

Relational DB Design: Functional Dependencies / Normalization

- Schema & Query Language determine application programs
- Goal: Have a Schema that makes queries simpler, avoids redundancy, update anomalies, loss of info

Example:

Schema 1: Emp (eno, ename, byr, sal, dno)
 Dept (dno, dname, floor, mgr) ← eno of mgr

Q1: Find all employees who make more than their manager:

```

Select E.ename
From Emp E, Dept D, Emp E2
Where (D.mgr = E2.eno) and (E.dno = D.dno)
and (E.sal > E2.sal)
  
```

Q2: Find all ~~dept~~ depts with a max salary > 2x average salary

```

Select (d.dname)
From Emp e, Dept d
Where (d.dno = e.dno)
Group by d.dno, d.dname
Having (Max(e.sal) > 2 * Avg(e.sal))
  
```

Schema 2: ED (eno, ename, byr, sal, dno, dname, floor, mgr)

Q1: Select e.ename
 from ED e1, ED e2
 where (e1.mgr = e2.eno) and (e1.sal > e2.sal)

Q2: Select e.dname
 From ED e
 Group by e.dno, e.dname
 Having MAX(e.sal) > 2 * AVG(e.sal)

* We get simpler queries, but we also get into trouble:

- Redundancy:

- each dpt is repeated once for each employee
- there is potential inconsistency (update anomalies).
We may change the location of a dpt in one tuple leaving it with the old value in another
- A simple change is translated into multiple replacements.
(e.g. change the manager of the "shoe" department)

- No independent existence:

- A dpt can not exist without employees
- when we delete the last employee of a dpt, we automatically lose track of the dept.

Objective of DB Design:

- No redundancy for space efficiency
- update integrity
- Semantic clarity
- Linguistic efficiency (the simpler the queries the better for both the user and query optimizer)
- Performance (binary relations imply most queries will have a large number of joins)

Tools for DB design : Functional Dependencies + Normalization

Functional Dependency : Let A, B be sets of attributes of R .

The FD $A \rightarrow B$ holds if a value for A uniquely determines a value for B .

More formally : $A \rightarrow B$ if :

\forall pairs (t_1, t_2) of tuples in relation r such that

$t_1[A] = t_2[A]$, it is also the case $t_1[B] = t_2[B]$

Note let K be a super-key of R , then

$$K \rightarrow R$$

Examples

Emp : $e_{no} \rightarrow e_{name}$
 $e_{no} \rightarrow d_{no}$

$e_{ssl} \not\rightarrow e_{name}$

Dept : $d_{no} \rightarrow mgr$
 $d_{no} \rightarrow floor$

FD : $d_{no} \rightarrow floor$
 $e_{no} \rightarrow ssl$
 $e_{no} \rightarrow mgr$

Various types of FD's help identify bad designs

- Trivial Dependency : $(A, B) \rightarrow A$ (Identity)
- Partial dependency : (A, B) is a key and $A \rightarrow C$

Ex: supply ($\overset{\text{supplier \#}}{\downarrow} s_{no}, \overset{\text{part \#}}{\downarrow} p_{no}, proj_{no}, scity, projcity, qty$)

key is $(s_{no}, p_{no}, proj_{no})$

$s_{no} \rightarrow scity$
 $proj_{no} \rightarrow projcity$ } partial dependencies

- Transitive Dependency : A is a key, B is not, and $A \rightarrow B \rightarrow C$

Ex : ED($e_{no}, \dots, d_{no}, \dots, mgr$) where e_{no} is the key.

$e_{no} \rightarrow d_{no} \rightarrow mgr$
 $e_{no} \rightarrow mgr$ } transitive dependency

Depending on the existence of various FDs, schemes are classified in various Normal Forms

1NF : Every attribute has an atomic value (as opposed to set value) Note, all relational model schemas are by definition 1NF

2NF : 1NF and no partial dependencies (partial dependencies ⇒ redundancy)

3NF : 2NF and no transitive dependencies to attributes that are not part of a key. (transitive dependencies ⇒ redundancy)

Equivalently, if $X \rightarrow A$ is a FD then either:
a) it is trivial, or
b) X is a superkey, or
c) A is a subset of a candidate key.

