

Database design-

You are given a

big relation schema R

and a set of fds F involving attributes of R .

Need to come up with a set of probably smaller schemas that can store the same information.

Consider:

Enrol = (s-id, c-id, grade, gpa)

$F = \{ \begin{array}{l} s\text{-id}, c\text{-id} \rightarrow \text{grade} \\ s\text{-id} \rightarrow \text{gpa} \end{array} \}$

What is the problem with this schema?

If a student is enrolled in two or more courses, student's gpa is repeated.

"Duplication of information, i.e. redundancy"

Bad for updates - Need to keep them consistent

If a student is not enrolled in any course but has a 'gpa', how to store this information?

"Inability to represent certain info" ?

Normal Forms:

Tell you when a relation schema is "good"
with respect to a set of fds.

BCNF : Boyce Codd Normal Form:

Def: A relation schema R is in BCNF with respect to
a set of fds F if \forall ^{one of the} following conditions
satisfied for every fd $\alpha \rightarrow \beta$ in F^+ :

(a) $\alpha \rightarrow \beta$ is a trivial fd (i.e., $\beta \subseteq \alpha$)

(b) α is a super key of R ,
i.e., $\alpha \rightarrow R$ is in F^+

examples: ex 1: Accounts = (acct#, bal, br-name)

$F = \{ \text{acct\#} \rightarrow \text{bal, br-name} \}$

ex 2:

Accounts = (acct#, bal, br-name, cust-name)

$F = \{ \text{acct\#} \rightarrow \text{bal, br-name, cust-name} \}$

ex 3: $\text{Enrol} = (S\text{-id}, C\text{-id}, \text{grade}, \text{gpa})$

$F = \{ S\text{-id}, C\text{-id} \rightarrow \text{grade}, \\ S\text{-id} \rightarrow \text{gpa} \}$

not in BCNF

ex 4: $\text{Account-Deposit} = (\text{Cust-name}, \text{acct\#}, \text{bal}, \text{br-name})$

$F = \{ \text{acct\#} \rightarrow \text{bal}, \text{br-name} \}$

not in BCNF

What should be done if not in BCNF?

Split the relation schema, i.e. decompose
into smaller schemas.

Decomposition: A decomposition of R is a set of relation schemas (R_1, R_2, \dots, R_k) such that

$$R = R_1 \cup R_2 \cup \dots \cup R_k.$$

Given a relation schema R and a set of fds F ,

when is a decomposition good?

Ex: Accounts-deposit = (cust-name, Acct#, bal, br-name)

$$F = \{ \text{Acct\#} \rightarrow \text{bal, br-name} \}$$

Consider the decomposition:

$$R_1 = (\text{cust-name, br-name}), \quad R_2 = (\text{Acct\#, bal, br-name})$$

is this a good decomposition?

No.

Why?

Consider the instance of Accounts-deposit

cust-name	Acct#	bal	br-name
John	100	10K	UIC
Mary	100	10K	UIC
Jeff	200	5K	UIC

instance of R_1

cust-name	br-name
John	UIC
Mary	UIC
Jeff	UIC

instance of R_2

Acct#	bal	br-name
100	10K	UIC
200	5K	UIC

The only way to get the original instance of Accounts-deposit is to take the join of the two instances of

R_1, R_2 .

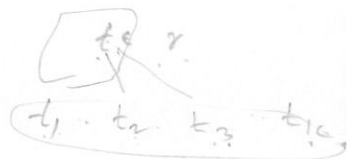
when The result of join:

	Cust-name	Acct#	bal	br-name
	John	100	10K	UIC
→	John	200	5K	UIC
	Mary	100	10K	UIC
	Mary	200	5K	UIC
→	Jeff	100	10K	UIC
	Jeff	200	5K	UIC

get additional tuples.

Means loss of information!!

Loss-less Join decomposition:



Given R, F .

