

Introduction to Data Management

CSE 344

Lectures 18: BCNF

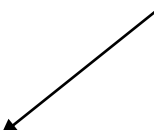
Announcements

- Midterm is posted and grades are up
 - Have until 8pm tonight for regrade requests
- WQ6 - 11pm July 31 (Monday)
- HW6 is up due Aug 1 (Tuesday)
 - Database design
- Will drop HW8 - not enough time to get to Spark lecture before it is due.

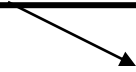
Review: Relation Decomposition

Break the relation into two:

Name	SSN	PhoneNumber	City
Fred	123-45-6789	206-555-1234	Seattle
Fred	123-45-6789	206-555-6543	Seattle
Joe	987-65-4321	908-555-2121	Westfield



Name	<u>SSN</u>	City
Fred	123-45-6789	Seattle
Joe	987-65-4321	Westfield



<u>SSN</u>	<u>PhoneNumber</u>
123-45-6789	206-555-1234
123-45-6789	206-555-6543
987-65-4321	908-555-2121

Anomalies have gone:

- No more repeated data
- Easy to move Fred to “Bellevue” (how ?)
- Easy to delete all Joe’s phone numbers (how ?)

Review: Functional Dependencies (FDs)

Definition $\forall t, t' \in R, A_1, \dots, A_m \rightarrow B_1, \dots, B_n$ holds in R if:
 $(t.A_1 = t'.A_1 \ \dots \ t.A_m = t'.A_m \rightarrow t.B_1 = t'.B_1 \ \dots \ t.B_n = t'.B_n)$

R	A_1	...	A_m		B_1	...	B_n		
t									
t'									

if t, t' agree here then t, t' agree here

Closure Algorithm

$X = \{A_1, \dots, A_n\}$.

Repeat until X doesn't change do:
if $B_1, \dots, B_n \rightarrow C$ is a FD **and**
 B_1, \dots, B_n are all in X
then add C to X.

Example:

1. name \rightarrow color
2. category \rightarrow department
3. color, category \rightarrow price

$\{\text{name, category}\}^+ =$
 $\{\text{name, category, color, department, price}\}$

Hence: $\text{name, category} \rightarrow \text{color, department, price}$

Keys

- A **superkey** is a set of attributes A_1, \dots, A_n s.t. for any other attribute B in the same relation, we have $A_1, \dots, A_n \rightarrow B$
- A **key** is a minimal superkey
 - A superkey and for which no subset is a superkey

Computing (Super)Keys

- For all sets X , compute X^+
- If $X^+ = [\text{all attributes}]$, then X is a superkey
- Try reducing to the minimal X 's to get the key

Example

Product(name, price, category, color)

name, category → price
category → color

What is the key ?

Example

Product(name, price, category, color)

name, category \rightarrow price
category \rightarrow color

What is the key ?

$(\text{name, category})^+ = \{ \text{name, category, price, color} \}$

Hence (name, category) is a key

How To Eliminate Anomalies

Main idea:

- $X \rightarrow A$ is OK if X is a (super)key
- $X \rightarrow A$ is not OK otherwise
 - Need to decompose the table, but how?

Boyce-Codd Normal Form

Boyce-Codd Normal Form

There are no
“bad” FDs:

Definition. A relation R is in BCNF if:

Whenever $X \rightarrow B$ is a non-trivial dependency,
then X is a superkey.

Equivalently:

Definition. A relation R is in BCNF if:

$\forall X$ in $X \rightarrow B$,
either $X^+ = X$ or $X^+ = [\text{all attributes}]$

BCNF Decomposition Algorithm

Normalize(R)

