

# Functional dependencies

- Given a schema  $R(t)$  and  $X, Y \subseteq t$ ,
  - FD is a constraint on the instances of  $R$ .

$X \rightarrow Y$  ( $X$  functionally determines  $Y$ ) iff

$\forall r$  VALID instance of  $R$

$\forall t_1, t_2 \in r$

if  $t_1[X] = t_2[X]$  then  $t_1[Y] = t_2[Y]$

NB Using De Morgan we can say

$\exists t_1, t_2 \in r$  if  $t_1[X] = t_2[X]$  AND  $t_1[Y] \neq t_2[Y]$

$\Rightarrow \forall r$  NOT a valid instance of  $R$

? Not sure

NB Whenever you introduce a variable,  
always say if it  $\exists$  or  $\forall$

- What does it mean for a specific instance  $t$  to satisfy a FD?

Given an instance  $r_0$  of  $R$

it is said to satisfy the FD  $X \rightarrow Y$  ( $r_0 \models X \rightarrow Y$ )

if the property holds for  $r_0$

(in formula)  $r_0 \models X \rightarrow Y$  iff  $\forall t_1, t_2 \in r_0, t_1[X] = t_2[X] \Rightarrow t_1[Y] = t_2[Y]$

## Expressing FD

'Book Code  $\rightarrow$  Title' is an indication of a BAD schema

Consider Book Code  $\rightarrow$  Author

1. direct expr.  $\leadsto BC = \Rightarrow A =$

"if BC is the same, then A is also same"

2. contraposition  $\leadsto A \neq \Rightarrow B \neq$

"if A different then BC different"

3. absurd  $\leadsto$

"it is impossible to have two lines with the same BC in different A"

?  $\mid \leadsto \neg (BC = \wedge A \neq)$   
. $\mid \leadsto BC = \wedge A \neq \Rightarrow \text{false}$

$A \Rightarrow B$   
is the same as  
 $\neg A \Rightarrow \neg B$

$A \Rightarrow B$   
 $\neg A \vee B$   
 $\neg (A \wedge \neg B)$   
 $(A \wedge \neg B) \Rightarrow \text{false}$   
nothing can imply false  
so this is the way to  
express impossibility in logic

the RULE

True  $\wedge A \wedge B \Rightarrow C \vee D \vee \text{False}$

True  $\wedge A \wedge \neg D \Rightarrow C \vee \neg B \vee \text{False}$

DeMorgan

$\neg (A \wedge B) \equiv \neg A \vee \neg B$

$\neg (A \vee B) \equiv \neg A \wedge \neg B$

## Examples as FD

1) At any given time, a teacher is at most in a classroom

$$\leadsto \text{time} = 1 \text{ teacher} = \Rightarrow \text{room} =$$

(also) It is impossible to have the same teacher at the same time in different classrooms

$$\leadsto \text{time} = 1 \text{ teacher} = 1 \text{ room} \neq \Rightarrow \text{false}$$

(then we set all equals to write the FD)

2) It is not possible for two different teachers to be in the same classroom at the same time

$$\leadsto \text{teacher} \neq 1 \text{ room} = 1 \text{ time} = \Rightarrow \text{False}$$

$$\text{room} = 1 \text{ time} = \Rightarrow \text{teacher} =$$

3) If two lessons take place on different floors, they belong to two different degree courses

$$\leadsto \text{lesson} \neq 1 \text{ floor} \neq \Rightarrow \text{degree} \neq$$

$$\leadsto \text{lesson} \neq 1 \text{ floor} \neq 1 \text{ degree} = \Rightarrow \text{false}$$

$$\leadsto \text{degree} \Rightarrow \text{lesson} = \vee \text{floor} =$$

(correct solution is)

$$\text{degree} = \Rightarrow \text{floor} =$$

4) If two different lessons take place on the same day for the same subject, they belong to two different CDE

$$\leadsto \text{lesson} \neq 1 \text{ day} = 1 \text{ subject} = \Rightarrow \text{degree} \neq$$

(imp.)  $\leadsto \text{lesson} \neq 1 \text{ day} = 1 \text{ subject} = 1 \text{ degree} = \Rightarrow \text{false}$

$$(\text{FD}) \leadsto \text{day} = 1 \text{ subject} = 1 \text{ degree} = \Rightarrow \text{lesson} =$$

## FD Notation

- $R \langle T, F \rangle$   
denotes a schema with attributes  $T$  and f.d.  $F$
- a FD is COMPLETE (bad name for "left side minimal")  
when the FD does not hold if we remove attributes  
 $\leadsto X \rightarrow Y$  is complete  $\Leftrightarrow \forall W \subset X, X \rightarrow Y$  is NOT valid  
(example)  $\text{fiscal\_code}, \text{surname} \rightarrow \text{name}$  NOT complete because  
fc directly has name
- a set of attribute  $X$  is a SUPERKEY (superset of a key  $X \subseteq X$ )  
when  $X$  determines every other attribute of the relation  $X \rightarrow T$   
(example)  $\text{student\_id} \rightarrow T$   
 $\text{fiscal\_code} \rightarrow T$   
 $\text{student\_id}, \text{fiscal\_code} \rightarrow T$   
given  $T$  as 'every other attribute',  
all these set of attributes  $T$   
are superkeys
- an attribute is a KEY if it is a SUPERKEY and it is COMPLETE,  
meaning that  $X \rightarrow T$  and is left side minimal

## More examples on FD

1) When two rooms are in a different floor,  
then they have a different number of seats

(wrong)  $\leadsto$  room  $\neq$   $\wedge$  floor  $\neq \Rightarrow$  seats  $\leadsto$  room  $\neq$  is redundant because  
it is already in floor  $\neq$

(correct)  $\leadsto$  floor  $\neq \Rightarrow$  seats  $\neq$

(NEVER BE REDUNDANT)

(imp)  $\leadsto$  floor  $\neq \wedge$  seats  $= \Rightarrow$  False

2) It is not possible to have the same teacher  
teaching the same subject in two different Degree Courses

(imp)  $\leadsto$  teacher =  $\wedge$  subject =  $\wedge$  degree  $\neq \Rightarrow$  false

(FD)  $\leadsto$  teacher =  $\wedge$  subject  $\Rightarrow$  degree

## FD implications

There are FDs that implies other FDs.

