

BD 2 - Forme Normali

Luca Cosmo

Università Ca' Foscari Venezia



Università
Ca' Foscari
Venezia

Introduzione

L'obiettivo delle forme normali è garantire che uno schema sia di buona qualità e viene spesso ottenuto tramite un processo di **normalizzazione** basato su una decomposizione dello schema di partenza.

Proprietà Desiderabili

- 1 Uno schema in forma normale non deve contenere **anomalie**
- 2 Il processo di normalizzazione deve **preservare i dati**
- 3 Il processo di normalizzazione deve **preservare le dipendenze**

Le due forme normali di maggior interesse sono:

- Boyce-Codd Normal Form (BCNF)
- Terza Forma Normale (3FN)

Vedremo che in generale **non possiamo** garantire tutte e tre le proprietà!

Forma Normale di Boyce-Codd (BCNF)

Definition (BCNF)

Uno schema di relazione $R(T, F)$ è in **BCNF** sse per ogni dipendenza funzionale $X \rightarrow Y \in F$ tale che $Y \not\subseteq X$ si ha che X è una superchiave.

Verificare se uno schema è in BCNF ha perciò costo **polinomiale**.

Example

| Articolo | Magazzino | Quantità | Indirizzo |
|----------|-----------|----------|-----------------|
| Flauto | Roma | 10 | Via Cavour, 7 |
| Oboe | Roma | 5 | Via Cavour, 7 |
| Arpa | Torino | 1 | Via Mazzini, 11 |

$R = \{\text{Articolo, Magazzino, Quantità, Indirizzo}\}$

$F = \{\text{Articolo} \rightarrow \text{Quantità}, \text{Magazzino} \rightarrow \text{Indirizzo}\}$

Dipendenze Anomale

Una dipendenza che viola BCNF è detta **anomala**. Nel nostro esempio abbiamo una dipendenza anomala Magazzino → Indirizzo.

| Articolo | Magazzino | Quantità | Indirizzo |
|----------|-----------|----------|-----------------|
| Flauto | Roma | 10 | Via Cavour, 7 |
| Oboe | Roma | 5 | Via Cavour, 7 |
| Arpa | Torino | 1 | Via Mazzini, 11 |

Tale dipendenza anomala evidenzia che lo schema mescola informazioni relative ai magazzini con altre **indipendenti** relative agli articoli.

Riuscite ad individuare delle anomalie pericolose in questo esempio?

Conversione in BCNF

L'algoritmo di conversione in BCNF è anche detto algoritmo di **analisi**, perché scompone lo schema originale fino a normalizzazione.

Sia $R(T, F)$ lo schema di partenza:

- 1 Se $R(T, F)$ è già in BCNF, ritorna $\{R(T, F)\}$
- 2 Altrimenti seleziona $X \rightarrow Y \in F$ che viola BCNF. Calcola gli insiemi di attributi $T_1 = X^+$ e $T_2 = X \cup (T \setminus T_1)$
- 3 Calcola le proiezioni $F_1 = \pi_{T_1}(F)$ e $F_2 = \pi_{T_2}(F)$
- 4 Decomponi ricorsivamente $R_1(T_1, F_1)$ e $R_2(T_2, F_2)$ in ρ_1 e ρ_2 .
- 5 Ritorna la loro unione $\rho_1 \cup \rho_2$.

Conversione in BCNF: Esempio (1/3)

Si consideri Telefoni($\{\text{Prefisso, Numero, Località}\}$, F) con:

$$F = \{\text{Prefisso Numero} \rightarrow \text{Località}, \text{Località} \rightarrow \text{Prefisso}\}.$$

La dipendenza $\text{Località} \rightarrow \text{Prefisso}$ viola BCNF, dato che:

$$\{\text{Località}\}_F^+ = \{\text{Località, Prefisso}\}.$$

Applicando l'algoritmo di conversione in BCNF, abbiamo:

- $R_1(\{\text{Località, Prefisso}\}, F_1)$ con F_1 da calcolare per proiezione
- $R_2(\{\text{Località, Numero}\}, F_2)$ con F_2 da calcolare per proiezione

