



CHAPTER 9

NORMALIZATION

Chapter Outline

- ▶ 1. Further topics in Functional Dependencies
 - ▶ 1.1 Inference Rules for FDs
 - ▶ 1.2 Equivalence of Sets of FDs
 - ▶ 1.3 Minimal Sets of FDs
- ▶ 2. Nulls, Dangling Tuples, Alternative Relational Designs
- ▶ 3. 1NF, 2NF, 3NF
- ▶ 4. BCNF

Chapter Outline

- ▶ 5. Multivalued Dependencies and Fourth Normal Form - further discussion
- ▶ 6. Other Dependencies and Normal Forms

Defining Functional Dependencies

- ▶ $X \rightarrow Y$ holds if whenever two tuples have the same value for X , they *must have* the same value for Y
 - ▶ For any two tuples t_1 and t_2 in any relation instance $r(R)$: If $t_1[X]=t_2[X]$, *then* $t_1[Y]=t_2[Y]$
- ▶ $X \rightarrow Y$ in R specifies a *constraint* on all relation instances $r(R)$
- ▶ Written as $X \rightarrow Y$; can be displayed graphically on a relation schema.
- ▶ FDs are derived from the real-world constraints on the attributes

1.1 Inference Rules for FDs

(1)

- ▶ **Definition:** An FD $X \rightarrow Y$ is **inferred from** or **implied by** a set of dependencies F specified on R if $X \rightarrow Y$ holds in *every* legal relation state r of R ; that is, whenever r satisfies all the dependencies in F , $X \rightarrow Y$ also holds in r .
- ▶ Given a set of FDs F , we can **infer** additional FDs that hold whenever the FDs in F hold

Inference Rules for FDs (2)

- ▶ Armstrong's inference rules:
 - ▶ IR1. (**Reflexive**) If $Y \text{ subset-of } X$, then $X \rightarrow Y$
 - ▶ IR2. (**Augmentation**) If $X \rightarrow Y$, then $XZ \rightarrow YZ$
 - ▶ (Notation: XZ stands for $X \cup Z$)
 - ▶ IR3. (**Transitive**) If $X \rightarrow Y$ and $Y \rightarrow Z$, then $X \rightarrow Z$
- ▶ IR1, IR2, IR3 form a **sound** and **complete** set of inference rules
 - ▶ These are rules hold and all other rules that hold can be deduced from these

Inference Rules for FDs (3)

- ▶ Some additional inference rules that are useful:
 - ▶ **Decomposition:** If $X \longrightarrow YZ$, then $X \longrightarrow Y$ and $X \longrightarrow Z$
 - ▶ **Union:** If $X \longrightarrow Y$ and $X \longrightarrow Z$, then $X \longrightarrow YZ$
 - ▶ **Pseudotransitivity:** If $X \longrightarrow Y$ and $WY \longrightarrow Z$, then $WX \longrightarrow Z$
- ▶ The last three inference rules, as well as any other inference rules, can be deduced from IR1, IR2, and IR3 (completeness property)

Closure

- ▶ **Closure** of a set F of FDs is the set F^+ of all FDs that can be inferred from F
- ▶ **Closure** of a set of attributes X with respect to F is the set X^+ of all attributes that are functionally determined by X
- ▶ X^+ can be calculated by repeatedly applying IR1, IR2, IR3 using the FDs in F

Example of Closure (1)

- For example, consider the following relation schema about classes held at a university in a given academic year.

CLASS (Classid, Course#, Instr_name, Credit_hrs, Text, Publisher, Classroom, Capacity).

- Let F , the set of functional dependencies for the above relation include the following f.d.s:

FD1: Classid \rightarrow Course#, Instr_name, Credit_hrs, Text, Publisher, Classroom, Capacity;

FD2: Course# \rightarrow Credit_hrs;

FD3: {Course#, Instr_name} \rightarrow Text, Classroom;

FD4: Text \rightarrow Publisher

FD5: Classroom \rightarrow Capacity

These f.d.s above represent the meaning of the individual attributes and the relationship among them and defines certain rules about the classes.