So, fix a set of FDs, other FDs are implied by  $F$ .

$$(def) \quad r \models F \Rightarrow r \models X \rightarrow Y \quad \forall r \leftarrow \begin{array}{l} \text{instance of} \\ \text{relation } R \end{array}$$

$\uparrow$  set of FDs

$\rightarrow$  What does it mean that

$F$  logically implies  $X \rightarrow Y$  (written  $F \models X \rightarrow Y$ )?

It means that for every instance of relation  $r$  ( $\forall r$ ),  
if  $r$  satisfies  $F$  ( $r \models F$ ) then  $r \models X \rightarrow Y$

• But we cannot prove that because we would need to check  $\forall r$

### Example

Let  $r$  be an instance of  $R \langle T, F \rangle$ , with  $F = \{X \rightarrow Y, X \rightarrow Z\}$   
and  $X, Y, Z \subseteq T$ . Let  $X' \subseteq X$ .

Then are other FDs satisfied by  $R$ , such as

- (trivial FD)  $X \rightarrow X'$
- $X \rightarrow Y, Z \quad \leadsto$  because  $X \rightarrow Y$  and  $X \rightarrow Z$  so  $X$  determines both  $Y, Z$   
so  $\{X \rightarrow Y, X \rightarrow Z\} \models X \rightarrow YZ$
- $\{X \rightarrow Y, Y \rightarrow Z\} \models X \rightarrow Z$

## Inference rules (or derivation rules)

These are the building blocks of problems.

We use a set of rules that are correct and complete.

Armstrong's "axioms" (actually "rules")

1. if  $Y \subseteq X$ , then  $X \rightarrow Y$  (reflexivity R)
2. if  $X \rightarrow Y$ ,  $Z \subseteq T$ , then  $XZ \rightarrow YZ$  (augmentation A)
3. if  $X \rightarrow Y$ ,  $Y \rightarrow Z$ , then  $X \rightarrow Z$  (transitivity T)

→ whenever you have an implication, you can add the same set left and right

## Derivation (or mechanical deduction)

(def) Let  $F$  be a set of FD,

$X \rightarrow Y$  is derivable from  $F$  ( $F \vdash X \rightarrow Y$ ),

if  $X \rightarrow Y$  can be inferred from  $F$  using Armstrong's axioms

We say  $F \vdash X \rightarrow Y$  means that starting from  $F$  we can derive the FD  $X \rightarrow Y$ .

Applying these rules, we can make our set of dependencies  $F$  bigger not bigger.

## Example

$R(A, B, C, D)$

$F = \{A \rightarrow B, BC \rightarrow D\}$

AC is a superkey? In other words,  $AC \rightarrow ABCD$ ?

→ Remember, a superkey is a set of attributes  $X$  that determines every other attribute.

It is very difficult to reason, so we use the axioms (or deduction rules).

Let's start from  $A \rightarrow B$  that is the first FD

using augmentation  $AC \rightarrow BC$  because  $C \subseteq T$  where  $T$  is every other attribute

because then I want to exploit the second FD  $BC \rightarrow D$ .

Now, we augment also the second FD  $BC \rightarrow D$  with  $Aug(BC)$

that becomes  $BC \rightarrow BCD$  because  $BC$  set union  $BC$  is just  $BC$ .

We know that  $AC \rightarrow BC$  and  $BC \rightarrow BCD$ ,

so using transitivity we get  $AC \rightarrow BCD$ .

In the end, we augment with  $A$  and get  $AC \rightarrow ABCD$ .

We proved that  $AC$  is a superkey without any reasoning,  
just with inference rules.



## Correctness and Completeness of Armstrong's axioms

(Theorem) Armstrong's axioms are correct and complete

The fact that by using these I will discover true FDs is obvious.

1. Correctness (easy)

$$\forall F, f \quad F \vdash f \Rightarrow F \models f$$

$\uparrow$  derives using rules       $\uparrow$  implies

a set of FDs  $F$   
can be used to  
derive with  $\vdash$  / imply by reasoning  $\models$   
another set of FDs  $f$

2. Completeness (difficult and more important)

$$\forall F, f \quad F \models f \Rightarrow F \vdash f$$

It means that everything that can be proven  
can be discovered also by applying deduction rules only.

(3. Decidability?)

We have an algorithm that can answer  
in polynomial time,  
if  $f$  can or cannot be derived from  $F$

## Closure of a set of dependencies F

Closure is a set operator that has a property such that  $A \subseteq A^+ = (A^+)^+$ .

It means take a set, make it bigger until you arrive to some form of boundary where you cannot make it bigger anymore.

(def) Given a set F of FD,

the closure of F, denoted by  $F^+$ , is

$$F^+ = \{X \rightarrow Y \mid F \vdash X \rightarrow Y\} \quad \leadsto \text{all the dependencies } X \rightarrow Y \text{ that you can deduce by inference from } F$$

For example.

