# Database Systems

**Chapter # 14**

**Database Design Theory and Normalization**

**eman.shahid@nu.edu.pk**

# Normalization

- Normalization is the process of reorganizing/restructuring data in a database with a series of so called normal-forms, so that it meets two basic requirements:
  - There is no redundancy of data (all data is stored in only one place)
  - data dependencies are logical (all related data items are stored together).
- Normalization is important for many reasons, but chiefly because it allows databases to take up as little disk space as possible, resulting in increased performance.

# Four informal guidelines

- Four informal guidelines used to determine the quality of relation schema design:
    - Making sure that the semantics of the attributes is clear in the schema
    - Reducing the redundant information in tuples
    - Reducing the NULL values in tuples
    - Disallowing the possibility of generating spurious tuples.

**EMPLOYEE**                                                                    F.K.

| Ename | Ssn | Bdate | Address | Dnumber |
|-------|-----|-------|---------|---------|

P.K.

**DEPARTMENT**                          F.K.

| Dname | Dnumber | Dmgr_ssn |
|-------|---------|----------|

P.K.

**DEPT_LOCATIONS**
F.K.

| Dnumber | Dlocation |
|---------|-----------|

P.K.

**PROJECT**                                          F.K.

| Pname | Pnumber | Plocation | Dnum |
|-------|---------|-----------|------|

P.K.

**WORKS_ON**
F.K.          F.K.

| Ssn | Pnumber | Hours |
|-----|---------|-------|

P.K.

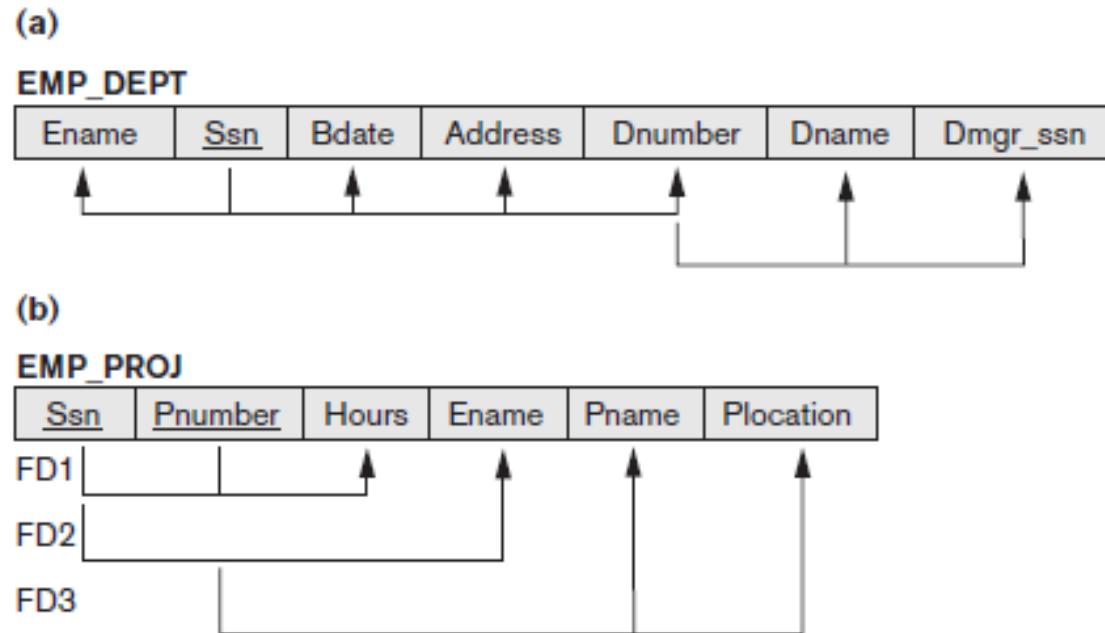**Figure 14.1** A simplified COMPANY relational database schema.

# Guideline 1

- Design a relation schema so that it is easy to explain its meaning.

- Do not combine attributes from multiple entity types and relationship types into a single relation.

- If the relation corresponds to a mixture of multiple entities and relationships, semantic ambiguities will result and the relation cannot be easily explained.

# Examples of Violating Guideline 1.

**Figure 14.3**

Two relation schemas suffering from update anomalies.
(a) EMP_DEPT and
(b) EMP_PROJ.

**(a)**

**EMP_DEPT**

| Ename | Ssn | Bdate | Address | Dnumber | Dname | Dmgr_ssn |
|-------|-----|-------|---------|---------|-------|----------|

**(b)**

**EMP_PROJ**

| Ssn | Pnumber | Hours | Ename | Pname | Plocation |
|-----|---------|-------|-------|-------|-----------|

FD1

FD2

FD3

They fare poorly
Against the above measure of design quality. They may be used as views, but they
Cause problems when used as base relations,

# Redundant Information in Tuples and Update Anomalies

- **Goal of schema design:** to minimize the storage space used by the base relations
  - Grouping attributes into relation schemas has a significant effect on storage space.

**EMPLOYEE**

| Ename | Ssn | Bdate | Address | Dnumber |
|-------|-----|-------|---------|---------|
| 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 | 291Berry, 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 |

**DEPARTMENT**

| Dname | Dnumber | Dmgr_ssn |
|-------|---------|----------|
| Research | 5 | 333445555 |
| Administration | 4 | 987654321 |
| Headquarters | 1 | 888665555 |

**Natural join of Employee and Department Tables**

Redundancy

**EMP_DEPT**

| Ename | Ssn | Bdate | Address | Dnumber | 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 FireOak, 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 |

# Anomalies

- Storing natural joins of base relations leads to an additional problem referred to as update anomalies.
    - Classified into:
        - **Insertion Anomalies**
        - **Deletion Anomalies**
        - **Modification Anomalies**

# Insertion Anomalies.

- Insertion anomalies can be differentiated into two types:

- To insert a new employee tuple into EMP_DEPT, we must include either the attribute values for the department that the employee works for, or NULLs (if the employee does not work for a department as yet).

