



Detecting Functional Dependencies

1.6.2017
Felix Naumann

Overview

1. Functional Dependencies

2. TANE

- Candidate sets
- Pruning Algorithm
- Dependency checking
- Approximate FDs

3. FD_Mine and DFD

4. Conditional FDs

5. fdep

6. HyFD

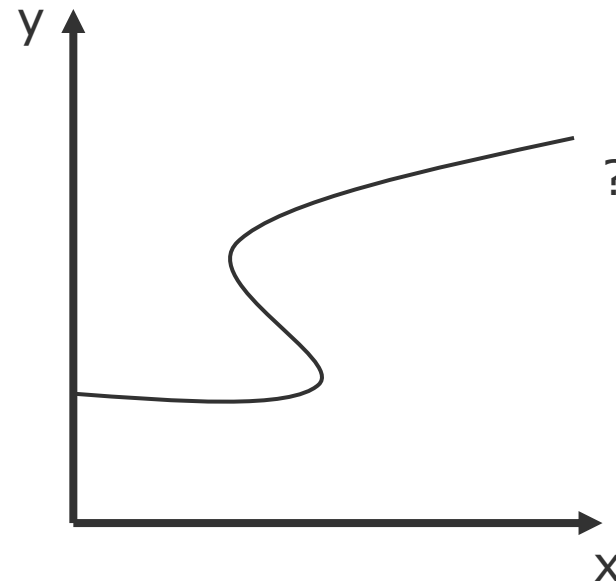
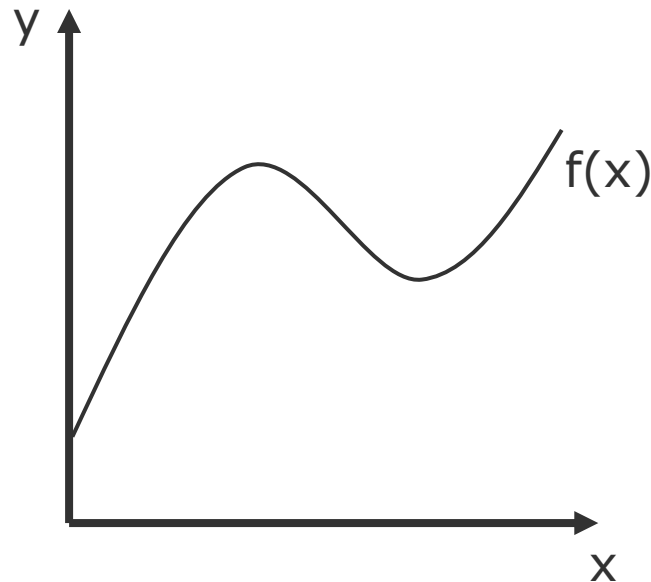
Thorsten Papenbrock



Felix Naumann
Data Profiling
Summer 2017

Definition – Functional Dependency

- „ $X \rightarrow A$ “ 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 $t_1, t_2 \in R$: $t_1[X] = t_2[X] \Rightarrow t_1[A] = t_2[A]$
- Generalization to sets: $X \rightarrow Y \Leftrightarrow X \rightarrow A$ for each $A \in Y$



Trivial and minimal FDs

- 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.
- Basis: Minimal set of FDs to derive all other FDs.
- Typical goal: Given a relation R, find all minimal completely non-trivial functional dependencies.

FD Inference Rules

- 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

- 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 partial/conditional FDs)
 - Anomaly detection
 - Data integrity constraints
 - Data curation rules
 - Query optimization: Independence of column attributes
 - Index selection

