

Limitations of DB Design Processes

- Provides a set of guidelines, does not result in a unique database schema
- Does not provide a way of evaluating alternative schemas
- Pitfalls:
 - Repetition of information
 - Inability to represent certain information
 - Loss of information
- Normalization theory provides a mechanism for analyzing and refining the schema produced by an E-R design

Redundancy and Other Problems

- Dependencies between attributes cause redundancy
 - Ex. All addresses in the same town have the same zip code
- Set valued attributes in the E-R diagram result in multiple rows in corresponding table
- Example: Person (PAN, Name, Address, Hobbies)
- A person entity with multiple hobbies yields multiple rows in table Person
 - Hence, the association between *Name* and *Address* for the same person is stored redundantly
- PAN is key of entity set, but (PAN, Hobby) is key of corresponding relation
- The relation Person can't describe people without hobbies

_	
ER Model Hobbies Person SIN Name Ad	<u>SSN Name Address Hobby</u> 1111 Joe 123 Main {biking, hiking}
Relational Model <u>SSN Name</u> 1111 Joe 1111 Joe	Address Hobby 123 Main biking 123 Main hiking Redundancy
	4

Redundancy leads to anomalies

- Update anomaly: A change in *Address* must be made in several places
- Deletion anomaly: Suppose a person gives up all hobbies. Do we:
 - Set Hobby attribute to null? No, since *Hobby* is part of key
 - Delete the entire row? No, since we lose other information in the row
- Insertion anomaly: *Hobby* value must be supplied for any inserted row since *Hobby* is part of key





- Result of E-R analysis need further refinement
- Appropriate decomposition can solve problems
- The underlying theory is referred to as *normalization* theory and is based on functional dependencies (and other kinds, like *multivalued dependencies*)



Functional Dependencies

- Definition: A *functional dependency* (FD) on a relation schema **R** is a constraint *X* → *Y*, where *X* and *Y* are subsets of attributes of **R**.
- <u>Definition</u>: An FD X → Y is *satisfied* in an instance r of R if for every pair of tuples, t and s: if t and s agree on all attributes in X then they must agree on all attributes in Y.
- <u>formally</u>: $\pi_X(t) = \pi_X(s) \Rightarrow \pi_y(t) = \pi_y(s)$ where $(\pi_X(t) \text{ is the projection of tuple t onto the attributes X})$





- Address → ZipCode
 Powai zip is 400076
- ArtistName \rightarrow BirthYear
 - Picasso was born in 1881
- VIN → Manufacturer, Engine type, ...
 VIN (vehicle information number) encodes all
 - VIN (vehicle information number) encodes all the information about manufacturer etc
- Author, Title → PublDate
 Shakespeare's Hamlet published in 1600





Entailment, Closure and Equivalence

- Definition: If F is a set of FDs on schema R and f is another FD on R, then F entails f if every instance r of R that satisfies every FD in F also satisfies f
 - Ex: $F = \{A \rightarrow B, B \rightarrow C\}$ and f is $A \rightarrow C$
 - If $Streetaddr \rightarrow Town$ and $Town \rightarrow Zip$ then $Streetaddr \rightarrow Zip$
 - Aka "Follows from"
- <u>Definition</u>: The *closure* of *F*, denoted *F*+, is the set of all FDs entailed by *F*
- **Definition**: *F* and *G* are *equivalent* if *F* entails *G* and *G* entails *F*









Soundness and Completeness

- Axioms are *sound*: If an FD *f*: X→ Y can be derived from a set of FDs *F* using the axioms, then *f* holds in every relation that satisfies every FD in *F*.
- Axioms are *complete*: If F entails f, then f can be derived from F using the axioms
- A consequence of completeness is the following (naïve) algorithm to determining if *F* entails *f*:
- <u>Algorithm</u>: Use the axioms in all possible ways to generate *F*+ (the set of possible FD's is finite so this can be done) and see if *f* is in *F*+

