



**Hasso  
Plattner  
Institut**

IT Systems Engineering | Universität Potsdam

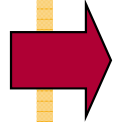
## Detecting Functional Dependencies

21.5.2013

Felix Naumann

# Overview

2



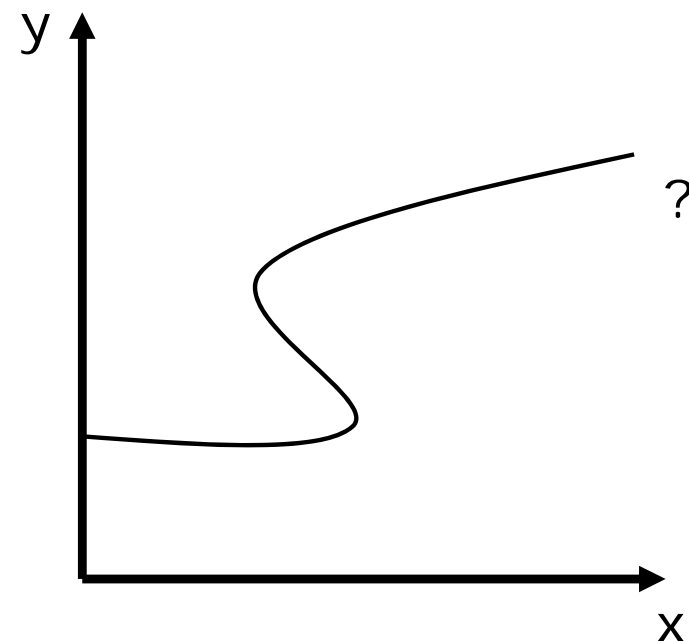
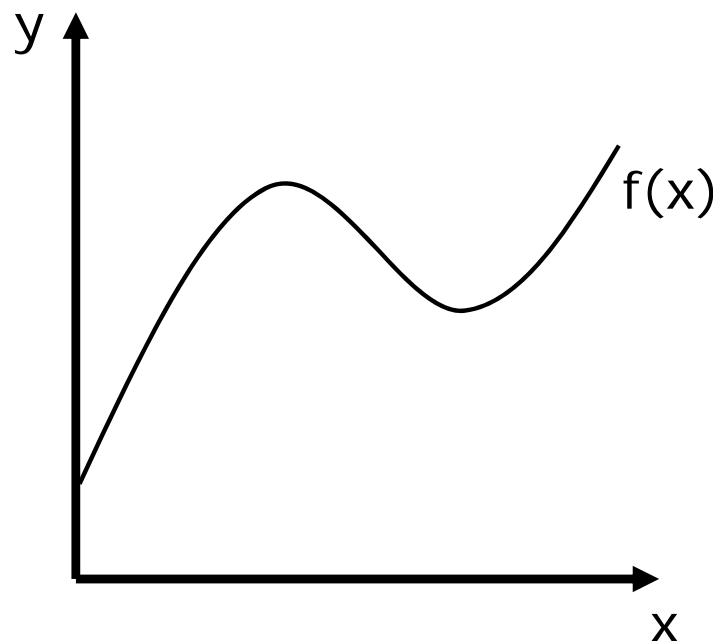
- Functional Dependencies
- TANE
  - Candidate sets
  - Pruning Algorithm
  - Dependency checking
  - Approximate FDs
- FD\_Mine
- Conditional FDs



# Definition – Functional Dependency

3

- „ $X \rightarrow A$ “ is a statement about a relation R: When two tuples have same value in attribute set X, they must have same values in attribute A.
- Formally:  $X \rightarrow A$  is an FD over R  $(R \models X \rightarrow A) \Leftrightarrow$  for all tuples  $t_1, t_2 \in R: t_1[X] = t_2[X] \Rightarrow t_1[A] = t_2[A]$
- Can generalize to sets:  $X \rightarrow Y \Leftrightarrow X \rightarrow A$  for each  $A \in Y$



# Trivial FDs

4

- Trivial: Attributes on RHS are subset of attributes on LHS
  - Street, City  $\rightarrow$  City
  - Any trivial FD holds
  
- Non-trivial: At least one attribute on RHS does not appear on LHS
  - Street, City  $\rightarrow$  Zip, City
- Completely non-trivial: Attributes on LHS and RHS are disjoint.
  - Street, City  $\rightarrow$  Zip
- Minimal FD: RHS does not depend on any subset of LHS.
  
- Typical goal: Given a relation  $R$ , find all minimal completely non-trivial functional dependencies.

# FD Inference Rules

5

- R1 Reflexivity  $X \supseteq Y \Rightarrow X \rightarrow Y$  (also  $X \rightarrow X$ )
  - Trivial FDs
- R2 Accumulation  $\{X \rightarrow Y\} \Rightarrow XZ \rightarrow YZ$ 
  - Aka: Augmentation
- R3 Transitivity  $\{X \rightarrow Y, Y \rightarrow Z\} \Rightarrow X \rightarrow Z$
  
- R1-R3 known as *Armstrong-Axioms*
  - Sound and complete
  
- R4 Decomposition  $\{X \rightarrow YZ\} \Rightarrow X \rightarrow Y$
- R5 Union  $\{X \rightarrow Y, X \rightarrow Z\} \Rightarrow X \rightarrow YZ$
- R6 Pseudotransitivity  $\{X \rightarrow Y, WY \rightarrow Z\} \Rightarrow WX \rightarrow Z$

# FD Discussion

6

- Schema vs. instance
- Keys as special case for FDs
  - $X$  is key of  $R$  if  $X \rightarrow R \setminus X$
- Uses for FDs
  - Schema design and normalization
  - Key discovery
  - Data cleansing (especially conditional FDs)

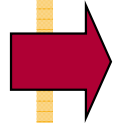
# Naive Discovery Approach

7

- Given relation  $R$ , detect all minimal, non-trivial FDs  $X \rightarrow A$ .
  
