# CSC 261/461 – Database Systems Lecture 10 (part 2)

Spring 2018

# Announcement

- Read Chapter 14 and 15
- You must self-study these chapters

  Too huge to cover in Lectures
- Project 2 Part 1 due tonight

# Agenda

- 1. Database Design
- 2. Normal forms & functional dependencies
- 3. Finding functional dependencies
- 4. Closures, superkeys & keys

# **Design Theory**

- Design theory is about how to represent your data to avoid *anomalies*.
- Achieved by Data Normalization, a process of analyzing a relation to ensure that it is well formed.
- Normalization involves decomposing relations with anomalies to produce smaller well structured relations.
- If a relation is normalized (or well formed), rows can be inserted, deleted and modified without creating anomalies.

# **Normalization Example**

- (Student ID)  $\rightarrow$  (Student Name, DormName, DormCost)
- However, if
  - $-(\text{DormName}) \rightarrow (\text{DormCost})$

Then, DormCost should be put into its own relation, resulting in:

 $\begin{array}{l} (\text{Student ID}) \not \rightarrow (\text{Student Name, DormName}) \\ (\text{DormName}) \not \rightarrow (\text{ DormCost}) \end{array}$ 

# **Normalization Example**

- (AttorneyID, ClientID) → (ClientName, MeetingDate, Duration)
- However, if
   − ClientID → ClientName
- Then: ClientName should be in its own relation:
- (AttorneyID, ClientID)  $\rightarrow$  (MeetingDate, Duration)
- (ClientID)  $\rightarrow$  (ClientName)

# **Normal Forms**

- $\underline{I^{st} Normal Form (INF)} = All tables are flat$
- <u>2<sup>nd</sup> Normal Form</u> = disused
- <u>3rd Normal Form (3NF)</u>
- **Boyce-Codd Normal Form** (BCNF)

DB designs based on *functional dependencies,* intended to prevent data *anomalies* 

• <u>4<sup>th</sup> and 5<sup>th</sup> Normal Forms</u> = see text books

# **Normalization Steps**



# 1<sup>st</sup> Normal Form (1NF)

Student	Courses
Mary	{CS145,CS229}
Joe	{CS145,CS106}

Student	Courses
Mary	CS145
Mary	CS229
Joe	CS145
Joe	CS106

#### Violates 1NF.

In 1<sup>st</sup> NF

#### **1NF Constraint:** Types must be atomic!

### **Data Anomalies & Constraints**

A poorly designed database causes *anomalies*:

Student	Course	Room
Mary	CSC261	101
Joe	CSC261	101
Sam	CSC261	101
••		

If every course is in only one room, contains <u>redundant</u> information!

# A poorly designed database causes *anomalies*:

Student	Course	Room
Mary	CSC261	101
Joe	CSC261	703
Sam	CSC261	101
••		

If we update the room number for one tuple, we get inconsistent data = an <u>update</u> <u>anomaly</u>

A poorly designed database causes *anomalies*:

Student	Course		Room	
••	••		•	

If everyone drops the class, we lose what room the class is in! = a <u>delete anomaly</u>

			А	poorly desig <b>a</b> l	gned datal <b>nomalies</b> :	base cause	2S	
				Student	Course	Room		
				Mary	CSC261	B01		Similarly, we
				Joe	CSC261	B01		can't reserve a room without
				Sam	CSC261	B01		students = an
CSC4	16	703		••				insert anomal
1				L	1	1	I	

Student	Course
Mary	CSC261
Joe	CSC261
Sam	CSC261
••	

Course	Room
CSC261	101
CSC257	601

Is this form better?

- Redundancy?
- Update anomaly?
- Delete anomaly?
- Insert anomaly?

Today: develop theory to understand why this design may be better **and** how to find this *decomposition*...

