



Schema Refinement and Normal Forms

Chapter 19

The Evils of Redundancy



- ❖ *Redundancy* is at the root of several problems associated with relational schemas

SSN	Name	Lot	Rate	W	H
123223666	Attishoo	48	8	10	40
23315368	Smiley	22	8	10	30
131243650	Smethurst	35	5	7	30
434263751	Guldu	35	5	7	32
612674134	Madayan	35	8	10	40

- ❖ *Functional dependency*: $R \rightarrow W$. *R determines W*
- ❖ Problems due to $R \rightarrow W$:
 - *Update anomaly*: Can we change W in just the first tuple
 - *Insertion anomaly*: What if we want to insert an employee and don't know the hourly wage for his rating?
 - *Deletion anomaly*: If we delete all employees with rating 5, we lose the information about the wage for rating 5!



Notations

- ❖ Consider relation obtained from Hourly_Emps:
 - Hourly_Emps (ssn, name, lot, rating, hrly_wages, hrs_worked)
- ❖ Notation: We will denote this relation schema by listing the attributes: **SNLRWH**
 - This is really the *set* of attributes {S,N,L,R,W,H}.
- ❖ Some FDs on Hourly_Emps:
 - *ssn is the key*: $S \rightarrow \text{SNLRWH}$
 - *rating determines hrly_wages*: $R \rightarrow W$



Example: Decomposition

S	N	L	R	W	H
123-22-3666	Attishoo	48	8	10	40
231-31-5368	Smiley	22	8	10	30
131-24-3650	Smethurst	35	5	7	30
434-26-3751	Guldu	35	5	7	32
612-67-4134	Madayan	35	8	10	40

Wages

R	W
8	10
5	7

Hourly_Emps2

S	N	L	R	H
123-22-3666	Attishoo	48	8	40
231-31-5368	Smiley	22	8	30
131-24-3650	Smethurst	35	5	30
434-26-3751	Guldu	35	5	32
612-67-4134	Madayan	35	8	40

Will 2 smaller
tables be better?



Functional Dependencies (FDs)

- ❖ A functional dependency $X \rightarrow Y$ holds over relation R if, for every allowable instance r of R:
 - $t1 \in r, t2 \in r, \pi_X(t1) = \pi_X(t2)$ implies $\pi_Y(t1) = \pi_Y(t2)$
 - i.e., given two tuples in r , if the X values agree, then the Y values must also agree. (X and Y are *sets* of attributes.)
- ❖ An FD is a statement about *all* allowable relations.
 - Must be identified based on semantics of application.
 - Given some allowable instance $r1$ of R, we *can* check if it violates some FD f , but we *cannot* tell if f holds over R!
- ❖ K is a candidate key for R means that $K \rightarrow R$



Reasoning About FDs

- ❖ Given some FDs, we can usually infer additional FDs:
 - $ssn \rightarrow did, did \rightarrow lot$ implies $ssn \rightarrow lot$
- ❖ An FD f is implied by a set of FDs F if f holds whenever all FDs in F hold.
 - $F^+ = \text{closure of } F$ is the set of all FDs that are implied by F .
- ❖ Armstrong's Axioms (X, Y, Z are sets of attributes):
 - Reflexivity: If $X \subseteq Y$, then $Y \rightarrow X$
 - Augmentation: If $X \rightarrow Y$, then $XZ \rightarrow YZ$ for any Z
 - Transitivity: If $X \rightarrow Y$ and $Y \rightarrow Z$, then $X \rightarrow Z$

Reasoning About FDs (Contd.)