  - For example, to insert a new tuple for an employee who works in department number 5, we must enter all the attribute values of department 5 correctly so that they are *consistent*

# Insertion Anomalies.

**EMPLOYEE**

| Ename | Ssn | Bdate | Address | Dnumber |
|---|---|---|---|---|
| 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 | 291Berry, 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 |

**DEPARTMENT**

| Dname | Dnumber | Dmgr_ssn |
|---|---|---|
| Research | 5 | 333445555 |
| Administration | 4 | 987654321 |
| Headquarters | 1 | 888665555 |

**EMP_DEPT**

| Ename | Ssn | Bdate | Address | Dnumber | 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 FireOak, 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 |

Redundancy

In the design of Figure 14.2, we do not have to worry about this consistency problem because we enter only the department number in the employee tuple; once in the database, as a single tuple in the DEPARTMENT relation.

For example, to insert a new tuple for an employee who works in department number 5, we must enter all the attribute values of department 5 correctly so that they are *consistent*

# Insertion Anomalies.

- **It is difficult to insert a new department that has no employees as yet in the EMP_DEPT relation.**
  - The only way to do this is to place NULL values in the

Redundancy

**EMP_DEPT**

| Ename | Ssn | Bdate | Address | Dnumber | 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 FireOak, 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 |

**This violates the entity integrity for EMP_DEPT because its primary key Ssn cannot be null.**

# EXAMPLE OF AN INSERT ANOMALY

- **Consider the relation:**
  - EMP_PROJ(Emp#, Proj#, Hours, Ename, Pname,Plocation)

- **Insert  Anomaly:**
  - Cannot insert a project unless an employee is assigned to it.

- **Conversely**
  - Cannot insert an employee unless an he/she is assigned to a project.

**EMP_PROJ**

|  |  |  | | Redundancy | Redundancy |
| --- | --- | --- | --- | --- | --- |
| Ssn | Pnumber | Hours | Ename | Pname | Plocation |
| 123456789 | 1 | 32.5 | Smith, John B. | ProductX | Bellaire |
| 123456789 | 2 | 7.5 | Smith, John B. | ProductY | Sugarland |
| 666884444 | 3 | 40.0 | Narayan, Ramesh K. | ProductZ | Houston |
| 453453453 | 1 | 20.0 | English, Joyce A. | ProductX | Bellaire |
| 453453453 | 2 | 20.0 | English, Joyce A. | ProductY | Sugarland |
| 333445555 | 2 | 10.0 | Wong, Franklin T. | ProductY | Sugarland |
| 333445555 | 3 | 10.0 | Wong, Franklin T. | ProductZ | Houston |
| 333445555 | 10 | 10.0 | Wong, Franklin T. | Computerization | Stafford |
| 333445555 | 20 | 10.0 | Wong, Franklin T. | Reorganization | Houston |
| 999887777 | 30 | 30.0 | Zelaya, Alicia J. | Newbenefits | Stafford |
| 999887777 | 10 | 10.0 | Zelaya, Alicia J. | Computerization | Stafford |
| 987987987 | 10 | 35.0 | Jabbar, Ahmad V. | Computerization | Stafford |
| 987987987 | 30 | 5.0 | Jabbar, Ahmad V. | Newbenefits | Stafford |
| 987654321 | 30 | 20.0 | Wallace, Jennifer S. | Newbenefits | Stafford |
| 987654321 | 20 | 15.0 | Wallace, Jennifer S. | Reorganization | Houston |
| 888665555 | 20 | Null | Borg, James E. | Reorganization | Houston |

# EXAMPLE OF AN UPDATE ANOMALY(Repeated Update)

- **Consider the relation:**

- EMP_PROJ(Emp#, Proj#, Hours Ename, Pname,Plocation)

- **Update Anomaly:**
  - Changing the name of project number P1 from "ProjectX" to "Customer-Accounting" may cause this update to be made for all 100 employees working on project P1.

| Ssn | Pnumber | Hours | Ename | Pname | Plocation |
|-----|---------|-------|-------|-------|-----------|
| 123456789 | 1 | 32.5 | Smith, John B. | ProductX | Bellaire |
| 123456789 | 2 | 7.5 | Smith, John B. | ProductY | Sugarland |
| 666884444 | 3 | 40.0 | Narayan, Ramesh K. | ProductZ | Houston |
| 453453453 | 1 | 20.0 | English, Joyce A. | ProductX | Bellaire |
| 453453453 | 2 | 20.0 | English, Joyce A. | ProductY | Sugarland |
| 333445555 | 2 | 10.0 | Wong, Franklin T. | ProductY | Sugarland |
| 333445555 | 3 | 10.0 | Wong, Franklin T. | ProductZ | Houston |
| 333445555 | 10 | 10.0 | Wong, Franklin T. | Computerization | Stafford |
| 333445555 | 20 | 10.0 | Wong, Franklin T. | Reorganization | Houston |
| 999887777 | 30 | 30.0 | Zelaya, Alicia J. | Newbenefits | Stafford |
| 999887777 | 10 | 10.0 | Zelaya, Alicia J. | Computerization | Stafford |
| 987987987 | 10 | 35.0 | Jabbar, Ahmad V. | Computerization | Stafford |
| 987987987 | 30 | 5.0 | Jabbar, Ahmad V. | Newbenefits | Stafford |
| 987654321 | 30 | 20.0 | Wallace, Jennifer S. | Newbenefits | Stafford |
| 987654321 | 20 | 15.0 | Wallace, Jennifer S. | Reorganization | Houston |
| 888665555 | 20 | Null | Borg, James E. | Reorganization | Houston |

EMP_PROJ — Redundancy, Redundancy

# EXAMPLE OF A DELETE ANOMALY

- **Consider the relation:**

- EMP_PROJ(Emp#, Proj#, Hours Ename, Pname,Plocation)

- **Delete Anomaly:**
  - When a project is deleted, it will result in deleting all the employees who work on that project.
  - Alternately, if an employee is the sole employee on a project, deleting that employee would result in deleting the corresponding project.

