

Chuẩn hóa dữ liệu (Normalization)

Các dạng chuẩn (1NF -> BCNF)
Thuật toán phân rã thành 3NF bảo toàn
phụ thuộc
Thuật toán phân rã thành BCNF có nổi
không mất



Chuẩn hóa (1)

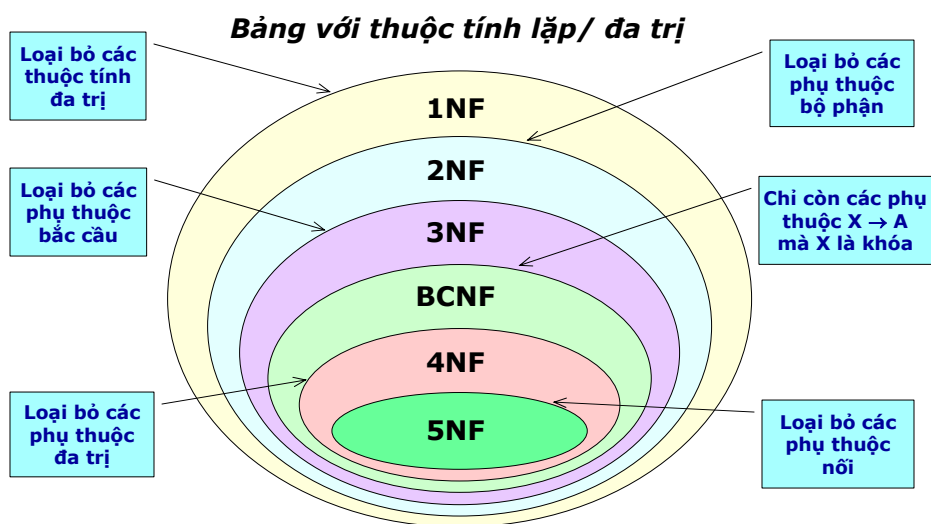
- Chuẩn hóa là quá trình tách bảng (phân rã) thành các bảng nhỏ hơn dựa vào các phụ thuộc hàm mà không làm mất thông tin.
- Các dạng chuẩn là các chỉ dẫn để thiết kế các bảng trong CSDL.
- Mục đích của chuẩn hóa là loại bỏ các dư thừa vì
 - Yêu cầu về không gian lưu trữ
 - Tránh các lỗi khi thao tác dữ liệu (Insert, Delete, Update).
- Chú ý: chuẩn hóa làm tăng thời gian truy vấn.



Chuẩn hoá (2)

- Các dạng chuẩn sử dụng phụ thuộc hàm
- Chuẩn 3 (3NF)
 - Luôn có được phân rã có nối không mất và bảo toàn phụ thuộc
 - Cho phép dư thừa 1 chút
- Chuẩn Boyce Codd (BCNF)
 - Tránh dư thừa
 - Luôn có phân rã có nối không mất
 - Không phải luôn có được phân rã bảo toàn phụ thuộc

Các dạng chuẩn (Normal forms)



Bảng R sau có phải là 1 quan hệ ?



MF	Tenfilm	NSX	Giathue	HSX	NPP	MaKH	TenKH	Điachi	Ngaydat
P001	Hồn yêu	1996	3.2	C 20	Fafim	K001	David	Paris	1/1/06
P002	Đảo vô hình	2006	5.0	C 20	Fafim	K001	David	Paris	1/2/06
						K002	Marie	Toulouse	5/1/06
						K003	John	Nice	1/2/06
P003	Hoa hậu FBI	2002	3.2	Wanner	Films	K003	John	Nice	1/2/06
P004	Taxi	2005	4.5	France	Fafim	K001	David	Paris	2/5/06
						K004	Helen	Bordeaux	1/2/06
P005	Đời cát	2003	2.5	VFC	VFC	K005	Taix	Toulouse	4/1/06

Còn bảng S sau là 1 quan hệ ?



