

DATA BASES DATA MINING

Fondements des Bases de Données : des Dépendances Fonctionnelles
aux Formes Normales

Équipe pédagogique BD



[http://liris.cnrs.fr/ecoquery/dokuwiki/doku.php?id=enseignement:
dbdm:start](http://liris.cnrs.fr/ecoquery/dokuwiki/doku.php?id=enseignement:dbdm:start)

Version du February 26, 2016

Example

Let $\mathcal{U} = \{id, name, address, cnum, desc, grade\}$ a set of attributes to model students and courses. We consider the following database schemas :

- ▶ $R1 = \{Data\}$ with $schema(Data) = \mathcal{U}^1$.
- ▶ $R2 = \{Student, Course, Enrollment\}$ avec
 - ▶ $schema(Student) = \{id, name, address\}$
 - ▶ $schema(Course) = \{cnum, desc\}$
 - ▶ $schema(Enrollment) = \{id, cnum, grade\}$

How to compare these schemas?

- ▶ Which one is the “best”?
- ▶ Why?

¹Similar to a spreadsheet.

Exemple

<i>Data</i>	<i>id</i>	<i>name</i>	<i>address</i>	<i>cnum</i>	<i>desc</i>	<i>grade</i>
	124	Jean	Paris	F234	Philo I	A
	456	Emma	Lyon	F234	Philo I	B
	789	Paul	Marseille	M321	Analyse I	C
	124	Jean	Paris	M321	Analyse I	A
	789	Paul	Marseille	CS24	BD I	B

Is there any problem here?

Exemple

<i>Data</i>	<i>id</i>	<i>name</i>	<i>address</i>	<i>cnum</i>	<i>desc</i>	<i>grade</i>
	124	Jean	Paris	F234	Philo I	A
	456	Emma	Lyon	F234	Philo I	B
	789	Paul	Marseille	M321	Analyse I	C
	124	Jean	Paris	M321	Analyse I	A
	789	Paul	Marseille	CS24	BD I	B

Is there any problem here?

Redundancies!

Redundancies

<i>Data</i>	<i>id</i>	<i>name</i>	<i>address</i>	<i>cnum</i>	<i>desc</i>	<i>grade</i>
	124	Jean	Paris	F234	Philo I	A
	456	Emma	Lyon	F234	Philo I	B
	789	Paul	Marseille	M321	Analyse I	C
	124	Jean	Paris	M321	Analyse I	A
	789	Paul	Marseille	CS24	BD I	B

Intuition on functional dependencies

- ▶ A student's *id* gives her/his name and address, so for each new enrollment, his/her name and address are duplicated!
- ▶ $\pi_{id,name,address}(Data)$ is the graph of a (partial) function $f : id \rightarrow name \times address$, similarly for $\pi_{cnum,desc}(Data)$
- ▶ $R2 = \{Student, Course, Enrollment\}$ is better than $R1 = \{Data\}$ because it avoids redundancies by keeping unrelated information (e.g., a student's name and a course's description) unrelated...

Functional is a theoretical tool to capture and reason on this phenomenon.

Functional Dependencies

Inference

Closure algorithm

Normalization

Functional dependencies : definition

Syntax

A *Functional Dependency (FD)* over a relation schema R is a formal expression of the form², with $X, Y \subseteq R$:

$$R : X \rightarrow Y$$

- ▶ $X \rightarrow Y$ is read “ X functionally determines Y ” or “ X gives Y ”
- ▶ A FD $X \rightarrow Y$ is **trivial** when $Y \subseteq X$
- ▶ A FD is **standard** when $X \neq \emptyset$.
- ▶ A set of attributes X is a **key** when $R : X \rightarrow R$

Semantics

Let r be a relation (a.k.a. *instance*) over R . The FD $R : X \rightarrow Y$ is *satisfied* by r , written $r \models R : X \rightarrow Y$, iff

$$\forall t_1, t_2 \in r. t_1[X] = t_2[X] \Rightarrow t_1[Y] = t_2[Y]$$

²We write $X \rightarrow Y$ when R is clear from the context.

What constraint is implied by a *non-standard* FD?
Why a *trivial* FD is said to be *trivial* ?