| | | | | Redundancy | Redundancy |
| --- | --- | --- | --- | --- | --- |
| **EMP_PROJ** | | | | | |
| Ssn | Pnumber | Hours | Ename | Pname | Plocation |
| 123456789 | 1 | 32.5 | Smith, John B. | ProductX | Bellaire |
| 123456789 | 2 | 7.5 | Smith, John B. | ProductY | Sugarland |
| 666884444 | 3 | 40.0 | Narayan, Ramesh K. | ProductZ | Houston |
| 453453453 | 1 | 20.0 | English, Joyce A. | ProductX | Bellaire |
| 453453453 | 2 | 20.0 | English, Joyce A. | ProductY | Sugarland |
| 333445555 | 2 | 10.0 | Wong, Franklin T. | ProductY | Sugarland |
| 333445555 | 3 | 10.0 | Wong, Franklin T. | ProductZ | Houston |
| 333445555 | 10 | 10.0 | Wong, Franklin T. | Computerization | Stafford |
| 333445555 | 20 | 10.0 | Wong, Franklin T. | Reorganization | Houston |
| 999887777 | 30 | 30.0 | Zelaya, Alicia J. | Newbenefits | Stafford |
| 999887777 | 10 | 10.0 | Zelaya, Alicia J. | Computerization | Stafford |
| 987987987 | 10 | 35.0 | Jabbar, Ahmad V. | Computerization | Stafford |
| 987987987 | 30 | 5.0 | Jabbar, Ahmad V. | Newbenefits | Stafford |
| 987654321 | 30 | 20.0 | Wallace, Jennifer S. | Newbenefits | Stafford |
| 987654321 | 20 | 15.0 | Wallace, Jennifer S. | Reorganization | Houston |
| 888665555 | 20 | Null | Borg, James E. | Reorganization | Houston |

# Guideline 2.

- **Guideline 2.** Design the base relation schemas so that no insertion, deletion, or modification anomalies are not present in the relations.

# Null Values

- If many of the attributes do not apply to all tuples in the relation, we end up with many NULLs in those tuples.

- This can waste space at the storage level and may also lead to problems with understanding the meaning of the attributes.

- Another problem with NULLs is how to account for them when aggregate operations such as COUNT or SUM are applied.

# Reasons for NULLS

- **Attribute not applicable or invalid** (e.g. Visa_Status may not apply to local students)

- **Attribute value unknown** (may exist) (e.g. Date_of_birth may be unknown for an employee)

- **Value known to exist, but unavailable** (e.g. Home_Phone_Number for an employee may exist, but may not be available and recorded yet.

- **Guideline 3.** As far as possible, avoid placing attributes in a base relation whose values may frequently be NULL.

- If NULLs are unavoidable, make sure that they apply in exceptional cases only and do not apply to a majority of tuples in the relation.

# Generation of Spurious Tuples

Additional tuples that are not present in original table are called spurious tuples because they represent spurious information that is not valid.

**(b)**

**EMP_LOCS**

| Ename | Plocation |
|---|---|
| Smith, John B. | Bellaire |
| Smith, John B. | Sugarland |
| Narayan, Ramesh K. | Houston |
| English, Joyce A. | Bellaire |
| English, Joyce A. | Sugarland |
| Wong, Franklin T. | Sugarland |
| Wong, Franklin T. | Houston |
| Wong, Franklin T. | Stafford |
| Zelaya, Alicia J. | Stafford |
| Jabbar, Ahmad V. | Stafford |
| Wallace, Jennifer S. | Stafford |
| Wallace, Jennifer S. | Houston |
| Borg, James E. | Houston |

**Natural Join**

**EMP_PROJ1**

| Ssn | Pnumber | Hours | Pname | Plocation |
|---|---|---|---|---|
| 123456789 | 1 | 32.5 | ProductX | Bellaire |
| 123456789 | 2 | 7.5 | ProductY | Sugarland |
| 666884444 | 3 | 40.0 | ProductZ | Houston |
| 453453453 | 1 | 20.0 | ProductX | Bellaire |
| 453453453 | 2 | 20.0 | ProductY | Sugarland |
| 333445555 | 2 | 10.0 | ProductY | Sugarland |
| 333445555 | 3 | 10.0 | ProductZ | Houston |
| 333445555 | 10 | 10.0 | Computerization | Stafford |
| 333445555 | 20 | 10.0 | Reorganization | Houston |
| 999887777 | 30 | 30.0 | Newbenefits | Stafford |
| 999887777 | 10 | 10.0 | Computerization | Stafford |
| 987987987 | 10 | 35.0 | Computerization | Stafford |
| 987987987 | 30 | 5.0 | Newbenefits | Stafford |
| 987654321 | 30 | 20.0 | Newbenefits | Stafford |
| 987654321 | 20 | 15.0 | Reorganization | Houston |
| 888665555 | 20 | NULL | Reorganization | Houston |

| | Ssn | Pnumber | Hours | Pname | Plocation | Ename |
|---|---|---|---|---|---|---|
| | 123456789 | 1 | 32.5 | ProductX | Bellaire | Smith, John B. |
| * | 123456789 | 1 | 32.5 | ProductX | Bellaire | English, Joyce A. |
| | 123456789 | 2 | 7.5 | ProductY | Sugarland | Smith, John B. |
| * | 123456789 | 2 | 7.5 | ProductY | Sugarland | English, Joyce A. |
| * | 123456789 | 2 | 7.5 | ProductY | Sugarland | Wong, Franklin T. |
| | 666884444 | 3 | 40.0 | ProductZ | Houston | Narayan, Ramesh K. |
| * | 666884444 | 3 | 40.0 | ProductZ | Houston | Wong, Franklin T. |
| * | 453453453 | 1 | 20.0 | ProductX | Bellaire | Smith, John B. |
| | 453453453 | 1 | 20.0 | ProductX | Bellaire | English, Joyce A. |
| * | 453453453 | 2 | 20.0 | ProductY | Sugarland | Smith, John B. |
| | 453453453 | 2 | 20.0 | ProductY | Sugarland | English, Joyce A. |
| * | 453453453 | 2 | 20.0 | ProductY | Sugarland | Wong, Franklin T. |
| * | 333445555 | 2 | 10.0 | ProductY | Sugarland | Smith, John B. |
| * | 333445555 | 2 | 10.0 | ProductY | Sugarland | English, Joyce A. |
| | 333445555 | 2 | 10.0 | ProductY | Sugarland | Wong, Franklin T. |
| * | 333445555 | 3 | 10.0 | ProductZ | Houston | Narayan, Ramesh K. |
| | 333445555 | 3 | 10.0 | ProductZ | Houston | Wong, Franklin T. |
| | 333445555 | 10 | 10.0 | Computerization | Stafford | Wong, Franklin T. |
| * | 333445555 | 20 | 10.0 | Reorganization | Houston | Narayan, Ramesh K. |
| | 333445555 | 20 | 10.0 | Reorganization | Houston | Wong, Franklin T. |