BCNF : (Boyce Codd Normal Form) 3NF and no transitive dependencies. Equivalently, if $X \rightarrow A$ is a FD then either:
a) it is trivial, or
b) X is a super key

4NF (5,6...) Multivalued dependencies (rarely used, mostly just theory for theories sake)

Distinguishing Example Schemas

1NF but not 2NF: Supply(sno, pno, projno, scity, projcity, qty)
Key = (sno, pno, projno)

2NF but not 3NF: ED(eno, ename, byr, sal, dno, dname, floor, mgr)
Key = (eno)

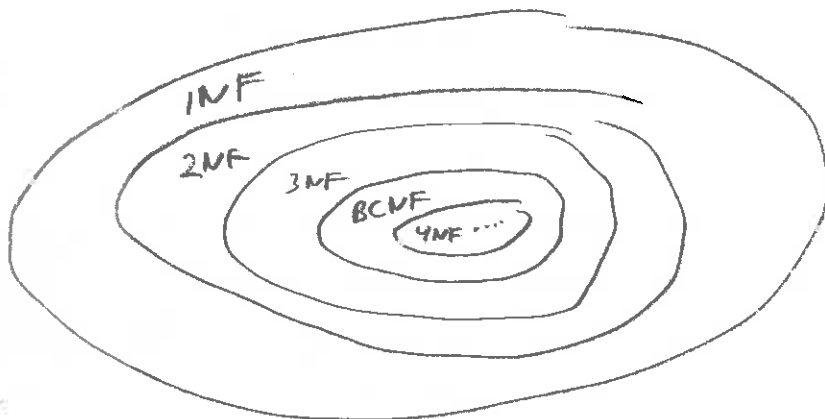
One receipt per person, food pair

3NF but not BCNF: Restaurant (person, food, rec-number)

Key = (person, food)

This is 3NF because it has transitive dependency to an attribute that is part of the key.

Each normal form is included in the next one higher up.



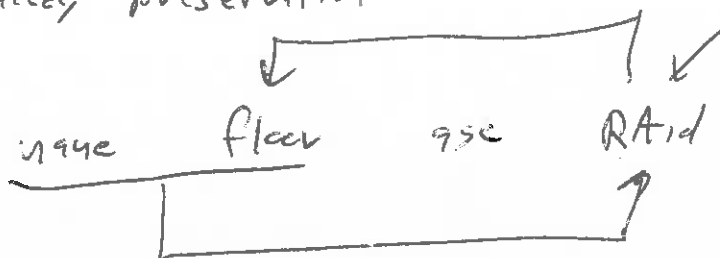
Ideally, a decomposition should satisfy:

- BCNF
- Lossless joins
- dependency preservation

If we can not get this, we settle for:

- 3NF
- Lossless joins
- dependency preservation

3NF example



Lossless joins: The decomposition should be done in a way such that no information is lost.

EX: $Relation = (name, loc, chanc, amount)$

↓ decompose

Ant-schema (amount, chanc)

Loan-schema (name, loc, amount)

~~loan~~

name	loc	amount
hampton	1	1000
norfolk	2	1500
wmbs	3	2000
NW	4	1500

~~amount~~

amount	chanc
1000	Sue
1500	Mike
2000	Biff
1500	Jane

amt \bowtie loan

=

hampton	1	1000	Sue
norfolk	2	1500	Mike
norfolk	2	1500	Jane
wmbs	3	2000	Biff
NW	4	1500	Mike
NW	4	1500	Jane

← loans

Gives us incorrect results! Hence we have lost information.