Example

<i>r</i>	A	B	C	D
<i>t</i> ₁	<i>a</i> ₁	<i>b</i> ₁	<i>c</i> ₁	<i>d</i> ₁
<i>t</i> ₂	<i>a</i> ₁	<i>b</i> ₁	<i>c</i> ₁	<i>d</i> ₂
<i>t</i> ₃	<i>a</i> ₁	<i>b</i> ₂	<i>c</i> ₂	<i>d</i> ₃
<i>t</i> ₄	<i>a</i> ₂	<i>b</i> ₂	<i>c</i> ₃	<i>d</i> ₄

- ▶ $r \models AB \rightarrow C$ (no counter-example)
- ▶ $r \models D \rightarrow ABCD$ (no counter-example)
- ▶ $r \not\models AB \rightarrow D$ (e.g., $t_1[AB] = t_2[AB]$ but $t_1[D] \neq t_2[D]$)
- ▶ $r \not\models A \rightarrow C$ (e.g., $t_2[A] = t_3[A]$ but $t_2[C] \neq t_3[C]$)

Functional Dependencies

Inference

Closure algorithm

Normalization

Logical implication

Definition

Let F be a set of FDs on a relation schema R and let f be a single FD on R . We overload \models for a **set** of FDs:

$$r \models F \text{ iff } \forall f \in F. r \models f$$

F **logical (semantically) implies** f , written

$$F \models f \text{ iff } \forall r. r \models F \Rightarrow r \models f$$

Example

With $F = \{A \rightarrow BCD, BC \rightarrow E\}$ and $r \models F$, the following hold as well:

- ▶ $r \models A \rightarrow CD$
- ▶ $r \models A \rightarrow E$

It can be proved using the definition of \models and basic reasoning on projection of tuples.

Armstrong's System for FD

Armstrong's System

The following rules constitute the so call *Armstrong's system* for FDs:

▶ Reflexivity

$$\frac{Y \subseteq X}{X \rightarrow Y}$$

▶ Augmentation

$$\frac{X \rightarrow Y}{WX \rightarrow WY}$$

▶ Transitivity

$$\frac{X \rightarrow Y \quad Y \rightarrow Z}{X \rightarrow Z}$$

Proof using Armstrong's system

Example

Let $\Sigma = \{A \rightarrow B, B \rightarrow C, CD \rightarrow E\}$ be a set of FDs on $\{A, B, C, D, E\}$.
We show that $\Sigma \vdash AD \rightarrow E$

$$\frac{\frac{\frac{A \rightarrow B \quad B \rightarrow C}{A \rightarrow C}}{AD \rightarrow CD} \quad CD \rightarrow E}{AD \rightarrow E}$$

Properties

Soundness and completeness

- ▶ The system is **sound** if $F \vdash f \Rightarrow F \models f$
if there is a proof, the proof is valid
- ▶ The system is **complete** if $F \models f \Rightarrow F \vdash f$
if it's valid, there is a proof

$$F \models \alpha \Leftrightarrow F \vdash \alpha$$

Soundness

Prove for every rule that, if its hypothesis are valid then its conclusion is valid as well.

Example : la transitivity

Let r be an instance on R s.t. $r \models X \rightarrow Y$ et $r \models Y \rightarrow Z$. Let $t_1, t_2 \in r$ be two tuples in r s.t. $t_1[X] = t_2[X]$, we have to show that $t_1[Z] = t_2[Z]$. Using $r \models X \rightarrow Y$ we deduce that $t_1[Y] = t_2[Y]$, then using $r \models Y \rightarrow Z$ we deduce that $t_1[Z] = t_2[Z]$. So the transitivity of FDs amounts to the transitivity of equality...

Additional rules

- ▶ Decomposition

$$\frac{X \rightarrow YZ}{X \rightarrow Y}$$

- ▶ Composition

$$\frac{X \rightarrow Y \quad X \rightarrow Z}{X \rightarrow YZ}$$

- ▶ Pseudo-transitivity

$$\frac{X \rightarrow Y \quad WY \rightarrow Z}{WX \rightarrow Z}$$

This rules are sound and can be (safely) added to Armstrong's system

Completeness

Preuve formelle

A (formal) proof of f from Σ using Armstrong' system written $\Sigma \vdash f$ is a *sequence* $\langle f_0, \dots, f_n \rangle$ of FDs s.t. $f_n = f$ et $\forall i \in [0..n]$:

- ▶ either $f_i \in \Sigma$;
- ▶ or f_i is the *conclusion* of a rule of which all its *antecedents* $f_0 \dots f_p$ appear before f_i in the sequence.