# Guideline 4

- Design relation schemas so that they can be joined with equality conditions on attributes that are appropriately related (primary key, foreign key) pairs in a way that guarantees that no spurious tuples are generated.

- Avoid relations that contain matching attributes that are not (foreign key, primary key) combinations because joining on such attributes may produce spurious tuples.
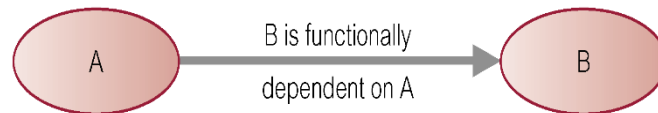
# Spurious Tuples

- There are two important properties of decompositions:
  - **Non-additive or losslessness of the corresponding join**
  - **Preservation of the functional dependencies.**

- **Note that:**
  - Property (a) is extremely important and *cannot* be sacrificed.
  - Property (b) is less stringent and may be sacrificed.

# Functional Dependencies

- A functional dependency is a constraint between two sets of attributes from the database.

- **Definition.**

- A functional dependency, denoted by $X \rightarrow Y$, between two sets of attributes X and Y that are subsets of R specifies a constraint on the possible tuples that can form a relation state r of R. The constraint is that, for any two tuples t1 and t2 in r that have

- t1[X] = t2[X], they must also have t1[Y] = t2[Y].

# Defining Functional Dependencies

- **If A and B are attributes of relation R, B is functionally dependent on A (denoted A → B), if each value of A in R is associated with exactly one value of B in R.**



- X → Y holds if whenever two tuples have the same value for X, they *must have* the same value for Y

  - For any two tuples t1 and t2 in any relation instance r(R): If t1[X]=t2[X], *then* t1[Y]=t2[Y]

- X → Y in R specifies a *constraint* on all relation instances r(R)

- Written as X → Y; can be displayed graphically on a relation schema as in Figures.  ( denoted by the arrow:  ).

- FDs are derived from the real-world constraints on the attributes

# Examples of FD constraints (1)

- Social security number determines employee name
    - SSN → ENAME

- Project number determines project name and location
    - PNUMBER → {PNAME, PLOCATION}

- Employee ssn and project number determines the hours per week that the employee works on the project
    - {SSN, PNUMBER} → HOURS

# Defining FDs from instances

- Note that in order to define the FDs, we need to understand the meaning of the attributes involved and the relationship between them.

- An FD is a property of the attributes in the schema R

- Given the instance (population) of a relation, all we can conclude is that an FD *may exist* between certain attributes.

- What we can definitely conclude is – that certain FDs *do not exist* because there are tuples that show a violation of those dependencies.

# Determining Functional Dependencies

**TEACH**

| Teacher | Course | Text |
|---------|--------|------|
| Smith | Data Structures | Bartram |
| Smith | Data Management | Martin |
| Hall | Compilers | Hoffman |
| Brown | Data Structures | Horowitz |

**Figure 14.7**
A relation state of TEACH with a *possible* functional dependency TEXT → COURSE. However, TEACHER → COURSE, TEXT → TEACHER and COURSE → TEXT are ruled out.

# Important Definitions

- **Determinant**
  - Refers to the attribute, or group of attributes, on the left-hand side of the arrow of a functional dependency.

- **Full Functional dependency:**
  - Indicates that if A and B are attributes of a relation, B is fully functionally dependent on A if B is functionally dependent on A, but not on any proper subset of A.
    - {StaffNo, StaffName} → BranchNo
    - StaffNo → BranchNo

- **Transitive Dependency**
  - A condition where A, B, and C are attributes of a relation such that if A → B and B → C, then C is transitively dependent on A via B (provided that A is not functionally dependent on B or C).

# Example Transitive Dependency

Staff Branch

| staffNo | sName | position | salary | branchNo | bAddress |
|---------|-------|----------|--------|----------|----------|
| SL21 | John White | Manager | 30000 | B005 | 22 Deer Rd, London |
| SG37 | Ann Beech | Assistant | 12000 | B003 | 163 Main St, Glasgow |
| SG14 | David Ford | Supervisor | 18000 | B003 | 163 Main St, Glasgow |
| SA9 | Mary Howe | Assistant | 9000 | B007 | 16 Argyll St, Aberdeen |
| SG5 | Susan Brand | Manager | 24000 | B003 | 163 Main St, Glasgow |
| SL41 | Julie Lee | Assistant | 9000 | B005 | 22 Deer Rd, London |

- **Consider functional dependencies in the Staff-Branch relation**

    **staffNo → sName, position, salary, branchNo, bAddress**

    **branchNo → bAddress**

- **Transitive dependency,**
    - **branchNo → bAddress exists on staffNo via branchNo**

29

# Normalization of Relations

- **Definition.** The normal form of a relation refers to the highest normal form condition that it meets, and hence indicates the degree to which it has been normalized.