Let $F \triangleq$ Set of n semantically obvious FD

More FD's can be inferred from F

$F^+ \triangleq$ set of all inferred FDs

Example

$$F = \left\{ \begin{array}{l} A \rightarrow B \\ A \rightarrow C \\ B \rightarrow K \end{array} \right\} \quad \text{then} \quad F^+ = \left\{ \begin{array}{l} A \rightarrow B \\ A \rightarrow C \\ A \rightarrow K \\ B \rightarrow K \end{array} \right\}$$

The closure of F can be systematically determined using the inference rules:

(Note: $F \models X \rightarrow Y$ denotes that $X \rightarrow Y$ can be inferred from F)

(1) reflexive rule: If $X \supseteq Y$, then $X \rightarrow Y$

(2) augmentation rule

$$a) \quad \{X \rightarrow Y\} \models XZ \rightarrow YZ$$

$$b) \quad \{X \rightarrow Y\} \models XZ \rightarrow Y$$

(3) transitive rule: $\{X \rightarrow Y, Y \rightarrow Z\} \models X \rightarrow Z$

(4) decomposition rule: $\{X \rightarrow YZ\} \models X \rightarrow Y$

(5) additive rule: $\{X \rightarrow Y, X \rightarrow Z\} \models X \rightarrow YZ$

(6) pseudotransitive rule: $\{X \rightarrow Y, WY \rightarrow Z\} \models WX \rightarrow Z$

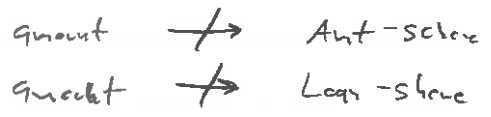
When decomposing a schema we must make sure there are no lossless joins. This can be done as follows:

Let $R \stackrel{\Delta}{=} \text{schema we are decompose}$
 $R_1, R_2 \stackrel{\Delta}{=} \text{two decomposed schemas (Note } R_1 \cup R_2 = R)$

The decomposition is a lossless-join decomposition if one of the two is true:

- a) $R_1 \cap R_2 \rightarrow R_1$
- b) $R_1 \cap R_2 \rightarrow R_2$

Note Ant-schema \wedge Loss-schema = quant



hence, decomposition may lose information & is rejected.

Dependency Preservation

The goal is to preserve functional dependencies without requiring a join. (Important if DBMS checks validity of dependencies after inserts, deletes, modifications)

- First, let $F = \{ \text{of all functional dependencies in } R \}$
- The closure of F , denoted F^+ , is the set of all logically implied dependencies. ~~(not by at back)~~

Ex: $F = \left\{ \begin{matrix} A \rightarrow B \\ A \rightarrow C \\ B \rightarrow H \end{matrix} \right\}$, then $F^+ = \left\{ \begin{matrix} A \rightarrow B \\ A \rightarrow C \\ A \rightarrow H \\ B \rightarrow H \end{matrix} \right\}$

- 1) reflexivity: $B \subset A, \text{ then } A \rightarrow B$
- 2) augmentation: if $A \rightarrow B, \text{ then } A \rightarrow AB$
- 3) transitivity: if $A \rightarrow B, \text{ and } B \rightarrow C, \text{ then } A \rightarrow C$

(1), (2), + (3) are known as Armstrong's axioms.

— Armstrong axioms are complete, i.e. by repeatedly applying (1)→(3) F^+ can be determined.

Algorithm for determining closure of attribute set α

result = α

while (changes to result)

{

for each FD $\beta \rightarrow \gamma$ in F

if $\beta \subseteq \text{result}$

then result = result $\cup \gamma$

}

Example:

assume $F =$ $A \rightarrow B$

$A \rightarrow C$

$CG \rightarrow H$

$CG \rightarrow I$

$B \rightarrow H$

A^+ (the closure of A under F) can be calculated

	result	
	<u>A</u>	
iter 1	A, B	$A \rightarrow B$
		$A \rightarrow C$
iter 1	A, B, C	$CG \rightarrow H$ no, $CG \rightarrow I$ no, $B \rightarrow H$ yes
iter 1	ABCH	thus $A^+ = \{A, B, C, H\}$

How about $\{AG\}^+$?