# **Functional Dependencies**

# Functional Dependencies for Dummies

- A relationship between attributes where one attribute (or group of attributes) determines the value of another attribute (or group of attributes) in the same table.
- Example:

SSN uniquely identify any Person

 $(SSN) \rightarrow (First Name, Last Name)$ 

Candidate Keys/Primary Keys and Functional Dependencies

- By definition:
- A candidate key of a relation functionally determines all other non key attributes in the row.
- Implies:
- A primary key of a relation functionally determines all other non key attributes in the row.

EmployeeID  $\rightarrow$  (EmployeeName, EmpPhone)

#### **Functional Dependency**

**Def:** Let A,B be *sets* of attributes We write  $A \rightarrow B$  or say A *functionally determines* B if, for any tuples  $t_1$  and  $t_2$ :  $t_1[A] = t_2[A]$  implies  $t_1[B] = t_2[B]$ and we call  $A \rightarrow B$  a *functional dependency* 

A->B means that

"whenever two tuples agree on A then they agree on B."



<u>Defn (again):</u> Given attribute sets  $A=\{A_1,...,A_m\}$ and  $B = \{B_1,...,B_n\}$  in R,



<u>Defn (again):</u> Given attribute sets  $A=\{A_1,...,A_m\}$ and  $B = \{B_1,...,B_n\}$  in R,

The *functional dependency*  $A \rightarrow B$ on **R** holds if for *any*  $t_i, t_j$  in **R**:



If t1,t2 agree here..

Defn (again):

Given attribute sets  $A = \{A_1, ..., A_m\}$ and  $B = \{B_1, ..., B_n\}$  in R,

The *functional dependency*  $A \rightarrow B$ on **R** holds if for *any*  $t_i, t_j$  in R:

$$\begin{split} t_i[A_1] &= t_j[A_1] \text{ AND } t_i[A_2] = t_j[A_2] \text{ AND } \dots \\ \text{AND } t_i[A_m] &= t_j[A_m] \end{split}$$



Defn (again):

Given attribute sets  $A = \{A_1, ..., A_m\}$ and  $B = \{B_1, ..., B_n\}$  in R,

The *functional dependency*  $A \rightarrow B$ on **R** holds if for *any*  $t_i, t_j$  in R:

 $\underline{if} t_i[A_1] = t_j[A_1] \text{ AND } t_i[A_2] = t_j[A_2] \text{ AND}$  $... \text{ AND } t_i[A_m] = t_j[A_m]$ 

 $\underline{\textbf{then}} t_i[B_1] = t_j[B_1] \text{ AND } t_i[B_2] = t_j[B_2]$ AND ... AND  $t_i[B_n] = t_j[B_n]$ 

## FDs for Relational Schema Design

- High-level idea: why do we care about FDs?
  - 1. Start with some relational schema
  - 2. Find out its functional dependencies (FDs)
  - 3. Use these to design a better schema
    - One which minimizes the possibility of anomalies

# **Functional Dependencies as Constraints**

# A **functional dependency** is a form of **constraint**

- *Holds* on some instances not others.
- Part of the schema, helps define a valid *instance*.

*Recall: an <i>instance* of a schema is a multiset of tuples conforming to that schema, *i.e. a table* 

Student	Course	Room
Mary	CS145	B01
Joe	CS145	B01
Sam	CS145	B01
	••	••

Note: The FD {Course} -> {Room} *holds on this instance* 

# Functional Dependencies as Constraints

Note that:

- You can check if an FD is violated by examining a single instance;
- However, you cannot prove that an FD is part of the schema by examining a single instance.
  - This would require checking every valid instance

Student	Course	Room
Mary	CS145	B01
Joe	CS145	B01
Sam	CS145	B01
		••

However, cannot *prove* that the FD {Course} -> {Room} is *part of the schema* 

## More Examples

An FD is a constraint which <u>holds</u>, or <u>does not hold</u> on an instance:

EmpID	Name	Phone	Position
E0045	Smith	1234	Clerk
E3542	Mike	9876	Salesrep
E1111	Smith	9876	Salesrep
E9999	Mary	1234	Lawyer