Completeness: $\Sigma \models X \rightarrow Y \Rightarrow \Sigma \vdash X \rightarrow Y$

We need a clear distinction between

- ▶ the **semantic** closure of X : $X^+ = \{A \mid \Sigma \models X \rightarrow A\}$
- ▶ the **syntactic** closure of X : $X^* = \{A \mid \Sigma \vdash X \rightarrow A\}$

Lemma: $\Sigma \vdash X \rightarrow Y \Leftrightarrow Y \subseteq X^*$

Completeness

$$\Sigma \models X \rightarrow Y \Rightarrow \Sigma \vdash X \rightarrow Y$$

$$\equiv \Sigma \not\models X \rightarrow Y \Rightarrow \Sigma \not\vdash X \rightarrow Y$$

$$\equiv \Sigma \not\models X \rightarrow Y \Rightarrow \exists r. (r \models \Sigma \wedge r \not\models X \rightarrow Y)$$

The crux is to find an instance r ,
with $X^* = X_1 \dots X_n$ et $Z_1 \dots Z_p = R \setminus X^*$

r	X_1	\dots	X_n	Z_1	\dots	Z_p
s	x_1	\dots	x_n	z_1	\dots	z_p
t	x_1	\dots	x_n	y_1	\dots	y_p

$$r \models \Sigma \text{ but } r \not\models X \rightarrow Y$$

Functional Dependencies

Inference

Closure algorithm

Normalization

Inference problem for FDs

Armstrong's system leads to a (inefficient) decision procedure for the *inference problem*.

Inference problem for FDs

Let F be a set of FDs and f a single FD, does $F \models f$ hold true?

Lemma: $F \models X \rightarrow Y$ iff $Y \subseteq X^+$

Thus, if we have an (efficient) algorithm to compute X^+ , we can (efficiently) solve the inference problem:

1. Given Σ and $X \rightarrow Y$, compute X^+ w.r.t. Σ
2. Return $Y \subseteq X^+$

Closure algorithm: $Closure(\Sigma, X)$

Data: Σ a set of FDs, X a set of d'attributes.

Result: X^+ , the closure of X w.r.t. Σ

```
1  $Cl := X$ 
2  $done := false$ 
3 while ( $\neg done$ ) do
4    $done := true$ 
5   forall the  $W \rightarrow Z \in \Sigma$  do
6     if  $W \subseteq Cl \wedge Z \not\subseteq Cl$  then
7        $Cl := Cl \cup Z$ 
8        $done := false$ 
9 return  $Cl$ 
```

How many times³ do we compute $W \subseteq CI \wedge Z \not\subseteq CI$ w.r.t. $|\Sigma| = n$?

³at worst, using a bad strategy at line 5.

Second algorithm

Data: Σ a set of FDs, X a set of d'attributes.

Result: X^+ , the closure of X w.r.t. Σ

```
10 for  $W \rightarrow Z \in F$  do
11    $count[W \rightarrow Z] := |W|$ 
12   for  $A \in W$  do
13      $list[A] := list[A] \cup W \rightarrow Z$ 
14  $closure := X, update := X$ 
15 while  $update \neq \emptyset$  do
16   Choose  $A \in update$ 
17    $update := update \setminus \{A\}$ 
18   for  $W \rightarrow Z \in list[A]$  do
19      $count[W \rightarrow Z] := count[W \rightarrow Z] - 1$ 
20     if  $count[W \rightarrow Z] = 0$  then
21        $update := update \cup (Z \setminus closure)$ 
22        $closure := closure \cup Z$ 
23 return  $closure$ 
```

Example : AE^+

$$\Sigma = \{A \rightarrow I; AB \rightarrow E; BI \rightarrow E; CD \rightarrow I; E \rightarrow C\}$$