Case 1 : A, B, C, H, I

Case 2 : no check,

Note $F^+ = \{A, B, C, G, H, I\}^+$

* In worst case this algorithm may take time quadratic in the size of F . There is a linear (but more complex) algorithm.

- Often FD are used to make sure the DB does not move into an inconsistent state by the insertion of a tuple or update.

- In order to minimize the number of FD that need to be tested, restrict F into a smaller (but just as powerful) set.

- (1) combine several FD into one where possible
- (2) remove extraneous attributes

- Attribute A is extraneous in α if $A \in \alpha$ and F logically implies $(F - \{\alpha \rightarrow \beta\}) \cup \{(\alpha - A) \rightarrow \beta\}$

- Attribute A is extraneous in β ($\alpha \rightarrow \beta$) if $A \in \beta$ and F logically implies $(F - \{\alpha \rightarrow \beta\}) \cup \{\alpha \rightarrow (\beta - A)\}$

- A canonical cover (or minimal cover) F_c is a set of FD such that F logically implies all FD in F_c , and F_c logically implies all FD in F . Furthermore, F_c must have the following properties:

- (i) no FD in F_c contains an extraneous attribute
- (ii) every FD has a single value on the R.H.S.
- (iii) we can not remove a FD and still be equivalent to F

Example: calculate F_c for $F =$

$$A \rightarrow BC$$

$$B \rightarrow C$$

$$A \rightarrow B$$

$$AB \rightarrow C$$

1) Make R.H.S have single values only

$$\hat{F} = \begin{array}{l} A \rightarrow B \\ A \rightarrow C \\ B \rightarrow C \\ \cancel{A \rightarrow B} \\ AB \rightarrow C \end{array}$$

2) A is extraneous in $AB \rightarrow C$ since $B \rightarrow C$ logically implies $AB \rightarrow C$,

thus $(\hat{F} - \{AB \rightarrow C\}) \cup \{B \rightarrow C\}$ logically implies F^+

removing A from $AB \rightarrow C$ gives $B \rightarrow C$ which is

already in result $\Rightarrow \begin{cases} A \rightarrow B \\ A \rightarrow C \\ B \rightarrow C \end{cases}$

3) we can remove $A \rightarrow C$ and still be equivalent to F

$$\hat{F} - \{A \rightarrow C\} \Rightarrow F^+$$

$$\Rightarrow F_c = \begin{cases} A \rightarrow B \\ B \rightarrow C \end{cases}$$

Minimal Cover

19.8e

Formal Algorithm: (note, is this also $A \rightarrow BC$
is replaced with $A \rightarrow B, A \rightarrow C, A \rightarrow B$)

1) Set $G = F$

2) replace each FD $X \rightarrow A_1, A_2, \dots, A_n$ in G by FDs $X \rightarrow A_1, X \rightarrow A_2, \dots, X \rightarrow A_n$

3) (4) For each FD $X \rightarrow A$ in G
= compute X^+ with respect to $\{G - (X \rightarrow A)\}$

if X^+ contains A , remove $X \rightarrow A$ from G

4) (3) for each remaining FD $X \rightarrow A \in G$

for each attribute B that is an element of X

~~compute $(X-B)^+$ with respect to $(G - (X \rightarrow A)) \cup ((X-B) \rightarrow A)$~~
~~if $(X-B)^+$ contains A , replace $X \rightarrow A$ with $(X-B) \rightarrow A$ in G~~

~~remove B if extraneous~~

These methods can be used to determine if a relation schema decomposition is dependency preserving.

