



Design Theory for Relational DBs: Functional Dependencies, Schema Decomposition, Normal Forms

EECS3421 - Introduction to Database Management Systems



# Database Design Theory

- Guides systematic improvements to database schemas
- General idea:
  - Express constraints on the data
  - Use these to decompose the relations
- Ultimately, get a schema that is in a "normal form"
  - guarantees certain desirable properties
  - "normal" in the sense of conforming to a standard
- The process of converting a schema to a normal form is called *normalization*

## Goal #1: remove redundancy

#### Consider this schema

Student Name	Student Email	Course	Instructor
Xiao	xiao@gmail	EECS3421	Smith
Xiao	xiao@gmail	EECS2031	Brown
Jaspreet	jaspreet@gmail	EECS3421	Smith

- What if... Xiao changes email addresses?
  - update anomaly: need to update more than one tuples
- What if... Xiao drops EECS2031?
  - deletion anomaly: loss of information that course is taught by Brown
- What if... We need to create a new course, EECS4411
  - *insertion anomaly*: how to fill rest of information (name, email, ...)

Multiple relations => exponentially worse

#### Goal #2: expressing constraints

#### • Consider the following schemata:

Students(yorkid, name, email) vs.

Students(yorkid, name) Emails(yorkid, address)

#### • Consider also:

House(street, city, value, owner, propertyTax)

VS.

House(street, city, value, owner) TaxRates(city, value, propertyTax) Maybe a student has more than one emails that we would like to register (in the first schema there will be a **redundancy**)

TaxRates are defined by city, so there is no need to repeat for each single House (in the first schema there will be a **redundancy**)

#### Dependencies, constraints are domain-dependent<sup>4</sup>

### Overview

- Part I: Functional Dependencies
- Part II: Schema Decomposition
- Part III: Normal Forms

## PART 1: FUNCTIONAL DEPENDENCIES



#### **Functional dependencies**

- Let X, Y be sets of attributes from relation R
- X -> Y (we say: "X functionally determines Y")
  - Any tuples in R which agree in all attributes of X must also agree in all attributes of Y
  - Or, "The values of attributes Y are a function of those in X"
  - Not necessarily an easy function to compute, mind you
  - => Consider X -> h, where h is the hash of attributes in X

#### Notational conventions

- "a", "b", "c" specific attributes
- "A", "B", "C" sets of (unnamed) attributes
- abc -> def same as {a,b,c} -> {d,e,f}

#### Most common to see <u>singletons</u> (X -> y or abc -> d)

# Rules and principles about FDs

- Rules
  - The splitting/combining rule
  - Trivial FDs
  - The transitive rule
- Algorithms related to FDs
  - the closure of a set of attributes of a relation
  - a minimal basis of a relation

## The Splitting/Combining rule of FDs

- Attributes on right independent of each other
  - Consider a,b,c -> d,e,f
  - "Attributes a, b, and c functionally determine d, e, and f"
    => No mention of d relating to e or f directly
- Splitting rule (useful to split up right side of FD)
   *abc -> def* becomes *abc -> d*, *abc -> e* and *abc -> f*
- No safe way to split left side
  - abc -> def is NOT the same as ab -> def and c -> def!
- Combining rule (useful to combine right sides):
  - *if abc -> d*, *abc -> e*, *abc -> f* holds, then *abc -> def* holds

## Splitting FDs – example

- Consider the relation and FD
  - EmailAddress(user, domain, firstName, lastName)
  - user, domain -> firstName, lastName
- The following hold
  - user, domain -> firstName
  - user, domain -> lastName
- The following do NOT hold!
  - user -> firstName, lastName
  - domain -> firstName, lastName

Gotcha: "doesn't hold" = "not all tuples" != "all tuples not"

## **Trivial FDs**

- Not all functional dependencies are useful
  - A -> A always holds
  - abc -> a also always holds (right side is subset of left side)
- FD with an attribute on both sides is "trivial"
  - Simplify by removing L ∩ R from R
     *abc* -> *ad* becomes *abc* -> *d*
  - Or, in singleton form, delete trivial FDs
     abc -> a and abc -> d becomes just abc -> d

### Transitive rule

- The transitive rule holds for FDs
  - Consider the FDs: *a* -> *b* and *b* -> *c*; then *a*->*c* holds
  - Consider the FDs: ad -> b and b -> cd; then ad->cd holds or just ad->c (because of the trivial dependency rule)