- **Normalization of data:** process of analyzing the given relation schemas based on their FDs and primary keys to achieve the desirable properties of
  - (1) minimizing redundancy
  - (2) minimizing the insertion, deletion, and update anomalies

- **Normal form:**
  - Condition using keys and FDs of a relation to certify whether a relation schema is in a particular normal form

# Normalization of Relations

- 2NF, 3NF, BCNF
  - based on keys and FDs of a relation schema
- 4NF
  - based on keys, multi-valued dependencies : MVDs;
- 5NF
  - based on keys, join dependencies : JDs
- Additional properties may be needed to ensure a good relational design (lossless join, dependency preservation)

# Practical Use of Normal Forms

- The database designers *need not* normalize to the highest possible normal form
  - (usually up to 3NF and BCNF. 4NF rarely used in practice.)
- **DE normalization**:
  - The process of storing the join of higher normal form relations as a base relation—which is in a lower normal form

# Definitions of Keys and Attributes Participating in Keys

- If a relation schema has more than one key, each is called a **candidate** key.
    - One of the candidate keys is *arbitrarily* designated to be the **primary key**, and the others are called **secondary keys**.
- A **Prime attribute** must be a member of *some* candidate key
- A **Nonprime attribute** is not a prime attribute— that is, it is not a member of any candidate key.

# Non-additive or Lossless Join Decomposition

- If we decompose a relation R into relations R1 and R2,
  - Decomposition is lossy if R1 ⋈ R2 ⊃ R
  - Decomposition is lossless if R1 ⋈ R2 = R

- To check for lossless join decomposition using FD set, following conditions must hold:
  - Union of Attributes of R1 and R2 must be equal to attribute of R. Each attribute of R must be either in R1 or in R2.
    - **Att(R1) U Att(R2) = Att(R)**
  - Intersection of Attributes of R1 and R2 must not be NULL.
    - **Att(R1) ∩ Att(R2) ≠ Φ**
  - Common attribute must be a key for at least one relation (R1 or R2)
  - **Att(R1) ∩ Att(R2) -> Att(R1) or Att(R1) ∩ Att(R2) -> Att(R2)**

# Example

- R (A, B, C, D)

- FD: {A->BC} is decomposed into R1(ABC) and R2(AD) which is a lossless join decomposition as:

- First condition holds true as Att(R1) U Att(R2) = (ABC) U (AD) = (ABCD) = Att(R).

- Second condition holds true as Att(R1) ∩ Att(R2) = (ABC) ∩ (AD) ≠ Φ

- Third condition holds true as Att(R1) ∩ Att(R2) = A is a key of R1(ABC) because A->BC is given.

# Dependency Preserving Decomposition

- If we decompose a relation R into relations R1 and R2, All dependencies of R either must be a part of R1 or R2 or must be *derivable* from combination of FD's of R1 and R2.

- **For Example**,

- **R (A, B, C, D)** with **FD set{A->BC}** is decomposed into **R1(ABC)** and **R2(AD)** which is dependency preserving because FD **A->BC** is a part of **R1(ABC)**.

# FIRST NORMAL FORM

It states that
the domain of an attribute
must include only *atomic*
(simple, indivisible) *values*
and
that the value of any
attribute in a tuple must be
a *single value* from the
domain of
that attribute.

**(a)**
**DEPARTMENT**

| Dname | Dnumber | Dmgr_ssn | Dlocations |
|---|---|---|---|

**(b)**
**DEPARTMENT**

| Dname | Dnumber | Dmgr_ssn | Dlocations |
|---|---|---|---|
| Research | 5 | 333445555 | {Bellaire, Sugarland, Houston} |
| Administration | 4 | 987654321 | {Stafford} |
| Headquarters | 1 | 888665555 | {Houston} |

**(c)**
**DEPARTMENT**

| Dname | Dnumber | Dmgr_ssn | Dlocation |
|---|---|---|---|
| Research | 5 | 333445555 | Bellaire |
| Research | 5 | 333445555 | Sugarland |
| Research | 5 | 333445555 | Houston |
| Administration | 4 | 987654321 | Stafford |
| Headquarters | 1 | 888665555 | Houston |

**Figure 14.9**
Normalization into 1NF. (a)
A relation schema that is not
in 1NF. (b) Sample state of
relation DEPARTMENT. (c)
1NF version of the same
relation with redundancy.

# Normalizing into 1NF.

**(a)**
**EMP_PROJ**

| | | Projs | |
|---|---|---|---|
| Ssn | Ename | Pnumber | Hours |

**(b)**
**EMP_PROJ**

| Ssn | Ename | Pnumber | Hours |
|---|---|---|---|
| 123456789 | Smith, John B. | 1 | 32.5 |
| | | 2 | 7.5 |
| 666884444 | Narayan, Ramesh K. | 3 | 40.0 |
| 453453453 | English, Joyce A. | 1 | 20.0 |
| | | 2 | 20.0 |
| 333445555 | Wong, Franklin T. | 2 | 10.0 |
| | | 3 | 10.0 |
| | | 10 | 10.0 |
| | | 20 | 10.0 |
| 999887777 | Zelaya, Alicia J. | 30 | 30.0 |
| | | 10 | 10.0 |
| 987987987 | Jabbar, Ahmad V. | 10 | 35.0 |
| | | 30 | 5.0 |
| 987654321 | Wallace, Jennifer S. | 30 | 20.0 |
| | | 20 | 15.0 |
| 888665555 | Borg, James E. | 20 | NULL |

**Figure 14.10**
Normalizing nested relations into 1NF. (a) Schema of the EMP_PROJ relation with a nested relation attribute PROJS. (b) Sample extension of the EMP_PROJ relation showing nested relations within each tuple. (c) Decomposition of EMP_PROJ into relations EMP_PROJ1 and EMP_PROJ2 by propagating the primary key.

**(c)**
**EMP_PROJ1**

| Ssn | Ename |
|---|---|

**EMP_PROJ2**

| Ssn | Pnumber | Hours |
|---|---|---|

# Second Normal Form

- Uses the concepts of **FDs, primary key**

- Definitions
  - **Prime attribute:** An attribute that is member of the primary key K
  - **Full functional dependency:** a FD  Y -> Z where removal of any attribute from Y means the FD does not hold any more

- Examples:
  - {SSN, PNUMBER} -> HOURS is a full FD since neither SSN -> HOURS nor PNUMBER -> HOURS hold
  - {SSN, PNUMBER} -> ENAME is not  a full FD (it is called a partial dependency ) since SSN -> ENAME also holds

# Second Normal Form

- **Definition.**

- In the 2NF, relational must be in 1NF.