# More Examples

EmpID	Name	Phone	Position
E0045	Smith	1234	Clerk
E3542	Mike	9876 ←	Salesrep
E1111	Smith	9876 ←	Salesrep
E9999	Mary	1234	Lawyer

{Position}  $\rightarrow$  {Phone}

# More Examples

EmpID	Name	Phone	Position
E0045	Smith	1234 →	Clerk
E3542	Mike	9876	Salesrep
E1111	Smith	9876	Salesrep
E9999	Mary	$1234 \rightarrow$	Lawyer

but *not* {Phone}  $\rightarrow$  {Position}

# ACTIVITY

А	В	С	D	E
1	2	4	3	6
3	2	5	1	8
1	4	4	5	7
1	2	4	3	6
3	2	5	1	8

Find at least *three* FDs which hold on this instance:

{	} → {	}	
{	} > {	}	
{	} > {	}	

# FINDING FUNCTIONAL DEPENDENCIES

## What you will learn about in this section

- 1. "Good" vs. "Bad" FDs: Intuition
- 2. Finding FDs
- 3. Closures

## "Good" vs. "Bad" FDs

#### We can start to develop a notion of **good** vs. **bad** FDs:

EmpID	Name	Phone	Position
E0045	Smith	1234	Clerk
E3542	Mike	9876	Salesrep
E1111	Smith	9876	Salesrep
E9999	Mary	1234	Lawyer

#### Intuitively:

EmpID -> Name, Phone, Position is "good FD"

 Minimal redundancy, less possibility of anomalies

## "Good" vs. "Bad" FDs

#### We can start to develop a notion of **good** vs. **bad** FDs:

EmpID	Name	Phone	Position
E0045	Smith	1234	Clerk
E3542	Mike	9876	Salesrep
E1111	Smith	9876	Salesrep
E9999	Mary	1234	Lawyer

#### Intuitively:

EmpID -> Name, Phone, Position is "good FD"

But Position -> Phone *is a* "bad FD"

 Redundancy!
 Possibility of data anomalies

### "Good" vs. "Bad" FDs

Student	Course	Room
Mary	CS145	B01
Joe	CS145	B01
Sam	CS145	B01
••	••	

Returning to our original example... can you see how the "bad FD" {Course} -> {Room} could lead to an:

- Update Anomaly
- Insert Anomaly
- Delete Anomaly
- ...

Given a set of FDs (from user) our goal is to:

- 1. Find all FDs, and
- 2. Eliminate the "Bad Ones".

## FDs for Relational Schema Design

- High-level idea: why do we care about FDs?
  - I. Start with some relational *schema*
  - 2. Find out its functional dependencies (FDs)
  - 3. Use these to **design a better schema** 
    - I. One which minimizes possibility of anomalies
- There can be a very large number of FDs...
   How to find them all efficiently?
- We can't necessarily show that any FD will hold **on all instances...** 
  - How to do this?

We will start with this problem: Given a set of FDs, F, what other FDs **must** hold?

Equivalent to asking: Given a set of FDs,  $F = \{f_{\scriptscriptstyle I}, \ldots f_{\scriptscriptstyle n}\},$  does an FD g hold?

**Inference problem**: How do we decide?

### Example:

#### Products

Name	Color	Category	Dep	Price
Gizmo	Green	Gadget	Toys	49
Widget	Black	Gadget	Toys	59
Gizmo	Green	Whatsit	Garden	99

Provided FDs:

- 1. {Name}  $\rightarrow$  {Color}
- 2. {Category}  $\rightarrow$  {Department}
- 3. {Color, Category}  $\rightarrow$  {Price}

Given the provided FDs, we can see that {Name, Category}  $\rightarrow$  {Price} must also hold on **any instance**...

Which / how many other FDs do?!?

Equivalent to asking: Given a set of FDs,  $F = \{f_{\scriptscriptstyle I}, \ldots f_{\scriptscriptstyle n}\},$  does an FD g hold?

**Inference problem**: How do we decide?

Answer: Three simple rules called Armstrong's Rules.

