        | Bianchi | Paolo | ...    | Tolfa | Roma | 10     | Physics   | Pippo | ...  | 26    |
| 01        | Rossi   | Mario | ...    | Roma  | Roma | 20     | Chemistry | Pluto | ...  | 27    |

È estremamente sbagliato salvare tutti i dati in una sola relazione, per ogni esame, devo salvare ogni volta tutti i dati di uno studente, creando ridondanza e spreco di memoria, inoltre se uno studente non ha sostenuto alcun esame, non apparirà nell'archivio, ed un esame che non è stato sostenuto da nessuno risulterà inesistente. È corretto suddividere tale relazione in 3 diverse tabelle :

| Studenti  |         |       |        |       |      | Corso  |           |       |
|-----------|---------|-------|--------|-------|------|--------|-----------|-------|
| Matricola | SurN    | Name  | BirthD | City  | Prov | ExCode | ExName    | Doc   |
| 01        | Rossi   | Mario | ...    | Roma  | Roma | 10     | Physics   | Pippo |
| 02        | Bianchi | Paolo | ...    | Tolfa | Roma | 20     | Chemistry | Pluto |
| 01        | Rossi   | Mario | ...    | Roma  | Roma |        |           |       |



| Esami     |        |      |       |
|-----------|--------|------|-------|
| Matricola | ExCode | Data | Grade |
| 01        | 10     | ...  | 28    |
| 01        | 20     | ...  | 27    |
| 02        | 10     | ...  | 26    |

## 4.1 Definizioni Formali

Uno **schema relazionale**, che si denota con  $R$  è un insieme di attributi del tipo:

$$R = \{A_1, A_2, A_3, \dots, A_k\}$$

Solitamente un sotto insieme di attributi di  $R$  viene indicato con le lettere  $X$  oppure  $Y$ , e l'unione si può denotare  $XY \equiv X \cup Y$ . Data una relazione  $R$ , una **tupla**  $r$  è una funzione che associa ad ogni attributo, un valore appartenente al suo dominio :

$$R = \{nome, cognome\} \quad r[nome] = marco, r[cognome] = rossi$$

Preso  $X \subset R$ , diciamo che due tuple  $r_1, r_2$  **coincidono** su  $X$  se :

$$\forall A \in X, r_1[A] = r_2[A]$$

| Corso |         |      |
|-------|---------|------|
| Nome  | Cognome | Voto |
| Paolo | Rossi   | 29   |
| Mario | Rossi   | 29   |

Le due tuple coincidono su  $\{\text{Cognome}, \text{Voto}\}$

Un **istanza** su uno schema relazionale  $R$  è l'insieme di tutte le tuple su  $R$ . Una **dipendenza funzionale** su  $R$  è una coppia ordinata  $X, Y$  di sotto-insiemi di  $R$ , ossia di attributi di  $R$ , che si denota con  $X \rightarrow Y$ , ed indica che :

$$\text{per ogni coppia di tuple } r_1, r_2, \text{ vale che } r_1[X] = r_2[X] \implies r_1[Y] = r_2[Y]$$

E si dice  $X$  *determina*  $Y$ , dove  $X$  è detto determinante ed  $Y$  dipendente. Le dipendenze funzionali esprimono dei *vincoli*. Uno schema  $R$  può avere un **insieme di dipendenze funzionali** denotato con  $F$  :

$$F = \{A \rightarrow B, D \rightarrow B, \dots, G \rightarrow H\}$$

se un istanza soddisfa ogni dipendenza di  $F$ , ossia tutte le dipendenze funzionali, è detta istanza **legale**. Le dipendenze godono di una relazione di *transitività*, ossia, se  $A \rightarrow B$  e  $B \rightarrow C$ , allora  $A \rightarrow C$ , non è però necessario aggiungere quest'ultima all'insieme  $F$ , dato che è implicita, e vogliamo mantenere  $F$  il più piccolo possibile. Ci sono quindi dipendenze funzionali che sono soddisfatte per ogni istanza della relazione, che però non è necessario includere in  $F$ , un altro esempio :

$$\begin{aligned} &\text{se } taxCode \rightarrow name, surname \\ &\text{è ovvio che } taxCode \rightarrow name \text{ e } taxCode \rightarrow surname \end{aligned}$$

Ma queste due ultime non saranno incluse in  $F$ . Tali dipendenze sono dette **banali**, ad esempio, se  $Y \subset X$ , se  $A \rightarrow X$  ovviamente  $A \rightarrow Y$ . Il numero delle dipendenze banali è *esponenziale* :

$$X \rightarrow Y \iff \forall A \in Y, X \rightarrow A$$

Tutte le dipendenze banali che non è necessario includere in  $F$ , si trovano in un insieme più grande detto **chiusura di  $F$** , che si denota con  $F^+$ , contenente tutte le dipendenze funzionali soddisfatte da un'istanza, è chiaro che  $F \subseteq F^+$ .