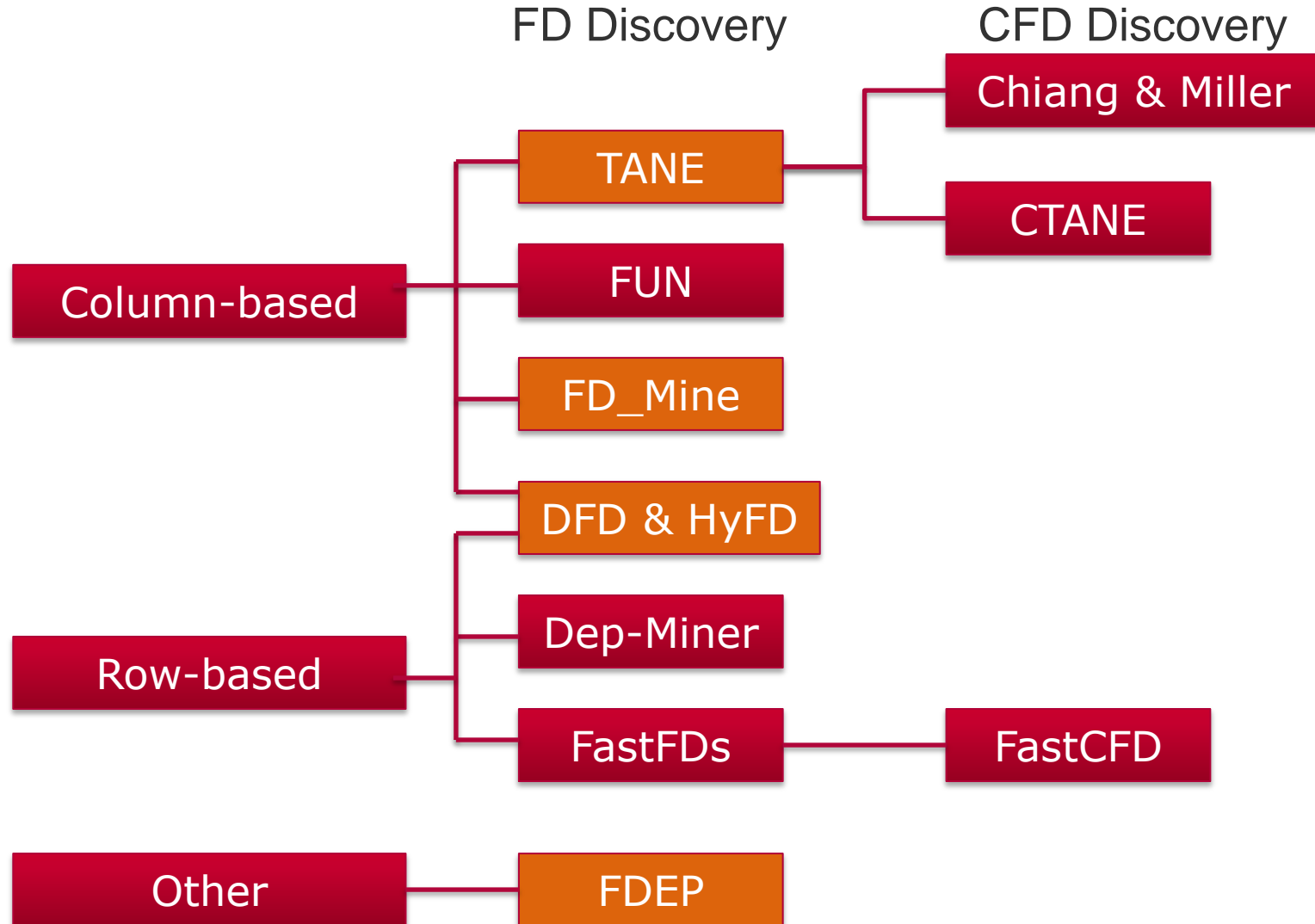
Naive Discovery Approach

- Task: Given relation R , detect all minimal, non-trivial FDs $X \rightarrow A$.

- For each column combination X
 - For each $A \in X$
 - 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
 - For each of the $|R|$ possibilities for RHS
 - check $2^{(|R|-1)}$ combinations for LHS

Existing Solutions



Felix Naumann
Data Profiling
Summer 2017

Overview

1. Functional Dependencies
2. **TANE**
 - Candidate sets
 - Pruning Algorithm
 - Dependency checking
 - Approximate FDs
3. FD_Mine and DFD
4. Conditional FDs