$A \rightarrow B$  in a relation  $R(A, B, C)$

the closure of  $A \rightarrow B$ , that is  $\{A \rightarrow B\}^+$

contains things like  $A \rightarrow A$ ,  $B \rightarrow B$ ,  $AB \rightarrow A$ ,  $ABC \rightarrow B$ ,  
 $AC \rightarrow B$ ,

We are writing everything that can be deduced from  $A \rightarrow B$  and  $R(A, B, C)$ .  
First, we write all the trivial FD, like

$$\begin{array}{llll} ABC \rightarrow AB & AC \rightarrow A & AB \rightarrow B & BC \rightarrow BC \\ ABC \rightarrow AC & AC \rightarrow C & AB \rightarrow A & BC \rightarrow B \quad \dots \\ ABC \rightarrow BC & AC \rightarrow AC & AB \rightarrow AB & BC \rightarrow C \end{array}$$

This is done using the rule of reflexivity,

$$\text{which is } Y \subseteq X \vdash X \rightarrow Y \quad \forall F$$

that is whenever  $X$  is a subset of  $Y$ ,  $X \rightarrow Y$  can be deduced from  $F$ .

Even starting from the empty set  $\{\}^+ = \{X \rightarrow Y \mid Y \subseteq X\}$

(continuation of closures of FDs)

Then, we have all the FDs that we get by augmentation and by Transitivity.

This is very interesting because we can say that

$F \vdash X \rightarrow Y \Leftrightarrow X \rightarrow Y \in F^+$   $\sim$   $F$  determines  $X \rightarrow Y$   
if and only if  
 $X \rightarrow Y$  belongs to  
the closure of  $F$

However, computing the closure of  $F$   
is expensive because grows exponentially.

So, it is useless in practice.

And we will use a new definition.

## Closure of a set of attributes X

(or  $X^+$  if  $F$  is clear from context)

(def) Given  $R \langle T, F \rangle$  and  $X \subseteq T$ ,  
the closure of  $X$  with respect to  $F$  (denoted by  $X_F^+$ ) is

$$X_F^+ = \{A_i \in T \mid F \vdash X \rightarrow A_i\}$$

$\leadsto$  an attribute  $A_i$  belongs to this set  $T$  if from  $X$  you can deduce that  $X$  determines  $A_i$

For example.

$$F = \{CF \rightarrow SID, SID \rightarrow N\}$$

$$(F_F^+ = \{SID, N\})$$

obvious  
 $CF \rightarrow SID$

this is because in the closure of  $CF$   
I can find also the name

So, the closure of an attribute are all those attributes that depend on this attribute directly or applying any of the Armstrong's Axioms.

## fundamental theorem

$$X \rightarrow Y \in F^+ \Leftrightarrow F \vdash X \rightarrow Y \Leftrightarrow Y \subseteq X_F^+$$

this was already stated in the previous pages

$\leadsto Y$  is included in the closure of  $X$

This means that if closing  $X$  I reach all attributes in  $Y$ , then  $X$  determines  $Y$ .

1.  $X \rightarrow Y \in F^+ \leadsto$  this is exponential
2.  $F \vdash X \rightarrow Y \leadsto$  we don't know how to solve it (all axioms randomly)  
(or try to apply)
3.  $Y \subseteq X_F^+ \leadsto$  this can be done in polynomial time

## Algorithm for computing the closures of $X$

Let's see an example.

$F = \{DB \rightarrow E, B \rightarrow C, A \rightarrow B\}$ . Compute  $(AD)^+$

$\leadsto X^+ = AD$  initialize the result variable with attributes  $A$  and  $D$ .

Then scan repeatedly the FDs in  $F$ .

Remember that  $X_F^+ = \{A_i \in T \mid F \vdash X \rightarrow A_i\}$

Assume that two lines are equal on  $A$  and  $D$ , look at FDs, what can you reach? Given  $A \rightarrow B$ , we can reach  $X^+ = ADB$ .

Now we restart the scan with  $ADB$  and

see what we can reach, that is  $X^+ = ADBE$ , given  $DB \rightarrow E$  by transitivity.

The same happens with  $B \rightarrow C$ , and we get  $X^+ = ADBEC$ .

We try once more to check that there's nothing left.

We observe that this algorithm is polynomial,

the max number of scans  $\leq$  number of attributes  $a$   
 $\leq$  number of dependencies  $p$

meaning that max scans  $\leq \min(a, p)$ .

The cost is given by  $p \cdot a \cdot \min(a, p) \approx n^3$ ,

so it is polynomial based on size of the input.

Certainly it is not exponential.

There is another algo, called 'fast closures' that is  $n^2$ , but we are not studying it.

## Slow Closure (pseudocode)

input  $R \langle T, F \rangle$ ,  $X \subseteq T$

output  $X^+$

algo  $\leadsto X^+ = X$

while  $X^+$  changes

for  $W \rightarrow V$  in  $F$  with  $W \subseteq X^+$  and  $V \notin X^+$

$X^+ = X^+ \cup V$

The question is:

$\rightarrow$  "is  $X \rightarrow Y$  derivable from  $F$  ( $F \vdash X \rightarrow Y$ )?"

just close  $X$  and see if it does reach all attributes in  $Y$  ( $Y \subseteq X_F^+$ )

$\rightarrow$  "is  $X$  a superkey?"

just compute the closure and see if you can reach all other attributes.

$\rightarrow$  "is  $X$  a key? (minimal superkey)"

check if there is a subset of  $X$  that is still a key,

eliminating one by one all the attributes in the set

and computing the closures.

Observe that this is still polynomial.

## Finding a Key

$F = \{DB \rightarrow E, B \rightarrow C, A \rightarrow B\}$  compute  $(AD)^+$

find a key.