## Identifying functional dependencies

- FDs are domain knowledge
  - Intrinsic features of the data you're dealing with
  - Something you know (or assume) about the data
- Database engine cannot identify FDs for you
  - Designer must specify them as part of schema
  - DBMS can only enforce FDs when told to
- DBMS cannot safely "optimize" FDs
  - It has only a finite sample of the data
  - An FD constrains the entire domain

# Coincidence or FD?

ID	Email	City	Country	Surname
1983	tom@gmail.com	Toronto	Canada	Fairgrieve
8624	mar@bell.com	London	Canada	Samways
9141	scotty@gmail.com	Winnipeg	Canada	Samways
1204	birds@gmail.com	Aachen	Germany	Lakemeyer

- What if we try to infer FDs from the data?
  - ID -> email, city, country, surname
  - email -> city, country, surname
  - city -> country
  - surname -> country

#### Domain knowledge required to validate FDs

# Keys and FDs

- Consider relation R with attributes A
- Superkey
  - Any  $\mathbf{S} \subseteq \mathbf{A}$  s.t.  $\mathbf{S} \rightarrow \mathbf{A}$

=> Any subset of **A** which determines all remaining attributes in **A** 

- Candidate key (or key)
  - $\mathbf{C} \subseteq \mathbf{A}$  s.t.  $\mathbf{C} \rightarrow \mathbf{A}$  and  $\mathbf{X} \rightarrow \mathbf{A}$  does not hold for any  $\mathbf{X} \subset \mathbf{C}$
  - => A superkey which contains no other superkeys
  - => Remove any attribute from **C** and you no longer have a key
- Primary key
  - The candidate key we use to identify the relation
  - => Always exists, only one allowed, doesn't matter which C we use
- Prime attribute
  - $\exists$  candidate key **C** s.t.  $x \in C$
  - => attribute that participates in at least one key

## FD: relaxes the concept of a "key"

- Superkey: X -> R
  - A superkey must include all remaining attributes of the relation on the RHS (Right-Hand-Side)
- Functional dependency: X -> Y
  - An FD can involve just a subset of them
- Example:

Houses(street, city, value, owner, tax)

- street, city -> value, owner, tax (both FD and key)
- city,value -> tax (FD only)

# Cyclic functional dependencies?

- Attributes on right side of one FD may appear on left side of another!
  - Simple example: assume relation (A, B) & FDs: A->B, B->A
  - What does this say about A and B?
- Example
  - studentID->email email->studentID

### Geometric view of FDs

- Let D be the domain of tuples in R
  - Every possible tuple is a point in **D**
- FD X on R restricts tuples in R to a subset of D
  - Points in **D** which violate **X** cannot be in **R**



## Inferring functional dependencies

- Problem
  - Given FDs  $X_1 \rightarrow a_1, X_2 \rightarrow a_2$ , etc.
  - Does some FD Y -> B (not given) also hold?
- Consider the dependencies
  - A -> B, B -> C

Does A -> C hold?

Intuitively, A -> C also holds The given FDs entail (imply) it (transitivity rule)

#### How to prove it in the general case?

#### **Closure test for FDs**

- Consider relation R
- Given attribute set  $A \subseteq R$  and FD set F
  - Denote A<sub>F</sub><sup>+</sup> as the closure of A relative to F

 $\Rightarrow A_{F}^{+} = \text{set of all FDs given or implied by } A$ 

• Computing the [transitive] closure of A

- While  $\exists X \in F'$  s.t. LHS(X)  $\subseteq A_F^+$ :

 $A_{F}^{+} = A_{F}^{+} U RHS(X)$ F' = F' - X

- At end: A -> B  $\forall B \in A_F^+$ 

#### Closure test – example

- Consider R(a,b,c,d,e,f) with FDs set F = {ab -> c, ac -> d, c -> e, ade -> f }
- Find  $A_F^+$  if A = ab or find  $\{a,b\}^+$



 ${a,b}^+={a,b,c,d,e,f}$  or  $ab \rightarrow cdef \rightarrow ab$  is a candidate key! <sup>22</sup>

#### Example : Closure Test

<b>R</b> (A, B, C, D, E)	X	<b>X</b> <sub>F</sub> <sup>+</sup>
<b>F</b> : AB -> C	A	{A, D, E}
A -> D	AB	{A, B, C, D, E}
D -> E	AC	{A, C, B, D, E}
AC -> B	В	{B}
	D	{D, E}

ls	AB -> E entailed by F?	Yes
ls	D -> C entailed by F?	No

Result:  $X_F^+$  allows us to determine all FDs of the form X -> Y entailed by F