- 1. Split/Combine,
- 2. Reduction, and
- 3. Transitivity... ideas by picture

### 1. Split/Combine (Decomposition & Union Rule)



$$A_1, ..., A_m \rightarrow B_1, ..., B_n$$

### 1. Split/Combine (Decomposition & Union Rule)



$$A_1, ..., A_m \rightarrow B_1, ..., B_n$$

... is equivalent to the following *n* FDs...

$$A_1, \dots, A_m \rightarrow B_i$$
 for i=1,...,n

### 1. Split/Combine (Decomposition & Union Rule)



And vice-versa,  $A_1, ..., A_m \rightarrow B_i$  for i=1,...,n

... is equivalent to ...

$$A_1, ..., A_m \rightarrow B_1, ..., B_n$$

## 2. Reduction/Trivial (Reflexive Rule)



$$A_1,...,A_m \rightarrow A_j$$
 for any j=1,...,m

# 3. Transitive Rule



$$A_1, ..., A_m \rightarrow B_1, ..., B_n$$
 and  
 $B_1, ..., B_n \rightarrow C_1, ..., C_k$ 

## 3. Transitive Rule



$$A_1, ..., A_m \rightarrow B_1, ..., B_n$$
 and  
 $B_1, ..., B_n \rightarrow C_1, ..., C_k$ 

implies

$$A_1, \dots, A_m \rightarrow C_1, \dots, C_k$$

## **Augmentation Rule**



 $A_1, ..., A_m \rightarrow B_1, ..., B_n$  implies

# **Augmentation Rule**



$$A_1, ..., A_m \rightarrow B_1, ..., B_n$$
  
implies  
$$X_1, A_1, ..., A_m \rightarrow B_1, ..., B_n$$

#### Example:

#### Products

Name	Color	Category	Dep	Price
Gizmo	Green	Gadget	Toys	49
Widget	Black	Gadget	Toys	59
Gizmo	Green	Whatsit	Garden	99

Provided FDs:

- 1. {Name}  $\rightarrow$  {Color}
- 2. {Category}  $\rightarrow$  {Department}
- 3. {Color, Category}  $\rightarrow$  {Price}

Which / how many other FDs hold?

Provided FDs:

- 1. {Name}  $\rightarrow$  {Color}
- 2. {Category}  $\rightarrow$  {Dept.}
- 3. {Color, Category}  $\rightarrow$  {Price}

### Example:

### **Inferred FDs:**

Inferred FD		Rule used
4. {Name, Category} ->	{Name}	?
5. {Name, Category} ->	{Color}	?
6. {Name, Category} ->	{Category}	?
7. {Name, Category -> { Category}	Color,	?
8. {Name, Category} ->	{Price}	?

Which / how many other FDs hold?

### Example:

### **Inferred FDs:**

Inferred FD		Rule used	1.	{Name
4. {Name, Category} -> {Name}		Trivial	2.	{Catego
5. {Name, Category} -> {Color}		Transitive (4 -> 1)	{C	)ept.}
6. {Name, Category} -> {Category}		Trivial	3. {Color, {Price}	
7. {Name, Category -> {Color, Category}		Split/combine (5 + 6)		
8. {Name, Category}	-> {Price}	Transitive (7 -> 3)		
	Can we find ar	n algorithmic way to do this?	0	
	Yes. But we ne	ed to learn about clos that!	sures	before

### Provided FDs:

1. {Name} $\rightarrow$ {Color}
2. {Category} →
{Dept.}
3. {Color, Category} $\rightarrow$
{Price}

## Closures

## Closure of a set of Attributes

**Given** a set of attributes  $A_1, ..., A_n$  and a set of FDs F: Then the <u>closure</u>,  $\{A_1, ..., A_n\}^+$  is the set of attributes **B** s.t.  $\{A_1, ..., A_n\} \rightarrow B$ 

<u>Example:</u>	F =	<pre>{name} → {color} {category} → {department} {color, category} → {price}</pre>
Example Closures:		<pre>{name}+ = {name, color} {name, category}+ = {name, category, color, dept, price} {color}+ = {color}</pre>