Example of Closure (2)

- The closures of attributes or sets of attributes for some example sets:

$\{ \text{Classid} \}^+ = \{ \text{Classid}, \text{Course\#}, \text{Instr_name}, \text{Credit_hrs}, \text{Text}, \text{Publisher}, \text{Classroom}, \text{Capacity} \} = \text{CLASS}$

$\{ \text{Course\#} \}^+ = \{ \text{Course\#}, \text{Credit_hrs} \}$

$\{ \text{Course\#}, \text{Instr_name} \}^+ = \{ \text{Course\#}, \text{Credit_hrs}, \text{Text}, \text{Publisher}, \text{Classroom}, \text{Capacity} \}$

Note that each closure above has an interpretation that is revealing about the attribute(s) on the left-hand-side. The closure of $\{ \text{Classid} \}^+$ is the entire relation CLASS indicating that all attributes of the relation can be determined from Classid and hence it is a key.

1.2 Equivalence of Sets of FDs

- ▶ Two sets of FDs F and G are **equivalent** if:
 - ▶ Every FD in F can be inferred from G , and
 - ▶ Every FD in G can be inferred from F
 - ▶ Hence, F and G are equivalent if $F^+ = G^+$
- ▶ Definition (**Covers**):
 - ▶ F **covers** G if every FD in G can be inferred from F
 - ▶ (i.e., if G^+ *subset-of* F^+)
- ▶ F and G are equivalent if F covers G and G covers F

1.3 Finding Minimal Cover of F.D.s (1)

- ▶ Just as we applied inference rules to expand on a set F of FDs to arrive at F^+ , its closure, it is possible to think **in the opposite direction** to see if we could shrink or reduce the set F to its *minimal form* so that the minimal set is still equivalent to the original set F .
- ▶ **Definition:** An attribute in a functional dependency is considered **extraneous attribute** if we can remove it without changing the closure of the set of dependencies. Formally, given F , the set of functional dependencies and a functional dependency $X \rightarrow A$ in F , attribute Y is extraneous in X if Y is a subset of X .

Minimal Sets of FDs (2)

- ▶ A set of FDs is **minimal** if it satisfies the following conditions:
 1. Every dependency in F has a single attribute for its RHS.
 2. We cannot remove any dependency from F and have a set of dependencies that is equivalent to F .
 3. We cannot replace any dependency $X \rightarrow A$ in F with a dependency $Y \rightarrow A$, where Y is a proper-subset-of X and still have a set of dependencies that is equivalent to F .

Problems with Null Values and Dangling Tuples (1)

4.1 Problems with NULL values

- ▶ when some tuples have NULL values for attributes that will be used to join individual relations in the decomposition that may lead to incomplete results.
- ▶ E.g., see Figure 15.2(a), where two relations EMPLOYEE and DEPARTMENT are shown. The last two employee tuples—‘Berger’ and ‘Benitez’—represent newly hired employees who have not yet been assigned to a department (assume that this does not violate any integrity constraints).
- ▶ If we want to retrieve a list of (Ename, Dname) values for all the employees. If we apply the NATURAL JOIN operation on EMPLOYEE and DEPARTMENT (Figure 15.2(b)), the two aforementioned tuples will *not* appear in the result.
- ▶ In such cases, LEFT OUTER JOIN may be used. The result is shown in Figure 15.2 (c).

.

Problems with Null Values and Dangling Tuples (2)

(a)

EMPLOYEE

Ename	<u>Ssn</u>	Bdate	Address	Dnum
Smith, John B.	123456789	1965-01-09	731 Fondren, Houston, TX	5
Wong, Franklin T.	333445555	1955-12-08	638 Voss, Houston, TX	5
Zelaya, Alicia J.	999887777	1968-07-19	3321 Castle, Spring, TX	4
Wallace, Jennifer S.	987654321	1941-06-20	291 Berry, Bellaire, TX	4
Narayan, Ramesh K.	666884444	1962-09-15	975 Fire Oak, Humble, TX	5
English, Joyce A.	453453453	1972-07-31	5631 Rice, Houston, TX	5
Jabbar, Ahmad V.	987987987	1969-03-29	980 Dallas, Houston, TX	4
Borg, James E.	888665555	1937-11-10	450 Stone, Houston, TX	1
Berger, Anders C.	999775555	1965-04-26	6530 Braes, Bellaire, TX	NULL
Benitez, Carlos M.	888664444	1963-01-09	7654 Beech, Houston, TX	NULL

DEPARTMENT

Dname	<u>Dnum</u>	Dmgr_ssn
Research	5	333445555
Administration	4	987654321
Headquarters	1	888665555

Figure 15.2
Issues with
NULL-value
joins. (a) Some
EMPLOYEE
tuples have
NULL for the
join attribute
Dnum.