→ start from a supkey.

like the set of all attributes, that is always a supkey because it trivially determines itself.

Let's start with ~~A~~BCDE

$(BCDE)^+ = BCDE$  because you cannot reach A from BCDE since A is not in the left side of any FD, so we need it

Try with A~~B~~CDE

$(ACDE)^+ = ACDEB$  we can reach B from  $A \rightarrow B$ , so it is a supkey

Now A~~B~~CDE, after removing B

$(ADE)^+ = ADEBC$  because of  $A \rightarrow B$  and  $B \rightarrow C$ , so supkey

Now A~~B~~CDE

$(AE)^+ = AEBC$  but we cannot reach D because it is never on the right side

And A~~B~~CDE

$(AD)^+ = ADBCE$  so AD is a key!

This algo is guaranteed to always find at least a key.

But in the worst case the complexity is  $n!$

(exponential complexity)

## Prime Attributes

An attribute is prime when it belongs to at least one key.

The problem of checking if an attribute is prime is NP-complete (expensive)

(def) Given a scheme  $R \langle T, F \rangle$ ,

we say that  $W \subseteq T$  is a candidate key of  $R$  if

- $W \rightarrow T \in F^+$   $\leadsto$   $W$  superkey
- $\forall V \subset W, V \rightarrow T \notin F^+$   $\leadsto$  if  $V \subset W$ ,  $V$  not a superkey



## Equivalence of FD. sets

It is possible to present the same set of dependencies in many similar ways.

$$F = \left\{ \begin{array}{l} FC \rightarrow SID \\ SID \rightarrow FC \\ FC \rightarrow N, S, A \end{array} \right\}$$

$$G = \left\{ \begin{array}{l} FC \rightarrow SID \\ SID \rightarrow FC \\ SID \rightarrow N, S, A \end{array} \right\}$$

They are equivalent!

Observe that  $F \vdash G \Leftrightarrow G \subseteq F^+$   
and  $G \vdash F \Leftrightarrow F \subseteq G^+$

(def) Two sets of FDs,  $F$  and  $G$ , are equivalent  $F \equiv G$  on the scheme  $R$ , iff  $F^+ = G^+$

So, from now on we focus on  $F^+$ , that is the reality (the complete truth) while  $F$  is just one way to tell the story.

## Choosing the best set of FDs (or cover)

And what if I want to make  $F$  as small as possible, going in the other direction with respect to closures.

$$\left\{ \begin{array}{l} FC, N \rightarrow S \\ FC \rightarrow N \end{array} \right\} \equiv \left\{ \begin{array}{l} FC \rightarrow N \\ FC \rightarrow S \end{array} \right\}$$

The smaller, the better. We prefer the one on the right.

When is a set of FDs minimal?

A cover is minimal when

1. there are no extraneous attributes that you may delete (like  $N$  before)
2. no redundant FDs that you may delete

An attribute  $A$  is extraneous when  $F \vdash (X - \{A\}) \rightarrow Y$ , meaning that it may be deleted and  $X$  still determines  $Y$

A FD  $X \rightarrow Y$  is redundant when  $F - \{X \rightarrow Y\} \vdash X \rightarrow Y$  meaning that it can be derived from  $F$

## Canonical Covers

A set of dependencies  $F$  is called canonical iff

1. the right part of FD in  $F$  is an attribute.

$$ABC \rightarrow DEF \rightsquigarrow \begin{array}{l} ABC \rightarrow D \\ ABC \rightarrow E \\ ABC \rightarrow F \end{array}$$

2. there are no extraneous attributes
3. no dependency in  $F$  is redundant

There is an algo to get the canonical form of  $F$

1. transform FD in the form  $X \rightarrow A$

2. delete extraneous attributes

$\rightsquigarrow$  for example, given  $ABC \rightarrow E$  how do I know if  $A$  is extraneous?  
meaning  $F \vdash ABC \rightarrow E$ ? compute  $BC_F^+$  and check if  $E \in BC_F^+$

3. eliminate redundant dependencies

$\rightsquigarrow$  take  $F - (AB \rightarrow C)$  and compute  $(AB)^+$   
is  $C \in AB^+$ ?

The canonical cover is not necessarily unique.

## Exercise

Consider  $R(AB \rightarrow D, AD \rightarrow B, A \rightarrow C, C \rightarrow A, BEC \rightarrow D, ABCDE)$

1) Are there any extraneous attributes in some dependencies?

$\leadsto$   $A$  is extraneous when  $F \vdash (X - \{A\}) \rightarrow Y$ ,  
 meaning that if we delete  $A$ , we still can get  $Y$ .  
 Compute  $(X - \{A\})^+$  and check if  $Y \in (X - \{A\})^+$ .  
 (Do it with all FDs?)

1.  $AB \leadsto B^+ = B$      $AB \leadsto A^+ = AC$  but  $D \notin A^+$
2.  $AD \leadsto D^+ = D$      $AD \leadsto A^+ = AC$  but  $B \notin A^+$
3.  $BEC \leadsto EC^+ = EC$      $BEC \leadsto BC^+ = BCA$  but  $D \notin BC^+$

No extraneous attributes

phelli said  
 $BC^+ = BCAD$   
 so  $D \in BC^+$   
 maybe because  
 $BEC \rightarrow D$

2) Are there any redundant dependencies?