Passiamo alla definizione di **chiave**. Un sotto-insieme di attributi  $K$  è detto chiave se :

- $K \rightarrow R \in F^+$
- $\nexists K' \subset K | K' \rightarrow R \in F^+$

Detto in maniera meno formale, non esisteranno due tuple di una relazione con gli stessi valori per le chiavi. Nei linguaggi come *SQL*, fra gli attributi della chiave se ne definisce uno in particolare detto **chiave primaria**, che non può assumere valore **NULL**.

## 4.2 Gli Assiomi di Armstrong

Dato uno schema  $R$ , un insieme di dipendenze funzionali  $F$ , ed un istanza legale  $r$ , se tale istanza soddisfa una dipendenza si scrive  $r \text{ SAT } X \rightarrow Y$ , abbiamo poi definito la chiusura di  $F$  come :

$$F^+ = \{X \rightarrow Y | \forall r \text{ legale } r \text{ SAT } X \rightarrow Y\}$$

Se voglio verificare che  $K$  sia una chiave, non mi è possibile eseguire manualmente il controllo su tutto  $F^+$ , che abbiamo visto essere di dimensioni esponenziali rispetto ad  $F$ , inoltre ci è impossibile controllare ogni istanza, dato che esse sono potenzialmente infinite. Definiamo adesso un nuovo insieme, a partire da  $F$ , che ci permette di esaminare lo schema relazionale più facilmente.

$F^A$  è un insieme di dipendenze funzionali su  $R$ , ottenuto dai seguenti *assiomi di Armstrong* :

- Se  $X \rightarrow Y \in F$  allora  $X \rightarrow Y \in F^A$
- Se  $Y \subseteq X \in R$  allora  $X \rightarrow Y \in F^A$  (**Riflessività**)
- Se  $X \rightarrow Y \in F^A$  allora  $XZ \rightarrow YZ \in F^A \forall Z \in R$  (**Aumento**)
- Se  $X \rightarrow Y \in F^A$  e Se  $Y \rightarrow Z \in F^A$  allora Se  $X \rightarrow Z \in F^A$  (**Transitività**)

Vediamo alcune proprietà che seguono dagli assiomi :

- $X \rightarrow Y \in F^A \wedge X \rightarrow Z \in F^A \implies X \rightarrow YZ \in F^A$  (Unione)
- $X \rightarrow Y \in F^A \wedge Z \subseteq Y \implies X \rightarrow Z \in F^A$  (Decomposizione)
- $X \rightarrow Y \in F^A \wedge WY \rightarrow Z \in F^A \implies XW \rightarrow Z \in F^A$  (Pseudotransitività)

*Dimostrazioni* delle proprietà : (Unione) Se ho  $X \rightarrow Y \in F^A$ , applicando *l'aumento*, utilizzando  $X$ , ottengo  $X \rightarrow XY \in F^A$ , Avendo poi  $X \rightarrow Z \in F^A$ , ed applicando *l'aumento* con  $Y$ , si ottiene  $XY \rightarrow YZ \in F^A$ . Ottengo quindi :

$$\begin{cases} X \rightarrow XY \in F^A \\ XY \rightarrow YZ \in F^A \end{cases} \implies \boxed{\text{per transitività}} \implies X \rightarrow YZ \in F^A \quad (18)$$

(Decomposizione) Se  $Z \subseteq Y$ , allora per *riflessività*, si ha  $Y \rightarrow Z \in F^A$ , dal momento che  $X \rightarrow Y \in F^A$ , si ha :

$$\begin{cases} X \rightarrow Y \in F^A \\ Y \rightarrow Z \in F^A \end{cases} \implies \boxed{\text{per transitività}} \implies X \rightarrow Z \in F^A \quad (19)$$

(Pseudotransitività) Se  $X \rightarrow Y \in F^A$ . applicando *l'aumento* con  $W$  si ottiene  $WX \rightarrow WY \in F^A$ , e ne consegue che :

$$\begin{cases} WX \rightarrow WY \in F^A \\ WY \rightarrow Z \in F^A \end{cases} \implies \boxed{\text{per transitività}} \implies WX \rightarrow Z \in F^A \quad (20)$$

*Osservazione* : Siano  $A_1, A_2, \dots, A_n$  un insieme di attributi, se  $\forall i, X \rightarrow A_i \in F^A$ , per unione,  $X \rightarrow A_1 A_2 \dots A_n \in F^A$ , analogamente, data tale condizione, per decomposizione :  $\forall i, X \rightarrow A_i \in F^A$ , è chiaro che : decomposizione  $\iff$  unione .