## Discarding redundant FDs

- Minimal basis: opposite extreme from closure
- Given a set of FDs F, want to find **minimal basis F'** s.t.
  - −  $F' \subseteq F$
  - F' entails  $X \forall X \in F$
- Properties of a minimal basis F'
  - RHS is always singleton
  - If any FD is removed from F', F' is no longer a minimal basis
  - If for any FD in F' we remove one or more attributes from the LHS of  $X \in F'$ , the result is no longer a minimal basis

## Constructing a minimal basis

Straightforward but time-consuming

- 1. Split all RHS into singletons
- 2.  $\forall X \in F'$ , test whether J = (F'-X)<sup>+</sup> is still equivalent to F<sup>+</sup>
- => Might make F' too small
- 3.  $\forall i \in LHS(X) \ \forall X \in F'$ , let LHS(X')=LHS(X)-i Test whether (F'-X+X')<sup>+</sup> is still equivalent to F<sup>+</sup>
- => Might make F' too big
- 4. Repeat (2) and (3) until neither makes progress

### Minimal Basis: Example

- Relation R: R(A, B, C, D)
- Defined FDs:
  - $F = \{A \rightarrow AC, B \rightarrow ABC, D \rightarrow ABC\}$

Find the minimal Basis M of F

## Minimal Basis: Example (cont.)

1<sup>st</sup> Step

- H = {A->A, A->C, B->A, B->B, B->C, D->A, D->B, D->C}

2<sup>nd</sup> Step

- A->A, B->B: can be removed as trivial
- A->C: can't be removed, as there is no other LHS with A
- B->A: can't be removed, because for J=H-{B->A} is B+=BC
- B->C: can be removed, because for J=H-{B->C} is B+=ABC
- D->A: can be removed, because for J=H-{D->A} is D+=DBA
- D->B: can't be removed, because for J=H-{D->B} is D+=DC
- D->C: can be removed, because for J=H-{D->C} is D+=DBAC Step outcome => H = {A->C, B->A, D->B}

## Minimal Basis: Example (cont.)

3<sup>rd</sup> Step

H doesn't change as all LHS in H are single attributes
 4<sup>th</sup> Step

- H doesn't change

Minimal Basis:  $M = H = \{A \rightarrow C, B \rightarrow A, D \rightarrow B\}$ 

Caveat: Different minimal bases are possible

### PART II: SCHEMA DECOMPOSITION



## FDs and redundancy

- Given relation R and FDs F
  - **R** often exhibits anomalies due to redundancy
  - F identifies many (not all) of the underlying problems
- Idea
  - Use **F** to identify "good" ways to split relations
  - Split **R** into 2+ smaller relations having less redundancy
  - Split up F into subsets which apply to the new relations (compute the projection of functional dependencies)

### Schema decomposition

- Given relation R and FDs F
  - Split **R** into  $\mathbf{R}_i$  s.t.  $\forall i \mathbf{R}_i \subset \mathbf{R}$  (no new attributes)
  - Split F into F<sub>i</sub> s.t. ∀i F entails F<sub>i</sub> (no new FDs)
  - Note: F<sub>i</sub> involves only attributes in R<sub>i</sub>
- Caveat: entirely possible to lose information
  - $F^+$  may entail FD X which is not in  $(U_i F_i)^+$
  - => Decomposition lost some FDs (*dependency not preserved*)
  - Possible to have  $\mathbf{R} \subset \bowtie_i \mathbf{R}_i$
  - => Decomposition lost some relationship (*lossy decomposition*)
- Goal: minimize anomalies without losing info

#### We'll revisit information loss later

## **Desired Properties of Decomposition**

- Lossless-join
- Dependency-preserving
- Anomaly-free (no redundancies)

#### This may be achieved through the use of Normal Forms <sup>37</sup>

## Splitting relations – example

• Consider the following relation **R**:

Student Name	Student Email	Course	Instructor
Xiao	xiao@gmail	eecs3421	Smith
Xiao	xiao@gmail	eecs4411	Brown
Jaspreet	jaspreet@gmaj	eecs3421	Smith

One possible decomposition of R

Students(name, email) Taking(email, course) Courses(course, instructor)

• Students ⋈ Taking ⋈ Courses reconstructs the right tuples!

# Gotcha: lossy join decomposition

• Consider a relation **R** with one more tuple

Student Name	Student Email	Course	Instructor
Xiao	xiao@gmail	eecs3421	Smith
Xiao	xiao@gmail	eecs4411	Brown
Jaspreet	jaspreet@gmail	eecs3421	Smith
Mary	mary@gmail	eecs4411	Rosenburg

- Students ⋈ Taking ⋈ Courses has bogus tuples!
  - Mary is not taking Brown's section of eecs4411
  - Xiao is not taking Rosenburg's section of eecs4411

Why did this happen? How to prevent it?

