

Design Theory

SID	Class	Room	Time	Lat	Lng
4749732	cs 145	Nvidia Aud	T/R 4:30-6	37.4277° N	122.1742° W
2720942	cs 145	Nvidia Aud	T/R 4:30-6	37.4277° N	122.1742° W
4823984	cs 145	Nvidia Aud	T/R 4:30-6	37.4277° N	122.1742° W
4287594	cs 145	Nvidia Aud	T/R 4:30-6
2984994	cs 145	Nvidia Aud	T/R 4:30-6
8472374	cs 145	Nvidia Aud	T/R 4:30-6
4723663	cs 145	Nvidia Aud	T/R 4:30-6
2478239	cs 145	Nvidia Aud	T/R 4:30-6
4763268	cs 145	Nvidia Aud	T/R 4:30-6
2364532	cs 145	Nvidia Aud	T/R 4:30-6
2364573	cs 145	Nvidia Aud	T/R 4:30-6
3476382	cs 145	Nvidia Aud	T/R 4:30-6
2347623	cs 145	Nvidia Aud	T/R 4:30-6
...
2364579	cs 245	Nvidia Aud	T/R 3-4:30	37.4277° N	122.1742° W
3476343	cs 245	Nvidia Aud	T/R 3-4:30	37.4277° N	122.1742° W
2322232	cs 245	Nvidia Aud	T/R 3-4:30	37.4277° N	122.1742° W

~375
cs145
students

~300
cs245
students

Problems:

- Lots of data repetition.
- A single change (for example, Room change) needs a lots of update operations.
- Deletion causes unexpected data lost problem.
- Incomplete data insertion causes problem.

Here **Lat, Lng** are dependent on **Room**; **Room, Time** are dependent on **Class**.
That is, **Room** → { **Lat, Lng** } and **Class** → { **Room, Time** }

Updated version: More efficient solution if you decompose the table into 3 different tables based on the dependencies.

SID	Class
4749732	cs 145
2720942	cs 145
4823984	cs 145
4287594	cs 145
2984994	cs 145
8472374	cs 145
4723663	cs 145
2478239	cs 145
4763268	cs 145
2364532	cs 145
2364573	cs 145
3476382	cs 145
2347623	cs 145
...	...
2364579	cs 245
3476343	cs 245
2322232	cs 245

Class	Room	Time
cs 145	Nvidia Aud	T/R 4:30-6
cs 245	Nvidia Aud	T/R 3-4:30
cs 246	Nvidia Aud	M/W 3-4:30

Room	Lat	Lng
Nvidia Aud	37.4277° N	122.1742° W

~375
cs145
students

~300
cs245
students

Data Anomalies >>

- Problems that occur when we try to cram too much into a single relation are called anomalies.

1. **Redundancy:** Information may be repeated unnecessarily in several tuples.

Student	Course	Room
Mary	CS145	B01
Joe	CS145	B01
Sam	CS145	B01
..

If every course is in only one room, contains **redundant** information!

2. **Update Anomaly:** We may change information in one tuple but leave the same information unchanged in another.

Student	Course	Room
Mary	CS145	B01
Joe	CS145	C12
Sam	CS145	B01
..

If we update the room number for one tuple, we get inconsistent data = an **update anomaly**

3. **Delete Anomaly:** If a set of values get deleted, we may lose other information as a side effect.

Student	Course	Room
..

If everyone drops the class, we lose what room the class was in! = a **delete anomaly**

4. **Insert Anomaly:** We can't insert a new row because of some missing value whose value can't be null.

Student	Course	Room
Mary	CS145	B01
Joe	CS145	B01
Sam	CS145	B01
..

...	CS229	C12
-----	-------	-----

 →

Similarly, we can't reserve a room without students = an **insert anomaly**

After decomposition (without anomalies):

Student	Course
Mary	CS145
Joe	CS145
Sam	CS145
..	..

Course	Room
CS145	B01
CS229	C12

Is this form better?

- Any Redundancy?
- Any Update anomaly?
- Any Delete anomaly?
- Any Insert anomaly?