#### 4.2.1 Chiusura di un Insieme di Attributi

Sia  $R$  uno schema relazionale, ed  $F$  l'insieme delle sue dipendenze funzionali. Preso un insieme di attributi  $X \subseteq R$ , definiamo come **la chiusura di  $X$  su  $F$** , che denotiamo con  $X_F^+$ , l'insieme di tutti gli attributi *determinati* da  $X$  dopo aver applicato gli assiomi di Armstrong.

$$X_F^+ = \{A | X \rightarrow A \in F^A\}$$

Banalmente, per riflessività,  $X \subseteq X_F^+$ .

Vediamo un importante *Lemma* : Sia  $R$  uno schema relazionale, ed  $F$  l'insieme delle sue dipendenze funzionali, vale che :  $X \rightarrow Y \in F^A \iff Y \subseteq X_F^+$ .

*Dimostrazione* : Sia  $Y = A_1, A_2, \dots, A_n$  un insieme di attributi, se  $Y \subseteq X_F^+$ , allora, per unione,  $\forall i, X \rightarrow A_i \in F^A$ . D'altro canto, se  $X \rightarrow Y \in F^A$ , per decomposizione,  $\forall i, X \rightarrow A_i \in F^A$ , quindi  $\forall i, A_i \in F^A$ , ne consegue che  $Y \subseteq X_F^+$ .

### 4.3 Teorema Fondamentale : $F^+ = F^A$

Segue uno dei teoremi più importanti del corso :

$$F^+ = F^A$$

Partendo da  $F$ , applicando ricorsivamente gli assiomi di Armstrong, si ottiene la chiusura di  $F$ .

#### 4.3.1 Dimostrazione

Procederemo col dimostrare tale teorema per doppia inclusione, dimostrando quindi la veridicità di  $F^+ \subseteq F^A$  e  $F^A \subseteq F^+$ .

$\boxed{F^A \subseteq F^+}$ : Si procede per induzione, come parametro da incrementare, si adopererà il numero di volte in cui applico un determinato assioma di Armstrong. Ossia, dimostreremo che, se una dipendenza è in  $F^+$ , è perchè si è applicato  $i$  volte un assioma di Armstrong su  $F$ . Si consideri quindi la seguente notazione :  $F_i^A$  è l'insieme delle dipendenze funzionali ottenute applicando  $i$  volte un assioma di Armstrong. Ne consegue ovviamente che :

$$F_0^A \subseteq F_1^A \subseteq F_2^A \subseteq \dots F_n^A \subseteq F^A$$

Lo scopo è quindi dimostrare, che applicando tali assiomi un numero di volte finito, si ottiene  $F^+$ .

**Caso Base** ( $i = 0$ ) :  $F_0^A = F \subseteq F^+$

**Passo Induttivo** : La tesi è, che  $F_i^A \subseteq F^+ \implies F_{i+1}^A \subseteq F^+$ . Distinguiamo i 3 casi di applicazione dei 3 assiomi.

- *Caso 1*  $X \rightarrow Y$  è stato ottenuto in  $F^A$  applicando la riflessività su  $Y \subseteq X$ . Se  $r$  è un'istanza di  $R$ , e  $t_1, t_2$  due tuple tali che  $t_1[X] = t_2[X]$ , ovviamente vale che  $t_1[Y] = t_2[Y]$ , quindi  $X \rightarrow Y \in F^+$ .
- *Caso 2*  $X \rightarrow Y$  è stato ottenuto in  $F^A$  applicando l'aumento su una dipendenza funzionale  $V \rightarrow W \in F^A$ , che è stato ottenuto applicando gli assiomi  $i - 1$  volte. Per ipotesi induttiva,  $V \rightarrow W \in F^+$ , vorrebbe dire che  $X = VZ$  e  $Y = WZ$  per qualche  $Z \subseteq R$ . Se  $r$  è un'istanza legale, e  $t_1, t_2$  due tuple tali che  $t_1[X] = t_2[X]$ , ne consegue che  $t_1[V] = t_2[V]$  e  $t_1[Z] = t_2[Z]$ , ma per ipotesi induttiva, essendo  $V \rightarrow W \in F^A$ , ne consegue che  $t_1[W] = t_2[W]$ , e da ciò ne consegue  $t_1[Z] = t_2[Z]$ , quindi  $t_1[Y] = t_2[Y]$ , ne consegue  $X \rightarrow Y \in F^+$ .
- *Caso 3*  $X \rightarrow Y$  è stato ottenuto in  $F^A$  applicando la transitività da due dipendenze  $X \rightarrow Z, Z \rightarrow Y \in F^A$ , ottenute ricorsivamente applicando gli assiomi  $i - 1$  volte. Se  $r$  è un'istanza legale, e  $t_1, t_2$  due tuple tali che  $t_1[X] = t_2[X]$ , ne consegue che  $t_1[Z] = t_2[Z]$ , ma per ipotesi induttiva allora  $t_1[Y] = t_2[Y]$ , quindi  $X \rightarrow Y \in F^+$ .