Felix Naumann
Data Profiling
Summer 2017

Tane – General Idea

- Two key ideas
 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

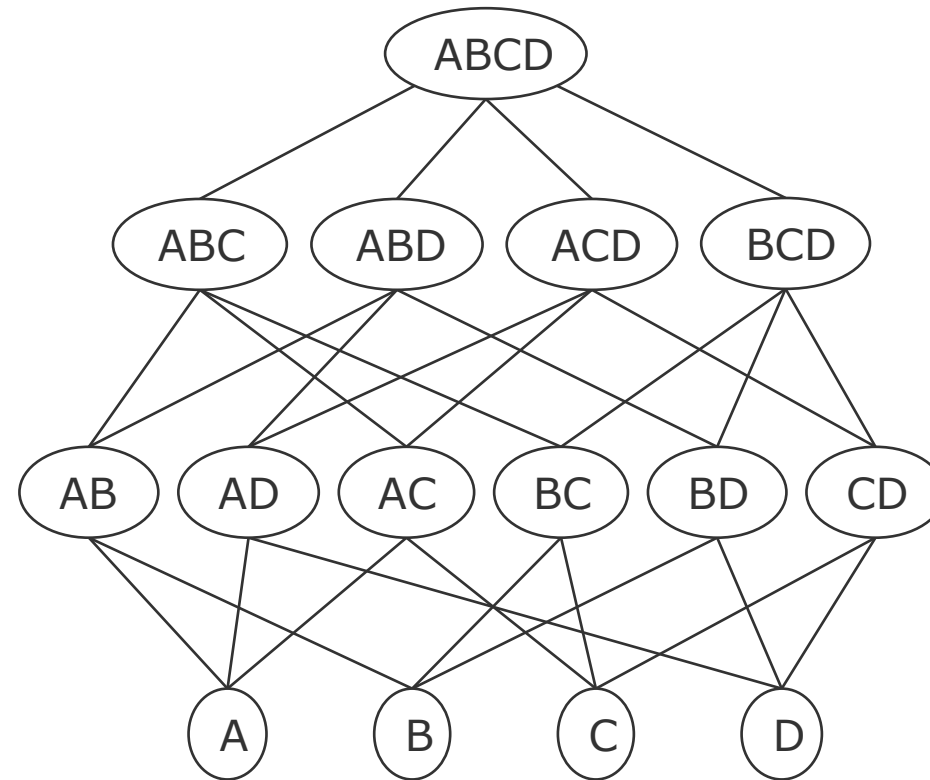


Hannu (T.T.)
Toivonen

Discovery strategy

- 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
 - Interpretation: Test each edge
 - Test on efficient data structure
 - "Stripped partitions" (=PLIs)



Overview

1. Functional Dependencies
2. TANE
 - **Candidate sets**
 - Pruning Algorithm
 - Dependency checking
 - Approximate FDs
3. FD_Mine and DFD
4. Conditional FDs



Felix Naumann
Data Profiling
Summer 2017

Candidate Sets

- 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

- 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: $A \notin C(AB) \cap C(AC) \cap C(BC)$
 - $AB \rightarrow C$ and $AC \rightarrow B$ need to be tested, because B and C appear in all candidate sets.

Improved RHS candidate pruning

- 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 all of X on LHS cannot be minimal.
 - Just remove B from candidate set.

- $C^+(X) = \{A \in R \mid \forall B \in X: X \setminus \{A, B\} \rightarrow B \text{ does not hold}\}$
 - Special case: $A = B$ corresponds to $C(X)$
 - Reminder: $C(X) = R \setminus \{A \in X \mid X \setminus A \rightarrow A \text{ holds}\}$

- 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

- $C^+(X) = \{A \in R \mid \forall B \in X: X \setminus \{A, B\} \rightarrow B \text{ does not hold}\}$
 - $C2 = \{R \setminus X\}$ 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: $C^+(ABC) = \{ABCD\} \setminus \{BD\} = \{AC\}$

Example for C3

- $C3 = \{A \in X \mid \exists B \in X \setminus A : X \setminus \{A, B\} \rightarrow B \text{ holds}\}$
 - 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

- 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 $X \rightarrow A$ (for any $A \notin X$).
 - $X \rightarrow A$ is always valid
- 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

1. Functional Dependencies
2. TANE
 - Candidate sets
 - **Pruning Algorithm**
 - Dependency checking
 - Approximate FDs
3. FD_Mine and DFD
4. Conditional FDs



Felix Naumann
Data Profiling
Summer 2017

TANE Base Algorithm

```
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
 - Without pruned nodes

Generating Lattice Levels

- $L_{l+1} = \{ X \mid |X| = l+1 \text{ and all subsets } Y \subset X \text{ of size } l \text{ are in } L_l \}$
 - This is the 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