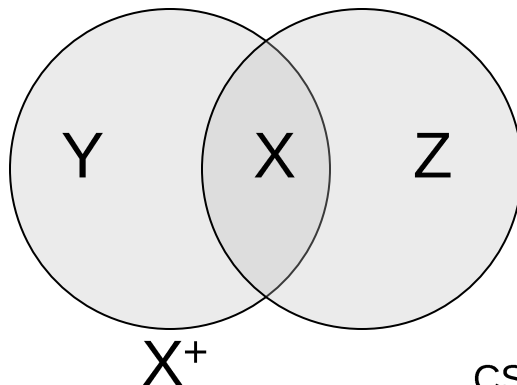
find X in $X \rightarrow B$ s.t.: $X \neq X^+$ and $X^+ \neq [\text{all attributes}]$

if (not found) **then** R is in BCNF \cup

let $Y = X^+ - X$; $Z = [\text{all attributes}] - X^+$

decompose R into $R_1(X \cup Y)$ and $R_2(X \cup Z)$

Normalize(R_1); Normalize(R_2);

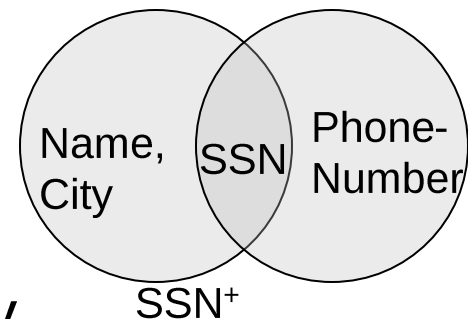


Want X in $X \rightarrow B$ s.t.: $X = X^+$ or $X^+ = [\text{all attributes}]$

Example

Name	SSN	PhoneNumber	City
Fred	123-45-6789	206-555-1234	Seattle
Fred	123-45-6789	206-555-6543	Seattle
Joe	987-65-4321	908-555-2121	Westfield
Joe	987-65-4321	908-555-1234	Westfield

$SSN \rightarrow \text{Name, City}$



The only key is: $\{SSN, \text{PhoneNumber}\}$

Hence $SSN \rightarrow \text{Name, City}$ is a “bad” dependency

In other words:

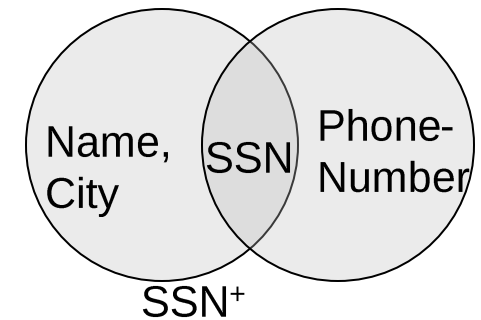
$SSN^+ = SSN, \text{Name, City}$ and is neither SSN nor All Attributes

Want X in $X \rightarrow B$ s.t.: $X = X^+$ or $X^+ = [\text{all attributes}]$

Example BCNF Decomposition

<u>Name</u>	<u>SSN</u>	<u>City</u>
Fred	123-45-6789	Seattle
Joe	987-65-4321	Westfield

$SSN \rightarrow \text{Name, City}$



<u>SSN</u>	<u>PhoneNumber</u>
123-45-6789	206-555-1234
123-45-6789	206-555-6543
987-65-4321	908-555-2121
987-65-4321	908-555-1234

Let's check anomalies:

- Redundancy ?
- Update ?
- Delete ?