- ❖ Couple of additional rules (that follow from AA):
 - *Union*: If $X \rightarrow Y$ and $X \rightarrow Z$, then $X \rightarrow YZ$
 - *Decomposition*: If $X \rightarrow YZ$, then $X \rightarrow Y$ and $X \rightarrow Z$
- ❖ Example: **Contracts**(*cid,sid,jid,did,pid,qty,value*), and:
 - C is the key: **$C \rightarrow CSJDPQV$**
 - Project purchases each part using single contract: **$JP \rightarrow C$**
 - Dept purchases at most one part from a supplier: **$SD \rightarrow P$**
- ❖ $JP \rightarrow C, C \rightarrow CSJDPQV$ imply $JP \rightarrow CSJDPQV$
- ❖ $SD \rightarrow P$ implies $SDJ \rightarrow JP$
- ❖ $SDJ \rightarrow JP, JP \rightarrow CSJDPQV$ imply $SDJ \rightarrow CSJDPQV$

Normal Forms



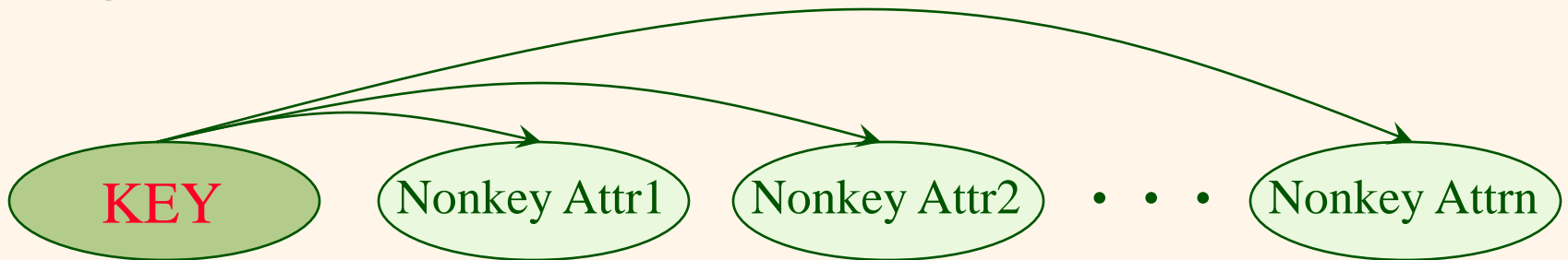
- ❖ Returning to the issue of schema refinement, the first question to ask is whether any refinement is needed!
- ❖ If a relation is in a certain *normal form*, then we know that certain kinds of problems cannot arise:
 - First normal form (1NF)
 - Second normal form (2NF)
 - Third normal form (3NF)
 - Boyce-Codd normal form (BCNF)
 - Fourth normal form (4NF)
 - Fifth normal form (5NF)

$$\text{BCNF} \subset 3\text{NF} \\ \subset 2\text{NF} \subset 1\text{NF}$$

Boyce-Codd Normal Form (BCNF)



- ❖ Reln R with FDs F is in **BCNF** if, for all possible $X \rightarrow A$
 - $A \in X$ (called a *trivial* FD), or
 - X contains a key for R.
- ❖ In other words, R is in **BCNF** if the only non-trivial FDs that hold over R are key constraints.
- ❖ BCNF ensures that no redundancy can be detected using FD information alone.



Third Normal Form (3NF)



- ❖ Reln R with FDs F is in **3NF** if, for all possible $X \rightarrow A$
 - $A \in X$ (called a *trivial* FD), or
 - X contains a key for R, or
 - A is part of some key for R.
- ❖ **Minimality** of a key is crucial in third condition above!
- ❖ If R is in BCNF, obviously in 3NF.
- ❖ If R is in 3NF, some redundancy is possible. It is a compromise, used when BCNF not achievable (e.g., no “good” decomposition, or performance considerations).

Decomposition of a Relation Scheme



- ❖ Suppose that relation R contains attributes $A_1 \dots A_n$.
A decomposition of R consists of replacing R by two or more relations such that:
 - Each new relation scheme contains a subset of the attributes of R (and no attributes that do not appear in R), and
 - Every attribute of R appears as an attribute of at least one of the new relations.
- ❖ Intuitively, decomposing R means we will store instances of the relation schemes produced by the decomposition, instead of instances of R .
- ❖ E.g., Can decompose **SNLRWH** into **SNLRH** and **RW**.