Def: A decomposition (R_1, R_2, \dots, R_k) of R is a loss less join decomposition ^{w.r.t. F} if for every legal instance r of R ^{w.r.t. F} the following condition holds:

$$\pi_{R_1}(r) \bowtie \pi_{R_2}(r) \bowtie \dots \bowtie \pi_{R_k}(r) = r. \quad (1)$$

(Note that the left side in (1) is always a Super Set of the right side)

Clearly, the above does not hold for the previous

example.

Consider the following decomposition of Accounts-deposit

$R_1 = (\text{Cust-name}, \text{Acct\#})$

$R_2 = (\text{Acct\#}, \text{bal}, \text{br-name})$

Recall $F = \{ \text{Acct\#} \rightarrow \text{bal}, \text{br-name} \}$

In this case, the decomposition is a loss less join decomposition.

why?

Sufficient Condition for loss-less join property:

Thm: Given R, F and given the decomposition (R_1, R_2) of R , this decomposition satisfies the loss-less join property if

$R_1 \cap R_2 \rightarrow R_1$ is in F^+ (i.e., $R_1 \cap R_2$ is a super key of R_1)

or
 $R_1 \cap R_2 \rightarrow R_2$ is in F^+ (i.e., $R_1 \cap R_2$ is a super key of R_2).

Proof: Suppose r is a legal instance of R and $R_1 \cap R_2 \rightarrow R_1$ is in F^+ , i.e., $R_1 \cap R_2$ is a super key of R_1 .

Now consider the join of the tables

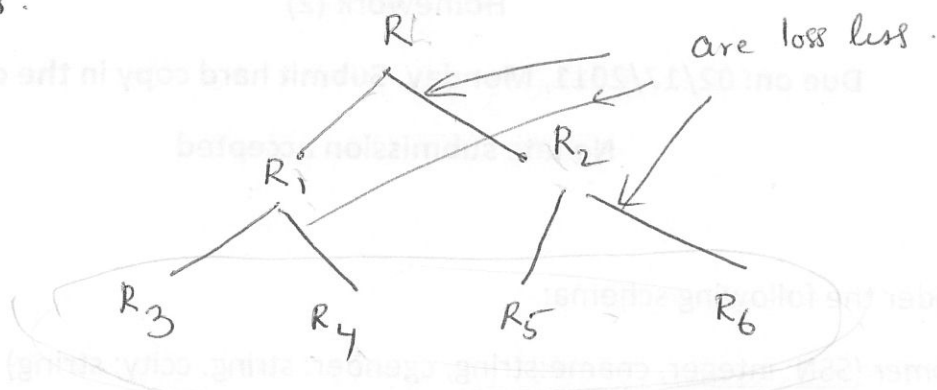
$\pi_{R_1}(r)$ and $\pi_{R_2}(r)$ there is exactly

For every tuple t_2 in $\pi_{R_2}(r)$ there is exactly one tuple t_1 in $\pi_{R_1}(r)$ with which t_2 joins.

Hence we get ~~there~~ ^{exactly r , i.e.,} will not be any

extra tuples in the result.