NGUOIQL	CAPDUOI1	CAPDUOI2	CAPDUOI3	CAPDUOI4
Bob	Jim	Mary	Beth	
Mary	Mike	Jason	Carol	Mark
Jim	Alan			

Cột lặp

- Jim chỉ có một cặp dưới ?
- Nếu Mary thêm một người cặp dưới ?
- Gộp các cột lặp thành một cột
- **Vấn đề gì xảy ra ???**



Dạng chuẩn 1 - 1NF

- Chuẩn 1 đảm bảo bảng đang xét là một quan hệ:
 - Giá trị của các cột là đơn trị.
 - Không chứa các cột lặp
- Cách chuyển sang dạng 1NF:
 - điền đủ dữ liệu vào các cột khác
 - Biến cột đa trị thành các dòng

MF	Tenfilm	NSX	Giathue	HSX	NPP	MaKH	TenKH	Điachi	Ngaydat
P001	Hồn yêu	1996	3.2	C 20	Fafim	K001	David	Paris	1/1/06
P002	Đảo vô hình	2006	5.0	C 20	Fafim	K001	David	Paris	1/2/06
P002	Đảo vô hình	2006	5.0	C 20	Fafim	K002	Marie	Toulouse	5/1/06
P002	Đảo vô hình	2006	5.0	C 20	Fafim	K003	John	Nice	1/2/06
P003	Hoa hậu FBI	2002	3.2	Wanner	Films	K003	John	Nice	1/2/06
P004	Taxi	2005	4.5	France	Fafim	K001	David	Paris	2/5/06
P004	Taxi	2005	4.5	France	Fafim	K004	Helen	Bordeaux	1/2/06
P005	Đời cát	2003	2.5	VFC	VFC	K005	Taix	Toulouse	4/1/06



Dạng chuẩn 2 - 2NF (1)

- Một quan hệ đạt 2NF nếu
 - Quan hệ đã đạt 1NF
 - Thuộc tính khác khóa phụ thuộc hoàn toàn vào khóa
- Một quan hệ đạt 2NF nếu thỏa mãn 1 trong các điều kiện sau
 - Khóa chính chỉ gồm 1 thuộc tính.
 - Bảng không có các thuộc tính không khóa.
 - Tất cả thuộc tính không khóa phụ thuộc hoàn toàn vào tập thuộc tính khóa chính.
- Gợi ý
 - Chỉ kiểm tra các quan hệ có đạt 2NF nếu quan hệ đó có khoá chính gồm 2 thuộc tính trở lên.
 - Để chuyển quan hệ từ dạng 1NF sang dạng 2NF, ta dùng phép chiếu.

Dạng chuẩn 2 - 2NF (2)

- Bảng R có các phụ thuộc hàm sau:
 - ① $MF \rightarrow \text{Tenfim, Giathue, NSX, HSX, N.ngu, NPP}$
 - ② $\text{MaKH} \rightarrow \text{TenKH, Diachi}$
 - ③ $MF, \text{MaKH} \rightarrow \text{Ngaythue}$
 - ④ $\text{HSX} \rightarrow \text{NPP}$
 - Khóa chính: MF, MaKH.
 - Các thuộc tính Tenfim, Giathue, TenKH, Diachi,...
 - là các thuộc tính không khóa
 - chỉ phụ thuộc vào một bộ phận của khóa
- R không đạt chuẩn 2

MF	Tenfim	NSX	Giathue	HSX	NPP	MaKH	TenKH	Điachi	Ngaydat
P001	Hồn yêu	1996	3.2	C 20	Fafim	K001	David	Paris	1/1/06
P002	Đảo vô hình	2006	5.0	C 20	Fafim	K001	David	Paris	1/2/06
P002	Đảo vô hình	2006	5.0	C 20	Fafim	K002	Marie	Toulouse	5/1/06
P002	Đảo vô hình	2006	5.0	C 20	Fafim	K003	John	Nice	1/2/06
P003	Hoa hậu FBI	2002	3.2	Wanner	Films	K003	John	Nice	1/2/06
P004	Taxi	2005	4.5	France	Fafim	K001	David	Paris	2/5/06
P004	Taxi	2005	4.5	France	Fafim	K004	Helen	Bordeaux	1/2/06
P005	Đời cát	2003	2.5	VFC	VFC	K005	Taix	Toulouse	4/1/06