### Information loss with decomposition

- Decompose R into S and T
  - Consider FD  $a \rightarrow b$ , with a only in S and b only in T
- FD loss
  - Attributes a and b no longer in same relation
  - => Must join T and S to enforce  $a \rightarrow b$  (expensive)
- Join loss
  - LHS and RHS no longer in same relation, no other connection
  - Neither  $(S \cap T) \rightarrow S$  nor  $(S \cap T) \rightarrow T$  in F<sup>+</sup>
  - => Joining T and S produces bogus tuples (irreparable)
- In our example:
  - ( $\{email, course\} \cap \{course, instructor\}$ ) =  $\{course\}$
  - course -/-> course, instructor and course -/-> email, course

# **Projecting FDs**

- Once we've split a relation we have to refactor our FDs to match
  - Each FDs must only mention attributes from one relation
- Similar to geometric projection
  - Many possible projections (depends on how we slice it)
  - Keep only the ones we need (minimal basis)



# FD projection algorithm

- Start with  $\mathbf{F}_{i} = \emptyset$
- For each subset X of R<sub>i</sub>
  - Compute X+
  - For each attribute a in X+
    - If a is in R<sub>i</sub>
      - $_{\circ}\,$  add X  $\rightarrow$  a to  $\mathbf{F}_{\mathrm{i}}$
- Compute the minimal basis of F<sub>i</sub>
- Projection is expensive
  - Suppose  $R_1$  has *n* attributes
  - How many subsets of R<sub>1</sub> are there?

# Making projection more efficient

- Ignore trivial dependencies
  - No need to add  $X \rightarrow A$  if A is in X itself
- Ignore trivial subsets
  - The empty set or the set of all attributes (both subsets of **X**)
- Ignore supersets of X if  $X^+ = R$ 
  - They can only give us "weaker" FDs (with more on the LHS)

## Example: Projecting FD's

- ABC with FD's A->B and B->C
  - *A*+=*ABC*; yields *A*->*B*, *A*->*C* 
    - We ignore A->A as trivial
    - We ignore the supersets of A, AB + and AC +, because they can only give us "weaker" FDs (with more on the LHS)
  - *B*<sup>+</sup>=*BC* ; yields *B*->*C*
  - $C^+=C$ ; yields nothing.
  - BC<sup>+</sup>=BC; yields nothing.

## **Example** -- Continued

- Resulting FD's: *A*->*B*, *A*->*C*, and *B*->*C*
- Projection onto AC : A->C
  - Only FD that involves a subset of  $\{A, C\}$
- Projection on BC: *B*->*C* 
  - Only FD that involves subset of {B, C}

### PART III: NORMAL FORMS



## Motivation for normal forms

- Identify a "good" schema
  - For some definition of "good"
  - Avoid anomalies, redundancy, etc.
- Many normal forms
  - 1<sup>st</sup>
  - 2<sup>nd</sup>
  - 3<sup>rd</sup>
  - Boyce-Codd
  - ... and several more we won't discuss...

#### $BCNF \subseteq 3NF \subseteq 2NF \subseteq 1NF$ (focus on 3NF/BCNF)

# 1<sup>st</sup> normal form (1NF)

- No multi-valued attributes allowed
  - Imagine storing a list/set of things in an attribute
  - => Not really even expressible in RA
- Counterexample
  - Course(name, instructor, [student,email]\*)
  - Redundancy in non-list attributes

Name	Instructor	Student Name	Student Email
eecs3421	Johnson	Xiao	xiao@gmail
		Jaspreet	jaspreet@utsc
		Mary	mary@utsc
eecs4411	Rosenburg	Jaspreet	jaspreet@utsc 50

## 2<sup>nd</sup> normal form (2NF)

- Non-prime attributes depend on candidate keys
  - Consider non-prime (ie. not part of a key) attribute 'a'
  - Then ∃FD X s.t. X -> a and X is a candidate key
- Counterexample
  - Movies(title, year, star, studio, studioAddress, salary)
  - FD: title, year -> studio; studio -> studioAddress; star->salary

Title	Year	Star	Studio	StudioAddr	Salary
Star Wars	1977	Hamill	Lucasfilm	1 Lucas Way	\$100,000
Star Wars	1977	Ford	Lucasfilm	1 Lucas Way	\$100,000
Star Wars	1977	Fisher	Lucasfilm	1 Lucas Way	\$100,000
Patriot Games	1992	Ford	Paramount	Cloud 9	\$2,000,000
Last Crusade	1989	Ford	Lucasfilm	1 Lucas Way	\$1,000,000