Minimal Cover Example

19.8p

$$F = \left. \begin{array}{l} ABH \rightarrow C \\ A \rightarrow D \\ C \rightarrow E \\ BGH \rightarrow F \end{array} \right\} \quad \left. \begin{array}{l} F \rightarrow AD \\ E \rightarrow F \\ BH \rightarrow E \end{array} \right\}$$

①

↓ simplify

STEP 1

① $ABH \rightarrow C$

⑤ $F \rightarrow A$

② $A \rightarrow D$

⑥ $F \rightarrow D$

③ $C \rightarrow E$

⑦ $E \rightarrow F$

④ $BGH \rightarrow F$

⑧ $BH \rightarrow E$

STEP 2

See if can remove FD,

1) ~~Does $(ABH)^+ \in C$ if ① is removed?~~

No → can not remove ①

2) ~~Does $A^+ \in D$ if remove ② → No~~

3) ~~Does $C^+ \in E$ if remove ③ → No~~

4) ~~Does $(BGH)^+ \in F$ if remove ④?~~

Yes, can remove ④

5) ~~Does $F^+ \in A$ if remove ⑤? No~~

6) ~~Does $F^+ \in D$ " " ⑥? yes can remove 6~~

7) ~~No~~

8) ~~No~~

STEP 2

Remove extraneous ATTR

① Is A extraneous in $ABH \rightarrow C$?~~If \exists $BH \rightarrow C$ by $ABH \rightarrow C$~~ Yes, if $BH \rightarrow C$ can be derivedfrom ~~$F \rightarrow (ABH \rightarrow C) \cup (BH \rightarrow C) F$~~ $BH \rightarrow E$ $E \rightarrow F \Rightarrow BH \rightarrow A \Rightarrow \text{~~ABH} \rightarrow A~~$ $F \rightarrow AD$ $\downarrow \checkmark$
 $ABH \rightarrow C$ ① $\Rightarrow BH \rightarrow C$ Is B extraneous in $BH \rightarrow C$? NO~~Yes~~Is H extraneous in $BH \rightarrow C$? NO~~①~~ ① $\Rightarrow BH \rightarrow C$ ④ Is G extraneous in $DBH \rightarrow F$?

$$\left. \begin{array}{l} BH \rightarrow E \\ E \rightarrow F \end{array} \right\} BH \rightarrow F$$

thus $BGH \rightarrow F$, the G is extra

$\Rightarrow BH \rightarrow F$

End of step 2

$BH \rightarrow C$

$A \rightarrow D$

$C \rightarrow E$

~~$BH \rightarrow F$~~ \Rightarrow

$F \rightarrow A$

~~$E \rightarrow D$~~ remove redundant

$E \rightarrow F$ FDs

~~$BH \rightarrow E$~~

Initial cover

\downarrow
 $BH \rightarrow C$

$A \rightarrow D$

$C \rightarrow E$

$F \rightarrow A$

$E \rightarrow F$

Now, how about schema decomposition:

Let R be decomposed into R_1, R_2, \dots, R_n

The restriction of F to R_i is the set F_i of all functional dependencies that include only attributes of R_i

Let $F' = F_1 \cup F_2 \cup \dots \cup F_n$

If $F'^+ = F^+$ the decomposition is dependency preserving

Example of how this is used:

Banker-schema = (branch-name, customer-name, banker-name)

$$F = \begin{cases} \text{banker} \rightarrow \text{branch} \\ \text{customer, branch} \rightarrow \text{banker} \end{cases}$$

Note, key = (customer, branch)

Banker-schema is not BCNF, why?

banker \rightarrow branch, banker is not a super key,
+ the dependency is key-trivial.

⇓ decompose

Banker-Branch-schema (banker, branch)

Customer-Banker-schema (customer, banker)

This is BCNF + lossless, but does it preserve dependency?