$\leadsto$   $X \rightarrow Y$  is redundant when  $F - \{X \rightarrow Y\} \vdash X \rightarrow Y$ ,  
 meaning that it can be derived from  $F$

Take  $F - (AB \rightarrow D)$  and see if  $D \in AB^+$   $\leftarrow$  we compute  $X^+$  on  $F - (X \rightarrow Y)$

Let's try eliminating FDs one by one and see if I can still reach the same attributes

1.  $F - (AB \rightarrow D) \leadsto AB^+ = ABCD$  because we have  $BEC \rightarrow D$ , so it is REDUNDANT
2.  $F - (AD \rightarrow B) \leadsto AD^+ = ADC$  NOT reaching  $B$
3.  $F - (A \rightarrow C) \leadsto A^+ = A$  NOT reaching  $C$
4.  $F - (C \rightarrow A) \leadsto C^+ = C$
5.  $F - (BEC \rightarrow D) \leadsto BEC^+ = BECAD$   $\leftarrow$  we can reach everything, in particular  $D$ , so this is REDUNDANT

but if we remove  $AB \rightarrow D$   
 then we cannot remove this  
 because otherwise we cannot  
 reach  $D$  anymore

Consider  $R(AB \rightarrow D, AD \rightarrow B, A \rightarrow C, C \rightarrow A, BEC \rightarrow D, ABCDE)$

3. How many keys can you find?

A key is basically a left side minimal superkey.

Let's start finding the superkeys, meaning the attributes  $X$  that determine every other attribute.

We start from the set of all attributes that is the trivial superkey.

$(ABCDE)^+ = ABCDE$  because of reflexivity axiom. (if  $Y \subseteq X$ )

Now we try to eliminate attributes one by one.

~~A~~  $\leadsto (BCDE)^+ = BCDEA$  superkey

~~AB~~  $\leadsto (CDE)^+ = CDEAB$  superkey  $\leftarrow$  also a key!

~~ABC~~  $\leadsto (DE)^+ = DE$  not a superkey

~~ABD~~  $\leadsto (CE)^+ = CEA$  not a superkey

~~ABE~~  $\leadsto (CD)^+ = CDAB$  not a superkey, because E is never on the right side of any FD, so it must be in every key to be resolvable

~~A~~~~B~~~~C~~~~D~~~~E~~  $\leadsto (ACDE)^+ = ACDEB$  superkey

~~A~~~~B~~~~C~~~~D~~~~E~~  $\leadsto (ADE)^+ = ADECB$  superkey

~~A~~~~B~~~~C~~~~D~~~~E~~  $\leadsto (AE)^+ = AEC$  NOT a superkey

Observe that we can find the key both on the original or the simplified set because they are equivalent.

## Decomposition of schemas

- As we know, to resolve redundancy in schemas, we decompose into smaller equivalent schemas.

It is a decomposition because the set union of the new attributes is the set of original attributes

(def) Given a schema  $R(T)$ ,

$P = \{R_1(T_1), \dots, R_k(T_k)\}$  is a decomposition of  $R$

iff  $\bigcup T_i = T$

- How do we choose the right decomposition?

We do not want to lose information (lossy decomposition)

In a lossless decomposition, we can get the original table with a natural join.

Use projection to create smaller tables.

What if there are already existing procedures, queries?

We just create a view of the original table from new tables.

- Lossless join decompositions (very important definition)

(def)  $P = \{R_1(T_1), \dots, R_k(T_k)\}$

is a lossless join decomposition of  $R(T)$

iff for each VALID instance  $r$  of  $R$

$$r = (\pi_{T_1} r) \bowtie \dots \bowtie (\pi_{T_k} r)$$

$\rightarrow$  doing projections and join we get the original table

Observe that this must be valid also for all the values

that I will insert in the future.

We need a theorem.

It is easy to prove that  $r \subseteq (\pi_{T_1} r) \bowtie \dots \bowtie (\pi_{T_k} r)$

$\rightarrow$  the way we lose information is by getting a bigger table with spurious tuples

## Binary decompositions

theorem of  
Binary Decomposition

There is a nice theorem which states that  
whenever there is a binary decomposition  
the decomposition is lossless

iff the intersection of the two schemas is a superkey  
either for  $T_1$  and  $T_2$ , meaning

$$(T_1 \cap T_2) \rightarrow t_1 \in F^+ \text{ or } (T_1 \cap T_2) \rightarrow T_2 \in F^+$$

It is essential that the column being used for connecting two tables  
is a superkey at least for one table.

This is a necessary and sufficient condition.

If your attribute is not a superkey for at least one table,  
then the decomposition is lossy.

How to check? We compute the closures  
contains all the attributes of  $T_1$  or  $T_2$  then it's okay.

The next step is to find one decomposition, not only check.

## Projection of FDs

This is necessary to go beyond binary decomposition.

(Def) Given the schema  $R \langle T, F \rangle$  and  $T_1 \subseteq T$ ,  
the projection of  $F$  over  $T_1$  is

$$\pi_{T_1}(F) = \{X \rightarrow Y \in F^+ \mid X, Y \subseteq T_1\}$$

Basically we want to project original FDs to the table  
where they belong to after decomposition.

First, consider their closures (also called meanings).

Second, if  $X, Y$  belongs to  $T_1$ , then they are projected into  $T_1$ .

Let's see an example.

Let  $R(A, B, C)$  and  $F = \{A \rightarrow B, B \rightarrow C, C \rightarrow A\}$

$$R_1(AB) \rightsquigarrow \pi_{AB}(F) \equiv \{A \rightarrow B, B \rightarrow A\}$$

$$R_2(AC) \rightsquigarrow \pi_{AC}(F) \equiv \{A \rightarrow C, C \rightarrow A\}$$

$B \rightarrow A$  is from  
 $B \rightarrow C$  and  $C \rightarrow A$   
(closures)

The steps are

1. compute the closures

2. syntactical projection

3. reduce to canonical form

synthetic projection means  
eliminating every dependency  
which is about attributes  
that are not in the table

remove redundant dependencies,  
like the trivial ones

How to calculate  $\pi_{T_1}(F)$ ?