Normalization >>

Normalization is a systematic approach of **decomposing** tables to eliminate data redundancy (repetition) and undesirable characteristics like insert, update and delete anomalies.

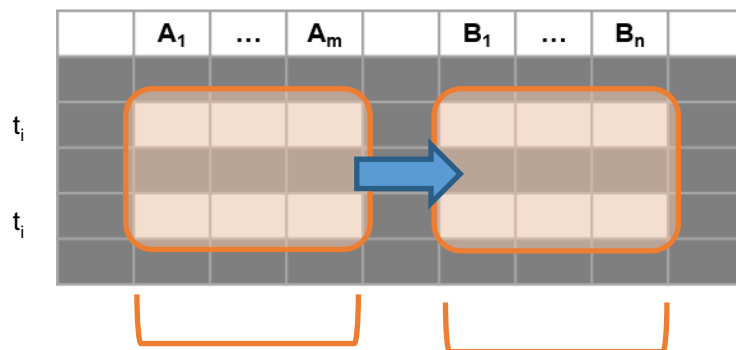
1.1) Functional Dependency:

- Let, $A = \{ A_1, A_2, \dots, A_m \}$ and $B = \{ B_1, B_2, \dots, B_n \}$ in R
- The functional dependency $A \rightarrow B$ on R holds if for any tuples t_i, t_j in R:
 $t_i[A] = t_j[A]$ implies $t_i[B] = t_j[B]$

that is whenever two or more tuples in R agree on all the attributes of A, they must also agree on all the attributes of B.



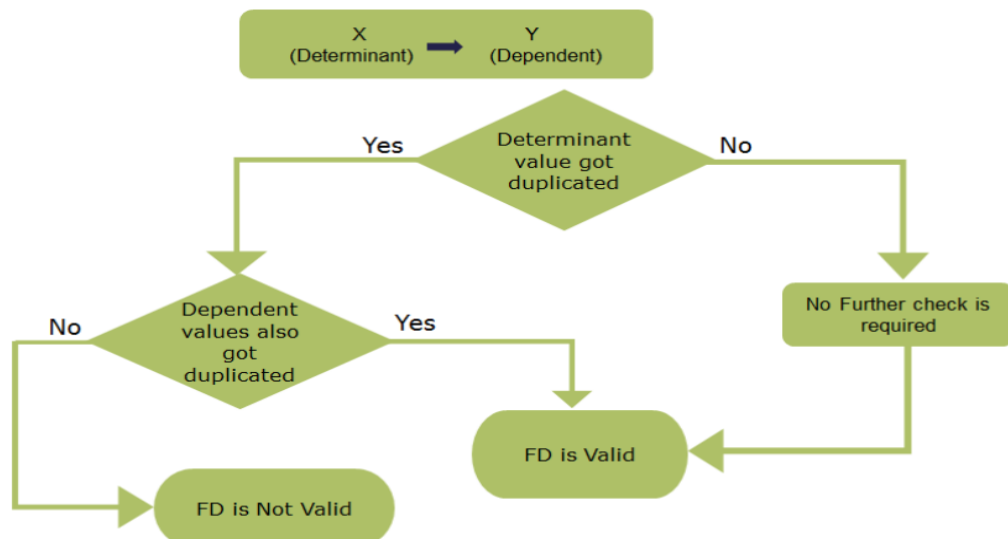
- if left side equals $t_i[A_1] = t_j[A_1], t_i[A_2]=t_j[A_2], \dots, t_i[A_m] = t_j[A_m]$
then right side also equals, $t_i[B_1] = t_j[B_1], t_i[B_2]=t_j[B_2], \dots, t_i[B_n] = t_j[B_n]$



If t_1, t_2 agree here.

...they also agree here!

- **Flow diagram:**



- FD is a constraint that *holds/does not hold* on an instance.
- A particular instance of R may coincidentally satisfy some FD but this FD may not hold for R in general.
- If the FD holds for every instances of Relation R, then FD becomes a part of the relational schema.
- *Example,*
 - {position} -> {phone} holds for this instance.

EmpID	Name	Phone	Position
E0045	Smith	1234	Clerk
E3542	Mike	9876 ←	Salesrep
E1111	Smith	9876 ←	Salesrep
E9999	Mary	1234	Lawyer

- {phone} -> {position} doesn't hold for this instance.