17

Befine the probability of the probability of

Augmentation

- If $X \rightarrow Y$, then $XZ \rightarrow YZ$ for any Z
- Let R=(A,B,C,D,E), $X=\{AB\} Y=\{CD\} Z=\{E\}$
- e.g.:
- t1 = (a1, b1, c1, d1, e1)
- t2 = (a2, b2, c2, d2, e2)
- $\pi_{XZ}(t1) = \pi_{XZ}(t2) \Rightarrow a1 = a2, b1 = b2, e1 = e2$

19

- Since $X \rightarrow Y$ and e1 = e2
- then c1 = c2, d1 = d2, e1 = e2
- $\blacksquare \Rightarrow \pi_{YZ}(t1) = \pi_{YZ}(t2)$

Transitivity If X -> Y, and Y -> Z then X -> Z Let R = (A,B,C,D,E), X={AB}, Y={CD} X={E} e.g.: t1 =(a1,b1,c1,d1,e1) t2 =(a2,b2,c2,d2,e2) assume X ->Y and Y -> Z $\pi_X(t1) = \pi_X(t2) => a1 = a2,b1 = b2$ Since X -> Y then c1 = c2,d1 = d2 $=> \pi_Y(t1) = \pi_Y(t2)$ Since Y -> Y then e1 = e2 $=> \pi_Z(t1) = \pi_Z(t2)$





Example $\blacksquare F: AB \to C, A \to D, D \to E, AC \to B$ Χ $X+_F$ Α $\{A, D, E\}$ AB $\{A, B, C, D, E\}$ (Hence AB is a key) В *{B}* D $\{D, E\}$ • Is $AB \rightarrow E$ entailed by F? Yes • Is $D \rightarrow C$ entailed by F? No ■ => $X+_F$ allows us to determine FDs of the form $X \rightarrow Y$ entailed by F23



Example

- **Problem**: Compute the attribute closure of *AB* with respect to the set of FDs :
 - $\blacksquare AB \rightarrow C \dots (a)$
 - $\blacksquare A \rightarrow D \dots (b)$
 - $\blacksquare D \longrightarrow E \dots (c)$
 - $\blacksquare AC \rightarrow B \dots (d)$
- Initially $closure = \{AB\}$
- Using (a) $closure = \{ABC\}$
- Using (b) $closure = \{ABCD\}$
- Using (c) *closure* = {*ABCDE*}



BCNF

- **Definition**: A relation schema **R** is in BCNF if for every FD *X*→ *Y* associated with **R** either
- 1. $Y \subseteq X$ (i.e., the FD is trivial) or
- 2. X is a superkey of **R**
- *Example*: Person1(*PAN*, *Name*, *Address*)
- The only FD is $PAN \rightarrow Name$, Address
- Since *PAN* is a key, Person1 is in BCNF

27

More examples

- Person (*PAN*, *Name*, *Address*, *Hobby*)
- The FD *PAN* → *Name, Address* does not satisfy requirements of BCNF since the key is (*PAN, Hobby*)
- HasAccount (*AccountNumber*, *ClientId*, *OfficeId*)
- The FD *AcctNum*→ *OfficeId* does not satisfy BCNF
- Requirements since keys are (*ClientId*, *OfficeId*) and (*AcctNum*, *ClientId*)





Decompositions

- Schema $\mathbf{R} = (R, F)$
 - R is set a of attributes
 - F is a set of functional dependencies over R
- The *decomposition of schema* **R** is a collection of schemas **R**i = (*Ri*, *Fi*) where
 - $R = \bigcup i Ri (no new attributes)$
 - *Fi* is a set of functional dependences involving only attributes of *Ri*
 - *F* entails *Fi* for all *i* (*no new FDs*)
- The *decomposition of an instance*, **r**, of **R** is a set of relations **r***i* = π_{Ri}(**r**) for all *i*

```
31
```