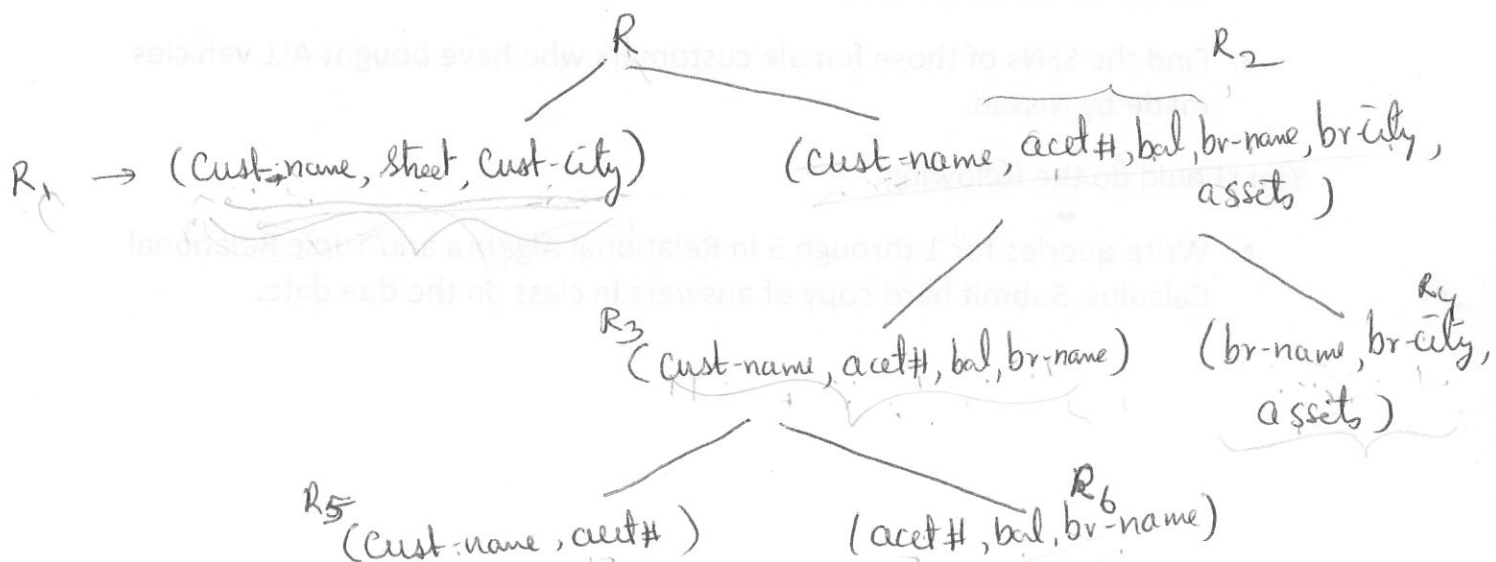
Suppose we perform a sequence of decompositions each of which is loss less than the final decomposition is loss less.



The final decomposition is (R_3, R_4, R_5, R_6)

Ex: $R = (\text{cust-name}, \text{street}, \text{cust-city}, \text{acct\#}, \text{bal}, \text{br-name}, \text{br-city}, \text{assets})$

$F = \left\{ \begin{array}{l} \text{cust-name} \rightarrow \text{street}, \text{cust-city} \\ \text{acct\#} \rightarrow \text{bal}, \text{br-name} \\ \text{br-name} \rightarrow \text{br-city}, \text{assets} \end{array} \right\}$



Final decomposition (R_1, R_5, R_6, R_4) is a loss-less join decomp.

"Loss-less join property" is critical and a must.

Another property called "dependency preservation" is desirable.

Given

R, F

consider a decomposition (R_1, R_2, \dots, R_k) of R .

Let F_1 be the set of all fds in F^+ whose attributes are from R_1

F_2 be the set of all fds in F^+ whose attributes are from R_2 ; i.e., $F_i = \{ \alpha \rightarrow \beta \in F^+ \mid \alpha, \beta \subseteq R_i \}$

Thus we have

F_1 on R_1

F_2 on R_2

F_k on R_k

The decomposition is dependency preserving if

$$F \subseteq (F_1 \cup F_2 \cup \dots \cup F_k)^+$$

Ex: $R = (\text{City}, \text{Street}, \text{ZipCode})$

$F = \{ \text{City}, \text{Street} \rightarrow \text{ZipCode} \\ \text{ZipCode} \rightarrow \text{City} \}$

$R_1 = (\text{Street}, \text{ZipCode})$, $R_2 = (\text{City}, \text{ZipCode})$

$F_2 = \{ \text{ZipCode} \rightarrow \text{City} \} \cup \text{Trivial fds.}$

$F_1 = \text{Trivial fds}$

$F \not\subseteq (F_1 \cup F_2)^+$. Is not dependency preserving?