```

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

```
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_1 ensures that supersets cannot be created during level generation (loops not executed on empty candidate sets)
- Lines 4-8: Key pruning

Example run

- $R = ABCD$, $C \rightarrow B$ 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

- $L_2 = AB, AC, AD, BC, BD, CD$
 - $C^+(X) = ABCD$ 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

- $L_3 = ABC, ABD, ACD, BCD$
 - $C^+(ABC) = C^+(BCD) = C$
 - $C^+(ABD) = C^+(ACD) = ABCD$
 - 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

- $L_4 = ABCD$
 - $C^+(ABCD) = C^+(ABC) \cap C^+(ABD) \cap C^+(ACD) \cap C^+(BCD) = \{\}$
 - 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

1. Functional Dependencies
2. TANE
 - Candidate sets
 - Pruning Algorithm
 - **Dependency checking**
 - Approximate FDs
3. FD_Mine and DFD
4. Conditional FDs



Felix Naumann
Data Profiling
Summer 2017

Notation and Partitions

- 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* (\approx PLI) is set of disjoint sets: $\pi_X = \{[t]_X \mid t \in R\}$
 - Each set has unique values for its X -attributes.
 - $|\pi|$ is number of equivalence classes in π .

Partitioning - Example

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

- Partition π *refines* partition π' if every equivalence class in π is a subset of some equivalence class in π' .
 - π has a “finer partitioning” than π' .
- $X \rightarrow A \Leftrightarrow \pi_X$ refines π_A
- π_X refines $\pi_A \Leftrightarrow |\pi_X| = |\pi_{X \cup A}|$
 - “ \Rightarrow ”: If π_X refines π_A then $\pi_{X \cup A} = \pi_X$.
 - “ \Leftarrow ”: $\pi_{X \cup A}$ always refines π_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 $\pi_{X \cup A} = \pi_X$.
- 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 π'

- Idea: Remove equivalence classes of size 1 from partitions.
 - Singleton equivalence class cannot violate any FD.
 - Same idea as for position lists (PLIs)
- Problem
 - $X \rightarrow A \Leftrightarrow |\pi_X| = |\pi_{X \cup A}|$ not true for stripped partitions π'
 - \Rightarrow : $|\pi_C| = |\pi_{C \cup A}| = 6$ and $|\pi'_C| = |\pi'_{C \cup A}| = 2$ (OK)
 - \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
 - $e(X) = 1 - |\pi_X| / |r|$
 - Fraction of distinct values is $|\pi_X| / |r|$, all others must be removed
 - $e(B) = 1 - |\pi_B|/r = 1 - 3/8 = 5/8$
 - $e(X) = (||\pi'_X|| - |\pi'_X|) / |r|$
 - $||\pi'_X||$ = sum of sizes of equivalence classes in π'
 - $e(B) = ((5+2) - 2) / |r| = 5/8$
 - $e(X) = e(Y) \Leftrightarrow |\pi_X| = |\pi_Y|$
 - $\Rightarrow X \rightarrow A \Leftrightarrow e(X) = e(X \cup A)$

	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

Felix Naumann
Data Profiling
Summer 2017

Computing partitions

- Compute partition π_A for each $A \in R$
 - Directly from database
 - Store only tuple's ID (Integers)
- Product $\pi_1 \cdot \pi_2$: Least refined partition that refines both π_1 and π_2
 - $\pi_X \cdot \pi_Y = \pi_{X \cup Y}$
 - Partitions π_X computed as product of two partitions of size $|X| - 1$.
 - Following bottom-up approach
 - 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 whether $e(X) = e(X \setminus A)$
- Also key pruning: X is key if $e(X) = 0$.

Partition Product

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

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

```

1   $\hat{\pi} := \emptyset$ 
2  for  $i := 1$  to  $|\hat{\pi}'|$  do
3      for each  $t \in c'_i$  do  $T[t] := i$ 
4       $S[i] := \emptyset$ 
5  for  $i := 1$  to  $|\hat{\pi}''|$  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  $\hat{\pi} := \hat{\pi} \cup \{S[T[t]]\}$ 
10              $S[T[t]] := \emptyset$ 
11 for  $i := 1$  to  $|\hat{\pi}'|$  do
12     for each  $t \in c'_i$  do  $T[t] := \text{NULL}$ 
13 return  $\hat{\pi}$ 

```

For each partition of first operand

For each tuple, remember in which partition it is

For each partition of second operand

Remember intersection

If intersection > 1, new element in π'

Clean up for re-use of T

	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

$\pi'_A = \{23, 45\}$
 $\pi'_B = \{16, 23, 45\}$
 $\pi'_C = \{126, 45\}$
 $\pi'_{A \cup B} =$
 $\pi'_{B \cup C} =$

Felix Naumann
Data Profiling
Summer 2017

Example run for π'_{BUC}

- Lines 2-4: $T[] = \langle 1, 2, 2, 3, 3, 1 \rangle$ and $S[] = \langle \perp, \perp, \perp \rangle$
 - Size of S: number of classes in π'_B .
- Lines 5-7 (for each class, for each tuple in π'_B)
 - $i=1$: $t = 1$ and $T[1]=1$: $S[] = \langle 1, \perp, \perp \rangle$
 - $i=1$: $t = 2$ and $T[2]=2$: $S[] = \langle 1, 2, \perp \rangle$
 - $i=1$: $t = 6$ and $T[6]=1$: $S[] = \langle 16, 2, \perp \rangle$
 - $i=1$: $t = 1$ and $|S[T[1]]| = 2$: New entry $\pi'_{BUC} = \{16\}$
 - $S[] = \langle \perp, 2, \perp \rangle$
 - $i=1$: $t = 2$ and $|S[T[2]]| = 1$: do nothing
 - $S[] = \langle \perp, \perp, \perp \rangle$
 - $i=1$: $t = 6$ and $|S[T[6]]| = 0$: do nothing
 - $S[] = \langle \perp, \perp, \perp \rangle$
 - $i=2$: $t = 4$ and $T[4]=3$: $S[] = \langle \perp, \perp, 4 \rangle$
 - $i=2$: $t = 5$ and $T[5]=3$: $S[] = \langle \perp, \perp, 45 \rangle$
 - $i=2$: $t = 4$ and $|S[T[4]]| = 2$: New entry $\pi'_{BUC} = \{16, 45\}$
 - $S[] = \langle \perp, \perp, \perp \rangle$
 - $i=2$: $t = 5$ and $|S[T[5]]| = 0$: do nothing
 - $S[] = \langle \perp, \perp, \perp \rangle$

	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

$$\pi'_A = \{23, 45\}$$

$$\pi'_B = \{16, 23, 45\}$$

$$\pi'_C = \{126, 45\}$$

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

$$\pi'_{BUC} = \{16, 45\}$$

Overview

1. Functional Dependencies
2. TANE
 - Candidate sets
 - Pruning Algorithm
 - Dependency checking
 - **Approximate FDs**
3. FD_Mine and DFD
4. Conditional FDs



Felix Naumann
Data Profiling
Summer 2017

Making TANE approximate

- 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 ε , 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 ε
- 3. Modify dependency checking algorithm
 - Efficient algorithm to compute error
 - Bounds to avoid error calculation

Approximate Dependency Checking