Problems with Null Values and Dangling Tuples (3)

(b)

Ename	<u>Ssn</u>	Bdate	Address	Dnum	Dname	Dmgr_ssn
Smith, John B.	123456789	1965-01-09	731 Fondren, Houston, TX	5	Research	333445555
Wong, Franklin T.	333445555	1955-12-08	638 Voss, Houston, TX	5	Research	333445555
Zelaya, Alicia J.	999887777	1968-07-19	3321 Castle, Spring, TX	4	Administration	987654321
Wallace, Jennifer S.	987654321	1941-06-20	291 Berry, Bellaire, TX	4	Administration	987654321
Narayan, Ramesh K.	666884444	1962-09-15	975 Fire Oak, Humble, TX	5	Research	333445555
English, Joyce A.	453453453	1972-07-31	5631 Rice, Houston, TX	5	Research	333445555
Jabbar, Ahmad V.	987987987	1969-03-29	980 Dallas, Houston, TX	4	Administration	987654321
Borg, James E.	888665555	1937-11-10	450 Stone, Houston, TX	1	Headquarters	888665555

(c)

Ename	<u>Ssn</u>	Bdate	Address	Dnum	Dname	Dmgr_ssn
Smith, John B.	123456789	1965-01-09	731 Fondren, Houston, TX	5	Research	333445555
Wong, Franklin T.	333445555	1955-12-08	638 Voss, Houston, TX	5	Research	333445555
Zelaya, Alicia J.	999887777	1968-07-19	3321 Castle, Spring, TX	4	Administration	987654321
Wallace, Jennifer S.	987654321	1941-06-20	291 Berry, Bellaire, TX	4	Administration	987654321
Narayan, Ramesh K.	666884444	1962-09-15	975 Fire Oak, Humble, TX	5	Research	333445555
English, Joyce A.	453453453	1972-07-31	5631 Rice, Houston, TX	5	Research	333445555
Jabbar, Ahmad V.	987987987	1969-03-29	980 Dallas, Houston, TX	4	Administration	987654321
Borg, James E.	888665555	1937-11-10	450 Stone, Houston, TX	1	Headquarters	888665555
Berger, Anders C.	999775555	1965-04-26	6530 Braes, Bellaire, TX	NULL	NULL	NULL
Benitez, Carlos M.	888665555	1963-01-09	7654 Beech, Houston, TX	NULL	NULL	NULL

Figure 15.2
Issues with NULL-value joins.
(b) Result of applying NATURAL JOIN to the EMPLOYEE and DEPARTMENT relations.
(c) Result of applying LEFT OUTER JOIN to EMPLOYEE and DEPARTMENT

Problems with Null Values and Dangling Tuples (4)

Problems with Dangling Tuples

- ▶ Consider the decomposition of EMPLOYEE into EMPLOYEE_1 and EMPLOYEE_2 as shown in Figure 15.3 (a) and 15.3 (b).
- ▶ Their NATURAL JOIN yields the original relation EMPLOYEE in Figure 15.2(a).
- ▶ We may use the alternative representation, shown in Figure 15.3(c), where we *do not include a tuple* in EMPLOYEE_3 if the employee has not been assigned a department (instead of including a tuple with NULL for Dnum as in EMPLOYEE_2).
- ▶ If we use EMPLOYEE_3 instead of EMPLOYEE_2 and apply a NATURAL JOIN on EMPLOYEE_1 and EMPLOYEE_3, the tuples for Berger and Benitez will not appear in the result; these are called **dangling tuples** in EMPLOYEE.