EmpID	Name	Phone	Position
E0045	Smith	1234 →	Clerk
E3542	Mike	9876	Salesrep
E1111	Smith	9876	Salesrep
E9999	Mary	1234 →	Lawyer

Practice:

A	B	C
1	2	3
2	2	3
3	2	3
4	3	2
5	2	3
6	3	2

A → A	Valid	AB → A	Valid
A → B	Valid	AB → B	Valid
A → C	Valid	AB → C	Valid
B → A	Invalid	BC → A	Invalid
B → B	Valid	BC → B	Valid
B → C	Valid	BC → C	Valid
C → A	Invalid	CA → A	Valid
C → B	Valid	CA → B	Valid
C → C	Valid	CA → C	Valid

- Why we need FDs in Database Design:
 - First we will start with some relational schema (received from ERD)
 - [Task 1]** Then we will find out its Functional Dependencies.
 - [Task 2]** Finally by using these FDs we will design a better schema that will minimize the possibility of anomalies.

1.2) Task 1 (Discover all FDs):

Armstrong's Axioms:

- Reflexivity rule:** If α is a set of attributes and $\beta \subseteq \alpha$, then $\alpha \rightarrow \beta$ holds.
Ex: $AB \rightarrow B$, here B is a subset of AB.
- Augmentation rule:** If $\alpha \rightarrow \beta$ holds and γ is a set of attributes, then $\gamma\alpha \rightarrow \gamma\beta$ holds.
Ex: if $AB \rightarrow C$ holds then, $ABD \rightarrow CD$ holds
- Transitivity rule:** If $\alpha \rightarrow \beta$ holds and $\beta \rightarrow \gamma$ holds, then $\alpha \rightarrow \gamma$ holds.
Ex: if $A \rightarrow B$ and $B \rightarrow C$, then $A \rightarrow C$

- Additional Rules:
 - i. **Union rule:** If $\alpha \rightarrow \beta$ holds and $\alpha \rightarrow \gamma$ holds, then $\alpha \rightarrow \beta\gamma$ holds.
 Ex: if $A \rightarrow B$ and $A \rightarrow C$, then $A \rightarrow BC$
 - ii. **Decomposition rule:** If $\alpha \rightarrow \beta\gamma$ holds, then $\alpha \rightarrow \beta$ holds and $\alpha \rightarrow \gamma$ holds.
 Ex: if $A \rightarrow BC$ then $A \rightarrow B$ and $A \rightarrow C$
 - iii. **Pseudo-transitivity rule:** If $\alpha \rightarrow \beta$ holds and $\gamma\beta \rightarrow \delta$ holds, then $\alpha\gamma \rightarrow \delta$ holds.
 Ex: if $A \rightarrow B$ holds and $CB \rightarrow D$ then $CA \rightarrow D$

- Let, $R = (A, B, C, G, H, I)$ and $F = \{A \rightarrow B, A \rightarrow C, CG \rightarrow H, CG \rightarrow I, B \rightarrow H\}$

Then:

- ▶ $A \rightarrow H$. Since $A \rightarrow B$ and $B \rightarrow H$ hold, we apply the transitivity rule.
- ▶ $CG \rightarrow HI$. Since $CG \rightarrow H$ and $CG \rightarrow I$, the union rule implies that $CG \rightarrow HI$
- ▶ $AG \rightarrow I$. Since $A \rightarrow C$ and $CG \rightarrow I$, the pseudo-transitivity rule implies that $AG \rightarrow I$ holds.

- Functional Dependency Closure: [out of syllabus]

```

F+ = F
repeat
  for each functional dependency f in F+
    apply reflexivity and augmentation rules on f
    add the resulting functional dependencies to F+
  for each pair of functional dependencies f1 and f2 in F+
    if f1 and f2 can be combined using transitivity
      add the resulting functional dependency to F+
until F+ does not change any further
  
```

Figure A procedure to compute F^+ .