```

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  $\mathcal{C}^+(X)$ 

```

```

5'      if  $e(X \setminus \{A\} \rightarrow A) \leq \varepsilon$  then
8'          if  $X \setminus \{A\} \rightarrow A$  holds exactly then
9'              remove all  $B$  in  $R \setminus X$  from  $\mathcal{C}^+(X)$ 

```

Exact version

Modifications

Computing error

- 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: $B \rightarrow A$

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

- Computing 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

Felix Naumann
Data Profiling
Summer 2017

Bounding on error

- 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
 - $e(X) - e(X \cup A) > \varepsilon$
 - $e(X) < \varepsilon$
- $e(B) = 4/8$
- $e(B \cup A) = 2/8$
- $e(B \rightarrow A) = 8/8 - (1+2+1+2)/8 = 2/8$

	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

1. Functional Dependencies
2. TANE
 - Candidate sets
 - Pruning Algorithm
 - Dependency checking
 - Approximate FDs
- 3. FD_Mine and DFD**
4. Conditional FDs



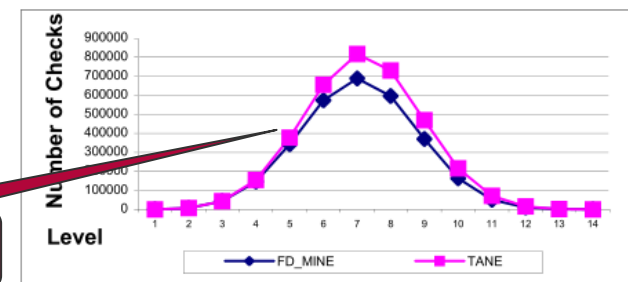
Felix Naumann
Data Profiling
Summer 2017

FD_Mine: Refinement of TANE

- 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

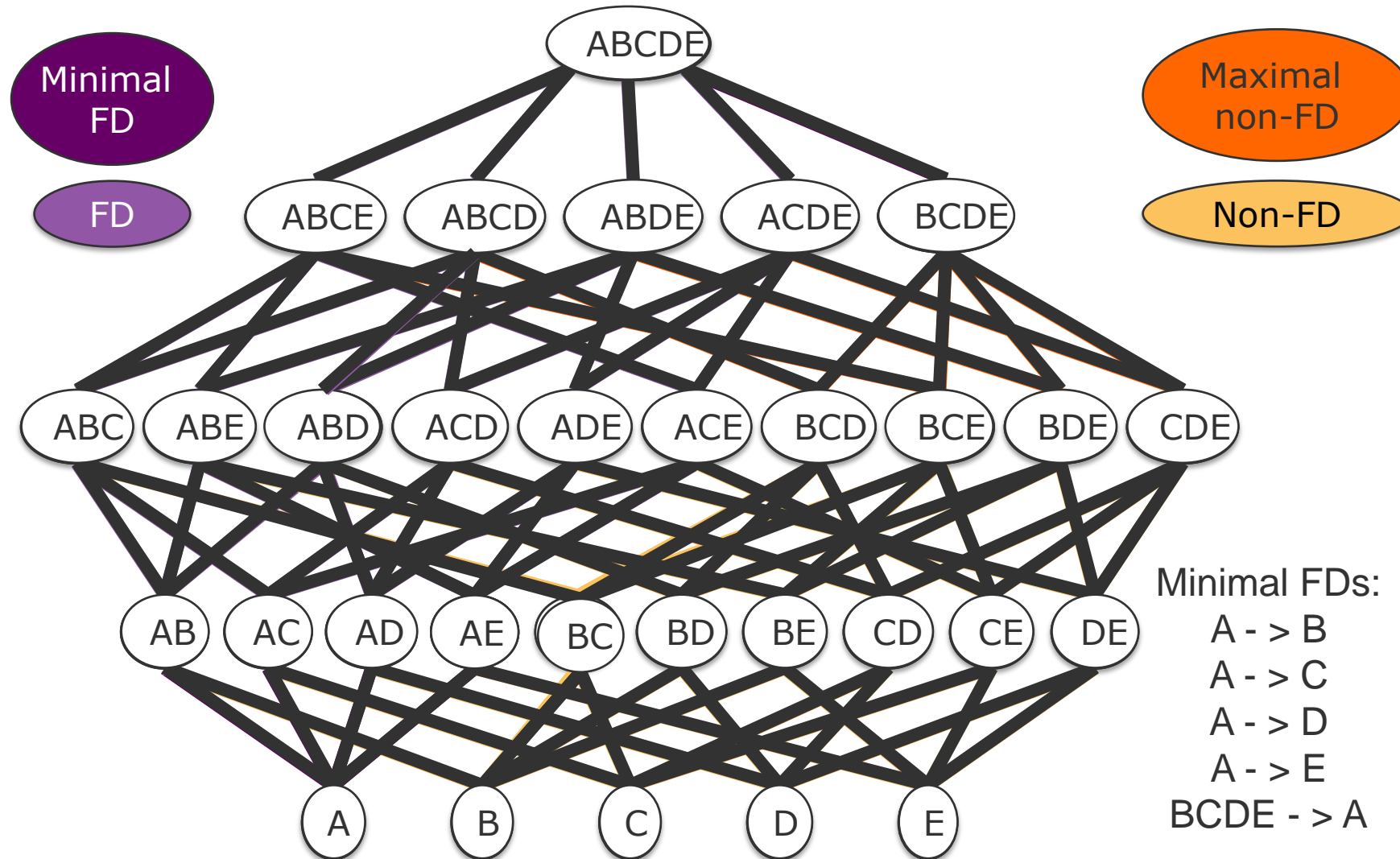
Pairs of UCCs