Algorithm for obtaining a loss-less join and BCNF

decomposition:

Input: R, F

Output: a loss-less join and BCNF decomposition.
(i.e., all relation schemas in the output are in BCNF).

final output in 'result'.

1. Compute F^+

2. $result := \{R\}$

3. $done := False;$

While not done

{ if $\exists R_i \in result$ such that R_i is not in BCNF

{ let $\alpha \rightarrow \beta \in F^+$ be such that
 $\alpha, \beta \subseteq R_i$, $\alpha \cap \beta = \emptyset$ and
 α is not a super key of R_i
(i.e., $\alpha \rightarrow R_i \notin F^+$);

$result := result - \{R_i\} \cup \{\alpha \beta\} \cup \{R_i - \beta\}$

else

$done := true;$

}

Features of the alg.:

Starts with $\text{result} = \{R\}$

- At each step chooses a relation's schema R_i in result that is not in BCNF.
- It chooses an "offending" fd $\alpha \rightarrow \beta \in F^+$ such that $\alpha \cap \beta = \emptyset$, α is not a super key of R_i ;

It replaces R_i by $(R_i - \beta)$ and $(\alpha\beta)$
(i.e., it decomposes R_i into $(R_i - \beta)$, $(\alpha\beta)$

(the decomposition is loss-less because α is a super key of $(\alpha\beta)$).

- It stops when all the relation schemas in result are in BCNF.

Ex:

$R = (\text{cust-name}, \text{acct\#}, \text{bal}, \text{br-name}, \text{br-city}, \text{assets}, \text{cust-city}, \text{street})$

$F = \{ \text{cust-name} \rightarrow \text{cust-city}, \text{street}$
 $\text{acct\#} \rightarrow \text{bal}, \text{br-name}$
 $\text{br-name} \rightarrow \text{br-city}, \text{assets} \}$

result = $\{ R \}$

R not in BCNF

1. Choose the fd: $\text{cust-name} \rightarrow \text{cust-city}, \text{street}$

\neq result = $\{ \overbrace{(\text{cust-name}, \text{acct\#}, \text{bal}, \text{br-name}, \text{br-city}, \text{assets})}^{R_1}, \underbrace{(\text{cust-name}, \text{cust-city}, \text{street})}_{R_2} \}$

The first relation R_1 not in BCNF;

Choose the fd: $\text{acct\#} \rightarrow \text{bal}, \text{br-name}$ $\text{br-name} \rightarrow \text{br-city}, \text{assets}$

2.

decompose R_1 into

$R_3 = (\text{cust-name}, \text{acct\#}, \text{bal}, \text{br-name})$

$R_4 = (\text{br-name}, \text{br-city}, \text{assets})$

R_3 not in BCNF

Choose $\text{acct\#} \rightarrow \text{bal}, \text{br-name}$

3.

decompose R_3 into

$R_5 = (\text{cust-name}, \text{acct\#})$

$R_6 = (\text{acct\#}, \text{bal}, \text{br-name})$

Final 'result': $\{ R_2, R_4, R_5, R_6 \}$

Does the decomposition satisfy
"dependency preservation"?

$$F_2 = \{ \text{cust-name} \rightarrow \text{cust-city, street} \} \cup \text{Trivial fds}$$

$$F_4 = \{ \text{br-name} \rightarrow \text{br-city, assets} \} \cup \text{Trivial fds}$$

$$F_5 = \text{Trivial fds}$$

$$F_6 = \{ \text{acct\#} \rightarrow \text{bal, br-name} \} \cup \text{Trivial fds}$$