$\pi_{\text{MaKH, TenKH, Diachi}}(R)$			$\pi_{\text{MF, Tenfim, NSX, Giathue, HSX, NPP}}(R)$						$\pi_{\text{MF, MaKH, Ngaydat}}(R)$		
MaKH	TenKH	Diachi	MF	Tenfim	NSX	Giathue	HSX	NPP	MF	MaKH	Ngaydat
K001	David	Paris	P001	Hồn yêu	1996	3.2	C 20	Fafim	P001	K001	1/1/06
K002	Marie	Toulouse	P002	Đảo vô hình	2006	5.0	C 20	Fafim	P002	K001	1/2/06
K003	John	Nice	P003	Hoa hậu FBI	2002	3.2	Wanner	Films	P002	K002	5/1/06
K004	Helen	Bordeaux	P004	Taxi	2005	4.5	France	Fafim	P002	K003	1/2/06
K005	Taix	Toulouse	P005	Đời cát	2003	2.5	VFC	VFC	P003	K003	1/2/06
									P004	K001	2/5/06
									P004	K004	1/2/06
									P005	K005	4/1/06



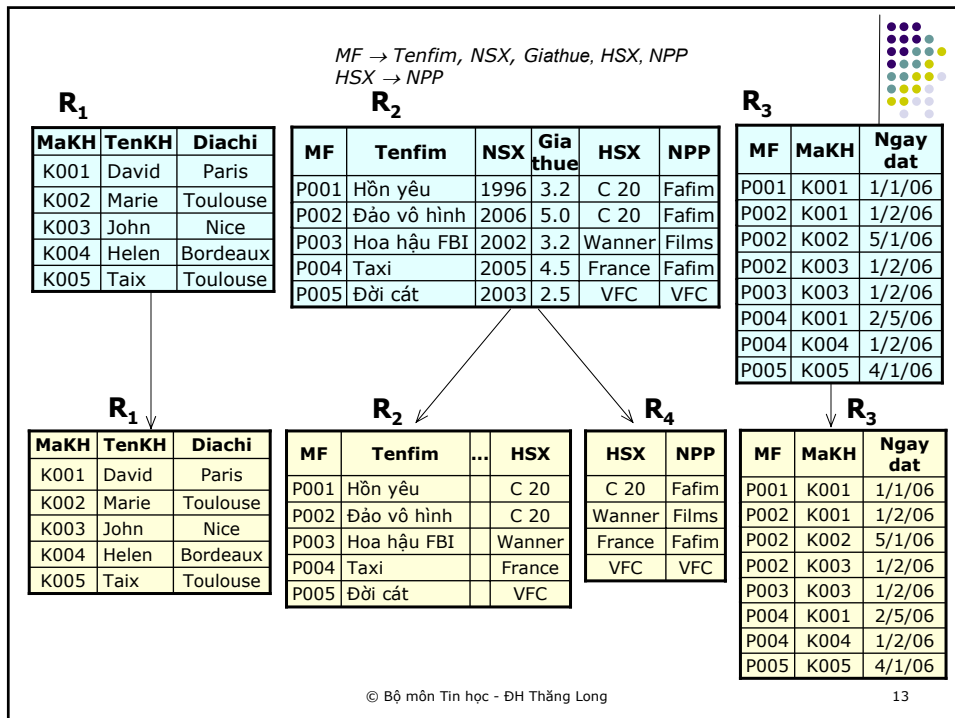
Các lỗi khi ở dạng 2NF

- Lược đồ $R_2(MF, Tenfim, NSX, Giathue, HSX, NPP)$ gồm các phụ thuộc hàm:
 - $MF \rightarrow Tenfim, NSX, Giathue, HSX, NPP$
 - $HSX \rightarrow NPP$
- R_2 đạt 2NF.
- Vấn đề
 - Có thể thêm nhà phân phối cho 1 hãng sản xuất nếu hãng đó chưa sản xuất bộ phim nào?
 - Nếu xoá thông tin 1 bộ phim có làm mất thông tin về hãng sản xuất và nhà phân phối không?
 - Khi thay đổi nhà phân phối cho 1 hãng sản xuất thì phải thực hiện thao tác đó bao nhiêu lần?