- For each column combination  $X$ 
  - For each pair of tuples  $(t_1, t_2)$ 
    - ◇ If  $t_1[X \setminus A] = t_2[X \setminus A]$  and  $t_1[A] \neq t_2[A]$ : Break
  
- Complexity
  - Exponential in number of attributes
  - times number of rows squared

# Overview

8



- Functional Dependencies
- TANE
  - Candidate sets
  - Pruning Algorithm
  - Dependency checking
  - Approximate FDs
- FD\_Mine
- Conditional FDs





# Tane – General Idea

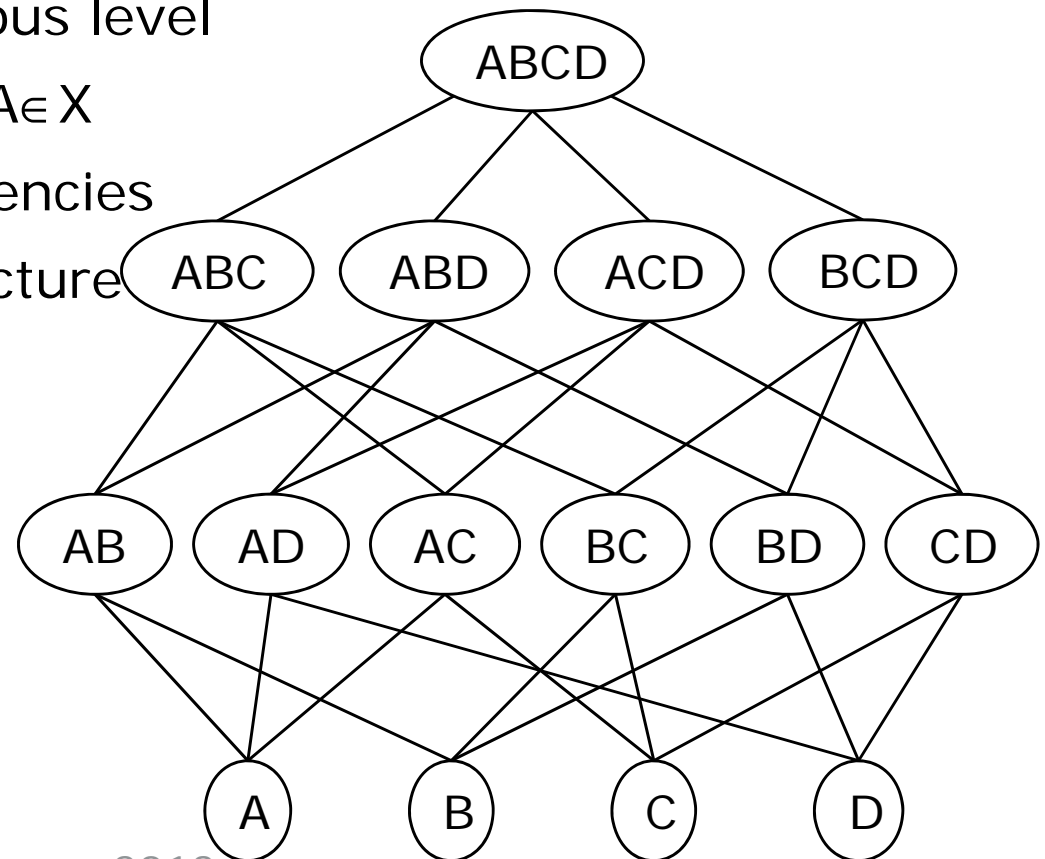
9

- Two elements of approach
  1. Reduce column combinations through pruning
    - ◇ Reasoning over FDs
  2. Reduce tuple sets through partitioning
    - ◇ Partition data according to attribute values
    - ◇ Level-wise increase of size of attribute set
      - Consider sets of tuples whose values agree on that set
  
- Huhtala, Y.; Kärkkäinen, J.; Porkka, P. & Toivonen, H.  
*TANE: An Efficient Algorithm for Discovering Functional and Approximate Dependencies*  
*Computer Journal*, **1999**, 42, 100-111

# Discovery strategy

10

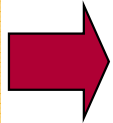
- Bottom up traversal through lattice
  - $\Rightarrow$  only minimal dependencies
  - Pruning
  - Re-use results from previous level
- For a set  $X$ , test all  $X \setminus A \rightarrow A, A \in X$ 
  - $\Rightarrow$  only non-trivial dependencies
  - Test on efficient data structure



# Overview

11

- Functional Dependencies
- TANE
  - Candidate sets
  - Pruning Algorithm
  - Dependency checking
  - Approximate FDs
- FD\_Mine
- Conditional FDs