Conversione in BCNF: Esempio (2/3)

Dato $F = \{\text{Prefisso Numero} \rightarrow \text{Località}, \text{Località} \rightarrow \text{Prefisso}\}$, calcoliamo la sua proiezione per $R_1(\{\text{Località}, \text{Prefisso}\})$:

- $\{\text{Località}\}_F^+ = \{\text{Località}, \text{Prefisso}\}$, da cui $\text{Località} \rightarrow \text{Prefisso} \in F_1$
- $\{\text{Prefisso}\}_F^+ = \{\text{Prefisso}\}$, da cui nessuna nuova dipendenza

Calcoliamo poi la sua proiezione per $R_2(\{\text{Località}, \text{Numero}\})$:

- $\{\text{Località}\}_F^+ = \{\text{Località}, \text{Prefisso}\}$, da cui nessuna nuova dipendenza
- $\{\text{Numero}\}_F^+ = \{\text{Numero}\}$, da cui nessuna nuova dipendenza

Abbiamo quindi $F_1 = \{\text{Località} \rightarrow \text{Prefisso}\}$ e $F_2 = \emptyset$.

Conversione in BCNF: Esempio (3/3)

Abbiamo decomposto Telefoni($\{\text{Prefisso, Numero, Località}\}$, F) con:

$$F = \{\text{Prefisso Numero} \rightarrow \text{Località}, \text{Località} \rightarrow \text{Prefisso}\}$$

nei seguenti schemi:

- $R_1(\{\text{Località, Prefisso}\}, \{\text{Località} \rightarrow \text{Prefisso}\})$
- $R_2(\{\text{Località, Numero}\}, \emptyset)$

Entrambi gli schemi sono in BCNF, ma è andata perduta la dipendenza funzionale $\text{Prefisso Numero} \rightarrow \text{Località}$!

Proprietà

La conversione in BCNF non garantisce la preservazione delle dipendenze!

Perdita di Dipendenze

Cosa succede se proviamo ad inserire (049, 513227, Este)?

Prefisso Numero → Località

Località → Prefisso

| Prefisso | Numero | Località |
|----------|--------|----------|
| 041 | 422865 | Venezia |
| 041 | 463212 | Venezia |
| 049 | 513227 | Padova |

Località → Prefisso

| Prefisso | Località |
|----------|----------|
| 041 | Venezia |
| 049 | Padova |

| Numero | Località |
|--------|----------|
| 422865 | Venezia |
| 463212 | Venezia |
| 513227 | Padova |

Correttezza della Conversione in BCNF

L'algoritmo di conversione in BCNF **termina** quando non ci sono più dipendenze anomale. Per garantire che ciò avverrà, dimostriamo che tutti gli schemi **con solo due attributi** sono in BCNF.

Consideriamo $R(\{A, B\}, F)$ e sia $X \rightarrow Y \in F$. Dimostriamo che in nessun caso viene violata BCNF:

- 1 Se $X = \{A\}$, ho due casi. Se $B \notin Y$, allora $Y \subseteq X$ e la dipendenza è banale. Se invece $B \in Y$, allora X è una superchiave.
- 2 Se $X = \{B\}$, il caso è simmetrico al precedente.
- 3 Se $X = \{A, B\}$, allora $Y \subseteq X$ e la dipendenza è banale.

La Conversione in BCNF Preserva i Dati

Preservazione dei Dati

La conversione in BCNF preserva i dati (segue dimostrazione).

Supponiamo che $R(T, F)$ sia decomposto in $\{R_1(T_1), R_2(T_2)\}$, allora deve esistere $X \rightarrow Y \in F$ che viola BCNF. Per costruzione $T_1 = X_F^+$ e $T_2 = X \cup (T \setminus T_1)$.

Osserviamo che $T_1 \cap T_2 = X$. Dato che $X \rightarrow X_F^+ \in F^+$, abbiamo che $T_1 \cap T_2 \rightarrow T_1 \in F^+$, quindi la decomposizione preserva i dati per il teorema visto nella lezione precedente¹.