Start with X = { $A_1$ , ...,  $A_n$ } and set of FDs F. **Repeat until** X doesn't change; **do**: **if** { $B_1$ , ...,  $B_n$ }  $\rightarrow$  C is in F **and** { $B_1$ , ...,  $B_n$ }  $\subseteq$  X **then** add C to X. **Return** X as X<sup>+</sup>

```
Start with X = \{A_1, ..., A_n\}, FDs F.
Repeat until X doesn't change;
do:
  if \{B_1, ..., B_n\} \rightarrow C is in F and \{B_1, ..., B_n\}
\ldots, B_n \subseteq X:
    then add C to X.
Return X as X<sup>+</sup>
\{name\} \rightarrow \{color\}
\{category\} \rightarrow \{dept\}
{color, category} \rightarrow
{price}
```

F =

{name, category}\* =
{name, category}

Start with  $X = \{A_1, ..., A_n\}$ , FDs F. **Repeat until** X doesn't change; do: if  $\{B_1, ..., B_n\} \rightarrow C$  is in F and  $\{B_1, ..., B_n\}$  $\ldots$ ,  $B_n$   $\subseteq X$ : then add C to X. Return X as X<sup>+</sup>  $\{name\} \rightarrow \{color\}$  $\{category\} \rightarrow \{dept\}$ {color, category}  $\rightarrow$ {price}

F =

{name, category}\* =
{name, category}

{name, category}\* =
{name, category, color}

Start with  $X = \{A_1, ..., A_n\}$ , FDs F. **Repeat until** X doesn't change; do: if  $\{B_1, ..., B_n\} \rightarrow C$  is in F and  $\{B_1, ..., B_n\}$  $\ldots$ ,  $B_n$   $\subseteq X$ : then add C to X. Return X as X<sup>+</sup>  $\{name\} \rightarrow \{color\}$  $\{category\} \rightarrow \{dept\}$ {color, category}  $\rightarrow$ {price}

F =

{name, category}\* =
{name, category}

{name, category}\* =
{name, category, color}

{name, category}\* =
{name, category, color, dept}

Start with  $X = \{A_1, ..., A_n\}$ , FDs F. **Repeat until** X doesn't change; do: if  $\{B_1, ..., B_n\} \rightarrow C$  is in F and  $\{B_1, ..., B_n\}$  $\ldots$ ,  $B_n$   $\subseteq X$ : then add C to X. Return X as X<sup>+</sup>  $\{name\} \rightarrow \{color\}$  $\{category\} \rightarrow \{dept\}$ {color, category}  $\rightarrow$ {price}

F =

{name, category}\* =
{name, category}

{name, category}\* =
{name, category, color}

{name, category}<sup>+</sup> =
{name, category, color, dept}

{name, category}\* =
{name, category, color, dept,
price}

# EXAMPLE

R(A,B,C,D,E,F)

$$\{A,B\} \rightarrow \{C\} \\ \{A,D\} \rightarrow \{E\} \\ \{B\} \rightarrow \{D\} \\ \{A,F\} \rightarrow \{B\}$$

}

}

Compute  $\{A,B\}^+ = \{A, B, B, A, B, A, B, B, A, B, A,$ 

Compute  $\{A, F\}^+ = \{A, F, F\}^+$ 

# EXAMPLE

R(A,B,C,D,E,F)

$$\{A,B\} \rightarrow \{C\} \\ \{A,D\} \rightarrow \{E\} \\ \{B\} \rightarrow \{D\} \\ \{A,F\} \rightarrow \{B\}$$

}

}

Compute {A,B}<sup>+</sup> = {A, B, C, D

Compute {A, F}<sup>+</sup> = {A, F, B

# EXAMPLE

R(A,B,C,D,E,F)

$$\{A,B\} \rightarrow \{C\} \\ \{A,D\} \rightarrow \{E\} \\ \{B\} \rightarrow \{D\} \\ \{A,F\} \rightarrow \{B\}$$