# Candidate Sets

12

- RHS candidate set  $C(X)$
- Stores only those attributes that might depend on **all** other attributes of  $X$ .
  - I.e., those that still need to be checked
  - If  $A \in C(X)$  then  $A$  does not depend on any proper subset of  $X$ .
- $C(X) = R \setminus \{A \in X \mid X \setminus A \rightarrow A \text{ holds}\}$
- Example:  $R = \{ABCD\}$ , and  $A \rightarrow C$  and  $CD \rightarrow B$ 
  - $C(A) = \{ABCD\} \setminus \{A\} = C(B) = C(C) = C(D)$
  - $C(AB) = \{ABCD\} \setminus \{A\}$
  - $C(AC) = \{ABCD\} \setminus \{C\} = \{ABD\}$
  - $C(CD) = \{ABCD\} \setminus \{C\}$
  - $C(BCD) = \{ABCD\} \setminus \{B\} = \{ACD\}$

# RHS candidate pruning

13

- For minimality it suffices to test  $X \setminus A \rightarrow A$  where
  - $A \in X$  and  $A \in C(X \setminus \{B\})$  for all  $B \in X$ .
  - I.e.,  $A$  is in **all** candidate sets of the subsets.
  
- Example
  - $X = \{ABC\}$ . Assume we know  $C \rightarrow A$  from previous step.
  - Need to test three dependencies:  $AB \rightarrow C$ ,  $AC \rightarrow B$ , and  $BC \rightarrow A$ 
    - ◇ We should not be testing  $BC \rightarrow A$ , because we know  $C \rightarrow A$
  - Candidate sets:
    - ◇  $C(AB) = \{ABC\}$ ,  $C(AC) = \{BC\}$ ,  $C(BC) = \{ABC\}$
  - E.g.  $BC \rightarrow A$  does not need to be tested for minimality, because  $A$  is not in all three candidate sets:  $A \notin C(AB) \cap C(AC) \cap C(BC)$
  - $AB \rightarrow C$ ,  $AC \rightarrow B$  need to be tested, because  $B$  and  $C$  appear in all candidate sets.

# Improved RHS candidate pruning

15

- Basis: Let  $B \in X$  and let  $X \setminus B \rightarrow B$  hold. If  $X \rightarrow A$ , then  $X \setminus B \rightarrow A$ .
  - Example:  $A \rightarrow B$  holds. If  $AB \rightarrow C$  holds, then also  $A \rightarrow C$ .
  - Use this to reduce candidate set: If  $X \setminus B \rightarrow B$  for some  $B$ , then any dependency with  $X$  on LHS cannot be minimal.
    - ◇ Just remove  $B$ .
  
- $C^+(X) = \{A \in R \mid \forall B \in X: X \setminus \{A, B\} \rightarrow B \text{ does not hold}\}$ 
  - Special case:  $A = B$  corresponds to  $C(X)$
  - $C(X) = R \setminus \{A \in X \mid X \setminus A \rightarrow A \text{ holds}\}$
  
- This definition removes three types of candidates.
  - $C1 = \{A \in X \mid X \setminus A \rightarrow A \text{ holds}\}$  (as before)
  - $C2 = \{R \setminus X\}$  if  $\exists B \in X: X \setminus B \rightarrow B$
  - $C3 = \{A \in X \mid \exists B \in X \setminus A : X \setminus \{A, B\} \rightarrow B \text{ holds}\}$

# Example for C2

16

- $C^+(X) = \{A \in R \mid \forall B \in X: X \setminus \{A, B\} \rightarrow B \text{ does not hold}\}$
- $C2 = \{R \setminus X\}$  if  $\exists B \in X: X \setminus B \rightarrow B$
  
- $R = ABCD, X = ABC$
- $C(X) = ABCD$  initially
- Discovery of  $C \rightarrow B$ 
  - Remove B from  $C(X)$
  - Additionally remove  $R \setminus X = D$
  - Ok, because remaining combination of LHS contains B and C.
    - ◇  $ABC \rightarrow D$  is not minimal because  $C \rightarrow B$
- Together:  $C^+(X) = \{AC\}$

# C3

17

- $C3 = \{A \in X \mid \exists B \in X \setminus A : X \setminus \{A, B\} \rightarrow B \text{ holds}\}$ 
  - Same idea as before, but for subsets
- Assume  $X$  has proper subset  $Y$  ( $X \supset Y$ ) such that  $Y \setminus B \rightarrow B$  holds for some  $B \in Y$ .
- Then we can remove from  $C(X)$  all  $A \in X \setminus Y$ .
  