<section-header>Example●. Schema (R, F) where●. A = {PAN, Name, Address, Hobby}●. B = {PAN → Name, Address}●. B a be decomposed into●. A = {PAN, Name, Address}●. B = {PAN → Name, Address}●. B = {PAN → Name, Address}●. B = {PAN, Hobby}●. F = { }

Decomposition into BCNF

- Consider relation R with FDs F.
- If X -> Y violates BCNF (and X \cap Y = \emptyset), decompose R into XY and R Y.
- Repeated application of this idea will give us a collection of relations that are in BCNF and guaranteed to terminate.
- e.g., CSJDPQV, key C, JP->C, SD->P, J->S
- To deal with SD -> P, decompose into SDP, CSJDQV.
- To deal with J -> S, decompose CSJDQV into JS and CJDQV
- In general, several dependencies may cause violation of BCNF.







Testing for LosslePANess

- A (binary) decomposition of $\mathbf{R} = (R, F)$ into $\mathbf{R}1 = (R1, F1)$ and $\mathbf{R}2 = (R2, F2)$ is lossless *if and only if*:
- either the FD $(R1 \cap R2) \rightarrow R1$ is in F+
- or the FD $(R1 \cap R2) \rightarrow R2$ is in F+
- Intuitively: the attributes common to R1 and R2 must contain a key for either R1 or R2.

37

Example

- Schema (R, F) where R = {PAN, Name, Address, Hobby}, F = {PAN → Name, Address} can be decomposed into
- $R1 = \{PAN, Name, Address\}, F1 = \{PAN \rightarrow Name, Address\} \&$
- $R2 = \{PAN, Hobby\}, F2 = \{\}$
- Since $R1 \cap R2 = PAN$ and $PAN \rightarrow R1$ the decomposition is lossless

The intuition behind this

Suppose R1 ∩ R2 → R2. Then a row of r1 can combine with exactly one row of r2 in the natural join (since in r2 a particular set of values for the attributes in R1 ∩ R2 defines a unique row)



Dependency Preservation Consider a decomposition of R = (R, F) into R1 = (R1, F1) and R2 = (R2, F2) An FD X → Y of F is in Fi iff X ∪ Y ⊆ Ri An FD, f ∈ F may be in neither F1, nor F2, nor even (F1 ∪ F2)+ Checking that f is true in r1 or r2 is (relatively) easy Checking f in r1 ∝ r2 is harder – requires a join *Ideally*: want to check FDs locally, in r1 and r2, and have a guarantee that every f ∈ F holds in r1 ∝ r2 The decomposition is dependency preserving iff the sets F and F1 ∪ F2 are equivalent: F+ = (F1 ∪ F2)+ Then checking all FDs in F, as r1 and r2 are updated, can be done by checking F1 in r1 and F2 in r2

- If f is an FD in F, but f is not in $F1 \cup F2$, there are two possibilities:
- 1. $f \in (F1 \cup F2)$ +
- If the constraints in *F1* and *F2* are maintained, *f* will be maintained automatically.
- 2. $f \notin (F1 \cup F2)$ +
- *f* can be checked only by first taking the join of r*l* and r2. This is costly.

41

Example

- Schema (R, F) where R = {SSN, Name, Address, Hobby} & F = {SSN → Name, Address} can be decomposed into
- $R1 = \{SSN, Name, Address\}, F1 = \{SSN \rightarrow Name, Address\} \&$
- $R2 = \{SSN, Hobby\}, F2 = \{\}$
- Since $F = F1 \cup F2$, it trivially holds that $F + = (F1 \cup F2) + \text{ i.e.}$, the decomposition is dependency preserving

Example

- Schema: (ABC; F), $F = \{A \rightarrow B, B \rightarrow C, C \rightarrow B\}$
- Decomposition:
- 1. $(AC, F1), F1 = \{A \rightarrow C\}$
- Note: $A \rightarrow C \notin F$, but in F +
- 2. $(BC, F2), F2 = \{B \rightarrow C, C \rightarrow B\}$
- $A \to B \notin (F1 \cup F2)$, but $A \to B \in (F1 \cup F2)$ +.
- So $F + = (F1 \cup F2) +$ and thus the decompositions is still dependency preserving
 ⁴³