- In the second normal form, all non-key attributes are fully functional dependent on the primary key

- OR there should be no partial dependency

- Partial dependency can be identified by the rule:
    - LHS(proper subset of candidate key/Primary key) && RHS (Non Prime Attribute)

# Second Normal Form

(a)

EMP_PROJ

| Ssn | Pnumber | Hours | Ename | Pname | Plocation |

FD1
FD2
FD3

## Solve this

# Second Normal Form

(a)

EMP_PROJ

| Ssn | Pnumber | Hours | Ename | Pname | Plocation |
|-----|---------|-------|-------|-------|-----------|

FD1

FD2

FD3

2NF Normalization

EP1

| Ssn | Pnumber | Hours |
|-----|---------|-------|

FD1

EP2

| Ssn | Ename |
|-----|-------|

FD2

EP3

| Pnumber | Pname | Plocation |
|---------|-------|-----------|

FD3

# Normalize this table upto 3nf

Staff Branch

| staffNo | sName | position | salary | branchNo | bAddress |
|---------|-------|----------|--------|----------|----------|
| SL21 | John White | Manager | 30000 | B005 | 22 Deer Rd, London |
| SG37 | Ann Beech | Assistant | 12000 | B003 | 163 Main St, Glasgow |
| SG14 | David Ford | Supervisor | 18000 | B003 | 163 Main St, Glasgow |
| SA9 | Mary Howe | Assistant | 9000 | B007 | 16 Argyll St, Aberdeen |
| SG5 | Susan Brand | Manager | 24000 | B003 | 163 Main St, Glasgow |
| SL41 | Julie Lee | Assistant | 9000 | B005 | 22 Deer Rd, London |

# Example Transitive Dependency

Staff Branch

| staffNo | sName | position | salary | branchNo | bAddress |
|---------|-------|----------|--------|----------|----------|
| SL21 | John White | Manager | 30000 | B005 | 22 Deer Rd, London |
| SG37 | Ann Beech | Assistant | 12000 | B003 | 163 Main St, Glasgow |
| SG14 | David Ford | Supervisor | 18000 | B003 | 163 Main St, Glasgow |
| SA9 | Mary Howe | Assistant | 9000 | B007 | 16 Argyll St, Aberdeen |
| SG5 | Susan Brand | Manager | 24000 | B003 | 163 Main St, Glasgow |
| SL41 | Julie Lee | Assistant | 9000 | B005 | 22 Deer Rd, London |

- **Consider functional dependencies in the Staff-Branch relation**

  **staffNo → sName, position, salary, branchNo, bAddress**

  **branchNo → bAddress**

- **Transitive dependency,**
  - **branchNo → bAddress exists on staffNo via branchNo**

# Third Normal Form

- **Definition.** a relation schema R is in 3NF if it satisfies 2NF and no nonprime attribute of R is transitively dependent on the primary key.

- **Definition of Transitive functional dependency:**
  - **Transitive functional dependency:** a FD  X -> Z that can be derived from two FDs   X -> Y and Y -> Z

- **Examples:**
  - SSN -> DMGRSSN is a **transitive** FD
    - Since SSN -> DNUMBER and DNUMBER -> DMGRSSN hold
  - SSN -> ENAME is **non-transitive**
    - Since there is no set of attributes X where SSN -> X and X -> ENAME

# Third Normal Form

**General Definition.** A relation schema R is in third normal form (3NF) if, whenever a nontrivial functional dependency X → A holds in R, either
(a) X is a super key/Candidate key of R, **OR**
(b) A is a prime attribute of R.