Dạng chuẩn 3 – 3NF

- Một quan hệ ở 3NF nếu
 - Quan hệ đã đạt 2NF
 - Không có các phụ thuộc bắc cầu.
- **Một lược đồ quan hệ R ở 3NF nếu với mọi phụ thuộc $X \rightarrow A$ đúng trong R và $A \notin X$ thì hoặc X là khóa bao hàm hoặc A là nguyên tố.**
- A là **thuộc tính nguyên tố (thuộc tính khóa)** nếu A là thuộc một khóa dự tuyển của R , ngược lại A là phi nguyên tố.
- VD: Cho lược đồ CSZ , $F = \{CS \rightarrow Z, Z \rightarrow C\}$, CS và SZ đều là khóa.
 \Rightarrow tất cả các thuộc tính là nguyên tố.



Các lỗi khi ở dạng 3NF (1)

- Ví dụ:
 - R(MF, Tenfim, NamSX, MaBS, Tinhtrang, HSX) và các mô tả như sau:
 - Mỗi fim có thể có nhiều bản sao
 - Mỗi bản sao của một phim có tình trạng khác nhau
 - Một phim chỉ có một tên, năm sản xuất và HSX.
- Từ thông tin trên ta có các phụ thuộc hàm sau:
 - MF → Tenfim, NamSX, HSX
 - MF, MaBS → Tinhtrang
- R có khoá (MF, MaBS)

Các lỗi khi ở dạng 3NF (2)



MF	Tenfilm	NSX	HSX	MBS	Tinhtrang
P001	Hồn yêu	1996	C 20	1	Mượn
P001	Hồn yêu	1996	C 20	2	Hỏng
P001	Hồn yêu	1996	C 20	3	Mượn
P002	Đảo vô hình	2006	C 20	1	Bán
P002	Đảo vô hình	2006	C 20	2	Rối
P003	Hoa hậu FBI	2002	Wanner	1	Rối

- R trên đạt 3NF.
- Lỗi phát sinh
 - Khi xóa bản sao 1 của phim P003 thì làm mất thông tin của phim này
 - Không thể thêm thông tin về 1 phim nếu phim đó chưa có bản sao nào

Dạng chuẩn Boyce-Codd - BCNF (1)



- Một lược đồ quan hệ R ở dạng BCNF nếu với mọi phụ thuộc $X \rightarrow A$ đúng trong R và $A \notin X$ thì X là khóa bao hàm.
- Một quan hệ ở BCNF thì cũng đạt 3NF.
- Trong thực hành các quan hệ đạt chuẩn 3NF là đủ. Tuy nhiên một quan hệ ở 3NF **không đảm bảo đã loại bỏ được tất cả các lỗi khi thao tác dữ liệu.**

Dạng chuẩn BCNF (2)



MF	Tenfim	NSX	HSX	MBS	Tinhtrang
P001	Hồn yêu	1996	C 20	1	Mượn
P001	Hồn yêu	1996	C 20	2	Hồng
P001	Hồn yêu	1996	C 20	3	Mượn
P002	Đảo vô hình	2006	C 20	1	Bán
P002	Đảo vô hình	2006	C 20	2	Rối
P003	Hoa hậu FBI	2002	Wanner	1	Rối

$\pi_{MF, Tenfim, NSX, HSX}(R)$

MF	Tenfim	NSX	HSX
P001	Hồn yêu	1996	C 20
P002	Đảo vô hình	2006	C 20
P003	Hoa hậu FBI	2002	Wanner

$\pi_{MF, MBS, Tinhtrang}(R)$

MF	MBS	Tinhtrang
P001	1	Mượn
P001	2	Hồng
P001	3	Mượn
P002	1	Bán
P002	2	Rối
P003	1	Rối

Thuật toán phân rã thành 3NF (1)



- Dùng để phân rã R thành (R_1, R_2, \dots, R_n) mà
 - R_i đạt 3NF
 - Phân rã bảo toàn phụ thuộc, và
 - Phân rã có nối không mất
- Thuật toán
 - IN: Lược đồ quan hệ R, tập phụ thuộc hàm F
 - OUT: Một phân rã bảo toàn phụ thuộc của R sao cho mỗi lược đồ quan hệ đều có dạng 3NF và ứng với hình chiếu của F trên lược đồ đó.