BCNF Decomposition Algo

Input: $\mathbf{R} = (R; F)$ *Decomp* := \mathbf{R} while there is $\mathbf{S} = (S; F') \in Decomp$ and \mathbf{S} not in BCNF do Find $X \rightarrow Y \in F'$ that violates BCNF // X isn't a superkey in \mathbf{S} Replace \mathbf{S} in *Decomp* with $\mathbf{S1} = (XY; F1)$, $\mathbf{S2} = (S - (Y - X); F2)$ // F1 = all FDs of F' involving only attributes of XY // F2 = all FDs of F' involving only attributes of S - (Y - X) end return *Decomp*

```
45
```



Properties of algorithm

- Let $X \rightarrow Y$ violate BCNF in $\mathbf{R} = (R,F)$ and $\mathbf{R1} = (R1,F1)$, $\mathbf{R2} = (R2,F2)$ is the resulting decomposition. Then:
- There are *fewer violations* of BCNF in **R1** and **R2** than there were in **R**
- $X \rightarrow Y$ implies X is a key of **R1**
- Hence $X \rightarrow Y \in F1$ does not violate BCNF in **R1** and, since
- $X \rightarrow Y \notin F2$, does not violate BCNF in **R2** either
- Suppose f is $X' \to Y'$ and $f \in F$ doesn't violate BCNF in **R**.
- If $f \in F1$ or F2 it does not violate BCNF in **R1** or **R2** either since X' is a superkey of **R** and hence also of **R1** and **R2**.

47

• The decomposition is *lossless* since $F1 \cap F2 = X$



More properties

- A BCNF decomposition is *not necessarily* dependency preserving
- But always lossless
- BCNF+lossless+dependency preserving is sometimes unachievable (recall HasAccount)



Example

- HasAccount (AcctNum, ClientId, OfficeId)
- ClientId, OfficeId \rightarrow AcctNum
- OK since LHS contains a key
- AcctNum \rightarrow OfficeId
- OK since RHS is part of a key
- HasAccount is in 3NF but it might still contain redundant information due to AcctNum → OfficeId (which is not allowed by BCNF)



Another Example

- Person (SSN, Name, Address, Hobby)
- (SSN, Hobby) is the only key.
- *SSN→ Name* violates 3NF conditions since *Name* is not part of a key and *SSN* is not a superkey



Minimal Cover

- A *minimal cover* of a set of dependencies, *T*, is a set of dependencies, *U*, such that:
- 1. *U* is equivalent to T(T + = U +)
- 2. All FDs in *U* have the form $X \rightarrow A$ where *A* is a single attribute
- 3. It is not possible to make *U* smaller (while preserving equivalence) by
 - Deleting an FD OR
 - Deleting an attribute from an FD (either from LHS or RHS)
- FDs and attributes that can be deleted in this way are called *redundant*





- Example: Can an attribute be deleted from $ABH \rightarrow C$?
- Compute AB + T, AH + T, BH + T.
- Since $C \in (BH)_{+_T}$, $BH \rightarrow C$ is entailed by T and A is redundant in $ABH \rightarrow C$.

57



Synthesizing 3NF Schemas

Starting with a schema $\mathbf{R} = (R, T)$ **step 1**: Compute a minimal cover, *U*, of *T*. The decomposition is based on *U*, but since U + = T + the same functional dependencies will hold A minimal cover for $T = \{ABH \rightarrow CK, A \rightarrow D, C \rightarrow E, BGH \rightarrow F, F \rightarrow AD, E \rightarrow F, BH \rightarrow E\}$ is

 $U = \{BH \rightarrow C, BH \rightarrow K, A \rightarrow D, C \rightarrow E, F \rightarrow A, E \rightarrow F\}$