$F^+ \subseteq F^A$ : Procediamo con una dimostrazione per assurdo. Supponiamo che esiste una dipendenza funzionale  $X \rightarrow Y \in F^+$ , tale che  $X \rightarrow Y \notin F^A$ . Utilizzeremo una particolare istanza  $r$  di  $R$  che ci aiuterà a giungere ad una contraddizione, tale istanza è :

| $X_F^+$ |    |    | $R - X_F^+$ |    |    |
|---------|----|----|-------------|----|----|
| A1      | A2 | A3 | B1          | B2 | B3 |
| 1       | 1  | 1  | 1           | 1  | 1  |
| 1       | 1  | 1  | 0           | 0  | 0  |

Tale istanza ha solamente due tuple, in rosso sono indicati gli attributi che appartengono ad  $X_F^+$ , in blu, gli attributi restanti di  $R$  che non appartengono ad  $X_F^+$ . Tali due tuple, sono identiche sugli attributi di  $X_F^+$ , ma differiscono sui restanti attributi. Dimostriamo adesso due fatti :

- **$r$  è un istanza legale** : Assumiamo che non lo sia, vuol dire che esiste almeno una dipendenza  $V \rightarrow W \in F$ , quindi anche in  $F^A$ , che non è soddisfatta da  $r$ , vuol dire che vi sono due tuple che sono uguali su  $V$  e differenti su  $W$ , questo implicherebbe che  $V \subseteq X_F^+$  e  $W \cap (R - X_F^+) = \emptyset$ , dal momento che  $V \subseteq X_F^+$ , per lemma è vero che  $X \rightarrow V \in F^A$ , per assioma di transitività, essendo che  $V \rightarrow W \in F$ , si ottiene che  $X \rightarrow W \in F^A$ , per lemma quindi  $W \subseteq X_F^+$ , ma ciò è una contraddizione, dato che abbiamo detto che  $W \cap (R - X_F^+) = \emptyset$ , quindi  $r$  è legale.
- **se  $X \rightarrow Y \in F^+$ , è impossibile che  $X \rightarrow Y \notin F^A$**  : Supponiamo che ciò sia vero, sappiamo che  $r$  è un'istanza legale, quindi se le dipendenze in  $F^+$  sono soddisfatte da ogni istanza legale, allora  $r$  soddisfa  $X \rightarrow Y$ . Sappiamo che  $X \subseteq X_F^+$ , dato che su  $r$  ci sono 2 tuple identiche su  $X$ , vi devono essere anche su  $Y$ , quindi  $Y \subseteq X_F^+$ , per lemma, ne consegue che  $X \rightarrow Y \in F^A$ , ma abbiamo detto inizialmente che  $X \rightarrow Y \notin F^A$ , quindi si ha una contraddizione.

Quindi, è necessariamente vero che  $X \rightarrow Y \in F^+ \implies X \rightarrow Y \in F^A$ . Quindi  $F^A = F^+$ . ■

## 4.4 La Terza Forma Normale

All'inizio di questo paragrafo 4, si è parlato di alcune caratteristiche che una base di dati deve avere per essere il più possibile "corretta", in modo da evitare la ridondanza o altre anomalie durante le interrogazioni. Formalizziamo l'insieme delle caratteristiche in quella che si chiama **terza forma normale (3NF)**. Ricordiamo prima una definizione :

Una dipendenza  $X \rightarrow A$  è detta *non banale* se  $A \notin X$ .

Definiamo adesso rigorosamente :

Una base di dati è in *Terza Forma Normale* se tutte le dipendenze non banali sono della forma  $K \rightarrow X$ , dove si verifica *almeno una* condizione fra :

- $K$  contiene una chiave (è superchiave).
- $X$  è attributo di una chiave/è contenuto in una chiave (è primo).

Si può verificare che una base di dati sia in 3NF controllando solo le dipendenze non banali ottenute decomponendo  $F$ .

*Esempio :* Si consideri la seguente base di dati :

$$R = \{A, B, C, D\} \quad F = \{A \rightarrow B, B \rightarrow CD\}$$

Sappiamo che la chiave è  $A$ , infatti si nota che  $A \rightarrow R \in F^+$  (più avanti definiremo come trovare in maniera algoritmica la chiave di uno schema relazionale). Controlliamo che  $R$  sia in 3NF, decomponendo  $F$ , otteniamo :

- $A \rightarrow B$  (che era già in  $F$ )
- $B \rightarrow C$
- $B \rightarrow D$

Per queste ultime due,  $B$  non è una superchiave, e sia  $C$  che  $D$  non sono elementi che appartengono ad una chiave, quindi tale schema *non* è in terza forma normale.