- Example  $X = ABCD$  and let  $C \rightarrow B$
- $X \supset Y = BC$  and  $X \setminus Y = AD$
- Thus can remove all  $AD$ .
  - Any remaining combination of LHS contains  $B$  and  $C$ .
    - ◇  $ABC \rightarrow D$  and  $BCD \rightarrow A$
  - Again, since  $C \rightarrow B$  any such FD is not minimal.
- Together:  $C^+(X) = \{C\}$



# More pruning of lattice: Key pruning

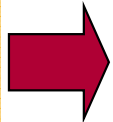
18

- Insight: If  $X$  is superkey and  $X \setminus B \rightarrow B$ , then  $X \setminus B$  is also a superkey.
- Case 1: If  $X$  is superkey, no need to test any  $X \rightarrow A$ .
- Case 2:
  - If  $X$  is superkey and not key, any  $X \rightarrow A$  is not minimal (for any  $A \notin X$ ).
  - If  $A \in X$  and  $X \setminus A \rightarrow A$  then  $X \setminus A$  is superkey, and no need to test.
- Summary: Can prune all keys and their supersets
- Later: Test for superkey-property based on “key-error” of partition

# Overview

19

- Functional Dependencies
- TANE
  - Candidate sets
  - Pruning Algorithm
  - Dependency checking
  - Approximate FDs
- FD\_Mine
- Conditional FDs



# TANE Base Algorithm

20

```
1   $L_0 := \{\emptyset\}$ 
2   $C^+(\emptyset) := R$ 
3   $L_1 := \{\{A\} \mid A \in R\}$ 
4   $\ell := 1$ 
5  while  $L_\ell \neq \emptyset$ 
6      COMPUTE_DEPENDENCIES( $L_\ell$ )
7      PRUNE( $L_\ell$ )
8       $L_{\ell+1} := \text{GENERATE\_NEXT\_LEVEL}(L_\ell)$ 
9       $\ell := \ell + 1$ 
```

- Each level  $L$  contains the corresponding nodes of the lattice

# Generating Lattice Levels

21

- $L_{l+1} = \{X \mid |X| = l+1 \text{ and all subsets } Y \subset X \text{ of size } l \text{ are in } L_l\}$ 
  - General apriori idea
  - Can use  $L_l$  to generate  $L_{l+1}$

```

1    $L_{l+1} := \emptyset$ 
2   for each  $K \in \text{PREFIX\_BLOCKS}(L_l)$  do
3       for each  $\{Y, Z\} \subseteq K, Y \neq Z$  do
4            $X := Y \cup Z$ 
5           if for all  $A \in X, X \setminus \{A\} \in L_l$  then
6                $L_{l+1} := L_{l+1} \cup \{X\}$ 
7   return  $L_{l+1}$ 

```

- Prefix blocks: disjoint sets from  $L_l$  with common prefix of size  $l-1$ 
  - All pairs for  $l = 1$
- Line 5: All subsets of new set must appear in lower level

# Dependency Computation

22

```

1  for each  $X \in L_\ell$  do
2       $\mathcal{C}^+(X) := \bigcap_{A \in X} \mathcal{C}^+(X \setminus \{A\})$ 
3  for each  $X \in L_\ell$  do
4      for each  $A \in X \cap \mathcal{C}^+(X)$  do
5          if  $X \setminus \{A\} \rightarrow A$  is valid then
6              output  $X \setminus \{A\} \rightarrow A$ 
7              remove  $A$  from  $\mathcal{C}^+(X)$ 
8              remove all  $B$  in  $R \setminus X$  from  $\mathcal{C}^+(X)$ 

```

} Trivial for L1:  
Nothing happens

- Line 2: Create candidate sets; each attribute must appear in all candidate sets of smaller size
- Line 4: Only test attributes from candidate set
- Line 5: Actual test on data
- Line 7: Reduce candidates by newly found dependent
- Line 8: Reduce candidates by all other attributes: cannot depend on all others, because any combination involving A and LHS is not minimal

# Pruning

23

```

1   for each  $X \in L_\ell$  do
2       if  $\mathcal{C}^+(X) = \emptyset$  do
3           delete  $X$  from  $L_\ell$ 
4       if  $X$  is a (super)key do
5           for each  $A \in \mathcal{C}^+(X) \setminus X$  do
6               if  $A \in \bigcap_{B \in X} \mathcal{C}^+(X \cup \{A\} \setminus \{B\})$  then
7                   output  $X \rightarrow A$ 
8               delete  $X$  from  $L_\ell$ 

```

- Line 3: Basic pruning
  - Deletion from  $L_\ell$  ensures that supersets cannot be created during level generation (loops not executed on empty candidate sets)
- Lines 4-8: Key pruning

# Example run

24

- $R = ABCD$ ,  $C \rightarrow B$  and  $AB \rightarrow D$  are to be discovered
  - Also:  $AC \rightarrow D$  through pseudo-transitivity
- $L_0 = \{\}$ ,
  - $C^+(\{\}) = ABCD$
  - nothing to do
- $L_1 = \{A\}, \{B\}, \{C\}, \{D\}$ 
  - $C^+(X) = ABCD$  for all  $X \in L_1$
  - Still nothing to do: No FDs can be generated from singletons
  - Thus, no pruning
- $L_2 = AB, AC, AD, BC, BD, CD$ 
  - E.g.  $C^+(AB) = C^+(AB \setminus A) \cap C^+(AB \setminus B) = ABCD \cap ABCD$
  - $C^+(X) = ABCD$  for all  $X \in L_2$
  - Dep. checks for  $AB$ :  $A \rightarrow B$  and  $B \rightarrow A$  Nothing happens