Compute {A,B}<sup>+</sup> = {A, B, C, D, E}

Compute  $\{A, F\}^+ = \{A, B, C, D, E, F\}$ 

# **3. CLOSURES, SUPERKEYS & KEYS**

## What you will learn about in this section

- 1. Closures
- 2. Superkeys & Keys

## Why Do We Need the Closure?

- With closure we can find all FD's easily
- To check if  $X \to A$ 
  - I. Compute X<sup>+</sup>
  - 2. Check if  $A \in X^+$

Note here that **X** is a *set* of attributes, but **A** is a *single* attribute. Why does considering FDs of this form suffice?

Recall the <u>Split/combine</u> rule:  $X \rightarrow A_1, ..., X \rightarrow A_n$  *implies*  $X \rightarrow \{A_1, ..., A_n\}$ 



We did not include {B,C}, {B,D}, {C,D}, {B,C,D} to save some space.

Step 1: Compute X<sup>+</sup>, for every set of attributes X:

Example: Given F =

 $\{A,B\} \rightarrow C$  $\{A,D\} \rightarrow B$  $\rightarrow D$ {B}

$${A}^{+} = {A}, {B}^{+} = {B,D}, {C}^{+} = {C}, {D}^{+} = {D}, {A,B}^{+} = {A,B,C,D}, {A,C}^{+} = {A,C}, {A,D}^{+} = {A,B,C,D}, {A,B,C}^{+} = {A,B,D}^{+} = {A,C,D}^{+} = {A,B,C,D}, {B,C,D}^{+} = {B,C,D}, {A,B,C,D}^{+} = {A,B,C,D}$$

Step 2: Enumerate all FDs X  $\rightarrow$  Y, s.t. Y  $\subseteq$  X<sup>+</sup> and X  $\cap$  Y =  $\emptyset$ :

 $\{A,B\} \rightarrow \{C,D\}, \{A,D\} \rightarrow \{B,C\}, \\ \{A,B,C\} \rightarrow \{D\}, \{A,B,D\} \rightarrow \{C\}, \\ \{A,C,D\} \rightarrow \{B\}$ 





Step 2: Enumerate all FDs X  $\rightarrow$  Y, s.t. Y  $\subset$  X<sup>+</sup> and X  $\cap$  Y =  $\varnothing$ :

 $\{A,B\} \rightarrow \{C,D\}, \{A,D\} \rightarrow \{B,C\},\$  $\{A,B,C\} \rightarrow \{D\}, \{A,B,D\} \rightarrow \{C\},\$  $\{A,C,D\} \rightarrow \{B\}$ 

The FD X  $\rightarrow$  Y is non-trivial

# Superkeys and Keys

## **Keys and Superkeys**

A <u>superkey</u> is a set of attributes  $A_1, ..., A_n$  s.t. for *any other* attribute **B** in R, we have  $\{A_1, ..., A_n\} \rightarrow B$ 

I.e. all attributes are functionally determined by a superkey

A **<u>key</u>** is a *minimal* superkey

Meaning that no subset of a key is also a superkey

# **Finding Keys and Superkeys**

- For each set of attributes X
  - 1. Compute  $X^+$
  - 2. If  $X^+$  = set of all attributes then X is a superkey
  - 3. If X is minimal, then it is a key

Do we need to check all sets of attributes?

# **Example of Finding Keys**

```
Product(name, price, category,
color)
```

```
{name, category} → price
{category} → color
```

What is a key?
## **Example of Keys**

Product(name, price, category, color)

{name, category} → price
{category} → color

{name, category}<sup>+</sup> = {name, price, category, color} = the set of all attributes ⇒ this is a superkey ⇒ this is a key, since neither name nor category alone is a superkey

## Acknowledgement

- Some of the slides in this presentation are taken from the slides provided by the authors.
- Many of these slides are taken from cs145 course offered by Stanford University.
- Thanks to YouTube, especially to <u>Dr. Daniel Soper</u> for his useful videos.