Il risultato si può quindi dimostrare **per induzione** sul numero di passi effettuati dall'algoritmo di conversione in BCNF.

¹Sia $\rho = \{R_1(T_1), R_2(T_2)\}$ una decomposizione di $R(T, F)$, si ha che ρ preserva i dati sse $T_1 \cap T_2 \rightarrow T_1 \in F^+$ oppure $T_1 \cap T_2 \rightarrow T_2 \in F^+$.

Proprietà di BCNF

Pregi

- + BCNF garantisce l'**assenza di anomalie** (no dipendenze anomale)
- + L'algoritmo di conversione in BCNF **preserva i dati**
- + Verificare se uno schema è in BCNF ha costo **polinomiale**

Difetti

- L'algoritmo di conversione in BCNF ha costo **esponenziale**, perché richiede di calcolare le proiezioni delle dipendenze²
- L'algoritmo di conversione in BCNF **non preserva le dipendenze** nel caso generale

²Esistono anche algoritmi di costo polinomiale, ma non sono usati in pratica perché producono schemi eccessivamente decomposti.

Terza Forma Normale (3NF)

Definition (3NF)

Uno schema di relazione $R(T, F)$ è in **3NF** sse per ogni dipendenza funzionale $X \rightarrow Y \in F$ tale che $Y \not\subseteq X$ si ha che X è una superchiave oppure tutti gli attributi di $Y \setminus X$ sono primi.

Verificare se uno schema è in 3NF ha costo **esponenziale**, perchè il calcolo degli attributi primi richiede di trovare tutte le chiavi.

Osservazione

Per definizione ogni schema in BCNF è anche in 3NF, ma non viceversa.

Esempio

Si consideri $\text{Telefoni}(\{\text{Prefisso}, \text{Numero}, \text{Località}\}, F)$ con:

$$F = \{\text{Prefisso} \text{ Numero} \rightarrow \text{Località}, \text{Località} \rightarrow \text{Prefisso}\}.$$

Calcoliamo le chiavi, osservando che Numero deve fare parte di tutte:

- $\{\text{Numero}\}_F^+ = \{\text{Numero}\}$
- $\{\text{Numero}, \text{Prefisso}\}_F^+ = \{\text{Numero}, \text{Prefisso}, \text{Località}\}$
- $\{\text{Numero}, \text{Località}\}_F^+ = \{\text{Numero}, \text{Località}, \text{Prefisso}\}$

Dato che $\{\text{Numero}, \text{Prefisso}\}$ e $\{\text{Numero}, \text{Località}\}$ sono chiavi, si ha che ogni attributo è primo e quindi lo schema è in 3NF.

Conversione in 3NF

L'algoritmo di conversione in 3NF è anche detto algoritmo di **sintesi**, perché basato sulla generazione di nuovi schemi più piccoli.

Sia $R(T, F)$ lo schema di partenza:

- 1 Costruisci G , una copertura canonica di F
- 2 Sostituisci in G ciascun insieme di dipendenze $X \rightarrow A_1, \dots, X \rightarrow A_n$ con una singola dipendenza $X \rightarrow A_1 \dots A_n$
- 3 Per ogni $X \rightarrow Y \in G$, crea un nuovo schema $S_i(XY)$
- 4 Elimina ogni schema contenuto in un altro schema
- 5 Se la decomposizione non contiene alcuno schema i cui attributi costituiscano una superchiave per R , aggiungi un nuovo schema $S(W)$ dove W è una chiave di R (garantisce la preservazione dei dati).

Conversione in 3NF: Esempio

Sia $R(\{A, B, C, D\}, \{AB \rightarrow C, C \rightarrow D, D \rightarrow B\})$, osserviamo che l'insieme delle dipendenze è già in forma canonica. Otteniamo quindi:

- $R_1(\{A, B, C\})$ tramite $AB \rightarrow C$
- $R_2(\{C, D\})$ tramite $C \rightarrow D$
- $R_3(\{B, D\})$ tramite $D \rightarrow B$

Nessuno schema è contenuto in un altro, quindi nessuno di essi viene eliminato. Poichè $\{A, B, C\}$ è una superchiave di R , non è necessario aggiungere altri schemi.