# Example run

25

- $L_2 = AB, AC, AD, BC, BD, CD$ 
  - $C^+(X) = ABCD$  for all  $X \in L_2$
  - Dep. checks for BC:  $B \rightarrow C$  (no!) and  $C \rightarrow B$  (yes!)
  - Output  $C \rightarrow B$
  - Delete B from  $C^+(BC)$ : ACD
  - Delete  $R \setminus BC$  from  $C^+(BC)$ : C
    - ◇ Note  $BC \rightarrow A$  and  $BC \rightarrow D$  are not minimal
- $L_3 = ABC, ABD, ACD, BCD$ 
  - $C^+(ABC) = C^+(AB) \cap C^+(AC) \cap C^+(BC) = C$
  - $C^+(BCD) = C^+(BC) \cap C^+(BD) \cap C^+(CD) = C$
  - $C^+(ABD) = C^+(ACD) = ABCD$  unchanged
  - Dep. check for ABC:  $ABC \cap C^+(ABC)$  are candidates
    - ◇  $AB \rightarrow C$  no! Did not check  $BC \rightarrow A$  and  $AC \rightarrow B$



# Example run

26

- $L_3 = ABC, ABD, ACD, BCD$ 
  - $C^+(ABC) = C^+(BCD) = C$
  - $C^+(ABD) = C^+(ACD) = ABCD$
  - Dep. check for ABD:  $ABD \cap C^+(ABD)$  are candidates
    - ◇  $AD \rightarrow B$  and  $BD \rightarrow A$ : no!
    - ◇  $AB \rightarrow D$ : yes! Output  $AB \rightarrow D$
    - ◇ Delete D from  $C^+(ABD)$ : ABC
    - ◇ Delete  $R \setminus ABD$  from  $C^+(ABD)$ : AB
  - Dep. check for BCD:  $BCD \cap C^+(BCD)$  are candidates
    - ◇ Only need to check  $BD \rightarrow C$ : no!
  - Dep. check for ACD:  $ACD \cap C^+(ACD)$  are candidates
    - ◇  $CD \rightarrow A$  and  $AD \rightarrow C$ : no!
    - ◇  $AC \rightarrow D$ : yes! Output  $AC \rightarrow D$
    - ◇ Delete D from  $C^+(ABD)$ : ABC
    - ◇ Delete  $R \setminus ACD$  from  $C^+(ABD)$ : AC

# Example run

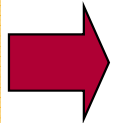
27

- $L_4 = ABCD$ 
  - $C^+(ABCD) = C^+(ABC) \cap C^+(ABD) \cap C^+(ACD) \cap C^+(BCD) = \{\}$
  - Nothing to check
  - Did not need to check
    - ◇  $BCD \rightarrow A$ : Not minimal because  $C \rightarrow B$
    - ◇  $ACD \rightarrow B$ : Not minimal because  $C \rightarrow B$
    - ◇  $ABD \rightarrow C$ : Not minimal because  $AB \rightarrow D$
    - ◇  $ABC \rightarrow D$ : Not minimal because  $AB \rightarrow D$

# Overview

28

- Functional Dependencies
- TANE
  - Candidate sets
  - Pruning Algorithm
  - Dependency checking
  - Approximate FDs
- FD\_Mine
- Conditional FDs



# Notation and Partitions

29

- Tuples  $t$  and  $u$  are equivalent wrt. attribute set  $X$  if  $t[A] = u[A]$  for all  $A \in X$ .
  
- Attribute set  $X$  partitions  $R$  into equivalence classes
  - Equivalence class of tuple  $t$ :
$$[t]_X = \{u \in R \mid \forall A \in X : t[A] = u[A]\}$$
  - Partition is set of disjoint sets:  $\pi_X = \{[t]_X \mid t \in R\}$ 
    - ◇ Each set has unique values for  $X$ -values.
  - $|\pi|$  is number of equivalence classes in  $\pi$ .

# Partitioning - Example

30

TupleID	A	B	C	D
1	1	a	\$	Flower
2	1	A		Tulip
3	2	A	\$	Daffodil
4	2	A	\$	Flower
5	2	b		Lily
6	3	b	\$	Orchid
7	3	C		Flower
8	3	C	#	Rose

- $[1]_A = [2]_A = \{1, 2\}$
- $\pi_A = \{\{1, 2\}, \{3, 4, 5\}, \{6, 7, 8\}\}$
- $\pi_{BC} = \{\{1\}, \{2\}, \{3, 4\}, \{5\}, \{6\}, \{7\}, \{8\}\}$
- $\pi_D = \{\{1, 4, 7\}, \{2\}, \{3\}, \{5\}, \{6\}, \{8\}\}$

# Partition refinement

31

- Partition  $\pi$  *refines* partition  $\pi'$  if every equivalence class in  $\pi$  is a subset of some equivalence class in  $\pi'$ .
  - $\pi$  has a finer partitioning than  $\pi'$ .
- $X \rightarrow A \Leftrightarrow \pi_X$  refines  $\pi_A$ 
  - If  $\pi_X$  refines  $\pi_A$  then  $\pi_{X \cup A} = \pi_A$ .
  - $\pi_{X \cup A}$  always refines  $\pi_A$ .
  - $\Rightarrow$  if  $\pi_{X \cup A} \neq \pi_A$  then  $|\pi_X| \neq |\pi_{X \cup A}|$ .
  - $\Rightarrow$  if  $|\pi_X| = |\pi_{X \cup A}|$  then  $\pi_{X \cup A} = \pi_A$ .
- Together:  $X \rightarrow A \Leftrightarrow \pi_X$  refines  $\pi_A \Leftrightarrow |\pi_X| = |\pi_{X \cup A}|$ 
  - This implies a simple check for an FD.
- $\pi_A = \{12, 345, 678\}$
- $\pi_B = \{12345, 678\}$
- $\pi_{AB} = \{12, 345, 678\}$