- Example:

Let, $R = (A, B, C, D)$ and $F = \{A \rightarrow B, B \rightarrow C\}$

$F^+ = \{$

$A \rightarrow A, A \rightarrow B, A \rightarrow C, B \rightarrow B, B \rightarrow C, C \rightarrow C, D \rightarrow D, AB \rightarrow A,$
 $AB \rightarrow B, AB \rightarrow C, AC \rightarrow A, AC \rightarrow B, AC \rightarrow C, AD \rightarrow A, AD \rightarrow B,$
 $AD \rightarrow C, AD \rightarrow D, BC \rightarrow B, BC \rightarrow C, BD \rightarrow B, BD \rightarrow C, BD \rightarrow D,$
 $CD \rightarrow C, CD \rightarrow D, ABC \rightarrow A, ABC \rightarrow B, ABC \rightarrow C, ABD \rightarrow A,$
 $ABD \rightarrow B, ABD \rightarrow C, ABD \rightarrow D, BCD \rightarrow B, BCD \rightarrow C, BCD \rightarrow D,$
 $ABCD \rightarrow A, ABCD \rightarrow B, ABCD \rightarrow C, ABCD \rightarrow D.$

$\}$

- Inefficient process!!!!

- **Closure of Attribute Set:**

Let α be a set of attributes. We call the set of all attributes functionally determined by α under a set F of functional dependencies the **closure** of α under F . we denote it by α^+ .

Algorithm:

```

result :=  $\alpha$ ;
while (changes to result) do
    for each functional dependency  $\beta \rightarrow \gamma$  in  $F$  do
        begin
            if  $\beta \subseteq \textit{result}$  then  $\textit{result} := \textit{result} \cup \gamma$ ;
        end

```

Figure An algorithm to compute α^+ , the closure of α under F .

Example:

Let, $R = (A, B, C, G, H, I)$ and $F = \{A \rightarrow B, A \rightarrow C, CG \rightarrow H, CG \rightarrow I, B \rightarrow H\}$

Now, attribute closure of AG that is $(AG)^+$:

Initially, $(AG)^+ = AG$
 = AG B [using $A \rightarrow B$ rule, as A is a part of AG]
 = AGB C [using $A \rightarrow C$ rule, as A is a part of AGB]
 = AGBC H [using $CG \rightarrow H$ rule, as CG is a part of AGBC]
 = AGBCH I [using $CG \rightarrow I$ rule, as CG is a part of AGBCH]
 = AGBCHI [using $B \rightarrow H$ rule, as B is a part of AGBCHI, no change]
 = **AGBCHI** [no more check is needed as every FDs is checked]

Now, $(AB)^+ = AB$
 = AB [using $A \rightarrow B$ rule, as A is a part of AB, no change]
 = AB C [using $A \rightarrow C$ rule, as A is a part of AB]
 = ABC H [using $B \rightarrow H$ rule, as B is a part of ABC]
 = ABCH [couldn't use $CG \rightarrow H$ rule, as CG is not a part of ABCH]
 = ABCH [couldn't use $CH \rightarrow I$ rule, as CG is not a part of ABCH]
 = **ABCH** [no more changes is possible]

Practice:

If $R = (A, B, C, D, E)$ and $F = \{B \rightarrow AC, C \rightarrow AB, ABC \rightarrow D, BD \rightarrow A, AD \rightarrow C, E \rightarrow D\}$

- a) Find all the attribute closures with single element of R.
- b) Find all the attribute closures for all the sets with two attributes from relation R.

Uses:

- ▶ **Superkey check:**
 To test if α is a superkey, we compute α^+ , and check if α^+ contains all attributes of R. **Ex:** $(AG)^+ = ABCGHI$
- ▶ **FD validity checking:**
 We can check if a functional dependency $\alpha \rightarrow \beta$ holds (or, in other words, is in F^+), by checking if $\beta \subseteq \alpha^+$ **Ex.** $AG \rightarrow I$ is valid as $(AG)^+ = ABCGHI$

► **Determine all FDs: [No need]**

For each $\gamma \subseteq R$, we find the closure γ^+ , and for each $S \subseteq \gamma^+$, we output a functional dependency $\gamma \rightarrow S$.