**Figure 14.11**
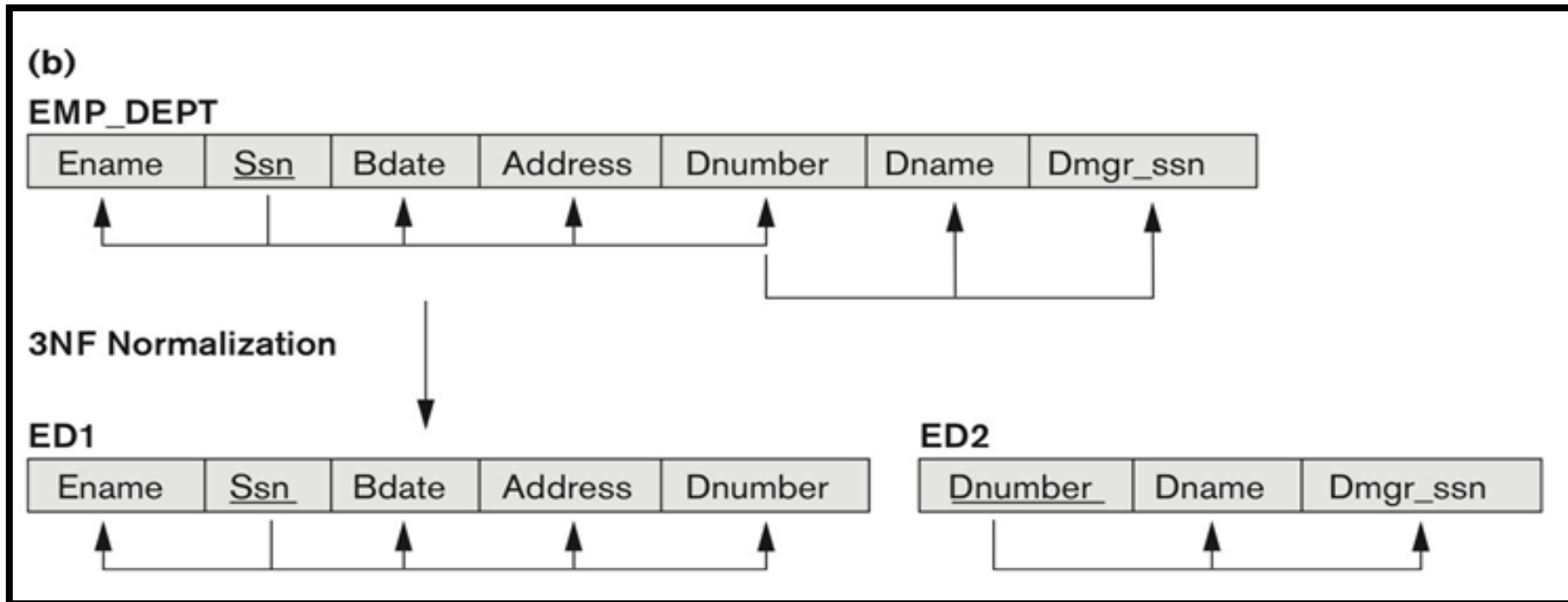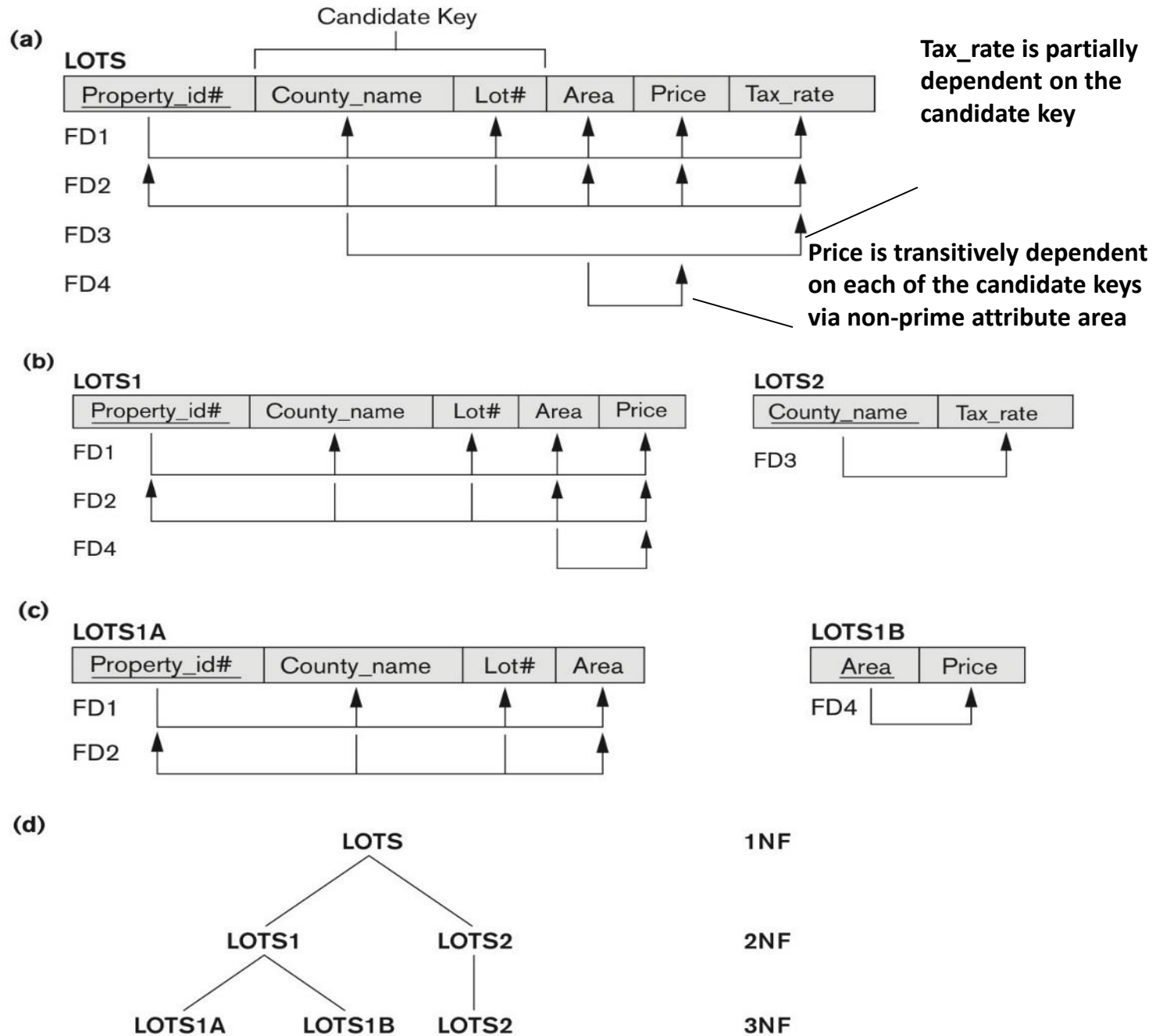Normalizing into 2NF and 3NF. (b) Normalizing EMP_DEPT into 3NF relations.

**Figure 14.12**
Normalization into 2NF and 3NF. (a) The LOTS relation with its functional dependencies FD1 through FD4. (b) Decomposing into the 2NF relations LOTS1 and LOTS2. (c) Decomposing LOTS1 into the 3NF relations LOTS1A and LOTS1B. (d) Progressive normalization of LOTS into a **3NF** design.

# Boyce-Codd Normal Form

- **Definition.** A relation schema R is in BCNF if whenever a nontrivial functional dependency X → A holds in R, **then X is a superkey of R.**
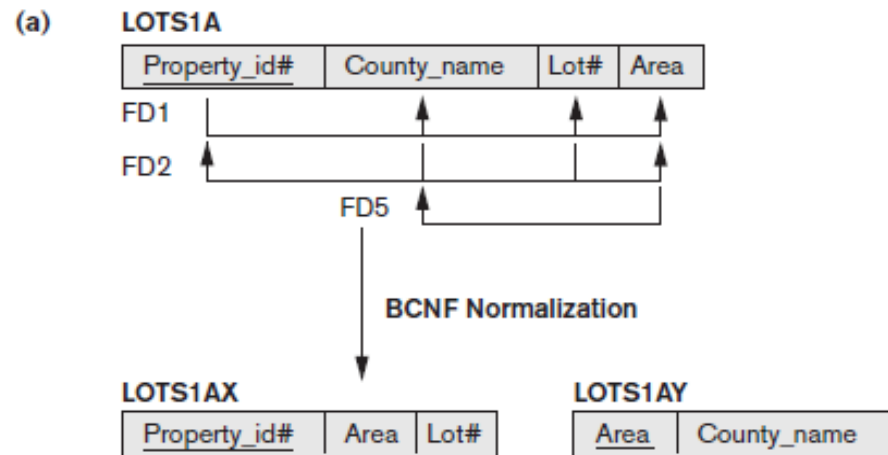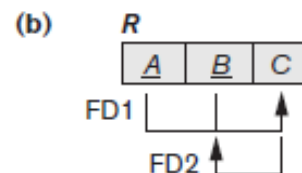


**Figure 14.13**
Boyce-Codd normal form. (a) BCNF normalization of LOTS1A with the functional dependency FD2 being lost in the decomposition. (b) A schematic relation with FDs; it is in 3NF, but not in BCNF due to the f.d. C → B.