	A	B
1	1	A
2	1	A
3	2	A
4	2	A
5	2	A
6	3	B
7	3	B
8	3	B

# Stripped partitions

32

- Idea: Remove equivalence classes of size 1 from partitions.
  - Singleton equivalence class cannot violate any FD.
  - Same idea as for position lists
- Problem
  - $X \rightarrow A \Leftrightarrow |\pi_X| = |\pi_{X \cup A}|$  not true for stripped partitions  $\pi'$
  - $\Rightarrow$ :  $|\pi_C| = |\pi_{C \cup A}| = 6$  and  $|\pi'_C| = |\pi'_{C \cup A}| = 2$
  - $\Leftarrow$ :  $|\pi'_B| = |\pi'_{B \cup C}| = 2$  but  $B \not\rightarrow C$
- Solution: Key error
  - $e(X)$  is minimum fraction of tuples to remove for  $X$  to be key
  - $e(X) = 1 - |\pi_X|/r$ 
    - ◇  $e(B) = 1 - |\pi_B|/r = 1 - 3/8 = 5/8$
  - $e(X) = (||\pi'_X|| - |\pi'_X|)/r$ 
    - ◇  $||\pi'_X|| =$  sum of sizes of equivalence classes in  $\pi'$
    - ◇  $e(B) = ((5+2) - 2)/r = 5/8$
  - $X \rightarrow A \Leftrightarrow e(X) = e(X \cup A)$

	A	B	C
1	1	A	a
2	1	A	b
3	2	A	c
4	2	A	c
5	2	A	d
6	3	B	e
7	3	B	e
8	3	D	f

# Computing partitions

33

- Compute partition  $\pi_A$  for each  $A \in R$ 
  - Directly from database
  - Only store tuples ID (Integers)
- Product  $\pi_1 \cdot \pi_2$ : Least refined partition that refines both  $\pi_1$  and  $\pi_2$ 
  - $\pi_X \cdot \pi_Y = \pi_{X \cup Y}$
  - Partitions  $\pi_X$  computed as product of two partitions of size  $|X|-1$ .
  - Algorithm on next slide
- Dependency checking:  $X \setminus A \rightarrow A$ 
  - Calculate  $e(X) = (||\pi'_X|| - |\pi'_X|)/r$  and  $e(X \setminus A) = \dots$
  - Check  $e(X) = e(X \setminus A)$
- Also key pruning:  $X$  is key if  $e(X) = 0$ .



# Partition Product

34

**Input:** Stripped partitions  $\widehat{\pi}' = \{c'_1, \dots, c'_{|\widehat{\pi}'|}\}$  and  $\widehat{\pi}'' = \{c''_1, \dots, c''_{|\widehat{\pi}''|}\}$ .

**Output:** Stripped partition  $\widehat{\pi} = \widehat{\pi}' \cdot \widehat{\pi}''$ .

```

1   $\widehat{\pi} := \emptyset$ 
2  for  $i := 1$  to  $|\widehat{\pi}'|$  do
3    for each  $t \in c'_i$  do  $T[t] := i$ 
4     $S[i] := \emptyset$ 
5  for  $i := 1$  to  $|\widehat{\pi}''|$  do
6    for each  $t \in c''_i$  do
7      if  $T[t] \neq \text{NULL}$  then  $S[T[t]] := S[T[t]] \cup \{t\}$ 
8    for each  $t \in c''_i$  do
9      if  $|S[T[t]]| \geq 2$  then  $\widehat{\pi} := \widehat{\pi} \cup \{S[T[t]]\}$ 
10      $S[T[t]] := \emptyset$ 
11  for  $i := 1$  to  $|\widehat{\pi}'|$  do
12    for each  $t \in c'_i$  do  $T[t] := \text{NULL}$ 
13  return  $\widehat{\pi}$ 

```

	A	B	C
1	0	4	7
2	1	5	7
3	1	5	8
4	2	6	9
5	2	6	9
6	3	4	7

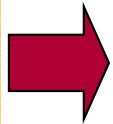
$\pi'_{A \cup B} = \{23, 45\}$

$\pi'_{B \cup C} = \{16, 45\}$

# Overview

35

- Functional Dependencies
- TANE
  - Candidate sets
  - Pruning Algorithm
  - Dependency checking
  - Approximate FDs
- FD\_Mine
- Conditional FDs





# Making TANE approximate

36

- Definition based on minimum number of tuples to be removed from  $R$  for  $X \rightarrow A$  to hold in  $R$ .
- Discovery problem:
  - Given relation  $R$  and threshold  $\varepsilon$ , find all minimal non-trivial FDs  $X \rightarrow A$  such that  $e(X \rightarrow A) \leq \varepsilon$ .
- 1. Define error: Fraction of tuples causing FD violation
  - Error  $e(X \rightarrow A) = \min\{|S| \mid S \subseteq R, R \setminus S \models X \rightarrow A\} / |R|$
- 2. Specify error threshold  $\varepsilon$
- 3. Modify dependency checking algorithm
  - Efficient algorithm to compute error
  - Bounds to avoid error calculation