Thuật toán phân rã thành 3NF (2)



- Thuật toán phân rã thành 3NF bảo toàn phụ thuộc
- Tính F_C
 $m = 0$
IF $(\alpha \rightarrow Y \in F_C)$ AND $(\alpha \rightarrow Z \in F_C)$
 $\{F_C = F_C \setminus \{\alpha \rightarrow Y, \alpha \rightarrow Z\} \cup \{\alpha \rightarrow YZ\}$
 $\} /* \text{Hợp các phụ thuộc có cùng vế trái} */$
FOR $(X \rightarrow Y \in F_C)$ DO
 $\{ m = m+1$
 $R_m = XY$
 $\}$

Thuật toán phân rã thành 3NF (3)



- Ví dụ 1**
 - $R = ABCD, F = \{AB \rightarrow C, C \rightarrow D, BC \rightarrow D, CD \rightarrow B\}$
 - Phân rã thành 3NF
- Giải quyết
 - Tính $F_C = \{AB \rightarrow C, C \rightarrow D, C \rightarrow B\}$
 - $F_C = \{AB \rightarrow C, C \rightarrow BD\}$
 - Ta có $R_1 = ABC, R_2 = BCD$
 - Phân rã $\rho = (ABC, BCD)$ bảo toàn phụ thuộc

Thuật toán phân rã thành 3NF (4)



- Thuật toán cải tiến để phân rã có nối không mất

Tính F_C

$m = 0$

IF $(\alpha \rightarrow Y \in F_C)$ AND $(\alpha \rightarrow Z \in F_C)$

$\{F_C = F_C \setminus \{\alpha \rightarrow Y, \alpha \rightarrow Z\} \cup \{\alpha \rightarrow YZ\}$

$\}$ /* Hợp các phụ thuộc có cùng vế trái */

FOR $(X \rightarrow Y \in F_C)$ DO

$\{ m = m + 1$

$R_m = XY \}$

IF ($\exists R_i$ với $1 \leq i \leq m$ chứa 1 khoá dự tuyển của R)

$\{ m = m + 1$

$R_m = 1$ khoá dự tuyển bất kỳ của R }

Thuật toán phân rã thành 3NF (5)



• Ví dụ 2

- $R = CTHRSG$, $F = \{C \rightarrow T, HR \rightarrow C, HT \rightarrow R, CS \rightarrow G, HS \rightarrow R\}$
- Tìm phân rã đạt 3NF bảo toàn phụ thuộc và có nối không mất.
- Giải quyết
 - Ta có $F_C = F$
 - $K = HS$
 - $\rho = (CT, HRC, HTR, CSG, HSR)$ đạt 3NF
 - Thêm lược đồ HS vào phân rã ρ ; có HSR chứa khóa
→ Phân rã ρ bảo toàn phụ thuộc và có nối không mất

Thuật toán phân rã thành 3NF (6)



- **Ví dụ 3**

- $F = \{MF \rightarrow \text{Tenfim}, \text{NSX}, \text{Giathue}, \text{HSX}, \text{NPP}; \text{HSX} \rightarrow \text{NPP}\}$
- Tìm phân rã đạt 3NF bảo toàn phụ thuộc và có nối không mất.

- **Giải quyết**

- $F_C = \{MF \rightarrow \text{Tenfim}, \text{NSX}, \text{Giathue}, \text{HSX}; \text{HSX} \rightarrow \text{NPP}\}$
- $K = (MF, \text{HSX})$
- Phân rã gồm $R_1 = (MF, \text{Tenfim}, \text{NSX}, \text{Giathue}, \text{HSX})$ và $R_2 = (\text{HSX}, \text{NPP})$ đạt 3NF

Phân rã thành BCNF (1)



- **Bổ đề:**
- 1. Mỗi lược đồ có 2 thuộc tính đều có dạng BCNF.
- 2. Nếu R không có dạng BCNF thì ta có thể tìm được các thuộc tính A và B trong R sao cho $(R - AB) \rightarrow A$ đúng (có thể $(R - AB) \rightarrow B$ cũng đúng). Điều ngược lại chưa chắc đã đúng.