NO!

$$F_1 = \{ \text{banker} \rightarrow \text{branch} \}$$

$$F_2 = \emptyset \text{ (only trivial dependencies)}$$

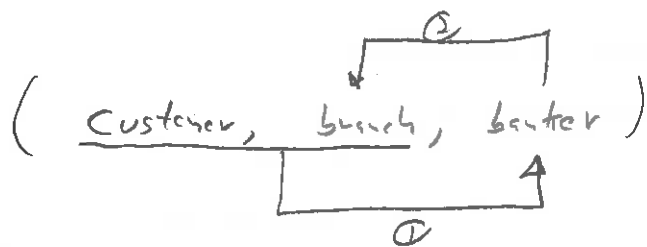
$$\Rightarrow F_1^+ \neq F_2^+$$

Thus, we could not satisfy

- a) BCNF, and
- b) lossless, and
- c) dependency preserving

But, we can get 3NF, lossless, dependency-preserving since

Banker-Scheme IS 3NF



① $X \rightarrow B$, X is a superkey

② $B \rightarrow A$, A is a subset of a candidate key

Advantage of Normal Form approach:

- It is a way to formalize a process that is usually only done using intuition

① Try to decompose R into BCNF (lossless)

Algorithm:

- 1) set $D = \{R\}$
- 2) while \exists schema $Q \in D$ that is not BCNF
 - a) choose $Q \in D$ that is not BCNF
 - b) find FD $X \rightarrow Y$ that violates BCNF
 - c) replace Q by two schemas $(Q - Y)$ and $(X \cup Y)$

Note, the final decomposition may not be ~~BCNF~~ dependency preserving,

② \Rightarrow test to see if dependency preserving,

③ If not dependency preserving, decompose into lossless & dependency-preserving 3NF

Algorithm

- 1) Find minimal cover F_c for F
- 2) for each left hand side $X \in F_c$ create a relation schemas $(X \cup A_1, X \cup A_2, \dots, X \cup A_n)$ where $X \rightarrow A_1, X \rightarrow A_2, \dots, X \rightarrow A_n$ are all the dependencies $\in F_c$ with X in left hand side
- 3) replace all remaining (unused) attributes into a single relation
- 4) If none of the relation schemas contains a key of R , create one more relation schemas that contains attributes that form a key

Example:

Lots (property-id, county-name, lot#, area, price, tax-mile)

⇓ abbreviate

Lots (pid, cname, lot#, area, price, tax)

F { $\begin{array}{l} \text{pid} \rightarrow \text{cname, lot\#, area, price, tax} \\ \text{cname, lot\#} \rightarrow \text{pid, area, price, tax} \\ \text{cname} \rightarrow \text{tax} \\ \text{area} \rightarrow \text{price} \end{array} \right\}$ two candidate keys!

Note $F^+ = F$

Is Lots BCNF?

No: $\text{cname} \rightarrow \text{tax}$ is a partial dependency, not even 3NF!

⇓ Decompose (BCNF decompose algorithm)

lots1 (pid, cname, lot#, area, price) lots2 (cname, tax)

Are both of these BCNF?

No: $\text{area} \rightarrow \text{price}$ is part of a transitive dependency!

⇓

lots1a (pid, cname, lot#, area) lots1b (area, price) lots2 (cname, tax)

BCNF? Yes!

But is it Dependency Preserving?

$$F_1 \stackrel{(lots!)}{=} \left\{ \begin{array}{l} pid \rightarrow cname, lot\#, area \\ cname, lot\# \rightarrow pid, area \end{array} \right.$$

$$F_2 \stackrel{(lots!)}{=} \left\{ area \rightarrow price \right.$$

$$F_3 \stackrel{!}{=} \left\{ cname \rightarrow tax \right.$$

$$F_1^+ \stackrel{!}{=} \left\{ \begin{array}{l} pid \rightarrow cname, lot\#, area \\ cname, lot\# \rightarrow pid, area \\ area \rightarrow price \\ cname \rightarrow tax \end{array} \right. +$$

$$= F^+$$

So Yes! it is dependency preserving

Lets add a 5th FD to our example:

area → change

In this case lots1 is not BCNF

∇ Decompose

- ① lots1x (pid, area, lot#)
- ② lots1y (area, change)
- ③ lots1z (area, price)
- ④ lots2 (change, tax)

Is it BCNF? Yes

Is it dependency preserving?