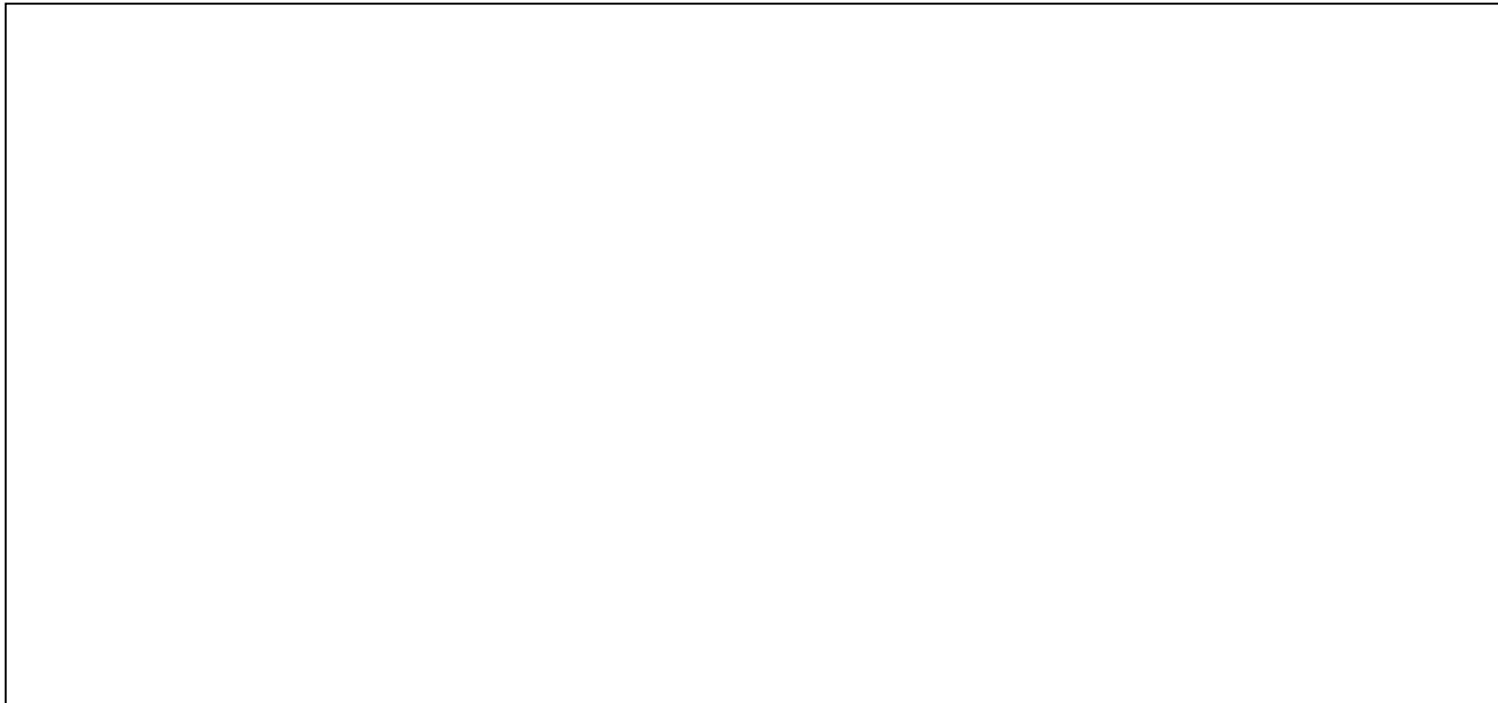
Find X in $X \rightarrow B$ s.t.: $X \neq X^+$ and $X^+ \neq [\text{all attributes}]$

Example BCNF Decomposition

Person(name, SSN, age, hairColor, phoneNumber)

SSN \rightarrow name, age

age \rightarrow hairColor



Find X in $X \rightarrow B$ s.t.: $X \neq X^+$ and $X^+ \neq [\text{all attributes}]$

Example BCNF Decomposition

Person(name, SSN, age, hairColor, phoneNumber)

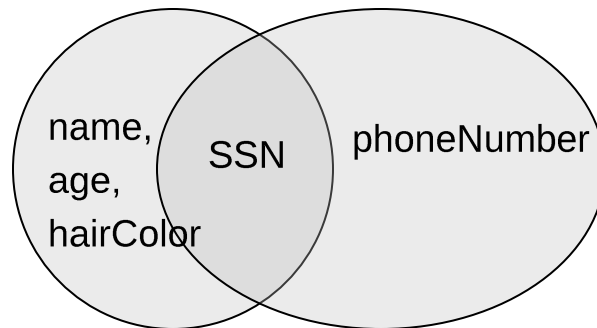
SSN \rightarrow name, age

age \rightarrow hairColor

Iteration 1: **Person**: SSN⁺ = SSN, name, age, hairColor

Decompose into: **P**(SSN, name, age, hairColor)

Phone(SSN, phoneNumber)



Find X in $X \rightarrow B$ s.t.: $X \neq X^+$ and $X^+ \neq [\text{all attributes}]$

Example BCNF Decomposition

Person(name, SSN, age, hairColor, phoneNumber)

SSN \rightarrow name, age

age \rightarrow hairColor

What are
the keys ?