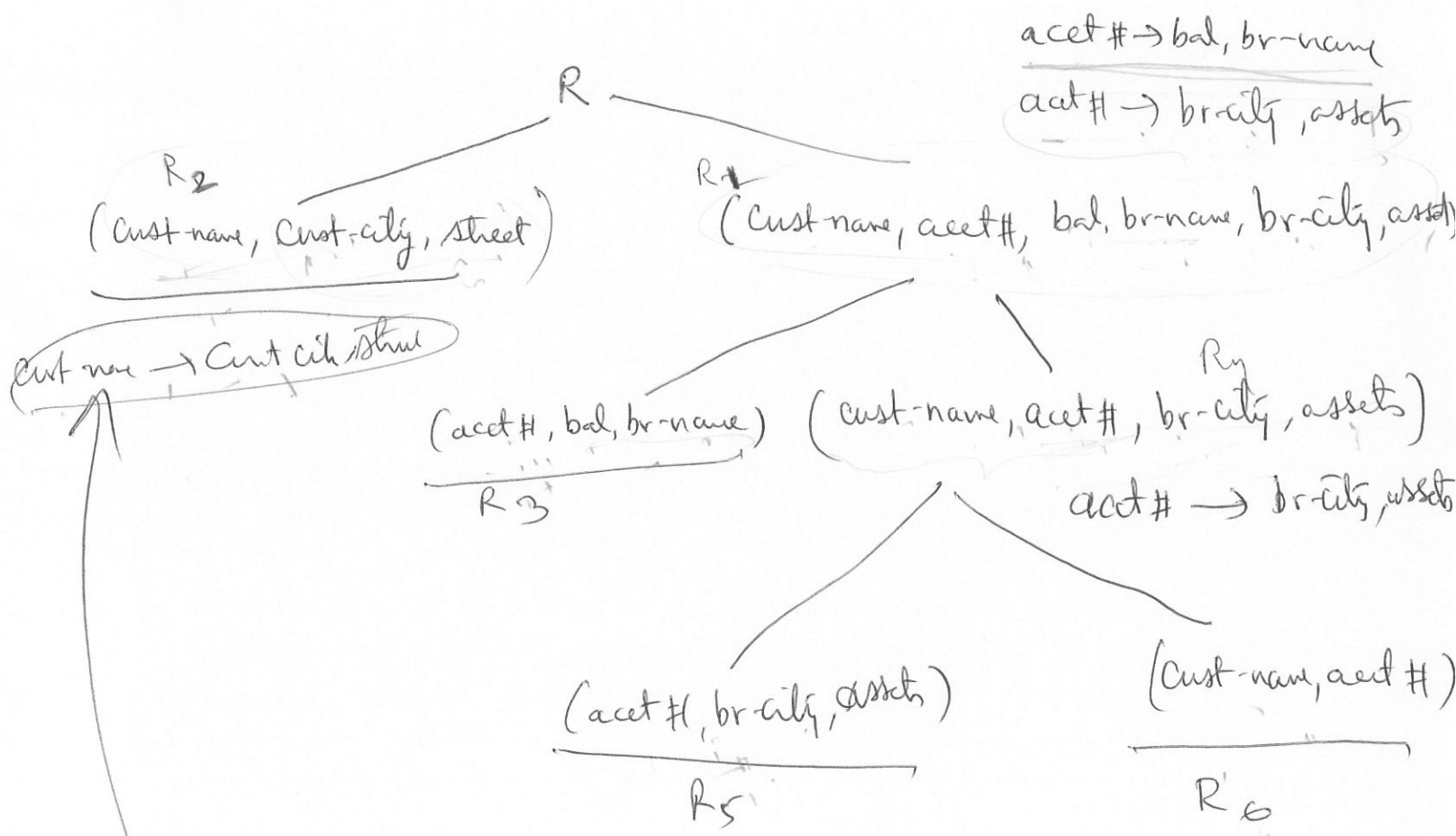
Clearly $F \subseteq (F_1 \cup F_2 \cup F_3 \cup F_4)^+$

Hence dependency preserving.