### 4.4.1 Dipendenze Parziali e Transitive

Sia  $R$  uno schema relazionale ed  $F$  il suo insieme delle dipendenze funzionali, si ha che :

- $X \rightarrow A \in F^+ | A \notin X$  è una **dipendenza parziale** se  $A$  non è primo ed  $X$  invece è contenuto nella chiave, si dice che  $A$  **dipende parzialmente** dalla chiave.
- $X \rightarrow A \in F^+ | A \notin X$  è una **dipendenza transitiva** se  $A$  non è primo e per ogni chiave  $K$  su  $R$ , si ha che  $X$  non è propriamente contenuto in  $K$  e  $K - X = \emptyset$ . Si dice che  $A$  **dipende transitivamente** da  $K$ (chiave) se  $K \rightarrow X \in F^+$  e  $X \rightarrow A \in F^+$  dove  $X$  non è chiave ed  $A$  non è parte della chiave.

Dato uno schema  $R$ , è in terza forma normale se e solo se **non** presenta *Dipendenze Parziali* o *Transitive* (analogamente, non presenta attributi che dipendono parzialmente o transitivamente da una chiave).

#### 4.4.2 Decomposizione in più Schemi

Vogliamo saper creare delle basi di dati che siano in terza forma normale, ossia che, ognuna delle sue relazioni sia in 3NF. Avendo uno schema che viola tale forma, è possibile *decomporlo*, creando più relazioni che a loro volta non violino la 3NF.

**Attenzione!** non tutte le decomposizioni in 3NF di uno schema che inizialmente non era in 3NF sono necessariamente corrette, esistono diverse scomposizioni in 3NF, ma alcune violano le dipendenze funzionali che vi erano in origine.

Vediamo un *Esempio*, ho un insieme di dipendenze  $F = \{A \rightarrow B, B \rightarrow C\}$ , dove  $A$  è la chiave. tale insieme viola la terza forma normale ( $B \rightarrow C$  è una dipendenza transitiva), vediamo due ipotesi di scomposizione :

$$(1) \begin{cases} R1 = \{A, B\} \text{ con } F = \{A \rightarrow B\} \\ R2 = \{B, C\} \text{ con } F = \{B \rightarrow C\} \end{cases}$$

$$(2) \begin{cases} R1 = \{A, B\} \text{ con } F = \{A \rightarrow B\} \\ R2 = \{A, C\} \text{ con } F = \{A \rightarrow C\} \end{cases}$$

Presa una qualsiasi istanza legale della relazione originale, se venisse decomposta in (1) o (2), il **JOIN** delle due relazioni scomposte dovrebbe creare una nuova istanza che sia ancora legale come quella originale, ciò non è soddisfatto dalla scomposizione (2), vediamo un esempio :

| R1 |    | ⋈ | R2 |    | = | R1⋈R2 |    |    |
|----|----|---|----|----|---|-------|----|----|
| A  | B  |   | A  | C  |   | A     | B  | C  |
| a1 | b1 |   | a1 | c1 |   | a1    | b1 | c1 |
| a2 | b2 |   | a2 | c2 |   | a2    | b1 | c2 |

Anche se sia  $R1$  che  $R2$  sono in 3NF, il risultato del loro **JOIN** non soddisfa la dipendenza funzionale  $B \rightarrow C$ .

Si noti come, anche se prendessimo in considerazione la decomposizione (1), essa soddisfa le dipendenze funzionali iniziali, ma non andrebbe comunque bene perché il **JOIN** genererebbe delle tuple che non vi erano in origine (perdita di informazioni).

In conclusione, la decomposizione di uno schema non in 3NF in più schemi in 3NF deve :

- Preservare le dipendenze funzionali dello schema originale.
- Poter ricostruire lo schema originale tramite un **JOIN** con ogni istanza legale senza aggiungere informazioni aggiuntive.

#### 4.4.3 Forma Normale Boyce-Codd

La terza forma normale non è la forma normale più restrittiva, ce ne sono altre, inclusa la seguente :

Una relazione è in **forma normale Boyce-Codd**, se, per ogni singola dipendenza funzionale, il determinante è una superchiave (o ovviamente, chiave).

Sia  $K$  la chiave :  $\forall A \rightarrow B \in F^+, K \subseteq A$

Una relazione che è in forma normale Boyce-Codd è anche in terza forma normale, non è vero però il contrario, se una relazione è in terza forma normale, non è necessariamente in forma normale Boyce-Codd. Non è sempre possibile decomporre uno schema in sotto-schemi in forma normale Boyce-Codd, è però sempre possibile farlo in sotto-schemi in terza forma normale, per questo, considereremo quest'ultima durante il corso.