Iteration 1: **Person**: SSN⁺ = SSN, name, age, hairColor

Decompose into: **P**(SSN, name, age, hairColor)

Phone(SSN, phoneNumber)

Iteration 2: **P**: age⁺ = age, hairColor

Decompose: **People**(SSN, name, age)

Hair(age, hairColor)

Phone(SSN, phoneNumber)

Find X in $X \rightarrow B$ s.t.: $X \neq X^+$ and $X^+ \neq [\text{all attributes}]$

Example BCNF Decomposition

Person(name, SSN, age, hairColor, phoneNumber)

SSN \rightarrow name, age

age \rightarrow hairColor

Note the keys!

Iteration 1: **Person**: SSN⁺ = SSN, name, age, hairColor

Decompose into: **P**(SSN, name, age, hairColor)

Phone(SSN, phoneNumber)

Iteration 2: **P**: age⁺ = age, hairColor

Decompose: **People**(SSN, name, age)

Hair(age, hairColor)

Phone(SSN, phoneNumber)

$R(A,B,C,D)$

$A \rightarrow B$
 $B \rightarrow C$

Example: BCNF

Recall: find X s.t.
 $X \neq X^+ \neq [\text{all-attrs}]$

$R(A,B,C,D)$

$R(A,B,C,D)$

$A \rightarrow B$
 $B \rightarrow C$

Example: BCNF

Recall: find X s.t.
 $X \neq X^+ \neq [\text{all-attrs}]$

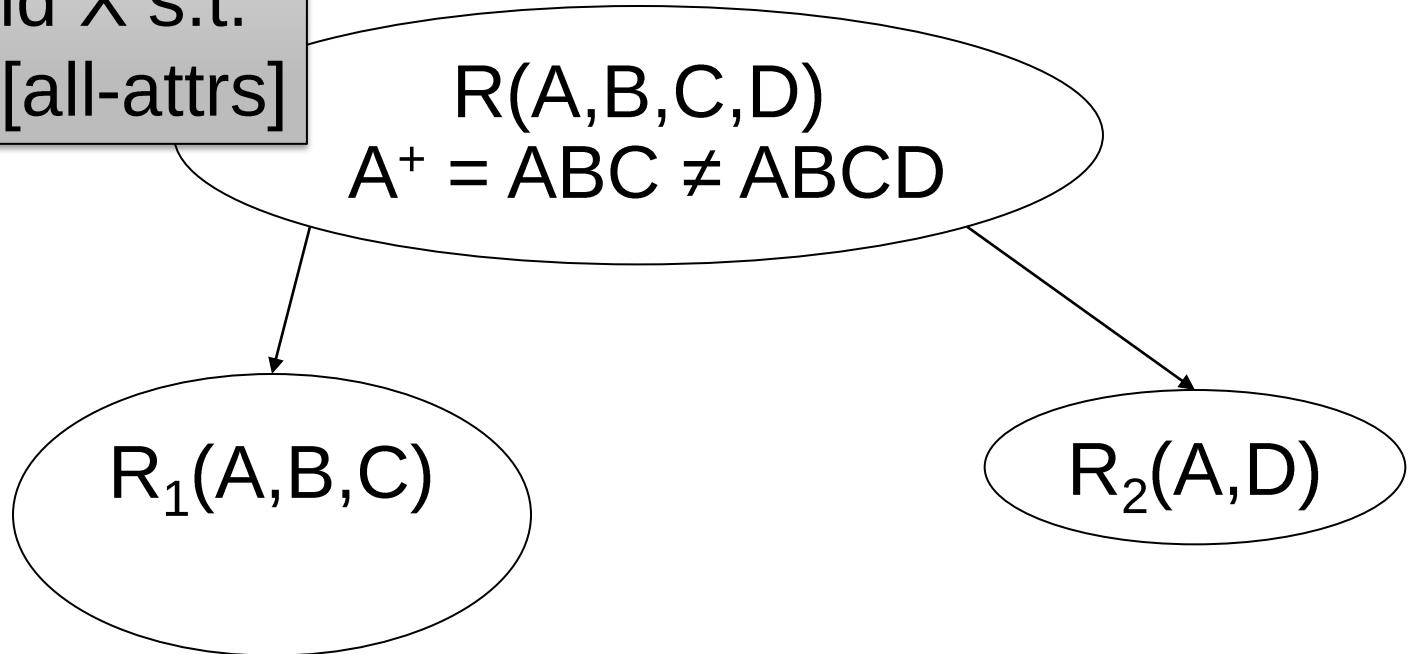