# Decomposition of Relations not in BCNF

FD1: {Student, Course} ⟶ Instructor
FD2:14 Instructor ⟶ Course

**TEACH**

| Student | Course | Instructor |
|---------|--------|------------|
| Narayan | Database | Mark |
| Smith | Database | Navathe |
| Smith | Operating Systems | Ammar |
| Smith | Theory | Schulman |
| Wallace | Database | Mark |
| Wallace | Operating Systems | Ahamad |
| Wong | Database | Omiecinski |
| Zelaya | Database | Navathe |
| Narayan | Operating Systems | Ammar |

# Decomposition of Relations not in BCNF

FD1: {Student, Course} → Instructor

FD2:14 Instructor → Course

**Normalized Relations:**

R1 (Instructor, Course) and R2(Instructor, Student)

**TEACH**

| Student | Course | Instructor |
|---------|--------|-----------|
| Narayan | Database | Mark |
| Smith | Database | Navathe |
| Smith | Operating Systems | Ammar |
| Smith | Theory | Schulman |
| Wallace | Database | Mark |
| Wallace | Operating Systems | Ahamad |
| Wong | Database | Omiecinski |
| Zelaya | Database | Navathe |
| Narayan | Operating Systems | Ammar |

## StaffPropertyInspection

| propertyNo | iDate | iTime | pAddress | comments | staffNo | sName | carReg |
|---|---|---|---|---|---|---|---|
| PG4 | 18-Oct-12 | 10.00 | 6 Lawrence St, Glasgow | Need to replace crockery | SG37 | Ann Beech | M231 JGR |
| PG4 | 22-Apr-13 | 09.00 | 6 Lawrence St, Glasgow | In good order | SG14 | David Ford | M533 HDR |
| PG4 | 1-Oct-13 | 12.00 | 6 Lawrence St, Glasgow | Damp rot in bathroom | SG14 | David Ford | N721 HFR |
| PG16 | 22-Apr-13 | 13.00 | 5 Novar Dr, Glasgow | Replace living room carpet | SG14 | David Ford | M533 HDR |
| PG16 | 24-Oct-13 | 14.00 | 5 Novar Dr, Glasgow | Good condition | SG37 | Ann Beech | N721 HFR |

The First Normal Form(1NF) StaffPropertyInspection relation.

StaffPropertyInspection    (propertyNo, iDate, iTime, pAddress, comments, staffNo, sName, carReg)

## StaffPropertyInspection

| propertyNo | iDate | iTime | pAddress | comments | staffNo | sName | carReg |
|---|---|---|---|---|---|---|---|
| PG4 | 18-Oct-12 | 10.00 | 6 Lawrence St, Glasgow | Need to replace crockery | SG37 | Ann Beech | M231 JGR |
| PG4 | 22-Apr-13 | 09.00 | 6 Lawrence St, Glasgow | In good order | SG14 | David Ford | M533 HDR |
| PG4 | 1-Oct-13 | 12.00 | 6 Lawrence St, Glasgow | Damp rot in bathroom | SG14 | David Ford | N721 HFR |
| PG16 | 22-Apr-13 | 13.00 | 5 Novar Dr, Glasgow | Replace living room carpet | SG14 | David Ford | M533 HDR |
| PG16 | 24-Oct-13 | 14.00 | 5 Novar Dr, Glasgow | Good condition | SG37 | Ann Beech | N721 HFR |

The First Normal Form(1NF) StaffPropertyInspection relation.

**StaffPropertyInspection**

| propertyNo | iDate | iTime | pAddress | comments | staffNo | sName | carReg |
|---|---|---|---|---|---|---|---|

fd1 (Primary key)

fd2 (Partial dependency)

fd3 (Transitive dependency)

fd4

fd5 (Candidate key)

fd6 (Candidate key)

fd1   propertyNo, iDate → iTime, comments, staffNo,
      sName, carReg                                    (Primary key)
fd2   propertyNo → pAddress                            (Partial dependency)
fd3   staffNo → sName                                  (Transitive dependency)
fd4   staffNo, iDate → carReg
fd5   carReg, iDate, iTime → propertyNo, pAddress,
      comments, staffNo, sName                         (Candidate key)
fd6   staffNo, iDate, iTime → propertyNo, pAddress, comments  (Candidate key)

# Second Normal Form (2NF)

| | |
|---|---|
| Property | (propertyNo, pAddress) |
| PropertyInspection | (propertyNo, iDate, iTime, comments, staffNo, sName, carReg) |

# Third Normal Form (3NF)

| | |
|---|---|
| Property | (propertyNo, pAddress) |
| Staff | (staffNo, sName) |
| PropertyInspect | (propertyNo, iDate, iTime, comments, staffNo, carReg) |

# Boyce–Codd Normal Form (BCNF)

| | |
|---|---|
| StaffCar | (staffNo, iDate, carReg) |
| Inspection | (propertyNo, iDate, iTime, comments, staffNo) |

Property Relation
fd2    propertyNo → pAddress

Staff Relation
fd3    staffNo → sName

PropertyInspect Relation
fd1′    propertyNo, iDate → iTime, comments, staffNo, carReg
fd4    staffNo, iDate → carReg
fd5′    carReg, iDate, iTime → propertyNo, comments, staffNo
fd6′    staffNo, iDate, iTime → propertyNo, comments