## 4.5 Calcolo per la Chiusura di $X$

Nel capitolo precedente 4.2.1 abbiamo dato la definizione dell'insieme  $X^+$ , ossia la chiusura di  $X$  (Un insieme di attributi dello schema relazionale), l'insieme contenente tutti gli attributi  $A$  tali che  $X \rightarrow A \in F^A$ . Vedremo adesso un *algoritmo* in grado di calcolare, a partire da un attributi, la sua chiusura, ciò risulterà utile anche in funzione di stabilire quali sono le chiavi di una relazione.

### 4.5.1 L'Algoritmo

**Input** : Lo schema  $R$ , l'insieme delle dipendenze funzionali  $F$ , ed  $X \subseteq R$ .

**Output** : Denotato con  $Z_f$ , sarà la chiusura di  $X$ .

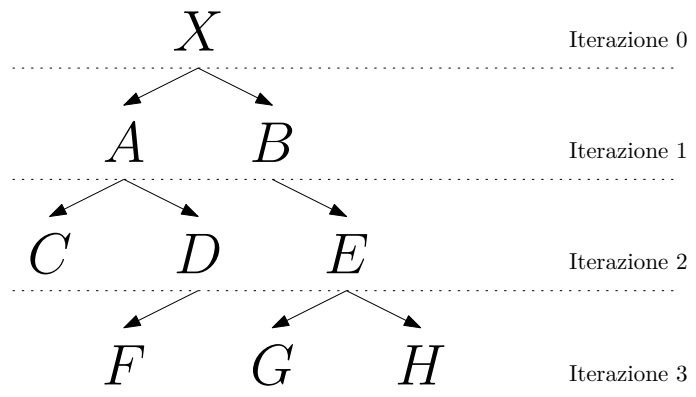
```

begin {
    Z = X
    S = {A |  $\exists Y \rightarrow V \in F, Y \subseteq Z \wedge A \in V$ }
    while S  $\not\subseteq$  Z {
        Z = X
        S = {A |  $\exists Y \rightarrow V \in F, Y \subseteq Z \wedge A \in V$ }
    }
    return Z
}

```

Ad ogni passo iterativo, ad  $S$  aggiungiamo tutti gli attributi determinati da  $X$ , direttamente o per transitività, si può immaginare un albero di "determinanza" che ad ogni passo cresce verso il basso.

In questo esempio,  $X$  determina  $A$  e  $B$ , quindi all'iterazione 1,  $A$  e  $B$  verranno inclusi in  $Z$ , per transitività, gli attributi  $C$  e  $D$  determinati da  $A$ , ed  $E$  determinato da  $B$ , sono anche determinati da  $X$ , quindi all'iterazione 2 verranno inclusi in  $Z$ , stesso procedimento per l'iterazione 3.



Denoteremo con  $Z_0, Z_1, Z_2 \dots, Z_n$  la sequenza dei valori che assume  $Z$ , dove  $Z_i$  è il valore di  $Z$  all' $i$ -esima iterazione, stessa cosa per la variabile  $S$  con la sequenza  $S_0, S_1, S_2 \dots, S_n$ .

L'algoritmo termina quando  $S$  sarà incluso in  $Z$ , finché l'algoritmo continuerà l'esecuzione, è chiaro che  $Z \cup S \neq Z$ . Per ogni iterazione, sarà vero che  $Z_{i+1} = Z_i \cup S_i$ , ciò implica ovviamente che  $Z_0 \subset Z_1 \subset Z_2 \dots \subset Z_n$ , quindi la sequenza dei valori di  $Z$ , è strettamente

crescente ad ogni iterazione, ed è ovviamente "limitata" da  $R$ , è ovvio che  $\exists f < \infty | Z_f = Z_{f+1} \implies S_f \subset Z_f \implies$  l'algoritmo, ad un certo punto termina. Abbiamo quindi denotato con  $Z_f$  il valore che assume  $Z$  al termine dell'algoritmo.