Problems with Null Values and Dangling Tuples (5)

(a) EMPLOYEE_1

Ename	<u>Ssn</u>	Bdate	Address
Smith, John B.	123456789	1965-01-09	731 Fondren, Houston, TX
Wong, Franklin T.	333445555	1955-12-08	638 Voss, Houston, TX
Zelaya, Alicia J.	999887777	1968-07-19	3321 Castle, Spring, TX
Wallace, Jennifer S.	987654321	1941-06-20	291 Berry, Bellaire, TX
Narayan, Ramesh K.	666884444	1962-09-15	975 Fire Oak, Humble, TX
English, Joyce A.	453453453	1972-07-31	5631 Rice, Houston, TX
Jabbar, Ahmad V.	987987987	1969-03-29	980 Dallas, Houston, TX
Borg, James E.	888665555	1937-11-10	450 Stone, Houston, TX
Berger, Anders C.	999775555	1965-04-26	6530 Braes, Bellaire, TX
Benitez, Carlos M.	888665555	1963-01-09	7654 Beech, Houston, TX

Figure 15.3

The dangling tuple problem. (a) The relation EMPLOYEE_1 (includes all attributes of EMPLOYEE from Figure 15.2(a) except Dnum). (b) The relation EMPLOYEE_2 (includes Dnum attribute with NULL values). (c) The relation EMPLOYEE_3 (includes Dnum attribute but does not include tuples for which Dnum has NULL values).

(b) EMPLOYEE_2

<u>Ssn</u>	Dnum
123456789	5
333445555	5
999887777	4
987654321	4
666884444	5
453453453	5
987987987	4
888665555	1
999775555	NULL
888664444	NULL

(c) EMPLOYEE_3

<u>Ssn</u>	Dnum
123456789	5
333445555	5
999887777	4
987654321	4
666884444	5
453453453	5
987987987	4
888665555	1

Example of Lossless-Join Decomposition

► **Lossless join decomposition**

► Decomposition of $R = (A, B, C)$

$$R_1 = (A, B) \quad R_2 = (B, C)$$

A	B	C
α	1	A
β	2	B

r

A	B
α	1
β	2

$\Pi_{A,B}(r)$

B	C
1	A
2	B

$\Pi_{B,C}(r)$

$\Pi_A(r) \bowtie \Pi_B(r)$

A	B	C
α	1	A
β	2	B

First Normal Form

- ▶ Domain is **atomic** if its elements are considered to be indivisible units
 - ▶ Examples of non-atomic domains:
 - ▶ Set of names, composite attributes
 - ▶ Identification numbers like CS101 that can be broken up into parts
- ▶ A relational schema R is in **first normal form** if the domains of all attributes of R are atomic
- ▶ Non-atomic values complicate storage and encourage redundant (repeated) storage of data
 - ▶ Example: Set of accounts stored with each customer, and set of owners stored with each account

First Normal Form (Cont'd)

- ▶ Atomicity is actually a property of how the elements of the domain are used.
 - ▶ Example: Strings would normally be considered indivisible
 - ▶ Suppose that students are given roll numbers which are strings of the form *CS0012* or *EE1127*
 - ▶ If the first two characters are extracted to find the department, the domain of roll numbers is not atomic.
 - ▶ Doing so is a bad idea: leads to encoding of information in application program rather than in the database.

Second Normal Form

Second Normal Form

A relation is in **2NF** if it is in 1NF, and every non-key attribute is fully dependent on each candidate key. (That is, we don't have any partial functional dependency.)

- 2NF (and 3NF) both involve the concepts of key and non-key attributes.
- A *key attribute* is any attribute that is part of a key; any attribute that is not a key attribute, is a *non-key attribute*.
- A relation in 2NF will not have any partial dependencies

Third Normal Form

- ▶ A relation schema R is in **third normal form (3NF)** if for all:

$$\alpha \rightarrow \beta \text{ in } F^+$$

at least one of the following holds:

- ▶ $\alpha \rightarrow \beta$ is trivial (i.e., $\beta \in \alpha$)
- ▶ α is a superkey for R
- ▶ Each attribute A in $\beta - \alpha$ is contained in a candidate key for R .

(NOTE: each attribute may be in a different candidate key)

- ▶ If a relation is in BCNF it is in 3NF (since in BCNF one of the first two conditions above must hold).
- ▶ Third condition is a minimal relaxation of BCNF to ensure dependency preservation.