59

step 2: Partition *U* into sets *U1*, *U2*, ... *Un* such that the LHS of all elements of *Ui* are the same $U1 = \{BH \rightarrow C, BH \rightarrow K\}, \\ U2 = \{A \rightarrow D\}, \\ U3 = \{C \rightarrow E\}, \\ U4 = \{F \rightarrow A\}, \\ U5 = \{E \rightarrow F\}$





BCNF Design Strategy

- The resulting decomposition, **R0**, **R1**, ... **Rn**, is
- 1. Dependency preserving (since every FD in U is a FD of some schema)
- 2. Lossless (although this is not obvious)
- 3. In 3NF (although this is not obvious)
- Strategy for decomposing a relation
- 1. Use 3NF decomposition first to get lossless, dependency preserving decomposition
- If any resulting schema is not in BCNF, split it using the BCNF algorithm (but this may yield a nondependency preserving result)



Denormalization

- Tradeoff: *Judiciously* introduce redundancy to improve performance of certain queries
- Example: Add attribute Name to Transcript (making it Transcript') SELECT T.Name, T.Grade FROM Transcript' T

WHERE T.CrsCode = 'CS317' AND T.Semester = 'F2007'

- Join is avoided
- If queries are asked more frequently than Transcript is modified, added redundancy might improve average performance
- But, Transcript' is no longer in BCNF since key is (*StudId*, *CrsCode*, *Semester*) and *StudId* → *Name*

65

Attribute Independence in **BCNF** BCNF schemas can have redundancy, e.g., when we force two or more manymany relationships in a single DeptName Number Textbook Instructor relation. FCDB 4604 CS Murali Consider the schema Courses(Number, DeptName, 4604 CS | SQL Made Easy Murali Textbook, Instructor). 4604 CS FCDB Ramakrishan • Each Course can have multiple 4604 CS | SQL Made Easy Ramakrishan required Textbooks. • Each Course can have multiple Instructors. Instructors can use any required textbook in a Course. 66





MVD

- A multi-valued dependency (MVD or MD) is an assertion that two sets of attributes are independent of each other.
- The multi-valued dependency A1 A2 . . .An ->> B1 B2 . . .Bm holds in a relation R if for every pair of tuples t and u in R that agree on all the A's, we can find a tuple v in R that agrees
- 1. with both t and u on A's,
- 2. with t on the B's, and
- 3. with u on all those attributes of R that are not A's or B's.

[Number	DeptName	Textbook	Instructor
ĺ	4604	CS	FCDB	Murali
ľ	4604	CS	SQL Made Easy	Murali
Ī	4604	CS	FCDB	Ramakrishan
	4604	CS	SQL Made Easy	Ramakrishan

Example

- Number DeptName ->>Textbook is an MD. For every pair of tuples t and u that agree on Number and DeptName, we can find a tuple v that agrees
- 1. with both t and u on Number and DeptName,
- 2. with t on Textbook, and with u on Instructor.
- Number DeptName ->> Instructor is an MD. For every pair of tuples t and u that agree on Number and DeptName, we can find a tuple v that agrees
- 1. with both t and u on Number and DeptName,
- 2. with t on Instructor, and with u on Textbook.







Fourth Normal Form

- A relation R is in fourth normal form (4NF) if for every non-trivial MD A1 A2 . . . An ->> B1B2 . . . Bm, {A1,A2, . . . , An} is a superkey.
- A relation in 4NF is also in BCNF since an FD is a special case of an MD.



Relationships am Forms	ongs	st No	rmal				
 4NF implies BCNF, i.e., if a relation is in 4NF, it is also in BCNF. BCNF implies 3NF, i.e., if a relation is in BCNF, it is also in 3NF. 							
Property	3NF	BCNF	4NF				
Eliminates redundancy due to FDs	Maybe	Yes	Yes				
Eliminates redundancy due to MDs	No	No	Yes				
Preserves FDs	Yes	Maybe	Maybe				
1165617651105	4						
Preserves MDs	Maybe	Maybe	Maybe				