# Approximate Dependency Checking

37

<p>1    <b>for each</b> <math>X \in L_\ell</math> <b>do</b></p> <p>2        <math>C^+(X) := \bigcap_{A \in X} C^+(X \setminus \{A\})</math></p> <p>3    <b>for each</b> <math>X \in L_\ell</math> <b>do</b></p> <p>4        <b>for each</b> <math>A \in X \cap C^+(X)</math> <b>do</b></p> <p>5            <b>if</b> <math>X \setminus \{A\} \rightarrow A</math> <b>is valid then</b></p> <p>6                <b>output</b> <math>X \setminus \{A\} \rightarrow A</math></p> <p>7                <b>remove</b> <math>A</math> <b>from</b> <math>C^+(X)</math></p> <p>8                <b>remove all</b> <math>B</math> <b>in</b> <math>R \setminus X</math> <b>from</b> <math>C^+(X)</math></p>	 <p>Exact version</p>
<hr/>	
<p>5'            <b>if</b> <math>e(X \setminus \{A\} \rightarrow A) \leq \varepsilon</math> <b>then</b></p> <p>8'                <b>if</b> <math>X \setminus \{A\} \rightarrow A</math> <b>holds exactly then</b></p> <p>9'                <b>remove all</b> <math>B</math> <b>in</b> <math>R \setminus X</math> <b>from</b> <math>C^+(X)</math></p>	 <p>Modifications</p>

# Computing error

38

- Error  $e(X \rightarrow A) = \min\{|S| \mid S \subseteq R, R \setminus S \models X \rightarrow A\} / |R|$
- Any equivalence class  $c \in \pi_X$  is the union of one or more equivalence classes  $c_1', c_2', \dots \in \pi_{X \cup A}$
- For each  $c \in \pi_X$  the tuples in all but one of the  $c_i$ 's must be removed for  $X \rightarrow A$  to hold.
- Minimum number to remove: Size of  $c$  minus size of largest  $c_i$ '.

- $$e(X \rightarrow A) = 1 - \frac{\sum_{c \in \pi_X} \max\{|c'| \mid c' \in \pi_{X \cup A} \wedge c' \subseteq c\}}{|R|}$$

- Example:  $B \rightarrow A$

- $\pi_A = \{12, 345, 678\}$
- $\pi_B = \{1, 234, 56, 78\}$
- $\pi_{AB} = \{1, 2, 34, 5, 6, 78\}$
- $|\pi_B| \neq |\pi_{BA}|$
- $e(B \rightarrow A) = 8/8 - (1+2+1+2)/8 = 2/8$

- Also possible on stripped partitions – not here.

	A	B
1	1	a
2	1	A
3	2	A
4	2	A
5	2	b
6	3	b
7	3	C
8	3	C

# Bounding on error

40

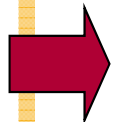
- Computing error is in  $O(|R|)$
- $e(X) - e(X \cup A) \leq e(X \rightarrow A) \leq e(X)$
- I.e., do not calculate FD error if
  - $e(X) - e(X \cup A) > \epsilon$
  - $e(X) < \epsilon$
- $e(B) = 4/8$
- $e(B \cup A) = 2/8$
- $e(B \rightarrow A) = 8/8 - (1+2+1+2)/8 = 2/8$

	A	B
1	1	a
2	1	A
3	2	A
4	2	A
5	2	b
6	3	b
7	3	C
8	3	C

# Overview

41

- Functional Dependencies
- TANE
  - Candidate sets
  - Pruning Algorithm
  - Dependency checking
  - Approximate FDs
- FD\_Mine
- Conditional FDs



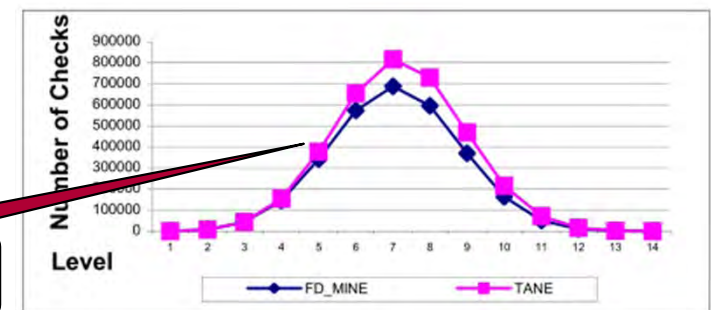
# FD\_Mine: Refinement of TANE

42

- If  $X \rightarrow Y$  and  $Y \rightarrow X$  hold, then  $X$  and  $Y$  are *equivalent candidates*, denoted as  $X \leftrightarrow Y$ .
- Make use of additional FD properties
  - $X \leftrightarrow Y$  and  $XW \rightarrow Z \Rightarrow YW \rightarrow Z$
  - $X \leftrightarrow Y$  and  $WZ \rightarrow X \Rightarrow WZ \rightarrow Y$
- Example
  - $A \rightarrow D$  and  $D \rightarrow A \Rightarrow A \leftrightarrow D$
  - Examination:  $AB \rightarrow C$  and  $BC \rightarrow A$
  - Inferred:
    - ◇  $BD \rightarrow C$  (property 1)
    - ◇  $BC \rightarrow D$  (property 2)
  - $D$  can be removed from table

