

Bases de données

Lecture 12

Nabil Mustafa

mustafa@lipn.univ-paris13.fr

new idea

LOSSLESS DECOMPOSITIONS

LOSSLESS DECOMPOSITIONS

Given a relation:

R (a, b, c, d, e)

say we decompose it into two tables:

R1 (a, b, c, d)

R2 (a, b, c, e)

Can we 'lose' information by doing this?
... what does that even mean?

LOSSLESS DECOMPOSITIONS

DECOMPOSING
TABLE:

R

a	b	c	d	e
a1	b1	c1	d1	e1
a1	b1	c1	d2	e2

R1

a	b	c	d
a1	b1	c1	d1
a1	b1	c1	d2

R2

a	b	c	e
a1	b1	c1	e1
a1	b1	c1	e2

RECONSTRUCTING
IT BACK:

```
SELECT R1.a, R1.b, R1.c, R1.d, R2.e
FROM R1, R2
WHERE R1.a = R2.a AND R1.b = R2.b AND R1.c = R2.c
```

JOINING **R1** and **R2** around a, b, c

Problem: failed to reconstruct
the original table!

a	b	c	d	e
a1	b1	c1	d1	e1
a1	b1	c1	d1	e2
a1	b1	c1	d2	e1
a1	b1	c1	d2	e2

≠ **R**

LOSSLESS DECOMPOSITIONS

Given a relation:

R (a, b, c, d, e)

say we decompose it into two tables:

R1 (a, b, c, d)

R2 (a, b, c, e)

Goal: (**R1** NATURAL JOIN **R2**) *must* give back exactly **R**

called a **LOSSLESS DECOMPOSITION**

LOSSLESS DECOMPOSITIONS

R

a	b	c	d	e
a1	b1	c1	d1	e1
a1	b1	c1	d2	e2



R1

a	b	c	d
a1	b1	c1	d1
a1	b1	c1	d2



R2

a	b	c	e
a1	b1	c1	e1
a1	b1	c1	e2

JOINING R1 and R2 around a, b, c

Why is this happening?

Reason: Multiple values
of d and e for fixed
values of a, b, c!

```
SELECT R1.a, R1.b, R1.c, R1.d, R2.e  
FROM R1, R2  
WHERE R1.a = R2.a AND R1.b = R2.b AND R1.c = R2.c
```

a	b	c	d	e
a1	b1	c1	d1	e1
a1	b1	c1	d1	e2
a1	b1	c1	d2	e1
a1	b1	c1	d2	e2

LOSSLESS DECOMPOSITIONS

R1						
<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>
a1	b1	c1			f1	g1
a2	b2	c2	d1	e1	f2	g2

R2

So the join gives 4 lines instead of 2!

LOSSLESS DECOMPOSITIONS

R1						
<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>
a1	b1	c1			f1	g1
a2	b2	c2	d1	e1	f2	g2

R2

$d, e \longrightarrow f, g$

LOSSLESS DECOMPOSITIONS

R1						
<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>
a1	b1	c1			f1	g1
a2	b2	c2	d1	e1		

R2			
d, e	→ f, g		

Now, each line of **R1** extends to a unique line in NATURAL JOIN

... the join gives 2 lines, correctly

LOSSLESS DECOMPOSITIONS

R1						
<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>
a1	b1	c1			f1	g1
a2	b2	c2	d1	e1	f2	g2

R2			
d, e	→ a, b, c		

LOSSLESS DECOMPOSITIONS

R1						
<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>
a1	b1	c1			f1	g1
			d1	e1	f2	g2

$\boxed{d, e \longrightarrow a, b, c}$

Now, each line of **R2** extends to a unique line in NATURAL JOIN
 ... the join gives 2 lines, correctly

LOSSLESS DECOMPOSITIONS

Given a relation **R**

say we decompose it into two tables:

R1

R2

CLAIM: If

$R1 \cap R2 \longrightarrow R1$

or

$R1 \cap R2 \longrightarrow R2$

then $(R1 \text{ NATURAL JOIN } R2)$ is equal to **R**

each row of $(R1 \text{ NATURAL JOIN } R2)$ gives *at least* one row of **R**

to prove: each row of $(R1 \text{ NATURAL JOIN } R2)$ gives *exactly one* row of **R**