Suppose we choose the fd

$$\text{acct\#} \rightarrow \text{bal, br-name}$$

in the 2nd step, we get a different decomp.



F_2

R_2, R_3, R_5, R_6

$F_3 = \{ \text{acct\#} \rightarrow \text{bal}, \text{br-name} \}$

$F_5 = \{ \text{acct\#} \rightarrow \text{br-city}, \text{assets} \}$

$F_6 = \{ \text{Cust-name} \rightarrow \text{acct\#} \}$

$\text{br-name} \rightarrow \text{br-city}, \text{assets}$

The previous algorithm does not guarantee
dependency preservation!

For a given R and F : Can we get a loss-less, BCNF and
dependency preserving decomposition

Consider the following example:

$R = (\text{City}, \text{Street}, \text{Zipcode})$

$F = \{ \text{City}, \text{Street} \rightarrow \text{Zipcode}$
 $\text{Zipcode} \rightarrow \text{City} \}$

Clearly R is not in BCNF.

Consider the following decomposition:

$R_1 = (\text{Zipcode}, \text{City})$ $R_2 = (\text{Street}, \text{Zipcode})$

$F_1 = \{ \text{Zipcode} \rightarrow \text{City} \}$ trivial $F_2 = \emptyset$ trivial

It is a loss-less join decomposition.

Both R_1, R_2 are depend in BCNF.

It is not dependency preserving.

BCNF ~~pro~~ normal form is too strong !!

3rd Normal Form:

Def: Relation Schema R is in 3NF with respect to a set F of fds, if for every $\alpha \rightarrow \beta$ in F^+ , one at least one of the following conditions is satisfied:

- (i) $\alpha \rightarrow \beta$ is a trivial fd (i.e., $\beta \subseteq \alpha$)
- (ii) α is a super key of R (i.e., $\alpha \rightarrow R \in F^+$)
- (iii) each attribute in $\beta - \alpha$ is contained in a candidate key of R .
(Called a prime attribute)

Ex1: The relation schema $R = (\text{City}, \text{Street}, \text{Zipcode})$
 $F = \{ \text{City}, \text{Street} \rightarrow \text{Zipcode}, \text{Zipcode} \rightarrow \text{City} \}$
is in 3NF with resp. F .

Ex2: $R = (\text{Cust-name}, \text{br-name}, \text{banker-name})$
 $F = \{ \text{Cust-name}, \text{br-name} \rightarrow \text{banker-name}, \text{banker-name} \rightarrow \text{br-name} \}$
is in 3NF (not in BCNF)

3NF allows some duplication of information

Any R that is in BCNF with resp. to F
is in 3NF with resp. to F .

The Converse is not always true!

We Given R, F ; we can always get
a decomposition that ~~is~~ satisfies loss-less join
property, that is dependency preserving
and all relations schemas are in 3NF.

~~Ideally~~
Given R, F , ideally we like to get a decomposition
Satisfying ~~the~~ (i.e., every relation schema is in
BCNF)

- (i) BCNF property
- (ii) loss-less join property
- (iii) dependency preserving.

If it is not possible we get a decomposition
Satisfying (i) 3NF property (i.e., every schema is in 3NF)
and Satisfying properties (ii) and (iii).

$R = (\text{cust-name}, \text{banker-name}, \text{br-name}, \text{cust-city}, \text{cust-street})$

$F = \{ \text{cust-name} \rightarrow \text{cust-city}, \text{cust-street}$
 $\text{banker-name} \rightarrow \text{br-name}$
 $\text{cust-name}, \text{br-name} \rightarrow \text{banker-name} \}$

$R_1 = (\text{cust-name}, \text{cust-city}, \text{cust-street})$ ✓

$R_2 = (\text{cust-name}, \text{br-name}, \text{banker-name})$

The above decomposition is a 3NF decomposition
satisfies loss-less join property and
dependency preservation.