**Teorema** :  $Z_f = X^+$

**Dimostrazione** : Facciamo una dimostrazione per doppia inclusione, partendo da

$\boxed{Z_f \subseteq X^+}$  : Utilizzeremo l'indice  $i$  delle iterazioni per una dimostrazione per induzione, ho il *caso base* :  $Z_0 \subseteq X^+$ , so che  $Z_0 = \{X\}$ , e la chiusura di  $X$  sicuramente contiene  $X$ , quindi  $Z_0 \subseteq X^+$  è verificato, passiamo al *passo induttivo* : la nostra ipotesi è che  $Z_i \subseteq X^+$  e vogliamo dimostrare che  $Z_{i+1} \subseteq X^+$ , prendo un qualsiasi  $A \in Z_{i+1}$ , sapendo che  $Z_{i+1} = Z_i \cup S_i$ , l'attributo  $A$  o appartiene a  $Z_i$  oppure a  $S_i$ , se  $A$  dovesse appartenere a  $Z_i$ , avremmo già risolto e dimostrato il passo induttivo, se invece  $A$  dovesse appartenere ad  $S_i$ , vuol dire che  $\exists Y \rightarrow V \in F$  tale che  $Y \subseteq Z_i \wedge A \in V$ . Utilizzo l'ipotesi che  $Z_i \subseteq X^+$ , quindi se  $Y \subseteq Z_i \implies Y \subseteq X^+$ , e ciò significa che  $X \rightarrow Y \in F^A$ , e sapendo che  $Y \rightarrow V \in F^A$ , per transitività ottengo che  $X \rightarrow V \in F^A$ , ricordando che  $A \in V$ , questo implica che  $X \rightarrow A \in F^A$ , ne concludiamo che  $A$ , che è un qualsiasi elemento di  $Z_{i+1}$  (in questo caso di  $S_i$ ), è anche in  $X^+$ , quindi la prima inclusione è verificata.

$\boxed{X^+ \subseteq Z_f}$  : Similmente alla dimostrazione del capitolo 4.3.1, ci serviremo di una particolare istanza  $r$ , composta da due tuple, in cui a sinistra compaiono tutti gli attributi contenuti in  $Z_f$ , e a destra quelli contenuti nel complementare.

$$r := \begin{array}{c} t_1 \\ t_2 \end{array} \begin{array}{|c|c|c|c|c|c|c|c|} \hline & \multicolumn{4}{Z_f} & \multicolumn{4}{R - Z_f} \\ \hline & 1 & 1 & 1 & \dots & 1 & 1 & \dots & 1 \\ \hline & 1 & 1 & 1 & \dots & 0 & 0 & \dots & 0 \\ \hline \end{array}$$

Dimostreremo che  $r$  sia legale, prendo una qualsiasi dipendenza funzionale  $Y \rightarrow V \in F$ ,  $t_1[Y] = t_2[Y]$ , so per certo che  $Y \subseteq Z_f$ , dato che, per come è strutturata l'istanza, non ci può essere un attributo in  $Z_f$  che determina un attributo in  $R - Z_f$  (Ciò è corretto in quanto, se un attributo fosse determinato da un elemento in  $Z_f$  (che sarebbe la chiusura di  $X$ ), ciò significherebbe che tale attributo è determinato da  $X$ , quindi dovrebbe esso stesso essere in  $Z_f$ ), quindi anche  $t_1[V] = t_2[V]$ , ed  $r$  è legale.

Sia  $A \in X^+$  un qualsiasi attributo della chiusura di  $X$ , allora  $X \rightarrow A \in F^A$ , ed  $r$  soddisfa tale dipendenza, sappiamo che  $X = Z_0 \subset Z_f \implies t_1[X] = t_2[X] \implies t_1[A] = t_2[A] \implies A \in Z_f$ .

$$\begin{cases} Z_f \subseteq X^+ \\ X^+ \subseteq Z_f \end{cases} \implies Z_f = X^+ \quad \blacksquare$$

*Esempio* : Si consideri il seguente schema relazionale :  $R = (A, B, C, D, E, H)$  con il seguente insieme di dipendenze funzionali :  $F = (AB \rightarrow CD, EH \rightarrow D, D \rightarrow H)$ , si calcoli  $AB^+$ .

- iterazione 0 -  $Z_0 = \{AB\}$ ,  $S_0 = \{CD\}$ , noto che  $S_0 \not\subseteq Z_0$ , quindi continuo.
- iterazione 1 -  $Z_1 = \{AB, CD\}$ ,  $S_1 = \{CD, H\}$ , noto che  $S_1 \not\subseteq Z_1$ , quindi continuo.
- iterazione 2 -  $Z_2 = \{AB, CD, H\}$ ,  $S_2 = \{CD, H\}$ , noto che  $S_2 \subseteq Z_2$ , quindi mi fermo.

**Output** :  $Z_2 = Z_f = AB^+ = \{AB, CD, H\}$



### 4.5.2 La Ricerca delle Chiavi

Sappiamo che, se un insieme di attributi  $X$  determina tutto lo schema relazionale  $R$ , allora tale insieme è una *chiave*. Se volessimo capire quali sono le chiavi di uno schema, potremmo quindi calcolare la chiusura di tutti i possibili sotto-insiemi di  $R$ , e compararli ad  $R$  stesso per vedere se sono equivalenti.

$$\forall X \in \mathcal{P}(R) : \text{ if } (\text{Algoritmo}(X) == R) \{ \text{return Algoritmo}(X) \}$$

Sappiamo però che, tutti i sotto insiemi di  $R$ , sono in numero  $2^{|R|}$ , e per ognuno di questi si dovrà eseguire l'algoritmo per calcolarne la chiusura 4.5.1, è chiaro che questo tipo di problema, necessita di un calcolo di una complessità *non polinomiale*, è quindi, per quanto corretto, impensabile di poter trovare la chiave di una relazione usando la "forza bruta".