$R(A,B,C,D)$
 $A^+ = ABC \neq ABCD$

$R(A,B,C,D)$

$A \rightarrow B$
 $B \rightarrow C$

Example: BCNF

Recall: find X s.t.
 $X \neq X^+ \neq [\text{all-attrs}]$

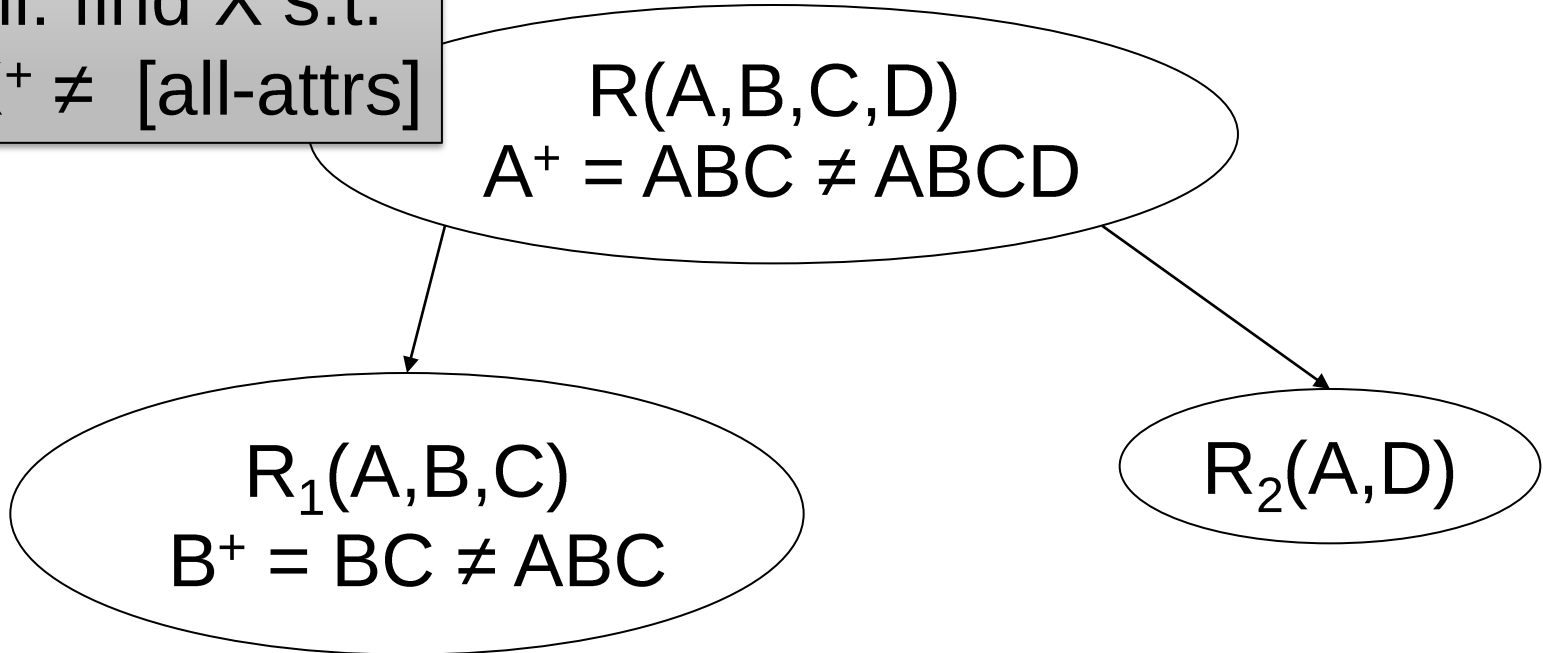


$R(A,B,C,D)$

$A \rightarrow B$
 $B \rightarrow C$

Example: BCNF

Recall: find X s.t.
 $X \neq X^+ \neq [\text{all-attrs}]$

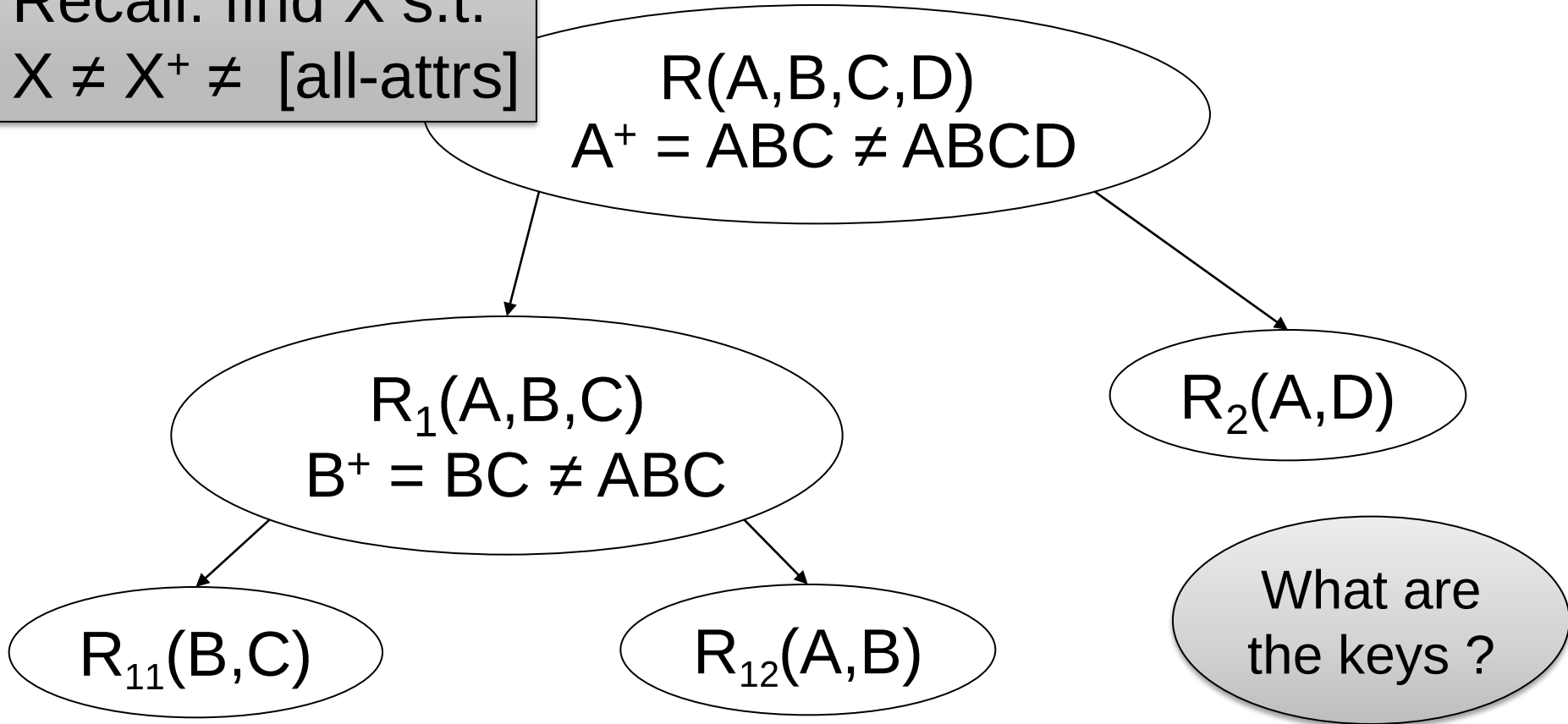


$R(A,B,C,D)$

$A \rightarrow B$
 $B \rightarrow C$

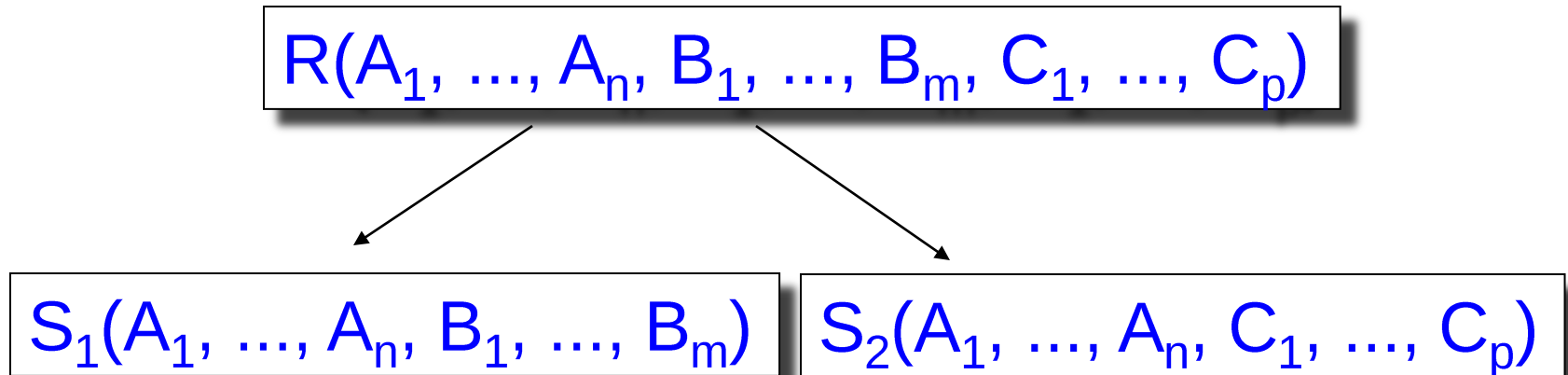
Example: BCNF

Recall: find X s.t.
 $X \neq X^+ \neq [\text{all-attrs}]$



What happens in R we first pick B^+ ? Or AB^+ ?

Decompositions in General



S_1 = projection of R on $A_1, \dots, A_n, B_1, \dots, B_m$

S_2 = projection of R on $A_1, \dots, A_n, C_1, \dots, C_p$

Not all decompositions are good - some are lossy

A Lossless Decomposition

Name	Price	Category
Gizmo	19.99	Gadget
OneClick	24.99	Camera
Gizmo	19.99	Camera

Name	Price
Gizmo	19.99
OneClick	24.99
Gizmo	19.99

Name	Category
Gizmo	Gadget
OneClick	Camera
Gizmo	Camera

A Lossy Decomposition

What is lossy here?