Calculate the determinants of  $T_1$ , meaning all the subsets of  $T_1$  ( $Y \subseteq T_1$ ).

Then for each det, compute the closures.

And in the end you remove redundancy.

$A \rightarrow \cancel{AB}$

$B \rightarrow \cancel{AB}$

$\cancel{AB \rightarrow AB}$



## Preserving Dependencies

(def) Given the schema  $R\langle T, F \rangle$ , the decomposition  $\rho = \{R_1, \dots, R_n\}$  preserves dependencies if the union of dependencies in  $\pi_i(F)$  is a cover of  $F$ .

for example.

$R(ABC, \{B \rightarrow C, C \rightarrow A\})$

$\swarrow \quad \searrow$   
 $R_1(AB) \quad R_2(BC)$   
 $\{B \rightarrow A\} \quad \{B \rightarrow C\}$

Is  $C \rightarrow A$  lost? Yes, because in the decomposition  $C^+ = C$

This can happen but it is not a problem.

It is a desirable property, but still not crucial.

(theorem) If there is an enemy decomposition that preserves dependencies and such that at least one  $T_j$  is a superkey, the  $\rho$  is a lossless join decomposition.

This is a generalization of the binary decomp. theorem.

It is sufficient but not necessary.

## Normal Forms

There are different types.

- 1NF  $\rightarrow$  each attribute has an elementary type (like number, date, string, ...)
  - $\hookrightarrow$  in NFNF (Not First Normal Form)  
also called object-relational system,  
where in a single cell you can put  
a list, array, ...
- 2NF, 3NF  $\leadsto$  impose restrictions on dependencies
- BCNF (Boyce-Codd Normal Form)  $\rightarrow$  it is the most natural and restrictive  
and the one we are studying

## BCNF

Intuition: if there exists in  $R$  a non-trivial  $X \rightarrow A$  dependency and  $X$  is not key,  
then  $X$  models the identity of an entity  
other than those modeled by the whole  $R$

In a table, given  $X \rightarrow A$ , if  $X$  is a superkey in that specific table,  
then there is no redundancy.

While, if  $X$  is not a superkey, then we have redundancy.

(Def)  $R \langle T, F \rangle$  is in BCNF  $\Leftrightarrow$  for every  $X \rightarrow A \in F^+$   
with  $A \notin X$  (meaning  $X$  is non-trivial)  
and  $X$  is a superkey

So  $X \rightarrow A$  must be trivial or, if not trivial,  $X$  must be a superkey.

It means that you cannot have two different lines with the same  $X$ .

Observe that we can express it in two alternative ways

1.  $A \in X$  or  $X$  superkey  $\leadsto$  True  $\Rightarrow A \in X$  or  $X$  superkey
2.  $A \notin X \Rightarrow X$  superkey

Since  $F^+$  is exponential, we cannot directly verify this property.

However, there is a theorem that we can use.

It states that if the BC property is true for every FD in  $F$ , it is also true for every FD in  $F^+$ .

Meaning that we can check in  $F$  and that's enough.

(Theorem)  $R \langle T, F \rangle$  is in BCNF  $\Leftrightarrow$  for every  $X \rightarrow A \in F$   
with  $A \notin X$  (non trivial),  
 $X$  is a superkey

It is the same as the definition, just with  $F$  instead of  $F^+$ .

Why the definition is with  $F^+$ ?

Because it is a property of  $F^+$ , depends on the meaning,  
not on a cover  $F$ .

But if we check a cover, we know that the property  
will be valid for all the covers.

Not clear...

If  $F$  is in canonical form, then each FD does not contain  
any extraneous attribute nor any redundant FD,  
so each superkey is also a key.

Meaning that in canonical cover we just check keys.

Observe that in the definition we have  $X \rightarrow A$   
that is an atomic FD because  $A$  is an attribute,  
while  $X$  is a set of attributes.

## Example

1. Teachers (Tax Code, Name, Dep, DepAddress)

→ We have a redundancy beca.  $dep \rightarrow dep, depAddress$   
meaning that each professor will have a DepAddress  
which is repeated for those of the same dep.  
This is not in BCNF.

2. Employees (Code, Qualification, ChildName)

→ Given  $code \rightarrow Qualification$ ,  
if an employee has 3 children it will appear 3 times.  
But  $(code)^t = qualification$ , so it is not a superkey,  
because we need also the children.  
This is not in BCNF

3. Libraries (BookCode, ShopName, ShopAddress, Title, Quantity)

→ Given  $BookCode \rightarrow Title$   
 $ShopName \rightarrow ShopAddress$   
 $BookCode, ShopName \rightarrow Quantity$  ← this is the only superkey  
But since the first two are FDs not trivial and not with superkeys,  
this cover will be redundant.

4. Telephones (AreaCode, Number, City, Subscriber, Street) ← not mobile phones

→ Given  $F = \{ AC, NU \rightarrow CI, SU, ST, CI \rightarrow AC \}$

But the city CI alone is not a superkey.  
This is not BCNF.

## The Analysis Algorithm (also the splitting algorithm)

Consider a schema where some FD violate the BCNF and display it.  
Let's see an example.

Teachers (TaxCode, Name, Dep, DepAddress)

$F = \{ \text{TaxCode} \rightarrow \text{Name}, \text{Dep}$   
 $\text{Dep} \rightarrow \text{DepAddress} \}$

The algo analyse each FD to see if it is or not in BCNF.