LOSSLESS DECOMPOSITIONS

Another way of looking at it:

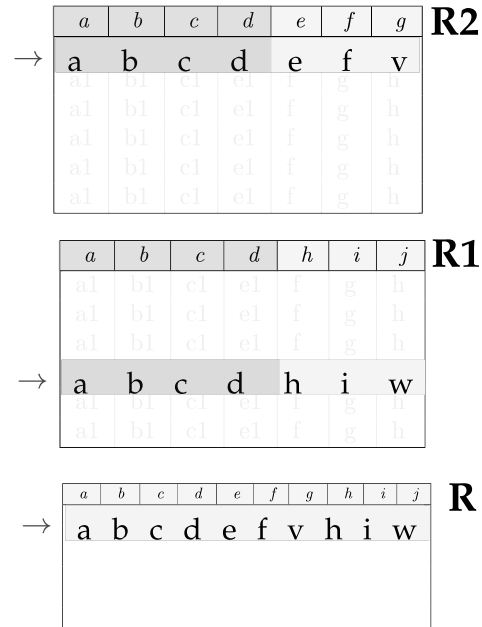
assume $R1 \cap R2 \longrightarrow R1$

consider the table **R2** line by line:

for each line of **R2**, exactly
one matching line of **R1**

these together give a line of **R**

So we cannot get anything 'extra' in **R**



LOSSLESS DECOMPOSITIONS

R

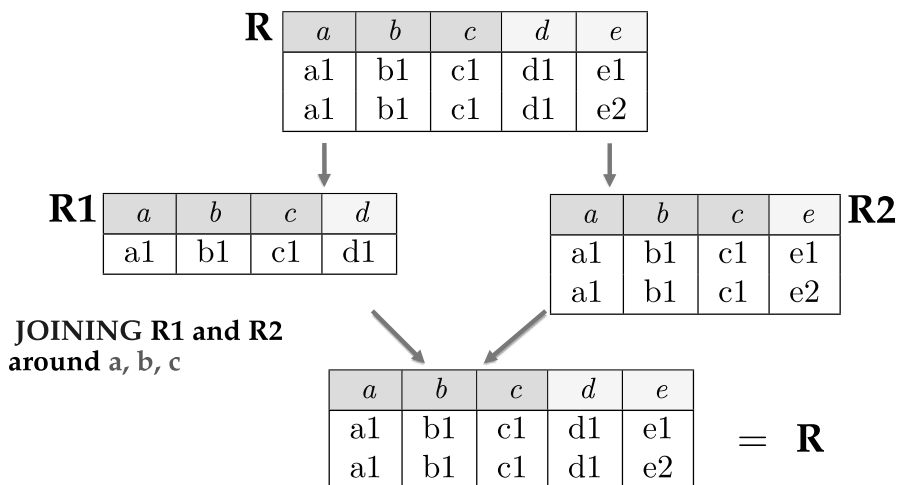
a	b	c	d	e
a1	b1	c1	d1	e1
a1	b1	c1	d2	e2

not possible with this dependency!

If we had this dependency

$a, b, c \longrightarrow d$

we get lossless decomposition!



LOSSLESS DECOMPOSITIONS

CLAIM: If

$$R1 \cap R2 \longrightarrow R1$$

or

$$R1 \cap R2 \longrightarrow R2$$

then $(R1 \text{ NATURAL JOIN } R2)$ is equal to R

Recall BCNF decomposition

If there is $\alpha \longrightarrow \beta$ in F^+ with

$(\beta - \alpha) \cap R$ not empty and $\alpha \subseteq R$ not a superkey of R

then

R1

$$\alpha \cup \beta$$

R2

$$R - \beta$$

Is this lossless? Note that $R1 \cap R2 = \alpha$ and $\alpha \longrightarrow R1$. Therefore,

\implies decompositions used in **BCNF** are lossless

new idea

3rd NORMAL FORM: DEFINITION

PRESERVING DEPENDENCIES

Given the relation:

$R(c, b, a, d)$

$a \rightarrow b$

$c, b \rightarrow a, d$

Is this in **BCNF** ? **NO**

$a \rightarrow b$

problem for R

$R1(c, a, d)$

$R2(a, b)$

BCNF decomposition:

$R2(a, b)$

$R1(c, a, d)$

Unpleasant side-effect: Difficult to verify the dependency

$c, b \rightarrow a, d$

in the two tables $R1$ and $R2$ without a join

PRESERVING DEPENDENCIES

Given a relation R and a set F of functional dependencies

we decompose R into tables $R1, \dots, Rk$

each table Ri has a list Fi of dependencies for it

$R(a, b, c, d, e, f)$

F $a, b \rightarrow f$

$d \rightarrow e, f$

$R1(a, b, c, d)$

$F1$

$R2(d, e, f)$

$F2$

$d \rightarrow e, f$

Over time, we insert more data into these k tables

each insert, we check new data satisfies Fi for table Ri

$d \rightarrow e, f$

for table $R2$

but unable easily to check dependencies *across* tables

$a, b \rightarrow f$

to avoid this problem, forced to check by **JOINS** *each* time we insert

3RD NORMAL FORMS

Given a relation \mathbf{R} and a set F of functional dependencies

\mathbf{R} is in BCNF if it has no redundancy with respect to F^+

\mathbf{R} is in BCNF if for all functional dependencies $\boxed{\alpha \longrightarrow \beta}$ in F^+
either $\boxed{\alpha \longrightarrow \beta}$ is *trivial* ... that is, $(\beta - \alpha) \cap \mathbf{R}$ is empty
or α is a superkey of \mathbf{R}

We now loosen the BCNF condition to allow some dependencies in \mathbf{R}

$\boxed{\alpha \longrightarrow \beta}$ can remain, even if α is not a superkey, **as long as**
each attribute in $\beta - \alpha$ is 'important'

3RD NORMAL FORMS

Given a relation \mathbf{R} and a set F of functional dependencies

\mathbf{R} is in 3NF if for all functional dependencies $\boxed{\alpha \longrightarrow \beta}$ in F^+
where $\alpha \subseteq \mathbf{R}$ and $\beta \subseteq \mathbf{R}$

either $\boxed{\alpha \longrightarrow \beta}$ is *trivial* ... that is, $(\beta - \alpha) \cap \mathbf{R}$ is empty
or α is a superkey of \mathbf{R}

or each attribute in $\beta - \alpha$ is part of some candidate key of \mathbf{R}



unintuitive condition, will become clear later!

note: all attributes in $\beta - \alpha$ need not be part of the *same* candidate key
different attributes can be part of different candidate keys

3RD NORMAL FORMS

Given the relations:

$R(c, b, a, d)$

$a \rightarrow b$

$c, b \rightarrow a, d$

Is this in **3NF** ? **YES**

$a \rightarrow b$

a problem for **R** for **3NF**?

Is b part of *some* candidate key for **R**? **yes**

$c, b \rightarrow a, d$

3RD NORMAL FORMS: ALGORITHM

Given a set of functional dependencies F

the algorithm for computing a **3NF** decomposition:

An Overview (details later)

Step 1. Simplify F to a 'minimal' list of dependencies F_c

Step 2. for each dependency $\alpha \rightarrow \beta$ in F_c , add relation

$\alpha \cup \beta$

Theorem: The above algorithm gives **3NF**, and preserves all dependencies!

new idea

EXTRANEOUS ATTRIBUTES

Step 1. Simplify F to a minimal list of dependencies F_c

remove all attributes from all dependencies in F that are *not necessary*
 \implies keep simplifying $F \rightarrow F'$ as long as $\text{closure}(F) = \text{closure}(F')$

F $a, c \rightarrow g$ $c, g \rightarrow b, d$ $a, c, d \rightarrow b$

$\downarrow ?$

F' $a, c \rightarrow g$ $c, g \rightarrow b, d$ $a, c \rightarrow b$

can do this if: $\text{closure}(F) = \text{closure}(F')$

F $a, c \rightarrow g$ $c, g \rightarrow b, d$ $a, c, d \rightarrow b$

$\downarrow ?$

F' $a, c \rightarrow g$ $c, g \rightarrow d$ $a, c \rightarrow b$

can do this if: $\text{closure}(F) = \text{closure}(F')$

Step 1. Simplify F to a minimal list of dependencies F_c

called extraneous

$$a, \alpha \longrightarrow \beta$$



Q: when can a be removed from this dependency ?