**Osservazione :** Se un insieme di attributi non compare mai come determinato in  $F$ , allora è parte della chiave. Con tale osservazione, si è può ridurre in maniera notevole il numero di attributi candidati ad essere chiave.

*Esempio :*  $R = (A, B, C, D, E, H)$   $F = (AB \rightarrow CD, C \rightarrow E, AB \rightarrow E, ABC \rightarrow D)$  Si vuole trovare la chiave di tale schema relazionale. Si facciano due osservazioni, l'attributi  $H$ , non compare in nessuna dipendenza funzionale, inoltre, gli attributi  $A$  e  $B$ , non compaiono mai come determinati in  $F$ , ma solo come determinanti. Si consideri quindi la chiusura di  $ABH$ .  $\text{Algoritmo}(ABH) = ABH^+ = (A, B, H, C, D, E) = R$ , quindi, per definizione 4.1, se non esiste nessun sottoinsieme di  $ABH$  che determina  $R$ , tale  $ABH$  è chiave. Ci si rende conto dopo pochi calcoli che  $ABH$  è chiave.

**Osservazione :** Gli attributi che compaiono esclusivamente come determinati, non sono chiavi. È consigliabile, iniziare a controllare le chiusure dei sottoinsiemi di cardinalità più alta, se essi non determinano  $R$ , non c'è bisogno di controllarne i sotto insiemi

*Remark :* Durante la prova d'esame, è necessario giustificare tutti i passaggi e le considerazioni che portano il candidato ad escludere un certo insieme di attributi durante la ricerca della chiave.

**Osservazione :** Le chiavi di uno schema relazionale possono essere molteplici, è quindi opportuno cominciare controllando tutti gli attributi che non compaiono mai come determinati. Se la chiusura di un insieme  $X$  è uguale ad  $R$ , bisogna verificare che non lo siano anche i suoi sottoinsiemi.

*Esempio :*  $R = (A, B, C, D, E, G, H)$   $F = (AB \rightarrow D, G \rightarrow A, G \rightarrow B, H \rightarrow E, H \rightarrow G, D \rightarrow H)$  Si vogliono trovare le 4 chiavi. Si cominci dagli attributi che determinano altri ma non sono determinati da nessuno, i candidati sono  $AB, G, D$  e  $H$ . Noto che  $C$  non compare mai in nessuna dipendenza funzionale, quindi sarà parte di tutti le chiavi. Comincio allora con  $ABC$ , trovo che  $\text{Algoritmo}(ABC) = R$ , controllo poi i sottoinsiemi e vedo che

$\text{Algoritmo}(AC) = (AB, BC^+)$  e che  $BC^+ = (BC)$ , quindi  $ABC$  è una chiave. Ne rimangono 3, gli altri 3 candidati erano  $G, D$  e  $H$ , utilizzo l'algoritmo e trovo che  $GC^+ = DC^+ = HC^+ = R$ , sono quindi tutte e 3 le chiavi, senza necessità di controllare i sottoinsiemi (dato che  $C$  da sola non può essere chiave, ma deve essere contenuto in ogni chiave).

Dato uno schema  $R$ , se viene trovata una chiave, spesso ci si chiede se tale chiave sia l'unica, è quindi opportuno fare un **test di unicità**, dato che, per controllare se uno schema sia in terza forma normale, necessitiamo di sapere quali siano tutte le chiavi.

Esiste un modo semplice per determinare se la chiave di uno schema  $R$  è unica :

$$\text{Sia } X = \bigcap_{V \rightarrow W \in F} R - (W - V)$$

Una volta aver trovato tale  $X$ , se  $X^+ = R$ , allora  $X$  è l'unica chiave di  $R$ , altrimenti, esistono più chiavi, ed  $X$  non è una superchiave di  $R$ .

*Esempio* :  $R = (A, B, C, D, E)$   $F = (AB \rightarrow C, AC \rightarrow B, D \rightarrow E)$ , voglio capire se la chiave sia unica, applico quindi il calcolo appena visto :

$$\begin{aligned} X &= (ABCDE - (C - AB) \cap ABCDE - (B - AC) \cap ABCDE - (E - D)) = \\ &= (ABCDE - C \cap ABCDE - B \cap ABCDE - E) = (ABDE \cap ACDE \cap ABCD) = AD \\ &\quad \text{calcolo } X^+ = AD^+ = AD \neq R, \text{ quindi ci sono più chiavi.} \end{aligned}$$