Goals of Normalization

- ▶ Let R be a relation scheme with a set F of functional dependencies.
- ▶ Decide whether a relation scheme R is in “good” form.
- ▶ In the case that a relation scheme R is not in “good” form, decompose it into a set of relation scheme $\{R_1, R_2, \dots, R_n\}$ such that
 - ▶ each relation scheme is in good form
 - ▶ the decomposition is a lossless-join decomposition
 - ▶ Preferably, the decomposition should be dependency preserving.

Boyce-Codd Normal Form

A relation schema R is in BCNF with respect to a set F of functional dependencies if for all functional dependencies in F^+ of the form

$$\alpha \rightarrow \beta$$

where $\alpha \subseteq R$ and $\beta \subseteq R$, at least one of the following holds:

- ▶ $\alpha \rightarrow \beta$ is trivial (i.e., $\beta \subseteq \alpha$)
- ▶ α is a superkey for R

Example schema *not* in BCNF:

instr_dept (ID, name, salary, dept_name, building, budget)

because $\text{dept_name} \rightarrow \text{building, budget}$
holds on *instr_dept*, but *dept_name* is not a superkey

BCNF and Dependency Preservation

- ▶ Constraints, including functional dependencies, are costly to check in practice unless they pertain to only one relation
- ▶ If it is sufficient to test only those dependencies on each individual relation of a decomposition in order to ensure that *all* functional dependencies hold, then that decomposition is *dependency preserving*.
- ▶ Because it is not always possible to achieve both BCNF and dependency preservation, we consider a weaker normal form, known as *third normal form*.

5. Multivalued Dependencies and Fourth Normal Form - Further Discussion (1)

Definition:

- ▶ A multivalued dependency (MVD) $X \twoheadrightarrow Y$ specified on relation schema R , where X and Y are both subsets of R , specifies the following constraint on any relation state r of R : If two tuples t_1 and t_2 exist in r such that $t_1[X] = t_2[X]$, then two tuples t_3 and t_4 should also exist in r with the following properties, where we use Z to denote $R - (X \cup Y)$:
 - ▶ $t_3[X] = t_4[X] = t_1[X] = t_2[X]$.
 - ▶ $t_3[Y] = t_1[Y]$ and $t_4[Y] = t_2[Y]$.
 - ▶ $t_3[Z] = t_2[Z]$ and $t_4[Z] = t_1[Z]$.
- ▶ An MVD $X \twoheadrightarrow Y$ in R is called a **trivial MVD** if (a) Y is a subset of X , or (b) $X \cup Y = R$.

Multivalued Dependencies and Fourth Normal Form (2)

- **Inference Rules for Functional and Multivalued Dependencies:**
 - IR1 (reflexive rule for FDs): If $X \supseteq Y$, then $X \rightarrow Y$.
 - IR2 (augmentation rule for FDs): $\{X \rightarrow Y\} \models XZ \rightarrow YZ$.
 - IR3 (transitive rule for FDs): $\{X \rightarrow Y, Y \rightarrow Z\} \models X \rightarrow Z$.
 - IR4 (complementation rule for MVDs): $\{X \twoheadrightarrow Y\} \models X \twoheadrightarrow (R - (X \cup Y))$.
 - IR5 (augmentation rule for MVDs): If $X \twoheadrightarrow Y$ and $W \supseteq Z$ then $WX \twoheadrightarrow YZ$.
 - IR6 (transitive rule for MVDs): $\{X \twoheadrightarrow Y, Y \twoheadrightarrow Z\} \models X \twoheadrightarrow (Z - Y)$.
 - IR7 (replication rule for FD to MVD): $\{X \rightarrow Y\} \models X \twoheadrightarrow Y$.
 - IR8 (coalescence rule for FDs and MVDs): If $X \twoheadrightarrow Y$ and there exists W with the properties that
 - (a) $W \cap Y$ is empty, (b) $W \rightarrow Z$, and (c) $Y \supseteq Z$, then $X \rightarrow Z$.