$(\text{TaxCode})^+ = \text{TaxCode}, \text{Name}, \text{Dep}, \text{DepAddress}$ ,  $\rightarrow$  this is a key

$(\text{Dep})^+ = \text{Dep}, \text{DepAddress}$   $\rightarrow$  Not a key

How to solve it? We create a table with the closure of the attribute that violates the BCNF, so  $(\text{Dep})^+$ , and another table with all the other attributes + X itself.

(algo) For every  $X \rightarrow Y$  that violates BCNF,  
 $R \langle T, F \rangle$  is recursively decomposed in  
 $R_1(X^+)$  and  $R_2(X, T - X^+)$

$\swarrow$  foreign  $K_{12}$

So  $R_1(\text{Dep}, \text{DepAddress})$  and  $R_2(\text{Dep}^+, \text{TaxCode}, \text{Name})$   
 $\downarrow$   $\downarrow$   
 $\text{Dep} \rightarrow \text{DepAddress}$   $\text{TaxCode} \rightarrow \text{Name}, \text{Dep}$

Now when we project the dependencies and Dep is a superkey of its table.  
This is a Binary Decomposition and it is guaranteed to be lossless because of the theorems we have seen.

The algorithm guarantees BCNF and lossless decomposition, but does not guarantee the FDs preservation.

## Steps of the algo

1. choose the FD that violates BCNF
2. split the schema
3. project the FD on the new tables
4. apply recursively

## Example.

Telephones (AreaCode, Number, City, Subscriber, Street)

$$F = \left\{ \begin{array}{l} AC, NU \rightarrow CI, SU, ST \\ CI \rightarrow AC \end{array} \right\} \rightarrow AC, NU \text{ is a superkey}$$
$$\sim \rightarrow CI \text{ is not a superkey}$$

We split based on CI.

$$R_1 (CI, AC) \quad R_2 (NU, CI^*, SU, ST)$$

We project the FD.

$$R_1(F) = \{CI \rightarrow AC\} \quad R_2(F) = \{\}$$

syntactic projection  
means take all the FD  
whose attributes are  
included in R2

Observe that if we calculate only the syntactic projection  $\rightarrow$

$R_2$  has no dependencies because they all contain AC  
which is not present in  $R_2$ .

But we should compute all the closures, included with 2, 3, ... terms

$$CI^+ = \dots, NU^+ = \dots, \dots$$

$$(CI, NU)^+ = \cancel{CI}, \cancel{NU}, \cancel{AC}, SU, ST \quad \text{and eliminating the trivial one,}$$

$$\text{we now can say } R_2(F) = \{CI, NU \rightarrow SU, ST\}$$

Observe that  $AC, NU \rightarrow CI$  has been lost, but it is not important.

## 3NF

It is weaker than BCNF.

(def)  $R\langle \Gamma, F \rangle$  is in 3NF if for every  $X \rightarrow A \in F^+$ ,  
with  $A \notin X$ ,  $X$  is a superkey or  $A$  is prime

(theo)  $\rightarrow$  same as (def) but with  $F$  instead of  $F^+$

Remember that  $A$  is prime when it belongs to at least one key.

Note that  $A$  is on the right side of the FD.

Since the BCNF might lose some FDs,

we might have to accept the 3NF with some redundancy  
to keep all the FDs.

To sum up, the 3NF admits a non-trivial and non-key dependency  
if the attributes on the right are prime;

while BCNF never admits any non-trivial and non-key dependence.

## The Synthesis Algorithm (to reach 3NF)

It groups attributes together, in the opposite direction with respect to the analysis.

It guarantees the lossless join property, the preservation of FDs, but not the BCNF.

However, it might happen that BCNF is reached.

Let's see an example.

$R(ABCDE, \{AD \rightarrow E, E \rightarrow B, E \rightarrow D, C \rightarrow E\})$

First, group all the dependencies that have the same left side

$G_1 = \{AD \rightarrow E, AD \rightarrow B\}$      $G_2 = \{E \rightarrow D\}$      $G_3 = \{C \rightarrow E\}$

Then for every group define a relation

$R_1(ADEB)$

$R_2(ED)$

$R_3(CE)$

Third, delete all relations that are included in others.

$R_1(ADEB)$

$R_3(CE)$

$\rightarrow$  we deleted  $R_2$  because  $R_2 \subseteq R_1$

Finally, check if one of these relations is a superkey.

This is because of the theorem that says

$\rightarrow$  whenever you have a decomposition which preserves any FDs if one of the relations of the decomposition is a superkey of the original table, then you are guaranteed to have a lossless join.

How can we be sure that no FDs are lost?

Because of the relation step.

Let's check the superkeys.  $R_1^+ = ADEB$      $R_3^+ = CED$     No superkey

We compute a key.

$ABCDE^+ = BCDE \rightarrow$  we need A, C

$ACDE^+ = \dots \rightarrow$  B not needed

$ACE^+ = ACEDB$

$AC^+ = ACEDB \rightarrow$  it is a key (not the only one)

So,  $R_1(ADEB), R_3(CE), R_4(AC)$



## Steps of the Synthesis's Algo

Let  $R \langle T, F \rangle$  with  $F$  canonical cover

1. partition into groups based on the determinant
2. define a relationship for each group with the same determinant
3. if a schema is contained into another, delete the smaller
4. check that at least one schema built is a superkey,  
otherwise find one arbitrary key and add that schema

Provided that you start with a set of FDs that are in canonical form,  
this algo is guaranteed to give a schema in 3NF.

Which algo is better?

(It depends on the goal you want to reach.

Which one is more computationally expensive?