Initialization

$List[A] = \{A \rightarrow D; AB \rightarrow E\}$	$count[A \rightarrow D] = 1$
$List[B] = \{AB \rightarrow E; BI \rightarrow E\}$	$count[AB \rightarrow E] = 2$
$List[C] = \{CD \rightarrow I\}$	$count[BI \rightarrow E] = 2$
$List[D] = \{CD \rightarrow I\}$	$count[CD \rightarrow I] = 2$
$List[E] = \{E \rightarrow C\}$	$count[E \rightarrow C] = 1$
$List[I] = \{BI \rightarrow E\}$	

Cover

Cover of a set of FDs

With $F^+ = \{f \mid F \models f\}$, let Σ et Γ be two sets of FDs,
 Γ is a cover of Σ iff $\Gamma^+ = \Sigma^+$

Data: F a set of FDs

Result: G a *minimal* cover of F

```
24  $G := \emptyset$ 
25 for  $X \rightarrow Y \in F$  do
26    $G := G \cup \{X \rightarrow X^+\}$ 
27 for  $X \rightarrow X^+ \in G$  do
28   if  $G - \{X \rightarrow X^+\} \vdash X \rightarrow X^+$  then
29      $G := G - \{X \rightarrow X^+\}$ 
30 return  $G$ 
```


Functional Dependencies

Inference

Closure algorithm

Normalization

Application of FD: Normalization

We write $\langle R, \Sigma \rangle$ with R a relation schema and Σ a set of FDs on R .
A set of attribute X is a *minimal key* of $\langle R, \Sigma \rangle$ iff:

- ▶ X is a key of R (i.e., $X \rightarrow R$ holds)
- ▶ X is *minimal* w.r.t. set inclusion: $\forall X' \subsetneq X \Rightarrow X' \not\rightarrow R$

Third Normal Form (3NF)

$\langle R, \Sigma \rangle$ is in **3NF** iff, for all *non-trivial* FD $X \rightarrow A$ of Σ^+ , one of the following conditions holds:

- ▶ X is a key of R
- ▶ A is a member of *at least* one minimal key of R^4

Boyce-Codd Normal Form (BCNF)

$\langle R, \Sigma \rangle$ is in **BCNF** iff, for all *non-trivial* $X \rightarrow A$ of Σ^+ , X is a key of R .

Informally, $\langle R, \Sigma \rangle$ is good when Σ is nothing but the key!

⁴An attribute that appears in *at least* one minimal key is said to be a *prime attribute*.

Example

3NF captures most of redundancies

- ▶ $\langle ABC, \{A \rightarrow B, B \rightarrow C\} \rangle$ is *not* in 3NF
A is the unique *minimal* key. Considering $B \rightarrow C$, C is *not* prime and B is *not* a key. Clearly, ABC should be divided into AB and BC
- ▶ $\langle ABC, \{AB \rightarrow C, C \rightarrow B\} \rangle$ is in 3NF
There are two *minimal* keys: AB and AC. Every attribute is prime so the 3NF condition holds. Unfortunately, some redundancies still hold but there is no way to decompose ABC into smaller relation without loss of FD!

BCNF captures all redundancies (expressed by FD)

- ▶ $\langle ABC, \{AB \rightarrow C, C \rightarrow B\} \rangle$ is *not* in BCNF
Considering $C \rightarrow B$, C alone is not a key.

Example

Back to the introductory example

With $\mathcal{U} = \{id, name, address, cnum, desc, grade\}$:

- ▶ the *natura*⁵ FDs are $f_1 = id \rightarrow name, address$, $f_2 = cnum \rightarrow desc$ and $f_3 = id, cnum \rightarrow grade$.
- ▶ The minimal key of \mathcal{U} is $\{id, cnum\}$.
- ▶ Without refinement, \mathcal{U} is not in 3NF, e.g., f_1 holds but id is not a key.
- ▶ The decomposition of \mathcal{U} into
 - ▶ $\langle \{id, name, address\}, \{f_1\} \rangle$
 - ▶ $\langle \{cnum, desc\}, \{f_2\} \rangle$
 - ▶ $\langle \{id, cnum, grade\}, \{f_3\} \rangle$

is *good* because the BCNF condition holds for each relation.

⁵Those which hold from the user's perspective, or alternatively, those that are true in the existing dataset.

End.