Name	Price	Category
Gizmo	19.99	Gadget
OneClick	24.99	Camera
Gizmo	19.99	Camera

Name	Category
Gizmo	Gadget
OneClick	Camera
Gizmo	Camera

Price	Category
19.99	Gadget
24.99	Camera
19.99	Camera

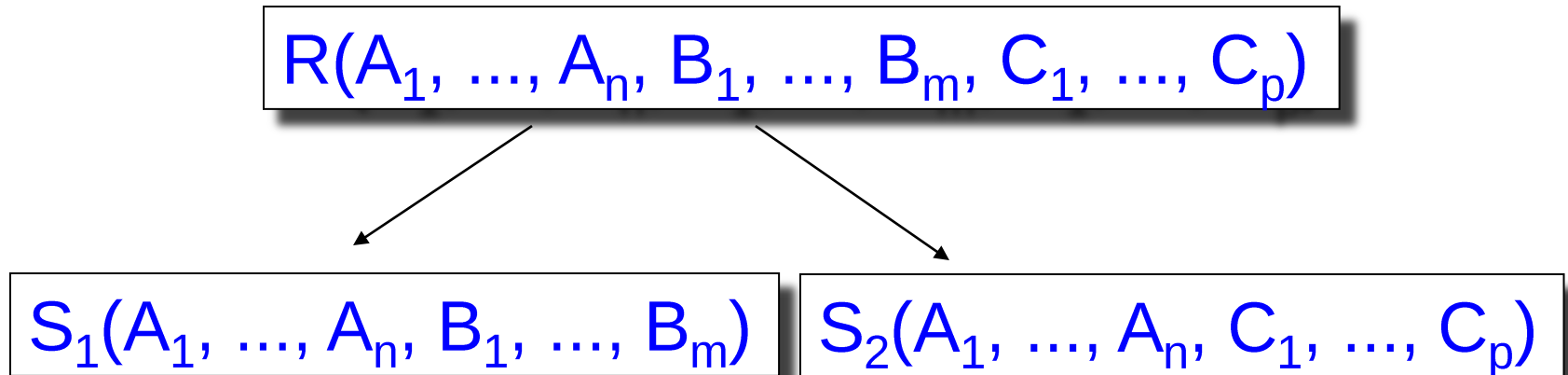
Lossy Decomposition

Name	Price	Category
Gizmo	19.99	Gadget
OneClick	24.99	Camera
Gizmo	19.99	Camera

Name	Category
Gizmo	Gadget
OneClick	Camera
Gizmo	Camera

Price	Category
19.99	Gadget
24.99	Camera
19.99	Camera

Decomposition in General



Let: S_1 = projection of R on $A_1, \dots, A_n, B_1, \dots, B_m$

S_2 = projection of R on $A_1, \dots, A_n, C_1, \dots, C_p$

The decomposition is called lossless if $R = S_1 \bowtie S_2$

Fact: If $A_1, \dots, A_n \rightarrow B_1, \dots, B_m$ then the decomposition is lossless

It follows that every BCNF decomposition is lossless

Schema Refinements = Normal Forms

- 1st Normal Form = all tables are flat
- 2nd Normal Form = obsolete
- Boyce Codd Normal Form = no bad FDs
- 3rd and 4th Normal Form = see book
 - BCNF is lossless but can cause loss of ability to check some FDs (see book 3.4.4)
 - 3NF fixes that (is lossless and dependency-preserving), but some tables might not be in BCNF – i.e., they may have redundancy anomalies

Normalization IRL

- It's good to know about BCNF (and maybe 3NF).
- Should every table always be in BCNF?

R(Item, OrderNo, Customer, Address, Zip Code, State)

What are the FDs:

Normalization IRL

R(Item, OrderNo, Customer, Address, Zip Code, State)

What are the FDs:

OrderNo -> Customer

Customer -> Address, Zip Code, State

Zip Code -> State

What is BCNF:

Normalization IRL

R(Item, OrderNo, Customer, Address, Zip Code, State)

What is BCNF:

R(Item, OrderNo, Customer, Address, Zip Code, State)
OrderNo⁺ = Customer, Address, Zip, Code, State

R₁(Item, OrderNo)

R₂(OrderNo, Customer, Address, Zip Code, State)

R₃(OrderNo, CustomerId)

R₄(CustomerId, Name, Address, Zip Code, State)

Are We Done?

Normalization IRL

R(Item, OrderNo, Customer, Address, Zip Code, State)

What is BCNF:

R(Item, OrderNo, Customer, Address, Zip Code, State)
OrderNo⁺ = Customer, Address, Zip, Code, State

R₁(Item, OrderNo)

R₂(OrderNo, Customer, Address, Zip Code, State)

R₃(OrderNo, CustomerId)

R₄(CustomerId, Name, Address, Zip Code)

R₅(Zip Code, State)

Is this a good idea?