Synthesis is polynomial, only the key is expensive (if needed).

Analysis is more difficult, because the projection  
of new tables is computationally expensive.

Deciding if a schema is in BCNF, it is polynomial,  
while deciding if it is in 3NF it is NP-hard.

## Example of Analysis Algo.

$R(ABCDE, \{AD \rightarrow E, AD \rightarrow B, E \rightarrow D, C \rightarrow E\})$

Scan all the FDs not check whether they violate the BC condition, by computing closures.

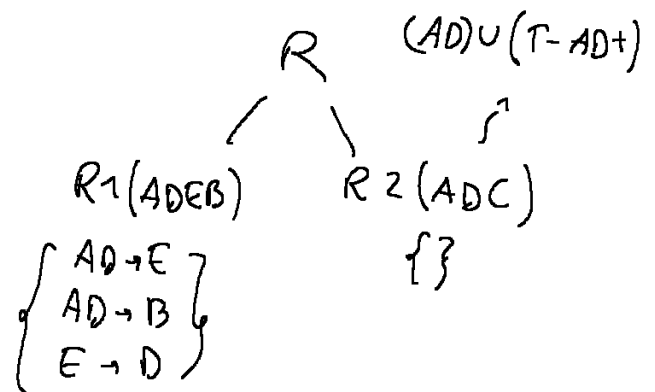
$AD^+ = ADEB$  not a superkey

We draw a tree because the operation is recursive and the result will be in the leaves of the tree.

Now we compute the syntactic projection and if no FD is lost, then because of a theorem, the syntactic projection coincide with the real projection.

If a FD is on one side, it cannot be on the other.

We observe from the tree that we have lost  $C \rightarrow E$ , so we need to compute the real projection.



with respect to  $R$   
✓

It is extremely expensive, compute the closure of every subset of  $R_1$

|                   |               |         |
|-------------------|---------------|---------|
| $A^+ = A$         | $AD^+ = ADEB$ | $ADE^+$ |
| $D^+ = D$         | $AE^+$        | $ADB^+$ |
| $E^+ = ED$        | $AB^+$        | $AEB^+$ |
| $B^+ = B$         | $DE^+$        | $DEB^+$ |
| $\vdots$          | $DB^+$        |         |
| <u>No new FDs</u> | $EB^+$        |         |

We basically look for the closure where we get a new FD because we used  $C \rightarrow E$ , that is the lost FD.

But we will never arrive anywhere because  $C$  is never on the right side. In this case, the two projection coincide with the syntactical projection.

$R(ABCDE, \{AD \rightarrow E, AD \rightarrow B, E \rightarrow D, C \rightarrow E\})$

Let's try now on the right side, with  $R2(ADC) \{\}$

$A^+ = A$

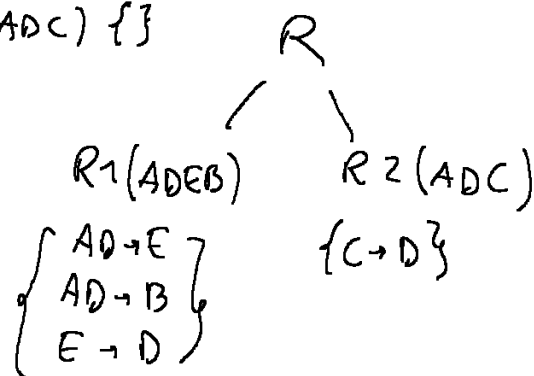
$AD^+ = ADDE$

$D^+ = D$

$AC^+ = ACD$

$C^+ = CED$

$DC^+ = DC$



this is interesting because, while

$C \rightarrow C$  is trivial and  $C \rightarrow E$  is useless because we do not have  $E$  in  $R2$ ,

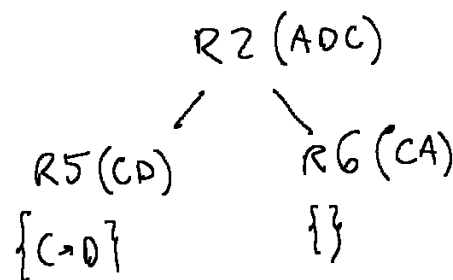
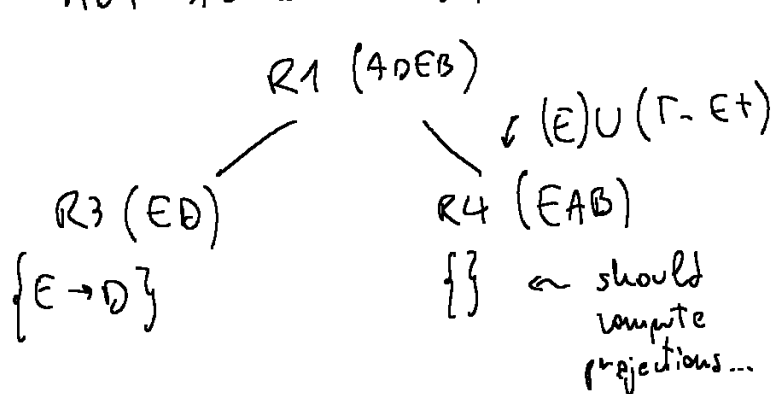
$C \rightarrow D$  can be projected in  $R2$ ,

meaning that we have learnt a new FD, since it was not in the syntactical projection.

TEST | Here you are authorized during test to skip the step of finding real projections, just mention it.

Now, is  $R1$  in BCNF? Let's compute the closures to check superkeys.

$AD^+ = ADEB$  superkey,  $E^+ = ED$  not superkey, then split



The leaves are in BCNF.

But we have lost a lot of FDs.