A: when β can be derived from just α in F

a can be removed if β is contained in the attribute closure of α in F

Example:

is this extraneous?



$$F \quad a, c \longrightarrow g \quad c, g \longrightarrow b, d \quad a, c, d \longrightarrow b$$

Yes

$$a, c \longrightarrow g \quad + \quad c, g \longrightarrow b, d \quad \implies \quad a, c \longrightarrow b$$

Step 1. Simplify F to a minimal list of dependencies F_c

$$\alpha \longrightarrow b, \beta$$

Q: when can b be removed from this dependency ?

A: when just the dependency $\alpha \longrightarrow \beta$ in F is sufficient to derive b

if the closure of α in $F - [\alpha \longrightarrow b, \beta] + [\alpha \longrightarrow \beta]$ contains b

Example:

is this extraneous?



$$F \quad a \longrightarrow b, c \quad b \longrightarrow c \quad a, b \longrightarrow c$$

Yes

$$a \longrightarrow b \quad b \longrightarrow c \quad \implies \quad a \longrightarrow c$$

```
def isAttributeExtraneous(F, alpha, beta, a):
    if a in beta:
        beta_prime = beta.difference({a})
        _F = list(F)
        _F.remove( [alpha, beta] )
        _F.append( [ alpha, beta_prime ] )
        if a in computeAttributeClosure(_F, alpha):
            return True, [ alpha, beta_prime ]
    if a in alpha:
        alpha_prime = alpha.difference({a})
        if beta.issubset( computeAttributeClosure(F, alpha_prime) ):
            return True, [ alpha_prime, beta ]
    return False, [ {}, {} ]
```

$$\alpha \longrightarrow b, \beta$$

To remove **b**, check if

the closure of α in $F - [\alpha \longrightarrow b, \beta] + [\alpha \longrightarrow \beta]$ contains **b**

$$a, \alpha \longrightarrow \beta$$

To remove **a**, check if

β is contained in the attribute closure of α in F

```
import itertools
def powerSet(inputset):
def computeAttributeClosure(F, X):
def isAttributeExtraneous(F, alpha, beta, A):
    if A in beta:
        beta_prime = beta.difference({A})
        _F = list(F)
        _F.remove( [alpha, beta] )
        _F.append( [ alpha, beta_prime ] )
        if A in computeAttributeClosure(_F, alpha):
            return True, [ alpha, beta_prime ]
    if A in alpha:
        alpha_prime = alpha.difference({A})
        if beta.issubset( computeAttributeClosure(F, alpha_prime) ):
            return True, [ alpha_prime, beta ]
    return False, [ {}, {} ]
d = [
    [ {'A'}, {'B', 'C'} ],
    [ {'B'}, {'C'} ],
    [ {'A', 'B'}, {'C'} ]
]
print( isAttributeExtraneous(d, *['A', 'B', 'C'], 'C') )
print( isAttributeExtraneous(d, *['A', 'B', 'C'], 'B') )
```



```
(True, [{'A'}, {'B'}])
(False, [{}, {}])
```

new idea

CANONICAL FORMS

ALGORITHM

Given a set of functional dependencies F

F_c is a **CANONICAL FORM** of F if :

- F_c is a minimal list of dependencies for F
- no two dependencies in F_c have same left side

In other words: from F , we want to construct F_c such that

$$\text{closure}(F) = \text{closure}(F_c)$$

no dependency in F_c has an extraneous attribute

no two dependencies in F_c have same left side

ALGORITHM

Given a set of functional dependencies F , algorithm to compute F_c

computeCanonicalCover (F)