Felix Naumann
Data Profiling
Summer 2017

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

DFD Explanation:
Tane visualized for R = (A,B,C,D,E)



DFD: Depth-first approach for functional dependency discovery

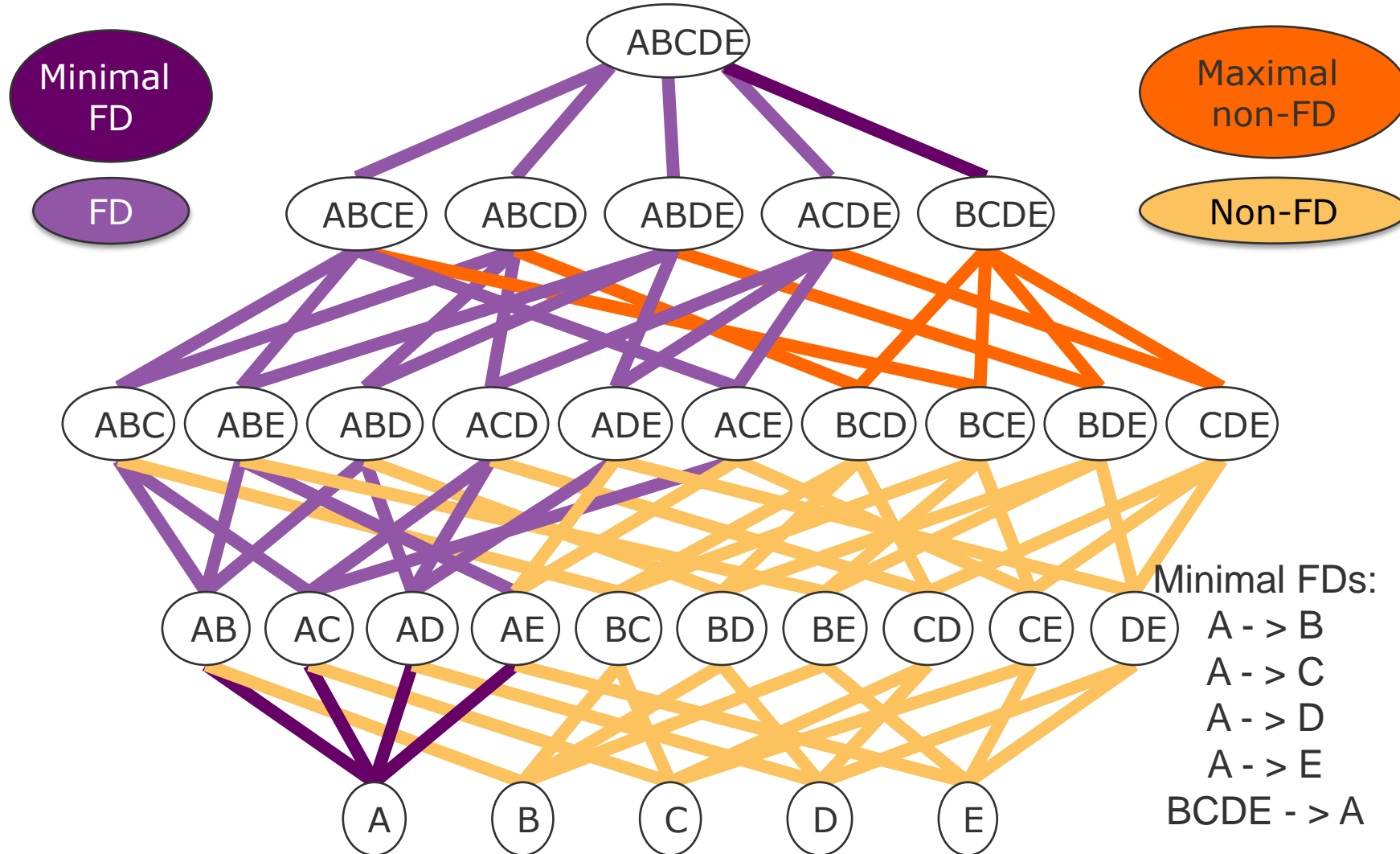
Idea: Aggressive Pruning

- Traverse depth-first and prune upwards and downwards
- Applied for key/unique discovery [DUCC, PVLDB 2013]
 - Key discovery is a subproblem of FD discovery
 - Adapt the concept of minimality in keys to LHS of FDs:

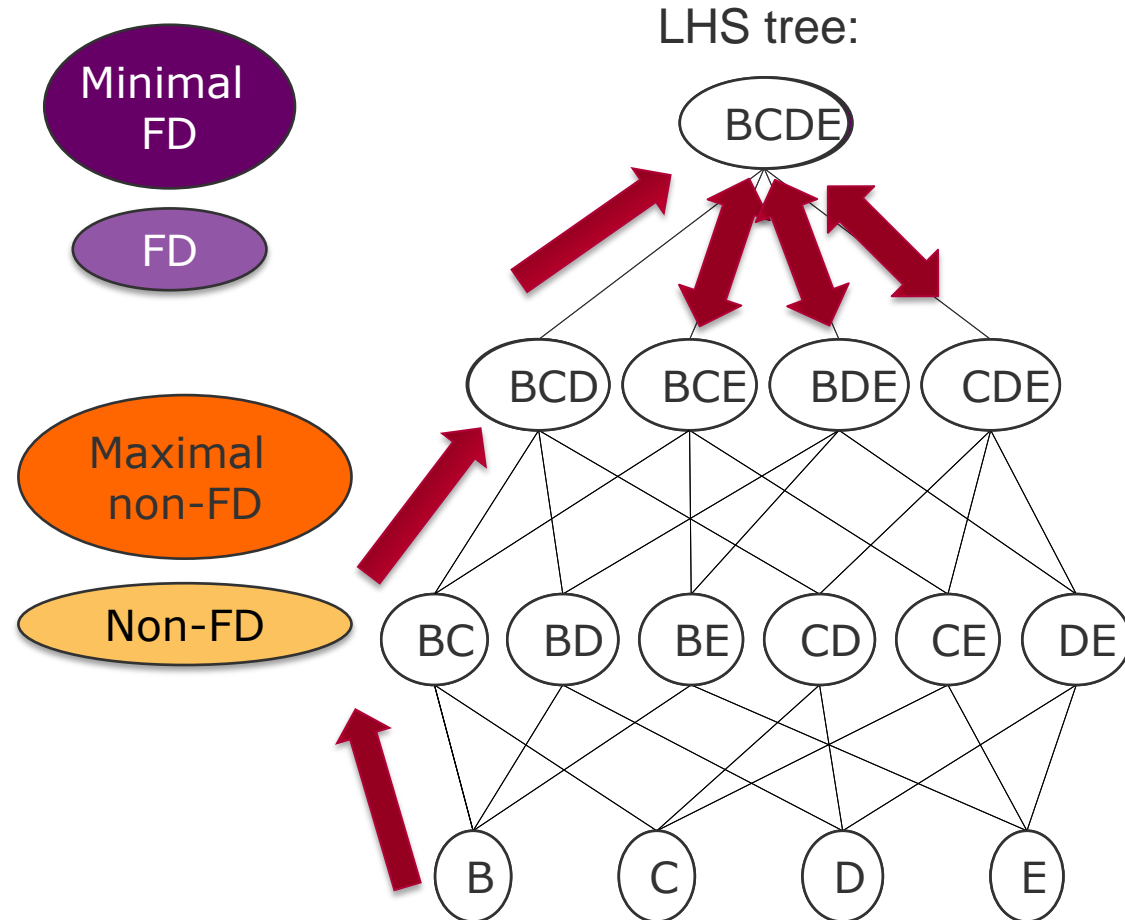
– An FD $X \rightarrow C$ is minimal if $\neg \exists X' \subset X : X' \xrightarrow{NOT} C$

– A non-dependency $X \not\rightarrow C$ is maximal if $\neg \exists X' \supset X : X' \rightarrow C$

Decompose Relation for each RHS



Decomposition for RHS=A



Checked combinations
Tane: 15
DFD: 7

Minimal FDs:
BCDE → A

Felix Naumann
Data Profiling
Summer 2017

Traversal Holes

- Aggressive traversal and pruning
 - Some nodes might never be reached.
- [GORDIAN, VLDB'04]:
 - Complement the set of **maximal non-keys** = set of **minimal keys**
- [DUCC, PVLDB'13]
 - Key observation: the **difference** of one set and the complement of its counterpart delivers the **unvisited nodes!**
- Hole discovery works for FDs:
 - Consider **minimal FD LHS** and **maximal non-FD LHS**

Functional Dependency Evaluation

dataSet	Columns	Rows	FDs	Tane	FUN	FD_Mine	Dep-Miner	FastFDs	FDep	DFD
iris	5	150	4	0.6s	0.1s	0.1s	0.1s	0.1s	0.1s	0.1s
balance-scale	5	625	1	0.9s	0.4s	0.3s	0.2s	0.5s	0.3s	0.2s
chess	7	28,056	1	2.0s	1.0s	3.0s	200.8s	200.1s	202.5s	0.9s
abalone	9	4,177	137	1.0s	0.3s	1.0s	2.9s	3.0s	4.1s	0.9s
nursery	9	12,960	1	3.1s	1.5s	6.0s	132.0s	131.9s	56.6s	1.1s
breast-cancer	11	699	46	1.4s	0.4s	1.5s	0.9s	1.0s	0.4s	0.9s
bridges	13	108	142	1.3s	0.5s	2.9s	0.2s	0.2s	0.2s	0.9s
echocardiogram	13	132	538	0.8s	0.1s	69.9s	0.1s	0.1s	0.1s	1.6s
adult	14	48,842	78	81.2s	150.2s	485.3s	5982s	5946s	760.7s	6.8s
letter	17	20,000	61	326s	553.9s	ML	865.4s	853.9s	292.3s	9.1s
hepatitis	20	155	8,250	10.9s	321.6s	TL	5363.1s	9.3s	0.5s	317.8s
horse	27	368	128,726	5451.s	TL	TL	TL	386.8s	15.7s	TL
fd-reduced-30	30	250,000	89,571	41.1s	78.4s	TL	391.9s	391.3s	TL	TL
flight	109	1,000	982,631	ML	TL	ML	TL	TL	213.5s	TL
plista	125	1,000	178,152	ML	TL	TL	TL	TL	26.4s	TL

Felix Naumann
Data Profiling
Summer 2017

Overview

1. Functional Dependencies
2. TANE
 - Candidate sets
 - Pruning Algorithm
 - Dependency checking
 - Approximate FDs
3. FD_Mine and DFD
4. **Conditional FDs**



Felix Naumann
Data Profiling
Summer 2017

Conditional Functional Dependencies

- Embedded (=partial) FD with its 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 $t_p[B]$ is...
 - a constant in $\text{dom}(B)$
 - an unnamed variable “_” for values in $\text{dom}(B)$
 - Special case: $T_p[B] = \text{“_”}$ for all B is equivalent to normal FD
- Idea similar to CINDs (later in lecture)

Semantics of CFDs

- $a \approx b$ (value a matches value 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 $t_1[Y] = t_2[Y] \approx t_p[Y]$
 - $t_p[X]$: identifies the set of tuples on which the constraint t_p applies:
 - $\{ t \mid t[X] \approx t_p[X] \}$
 - $t_1[Y] = t_2[Y] \approx t_p[Y]$: enforces the embedded FD, and the pattern of t_p

Example: Violation of CFDs

id	country	area-code	phone	street	city	zip
t1	44	131	1234567	Mayfield	NYC	EH4 8LE
t2	44	131	3456789	Crichton	NYC	EH4 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	—	—

Example: Violation of CFD by single tuple

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

- $\text{cust}([\text{country}, \text{area code}] \rightarrow [\text{city}], \text{Tp})$
- 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.

Felix Naumann
Data Profiling
Summer 2017

Further literature

- 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