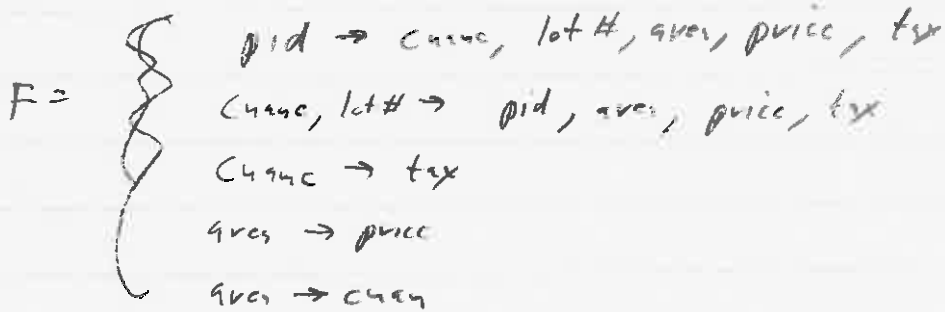
- F₁ = pid → area, lot#
- F₂ = area → change
- F₃ = area → price
- F₄ = change → tax

$$F^{1+} = \left\{ \begin{array}{l} pid \rightarrow change, lot\#, area, price, tax \\ change \rightarrow tax \\ area \rightarrow price, change, tax \end{array} \right.$$

but NOT change, lot# → pid, area, price, tax !

Not Dependency Preserving

Thus, we must settle for 3NF, depending previous, lessless and here use the decomposition algorithm:



Find minimal cover:

- 1 ~~pid → name~~
- 2 " → lot#
- 3 " → area
- 4 ~~" → price~~
- 5 ~~" → tax~~
- 6 name, lot# → pid
- 7 ~~" → area~~
- 8 ~~" → price~~
- 9 ~~" → tax~~
- 10 area → price
- 11 area → name
- 12 name → tax

pid⁺ (with respect to G - {pid → name})

= lot#, area, name, price, tax

this includes name, hence can remove ①

- can not remove 2

- lot 3
- yes 4
- yes 5
- lot 6
- yes 7
- yes 8
- yes 9
- lot 10
- lot 11
- lot 12

⇒

- pid → lot#
- " → area
- name, lot# → pid
- area → price
- area → name
- name → tax

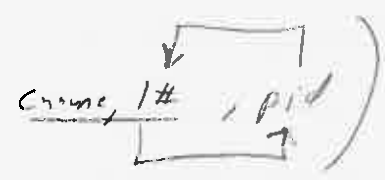
$$F^C \equiv \text{minimal cover} = \left\{ \begin{array}{l} pid \rightarrow \# \\ pid \rightarrow area \\ (name, \#) \rightarrow pid \\ area \rightarrow price \\ area \rightarrow crime \\ crime \rightarrow tax \end{array} \right.$$

3NF, lossless, preserves comp constraint:

- ① (pid, #, area)
- ② (area, price, crime)
- ③ (crime, #, pid)
- ④ (crime → tax)

It is 3NF!

(note, not BCNF because



① Check that it is dependency preserving:

$$F_1 = \text{pid} \rightarrow \text{I\#, gves}$$

$$F_2 = \text{gves} \rightarrow \text{price, cranc}$$

$$F_3 = \text{cranc, I\#} \rightarrow \text{pid}$$

$$F_4 = \text{cranc} \rightarrow \text{tax}$$

$$F^+ = \begin{cases} \text{pid} \rightarrow \text{I\#, gves, price, cranc, tax} \\ \text{cranc, I\#} \rightarrow \text{pid, gves, price, tax} \\ \text{cranc} \rightarrow \text{tax} \\ \text{gves} \rightarrow \text{price, cranc, tax} \end{cases}$$