$F_c = F$

While F_c keeps changing :

for each dependency $\alpha \longrightarrow \beta$ in F_c :

for each attribute a in $\alpha \cup \beta$:

if a is extraneous in $\alpha \longrightarrow \beta$:

delete it from $\alpha \longrightarrow \beta$

if there exist two dependencies $\alpha \longrightarrow \beta_1$ and $\alpha \longrightarrow \beta_2$ in F_c :

combine into one dependency $\alpha \longrightarrow \beta_1, \beta_2$

```
def computeCanonicalCover(F):
    OUT,changed = list(F),True
    while changed:
        changed = False
        for alpha_i, beta_i in OUT:
            for alpha_j, beta_j in OUT:
                if not changed and alpha_i == alpha_j and beta_i != beta_j:
                    OUT.append( [alpha_i, beta_i.union(beta_j)] )
                    OUT.remove( [alpha_i, beta_i] )
                    OUT.remove( [alpha_j, beta_j] )
                    changed = True
            for alpha, beta in OUT:
                for A in (alpha||beta):
                    to_update,[new_alpha,new_beta]=isAttributeExtraneous(OUT,*[alpha,beta],A)
                    if not changed and to_update:
                        if new_alpha != set() and new_beta != set():
                            OUT.append( [new_alpha, new_beta] )
                            OUT.remove( [alpha, beta] )
                            changed = True
                            break
    return OUT

d = [
    [ {'A'}, {'B'} ], #A->B
    [ {'A'}, {'C'} ], #A->C
    [ {'C', 'G'}, {'H'} ], #CG->H
    [ {'C', 'G'}, {'I'} ], #CG->I
    [ {'B'}, {'H'} ] #B->H
]
print( computeCanonicalCover(d) )
```



[[{'B'}, {'H'}], [{'A'}, {'B', 'C'}], [{'G', 'C'}, {'H', 'I'}]]

d	→	g, e
c, b	→	d
c, g	→	b
a, c	→	b
c, e	→	a

$a, c \longrightarrow g$ since can get from: $a, c \longrightarrow b$ $+$ $c, b \longrightarrow d$ $+$ $d \longrightarrow g, e$

$$c, b \longrightarrow d$$

$c, g \rightarrow b, d$ since can get from: $c, g \rightarrow b$ + $c, b \rightarrow d$

$a, c, d \longrightarrow b$ since can get from: $a, c \longrightarrow g + c, g \longrightarrow b, d$ $d \longrightarrow g, e$

$c, e \longrightarrow a, g$ since can get from: $c, e \longrightarrow a$ \vdash $a, c \longrightarrow b$ \vdash $c, b \longrightarrow d$

d	→	g, e
c, b	→	d
c, g	→	b
a, c	→	b
c, e	→	a

[[{'C': 'A'}, {'G'}], [{'D'}, {'E', 'G'}], [{'B', 'C'}, {'D'}], [{'C', 'G'}, {'D'}], [{'C', 'D'}, {'B'}], [{'E', 'C'}, {'A'}]]

```
[[{'C': 'A'}, {'G'}], [{'D'}, {'E'}, {'G'}], [{'B'}, {'C'}], [{'C'}, {'G'}], [{'C'}, {'D'}], [{"E"}, {"C"}], [{"A"}]]
```

Another example:

$a \longrightarrow b, c$ $b \longrightarrow a, c$ $c \longrightarrow a, b$

is this a valid canonical form?

$a \longrightarrow b$ $b \longrightarrow a, c$ $c \longrightarrow b$ **Yes**

is this a valid canonical form?

$a \longrightarrow c$ $c \longrightarrow b$ $b \longrightarrow a$ **Yes**

is this a valid canonical form?

$a \longrightarrow b$ $b \longrightarrow c$ $c \longrightarrow a$ **Yes**

is this a valid canonical form?

$a \longrightarrow c$ $b \longrightarrow c$ $c \longrightarrow a, b$ **Yes**

```
def isCanonicalCover(F, Fc):
    for alpha, beta in Fc:
        for A in alpha.union(beta):
            _isA_extra, X = isAttributeExtraneous(Fc, alpha, beta, A)
            if _isA_extra == True:
                return False

    R = set()
    for alpha, beta in F: R.update(alpha | beta)
    for K in powerSet(R):
        if computeAttributeClosure(F,K) != computeAttributeClosure(Fc,K):
            return False
    return True
```

new idea

3rd NORMAL FORM: ALGORITHM

FUNCTIONAL DEPENDENCY

Given a set of functional dependencies F

the algorithm for computing a **3NF** decomposition:

compute3NFDecomposition (F)

F_c = canonical cover of F

OUT = \emptyset

for each dependency $\boxed{\alpha \rightarrow \beta}$ in F_c :

if $\alpha \cup \beta$ not already part of some relation in OUT:

OUT += $\alpha \cup \beta$

Add a relation with a candidate key of all attributes (if required)

↑
to ensure that the decomposition is lossless

Theorem: The following algorithm gives **3NF** with **no** lost dependencies

```
def computeCanonicalCover(F):
def compute3NFDecomposition(F):
    OUT = []
    F_c = computeCanonicalCover(F)
    for alpha, beta in F_c:
        if ( [ R for R in OUT if (alpha).issubset( R ) ] == list() ):
            OUT.append( alpha beta )
r = [ {'A', 'B', 'C', 'D'}, {'A', 'B'} ]
d = [
    [ {'A'}, {'B'} ], #A->B
    [ {'C', 'B'}, {'A', 'D'} ] #CB->AD
]
print( computeBCNFDecomposition(d, r) )
print( compute3NFDecomposition(d) )
```



BCNF: [{'A', 'B'}, {'A', 'C', 'D'}]
 3NF: [{'A', 'B'}, {'A', 'C', 'D', 'B'}]

Theorem: The following algorithm gives **3NF** with **no** lost dependencies

Proof:

Given a set of dependencies F_c in canonical form

for each dependency $\alpha \longrightarrow \beta$ in F_c , we added the relation $\alpha \cup \beta$

no lost dependencies—added a relation for each dependency!

hard part: show that $\alpha \cup \beta$ is in 3NF

Have to show that it cannot happen that there is a table that is not in 3NF:

table	$\alpha \cup \beta$	from	$\alpha \longrightarrow \beta$	in F_c
with	$\gamma \longrightarrow a$ derivable from F_c $\gamma \cup \{a\} \subseteq \alpha \cup \beta$ yet γ is not a superkey of $\alpha \cup \beta$ and a not in any candidate key of $\alpha \cup \beta$			

$$\alpha \longrightarrow \beta$$

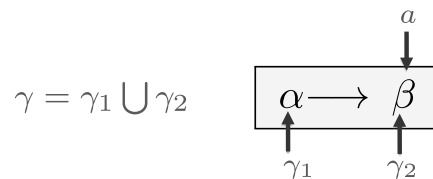
with

$\gamma \longrightarrow a$ derivable from F_c

$\gamma \cup \{a\} \subseteq \alpha \cup \beta$

yet γ is not a superkey of $\alpha \cup \beta$

and a not in any candidate key of $\alpha \cup \beta$



Key question: Did we derive $\gamma \longrightarrow a$ from F_c using $\alpha \longrightarrow \beta$?

yes \implies γ is a superkey in $\alpha \cup \beta$, since then all α is derivable from γ

no \implies a is redundant in $\alpha \longrightarrow \beta$, since knowing α , get γ_2 + knowing γ , get a

either way, we get a contradiction to our starting assumptions

$$\alpha \longrightarrow \beta$$

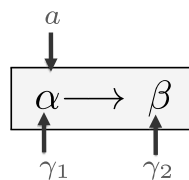
with

$\gamma \longrightarrow a$ derivable from F_c

$\gamma \cup \{a\} \subseteq \alpha \cup \beta$

yet γ is not a superkey of $\alpha \cup \beta$

and a not in any candidate key of $\alpha \cup \beta$



but then a is part of the candidate key α , so satisfies 3NF property

End of proof.

BCNF vs 3NF

In a perfect world, we want three things:

1. **No redundancy:** no non-trivial non-superkey dependency in any table
2. **Dependency preservation:** each dependency verifiable in some table
3. **Losslessness:** recover original table from JOINS of smaller tables



. . . but not always possible!

If we cannot have all three, then either:

- | |
|--|
| 3NF: satisfies 2. and 3., but can have redundancy <i>or</i>
BCNF: satisfies 1. and 3., but lose verifying some dependencies |
|--|

SUMMARY

1. Switched to a completely data dependencies point of view
2. Notion of functional dependency
3. Closure of functional dependencies
4. Closure of attributes
5. Boyce-Codd Normal Form (BCNF)
6. Lossless decompositions, BCNF is lossless
7. Dependency preservation when splitting tables
9. Canonical covers
8. Extraneous attributes
10. 3NF: definition and algorithm

**I gave you Python code
for each algorithm**

you should try examples
by running code.

test with your own
calculations and results.

The End