## 3<sup>rd</sup> normal form (3NF)

- Non-prime attr. depend *only* on candidate keys
  - Consider FD X -> a
  - Either  $a \in X \text{ OR } X$  is a superkey OR a is prime (part of a key)
  - => No transitive dependencies allowed
- Counterexample:
  - studio -> studioAddr

(studioAddr depends on studio which is not a candidate key)

Title	Year	Studio	StudioAddr
Star Wars	1977	Lucasfilm	1 Lucas Way
Patriot Games	1992	Paramount	Cloud 9
Last Crusade	1989	Lucasfilm	1 Lucas Way

## 3NF, dependencies, and join loss

- **Theorem**: always possible to convert a schema to lossless-join, dependency-preserving 3NF
- Caveat: always possible to create schemas in 3NF for which these properties do not hold
- FD loss example:
  - MovieInfo(title, year, studioName)
  - StudioAddress(title, year, studioAddress)
  - => Cannot enforce studioName -> studioAddress
- Join loss example:
  - Movies(title, year, star)
  - StarSalary(star, salary)
  - => Cannot enforce Movies ⋈ StarSalary yields bogus tuples (irreparable)

# **3NF Synthesis Algorithm**

**Objective**: Obtain a lossless and dependency-preserving decomposition of **R** 

1.Find a minimal cover  $\mathbf{F}_{\min}$  of  $\mathbf{F}$ 2.For each LHS(X):  $\forall X \in \mathbf{F}_{\min}$ , do:

- Create a relation schema using  $X_{F_{min}}$ +

3.Place any remaining attributes that have not been placed in any relations in step 2 in a single relation schema

4.If a key of **R** is not found in any relation, then add a trivial relation that consists of the key of **R** (if this trivial relation is useless, omit it)

# 3NF Synthesis Algorithm: Example

#### Question

Given:

- A Relation: **R**=(A, B, C, D, E, F, G, H)
- A Set of FDs F in R: F={A→CD, ACF→G, AD→BEF, BCG→D, CF→AH, CH→G, D→B, H→DEG}

Decompose R into a collection of relations  $R_i$  using the 3NF synthesis algorithm (which obtains a lossless and dependency-preserving decomposition of R)

# 3NF Synthesis Algorithm: Example

#### Answer

1. Find minimal basis of F:

 $\mathbf{F}_{min} = \{A \rightarrow C, A \rightarrow F, BCG \rightarrow D, CF \rightarrow A, CF \rightarrow H, D \rightarrow B, H \rightarrow D, H \rightarrow E, H \rightarrow G\}$ 

2. Create relation schemas based on  $X_{F_{min}}^{+}$ : *Closures*: { $A_{F_{min}}^{+} \rightarrow ACFH$ , BCG $_{F_{min}}^{+} \rightarrow BCGD$ , H $_{F_{min}}^{+} \rightarrow HDEG$ } *Relations*: **R1**(A, C, F, H), **R2**(B, C, G, D), **R3**(H, D, E, G)

3. No remaining attributes (of **R**), thus no need to place attributes in any of the available relations

4. A key of **R** was **A** which is already in relation **R1**, so no need to add a trivial relation that consists of the key of **R** 

#### Follow the detailed example online

## Boyce-Codd normal form (BCNF)

- One additional restriction over 3NF
  - All non-trivial FD have superkey LHS
- Counterexample
  - CanadianAddress(<u>street</u>, <u>city</u>, <u>province</u>, postalCode)
  - Candidate keys: {street, postalCode}, {street, city, province}
  - FD: postalCode -> city, province
  - Satisfies 3NF: city, province both non-prime
  - Violates BCNF: postalCode is not a superkey
  - => Possible anomalies involving postalCode

#### Do we care? How often do postal codes change?

## Limits of decomposition

- Pick two...
  - Lossless-join
  - Dependency-preservation
  - Anomaly-free
- 3NF
  - Always allows join lossless and dependency preserving
  - May allow some anomalies
- BCNF
  - Always excludes anomalies
  - May give up one of lossless-join or dependency-preserving

#### Use domain knowledge to choose 3NF vs. BCNF

### What is Next?

- Read Ullman & Widom's textbook (Chapter 3)
- Check detailed examples on Course's website
  - Sample 3NF Problem
  - Sample BCNF Problem
- Practice using online resources and examples