Thuật toán phân rã thành BCNF (1)

- Dùng để phân rã R thành (R_1, R_2, \dots, R_n) mà
 - R_i đạt BCNF
 - Phân rã có nối không mất
- Thuật toán:
 - IN: Lược đồ quan hệ R , tập phụ thuộc hàm F .
 - OUT: Một phân rã của R có nối không mất, sao cho mỗi lược đồ quan hệ trong phân rã có dạng BCNF ứng với hình chiếu của F trên lược đồ đó.



Thuật toán phân rã thành BCNF (2)

- Phương pháp: Phân rã lược đồ R thành 2 lược đồ:
 1. Lược đồ 1: có tập các thuộc tính XA , có dạng BCNF và phụ thuộc $X \rightarrow A$ đúng.
 2. Lược đồ 2: $R - A$
- Tiếp tục phân rã lược đồ $R - A$ theo các bước 1, 2 cho đến khi
 - không thể phân rã được nữa
 - lược đồ chỉ còn 2 thuộc tính

Thuật toán phân rã thành BCNF (3)



Chương trình chính:

```
BEGIN
  Z := R;
  REPEAT
    Phân rã Z thành Z-A và XA /*gọi TT phân rã D*/
    Thêm XA vào phân rã;
    Z := Z - A;
  UNTIL (không thể phân rã Z)
  Thêm Z vào phân rã
END.
```

Thuật toán phân rã thành BCNF (4)



• Thủ tục phân rã (D)

```
BEGIN
  IF Z không chứa A, B sao cho  $A \in (Z - AB)^+$  THEN
    return Z có dạng BCNF và không phân rã được
  ELSE
    BEGIN
      Tìm một cặp A và B; Y := Z;
      WHILE (Y chứa A và B sao cho  $(Y - AB) \rightarrow A$ ) DO
        Y := Y - B;
      return phân rã Z-A và Y /*Y là XA trong CT chính*/
    END;
  END;
```



Thuật toán phân rã thành BCNF (5)

- Cho lược đồ quan hệ CTHRSG và tập phụ thuộc $F = (C \rightarrow T, HR \rightarrow C, HT \rightarrow R, CS \rightarrow G, HS \rightarrow R)$
- Áp dụng thuật toán:
 - $Y = \text{CTHRSG}, A = C, B = T, C \in (\text{HRSG})^+ \Rightarrow Y = \text{CHRSRG}$
 - $Y = \text{CHRSRG}, A = R, B = C, R \in (\text{HSG})^+ \Rightarrow Y = \text{HRSRG}$
 - $Y = \text{HRSRG}, A = R, B = G, R \in (\text{HS})^+ \Rightarrow Y = \text{HRS}$
 - Không phân rã được nữa, **HRS** là một lược đồ trong phân rã.
 - $Z = \text{CTHRSG} - R = \text{CTHSG}$
 - Tiếp tục với $Y = \text{CTHSG}$



Thuật toán phân rã thành BCNF (6)

- $Y = \text{CTHSG}, A = T, B = H, T \in (\text{CSG})^+ \Rightarrow Y = \text{CTSG}$
- $Y = \text{CTSG}, A = T, B = S, T \in (\text{CG})^+ \Rightarrow Y = \text{CTG}$
- $Y = \text{CTG}, A = T, B = G, T \in (\text{C})^+ \Rightarrow Y = \text{CT}$
- CT không phân rã được nữa, **CT** là một lược đồ trong phân rã.
- $Z = \text{CHSG}$
- Tiếp tục với $Y = \text{CHSG}$
- ...
- Cuối cùng được phân rã (HRS, CT, CSG, CHS)
- **Chú ý:** Thứ tự chọn cặp A, B khác có thể thu được một phân rã khác, $\rho = (\text{CHS}, \text{THS}, \text{HSG}, \text{HSR})$.



Một vài chú ý về BCNF

- Với cùng lược đồ, có thể có nhiều phân rã khác nhau đạt BCNF
- Thuật toán chỉ sinh ra 1 trong các phân rã đó
- Phân rã BCNF có thể bảo toàn phụ thuộc
- Thứ tự chọn phụ thuộc hàm trong thuật toán có thể sinh ra phân rã khác nhau