La Conversione in 3NF Preserva i Dati e le Dipendenze

Preservazione delle Dipendenze

E' facile dimostrare che la conversione in 3NF preserva le dipendenze: poichè per ogni $X \rightarrow Y \in G$ viene creato uno schema $S_i(XY)$, abbiamo $X \rightarrow Y \in \pi_{XY}(G)$, quindi G è contenuto nell'unione delle proiezioni.

Preservazione dei Dati

L'ultimo passo della conversione in 3NF garantisce che la decomposizione contenga almeno uno schema i cui attributi formano una superchiave dello schema iniziale. Poichè la decomposizione preserva le dipendenze, essa deve preservare anche i dati per il teorema visto.

Correttezza dell'algoritmo di conversione in 3NF: non banale (vedi testo)

3NF ed Anomalie

Si consideri Telefoni($\{\text{Prefisso}, \text{Numero}, \text{Località}\}, F$) con:

$$F = \{\text{Prefisso Numero} \rightarrow \text{Località}, \text{Località} \rightarrow \text{Prefisso}\}.$$

Abbiamo già visto che lo schema è in 3NF, ma non garantisce l'assenza di anomalie. In particolare, si noti la replicazione del prefisso...

| Prefisso | Numero | Località |
|----------|--------|----------|
| 041 | 422865 | Venezia |
| 041 | 463212 | Venezia |
| 049 | 513227 | Padova |

Proprietà

La conversione in 3NF non garantisce l'assenza di anomalie!

Proprietà di 3NF

Pregi

- + L'algoritmo di conversione in 3NF **preserva i dati e le dipendenze**
- + L'algoritmo di conversione in 3NF ha costo **polinomiale**, perchè non richiede il calcolo delle proiezioni

Difetti

- Verificare se uno schema è in 3NF ha costo **esponenziale**, perchè richiede di identificare gli attributi primi
- Uno schema in 3NF può ancora contenere **anomalie**

Dipendenze Multivalore

Una nuova forma di anomalia non prevenuta neppure da BCNF si può verificare in presenza di **attributi multivalore indipendenti**. Per esempio la relazione sottostante non ha dipendenze funzionali non banali.

| Corso | LibroDiTesto | Docente |
|--------------|---------------------|----------|
| Basi di dati | Fondamenti di BD | Cosmo |
| Basi di dati | Web App Development | Cosmo |
| Basi di dati | Database Systems | Cosmo |
| Basi di dati | Fondamenti di BD | Raffaetà |
| Basi di dati | Web App Development | Raffaetà |
| Basi di dati | Database Systems | Raffaetà |

C'è però una forte ridondanza: se ci sono m docenti ed n libri di testo, si memorizzano $m \times n$ righe.

Dipendenze Multivalore

E' certamente possibile fare meglio, memorizzando solo $m + n$ righe.

| Corso | LibroDiTesto |
|--------------|---------------------|
| Basi di dati | Fondamenti di BD |
| Basi di dati | Web App Development |
| Basi di dati | Database Systems |

| Corso | Docente |
|--------------|----------|
| Basi di dati | Cosmo |
| Basi di dati | Raffaetà |

La teoria della normalizzazione è stata perciò generalizzata per rimuovere anche questo tipo di anomalie dovute alle **dipendenze multivalore** (4NF).

Checkpoint

Punti Chiave

- Conversione in BCNF come tecnica per l'eliminazione di anomalie, che però ha costo esponenziale e può perdere dipendenze
- Conversione in 3NF come alternativa più pratica in molti casi: costo polinomiale e preserva le dipendenze, ma può ammettere anomalie

Materiale Didattico

Fondamenti di Basi di Dati: Sezioni 5.4, 5.5 e 5.7