	A	B	C	D	E
1	0	0	0	1	0
2	0	1	0	1	0
3	0	2	0	1	2
4	0	3	1	1	0
5	4	1	1	2	4
6	4	3	1	2	2
7	0	0	1	1	0

Pairs of UCCs



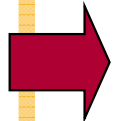
- Hong Yao, Howard J. Hamilton, Cory J. Butz: FD\_Mine: *Discovering Functional Dependencies in a Database Using Equivalences*. ICDM 2002: 729-732



# Overview

43

- Functional Dependencies
- TANE
  - Candidate sets
  - Pruning Algorithm
  - Dependency checking
  - Approximate FDs
- FD\_Mine
- Conditional FDs



# Conditional Functional Dependencies

44

- Idea similar to CINDs
- Embedded FD plus pattern tableau
- Definition CFD
  - Pair  $(X \rightarrow A, T_p)$ 
    - ◇  $X \rightarrow A$  is embedded FD
    - ◇  $T_p$  is pattern tableau, made up of pattern tuples  $t_p$
  - Pattern tuple with attributes  $B \in X \cup A$  where  $t_p[B]$ 
    - ◇ Constant in  $\text{dom}(B)$
    - ◇ Unnamed variable “\_” for values in  $\text{dom}(B)$
  - Special case:  $T_p[B] = \text{“_”}$  for all  $B$  is equivalent to normal FD

# Semantics of CFDs

45

- $a \approx b$  (a matches b) if
  - either a or b is “\_”
  - both a and b are constants and  $a = b$
  
- DB satisfies  $(R: X \rightarrow Y, T_p)$  iff
  - For any tuple  $t_p$  in the pattern tableau  $T_p$  and for any tuples  $t_1, t_2$  in DB:
    - ◇ If  $t_1[X] = t_2[X] \approx t_p[X]$ , then  $t_1[Y] = t_2[Y] \approx t_p[Y]$
  - $t_p[X]$ : identifying the set of tuples on which the constraint  $t_p$  applies:  $\{ t \mid t[X] \approx t_p[X] \}$
  - $t_1[Y] = t_2[Y] \approx t_p[Y]$ : enforcing the embedded FD, and the pattern of  $t_p$

Slide from Wenfei Fan

# Example: Violation of CFDs

46

id	country	area-code	phone	street	city	zip
t1	44	131	1234567	Mayfield	NYC	EH4 8LE
t2	44	131	3456789	Crichton	NYC	EH8 8LE
t3	01	908	3456789	Mountain Ave	NYC	07974
t4	01	908	9876543	Mainstreet	NYC	07974

- $\text{cust}([\text{country}, \text{zip}] \rightarrow [\text{street}], T_p)$

- Tuples t1 and t2 violate the CFD

- $t1[\text{country}, \text{zip}] = t2[\text{country}, \text{zip}] \approx t_p[\text{country}, \text{zip}]$

- But  $t1[\text{street}] \neq t2[\text{street}]$

- The CFD applies to t1 and t2 since they match  $t_p[\text{country}, \text{zip}]$

- Tuples t3 and t4 do not violate the CFD

- CFD does not apply to t3 and t4

country	zip	street
44	–	–

Slide from Wenfei Fan

# Example: Violation of CFD by single tuple

47

id	country	area-code	phone	street	city	zip
t1	44	131	1234567	Mayfield	NYC	EH4 8LE
t2	44	131	3456789	Crichton	NYC	EH8 8LE
t3	01	908	3456789	Mountain Ave	NYC	07974
t4	01	908	9876543	Mainstreet	NYC	07974

- $\text{cust}([\text{country}, \text{zip}] \rightarrow [\text{street}], \text{Tp})$
- Tuple t1 does not satisfy the CFD.

	country	zip	street
tp1	44	131	Edi
tp2	01	908	MH
tp3	–	–	–

- $t1[\text{country}, \text{area-code}] = t1[\text{country}, \text{area-code}] \approx tp1[\text{country}, \text{area-code}]$
- $t1[\text{city}] = t1[\text{city}]$ ; however,  $t1[\text{city}]$  does not match  $tp1[\text{city}]$
- In contrast to traditional FDs, a single tuple may violate a CFD.

# Further literature

48

- DepMiner
  - *Efficient Discovery of Functional Dependencies and Armstrong Relations*. Stéphane Lopes, Jean-Marc Petit, and Lotfi Lakhal. EDBT 2000
- CORDS
  - Ihab F. Ilyas, Volker Markl, Peter J. Haas, Paul Brown, Ashraf Aboulnaga: CORDS: Automatic Discovery of Correlations and Soft Functional Dependencies. SIGMOD Conference 2004: 647-658
- FastFDs
  - Catharine M. Wyss, Chris Giannella, Edward L. Robertson: *FastFDs: A Heuristic-Driven, Depth-First Algorithm for Mining Functional Dependencies from Relation Instances*. DaWaK 2001: 101-110
- CFDs
  - Loreto Bravo, [Wenfei Fan](#), [Shuai Ma](#): Extending Dependencies with Conditions. [VLDB 2007](#): 243-254