Example Decomposition

- ❖ Decompositions should be used only when needed.
 - SNLRWH has FDs $S \rightarrow \text{SNLRWH}$ and $R \rightarrow W$
 - Second FD causes violation of 3NF; W values repeatedly associated with R values. Easiest way to fix this is to create a relation RW to store these associations, and to remove W from the main schema:
 - i.e., we decompose SNLRWH into SNLRH and RW
- ❖ The information to be stored consists of SNLRWH tuples. If we just store the projections of these tuples onto SNLRH and RW , *are there any potential problems that we should be aware of?*

Problems with Decompositions



- ❖ There are three potential problems to consider:
 1. Some queries become more expensive.
 - e.g., How much did sailor Joe earn? (salary = $W \times H$)
 2. Given instances of the decomposed relations, we may not be able to reconstruct the corresponding instance of the original relation!
 - Fortunately, not in the SNLRWH example.
 3. Checking some dependencies may require joining the instances of the decomposed relations.
 - Fortunately, not in the SNLRWH example.
- ❖ Tradeoff: Must consider these issues vs. redundancy.

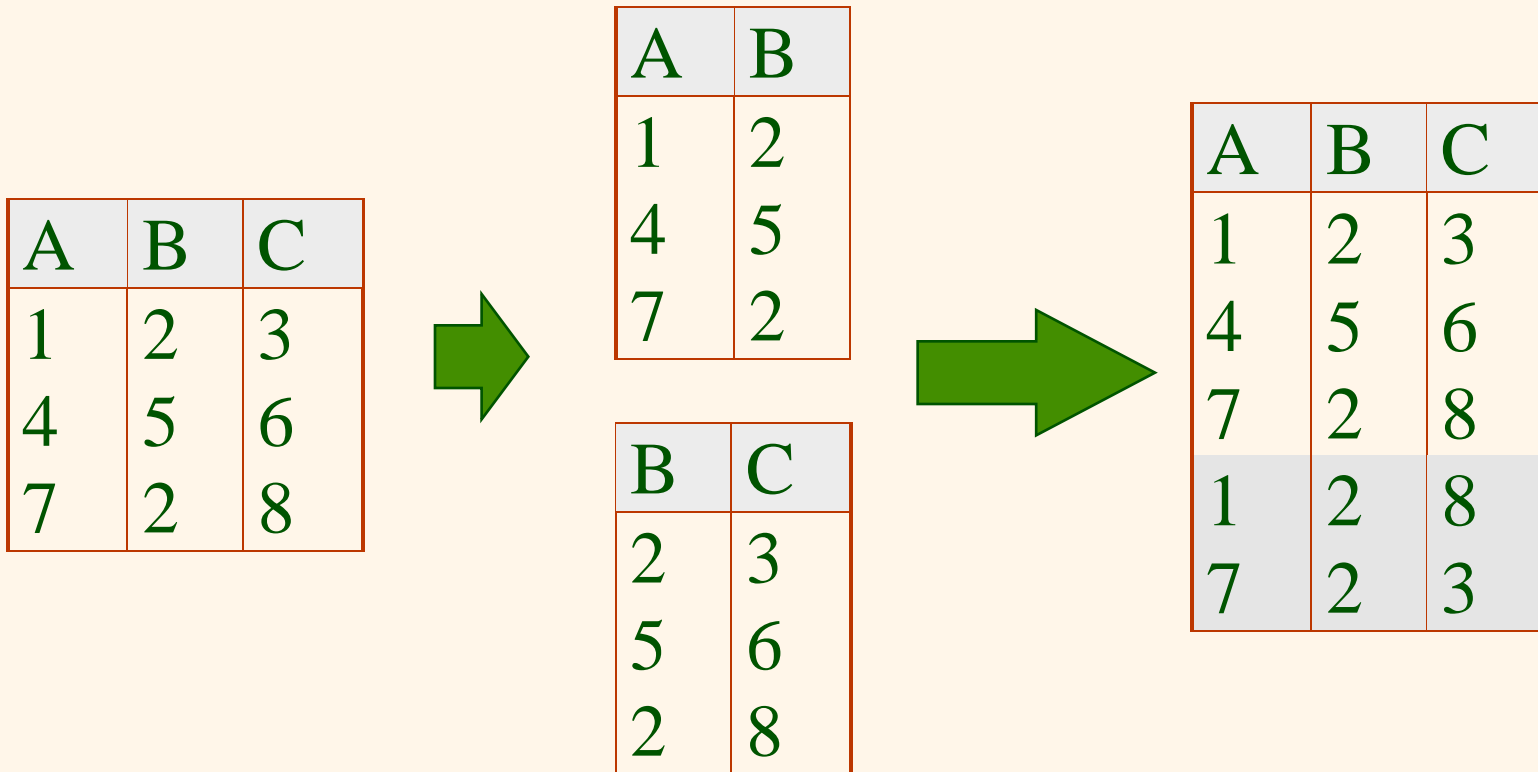


Lossless Join Decompositions

- ❖ Decomposition of R into X and Y is lossless-join w.r.t. a set of FDs F if, for every instance r that satisfies F :
 - $\pi_X(r) \bowtie \pi_Y(r) = r$
- ❖ It is always true that $r \subseteq \pi_X(r) \bowtie \pi_Y(r)$
 - In general, the other direction does not hold! If it does, the decomposition is lossless-join.
- ❖ Definition extended to decomposition into 3 or more relations in a straightforward way.
- ❖ *It is essential that all decompositions used to deal with redundancy be lossless! (Avoids Problem (2).)*



More on Lossless Join