Multivalued Dependencies and Fourth Normal Form (3)

Definition:

- ▶ A relation schema R is in **4NF** with respect to a set of dependencies F (that includes functional dependencies and multivalued dependencies) if, for every *nontrivial* multivalued dependency $X \twoheadrightarrow Y$ in F^+ , X is a superkey for R .
- ▶ Note: F^+ is the (complete) set of all dependencies (functional or multivalued) that will hold in every relation state r of R that satisfies F . It is also called the **closure** of F .

Multivalued Dependencies and Fourth Normal Form (4)

Fig. 15.4 Decomposing a relation state of EMP that is not in 4NF.

(a) EMP relation with additional tuples.

(b) Two corresponding 4NF relations EMP_PROJECTS and EMP_DEPENDENTS.

(a) EMP

<u>Ename</u>	<u>Pname</u>	<u>Dname</u>
Smith	X	John
Smith	Y	Anna
Smith	X	Anna
Smith	Y	John
Brown	W	Jim
Brown	X	Jim
Brown	Y	Jim
Brown	Z	Jim
Brown	W	Joan
Brown	X	Joan
Brown	Y	Joan
Brown	Z	Joan
Brown	W	Bob
Brown	X	Bob
Brown	Y	Bob
Brown	Z	Bob

(b) EMP_PROJECTS

<u>Ename</u>	<u>Pname</u>
Smith	X
Smith	Y
Brown	W
Brown	X
Brown	Y
Brown	Z

EMP_DEPENDENTS

<u>Ename</u>	<u>Dname</u>
Smith	Anna
Smith	John
Brown	Jim
Brown	Joan
Brown	Bob

6. Other Dependencies and Normal Forms

Join Dependency:

Definition:

- ▶ A join dependency (JD), denoted by $JD(R_1, R_2, \dots, R_n)$, specified on relation schema R , specifies a constraint on the states r of R .
- ▶ The constraint states that every legal state r of R should have a non-additive join decomposition into R_1, R_2, \dots, R_n ; that is, for every such r we have $*(\pi_{R_1}(r), \pi_{R_2}(r), \dots, \pi_{R_n}(r)) = r$

***Note:** an MVD is a special case of a JD where $n = 2$.*

- ▶ A join dependency $JD(R_1, R_2, \dots, R_n)$, specified on relation schema R , is a **trivial JD** if one of the relation schemas R_i in $JD(R_1, R_2, \dots, R_n)$ is equal to R .

Join Dependencies and Fifth Normal Form

Definition of 5NF:

- ▶ A relation schema R is in **fifth normal form (5NF)** (or **Project-Join Normal Form (PJNF)**) with respect to a set F of functional, multivalued, and join dependencies if,
 - ▶ for every nontrivial join dependency $JD(R_1, R_2, \dots, R_n)$ in F^+ (that is, implied by F),
 - ▶ every R_i is a superkey of R .
- ▶ Discovering join dependencies in practical databases with hundreds of relations is next to impossible. Therefore, 5NF is rarely used in practice.

Recap

- ▶ Functional Dependencies Revisited
- ▶ 1NF, 2NF, 3NF, BCNF
- ▶ Multivalued Dependencies and Fourth Normal Form, 5NF

Appendix

Properties of Relational Decompositions

- ▶ **Dependency Preservation Property of a Decomposition:**
 - ▶ **Dependency Preservation Property:**
 - ▶ A decomposition $D = \{R_1, R_2, \dots, R_m\}$ of R is **dependency-preserving** with respect to F if the union of the projections of F on each R_i in D is equivalent to F ; that is
$$((\pi_{R_1}(F)) \cup \dots \cup (\pi_{R_m}(F)))^+ = F^+$$

Properties of Relational Decompositions

2.3 Non-additive (Lossless) Join Property of a Decomposition:

- ▶ Definition: Lossless join property: a decomposition $D = \{R_1, R_2, \dots, R_m\}$ of R has the **lossless (nonadditive) join property** with respect to the set of dependencies F on R if, for *every* relation state r of R that satisfies F , the following holds, where $*$ is the natural join of all the relations in D :

$$* (\pi_{R_1}(r), \dots, \pi_{R_m}(r)) = r$$

- ▶ Note: The word loss in lossless refers to loss of information, not to loss of tuples. In fact, for “loss of information” a better term is “**addition of spurious information**”