- ❖ The decomposition of R into X and Y is **lossless-join wrt F** if and only if the closure of F contains:
- $X \cap Y \rightarrow X$, or
 - $X \cap Y \rightarrow Y$

Dependency Preserving Decomposition



- ❖ Consider CSJDPQV, C is key, $JP \rightarrow C$ and $SD \rightarrow P$.
 - BCNF decomposition: CSJDQV and SDP (lossless-join)
 - Problem: Checking (Enforcing) $JP \rightarrow C$ requires a join!
- ❖ **Dependency preserving decomposition** (Intuitive):
 - If R is decomposed into X, Y and Z, and we enforce the FDs on X, on Y and on Z, (*separately*) then all FDs that were given to hold on R must also hold. (Avoids Problem (3).)



Decomposition into BCNF

- ❖ Consider relation R with FDs F . If $X \rightarrow Y$ violates BCNF, decompose R into $R - Y$ and XY .
 - Repeated application of this idea will give us a collection of relations that are in BCNF; lossless join decomposition, and guaranteed to terminate.
 - e.g., CSJDPQV, key C , $JP \rightarrow C$, $SD \rightarrow P$, $J \rightarrow S$
 - To deal with $SD \rightarrow P$, decompose into SDP , CSJDQV.
 - To deal with $J \rightarrow S$, decompose CSJDQV into JS and $CJDQV$
- ❖ In general, several dependencies may cause violation of BCNF. The order in which we “deal with” them could lead to very different sets of relations!



Decomposition into 3NF

- ❖ Obviously, the algorithm for lossless join decomposition into BCNF can be used to obtain a lossless join decomposition into 3NF (typically, can stop earlier).
- ❖ To ensure dependency preservation, one idea:
 - If $X \rightarrow Y$ is not preserved, add relation XY .
 - Problem is that XY may violate 3NF! e.g., consider the addition of CJP to 'preserve' $JP \rightarrow C$. What if we also have $J \rightarrow C$?