1.3) Different types of Keys:

- **Superkey:**

Let R be a relation schema. A subset K of R is a **superkey** of R if, in any legal relation $r(R)$, for all pairs t_1 and t_2 of tuples in r such that $t_1 \neq t_2$, then $t_1[K] \neq t_2[K]$.

A set X of attributes in R is a superkey of R if and only if X^+ contains all attributes of R . In other words, X is a superkey if and only if it determines all other attributes.

- **Candidate key:**

X is a candidate key if and only if it is a superkey, but none of its proper subset is a superkey.

All candidate key finding **algorithm**:

Observation 1: any candidate key must contain attributes that have not appeared on the RHS of any functional dependency. (RHS keys are those keys that need help from others to be determined).

Observation 2: if an attribute has occurred on the RHS of some FD, but not on the LHS of any FD, then it cannot be in any candidate key. (These keys are determined by others and no other keys are dependent on them).

Final Algorithm:

- 1) Find all the attributes that have not appeared on the RHS of any FD. Denote this set by α
- 2) Denote the set of attributes that appear on the RHS of some FD, but not on the LHS of any FD by β
- 3) Compute the closure set α^+ , if $\alpha^+ = R$, then α is the **only candidate key**.
- 4) If $\alpha^+ \neq R$, then for each attribute x in $R - \beta$, test whether $\alpha \cup \{x\}$ is a candidate key. If not, try to add another attribute from $R - \beta$ to α and test whether it is candidate key.
- 5) Repeat step 4, until all candidate keys have been found.

Example 1:

If $R = (A, B, C, D, E)$ and $F = \{A \rightarrow C, CD \rightarrow B\}$

then, $\alpha = \{A, D, E\}, \beta = \{B\}$

Now $\alpha^+ = ABCDE = R$

So α is the only candidate key.

Example 2:

If $R = (A, B, C, D, E)$ and $F = \{A \rightarrow C, C \rightarrow BD, D \rightarrow A\}$

then, $\alpha = \{E\}, \beta = \{B\}$

Now $\alpha^+ = \{E\}$, not a superkey/candidate key. We will test each of $\{C, E\}, \{A, E\}, \{D, E\}$ next (not $\{B, E\}$).

$\{C, E\}^+ = \{C, E, B, D, A\}$. Therefore $\{C, E\}$ is a superkey. $\{C, E\}$ is also a candidate key since neither $\{E\}$ nor $\{C\}$ is a superkey.

$\{A, E\}^+ = \{A, E, C, B, D\}$. Similar to the above, $\{A, E\}$ is a candidate key.

Similarly we can verify $\{D, E\}$ is a candidate key.

Therefore $\{C, E\}$, $\{A, E\}$, $\{D, E\}$ are all of the candidate keys.

Practice 1:

If $R = (A, B, C, D, E)$ and $F = \{A \rightarrow BC, CD \rightarrow E, B \rightarrow D, E \rightarrow A\}$

- compute closure for each β in $\beta \rightarrow \gamma$ in F .
- List candidate keys of R .

Practice 2:

If $R = (A, B, C, D, E)$ and $F = \{A \rightarrow C, B \rightarrow D, AC \rightarrow D, CD \rightarrow E, E \rightarrow A\}$ then list the candidate keys of R

Practice 3:

If $R = (P, Q, R, S, T, U)$ and $F = \{PQ \rightarrow RTU, PR \rightarrow S, U \rightarrow P, R \rightarrow S, ST \rightarrow PU\}$ then list the candidate keys of R .

Practice 4:

If $R = (U, V, X, Y, Z)$ and $F = \{UV \rightarrow XZ, UX \rightarrow Y, X \rightarrow Y, VZ \rightarrow YX, Z \rightarrow UV\}$ then list the candidate keys of R .

1.4) Extraneous Attribute Detection:

An attribute of a functional dependency is said to be **extraneous** if we can remove it without changing the closure of the set of functional dependencies.

Let R be the relation schema, and let F be the given set of functional dependencies that hold on R . Consider an attribute A in a dependency $\alpha \rightarrow \beta$.

- If $A \in \beta$, to check if A is extraneous consider the set $F' = (F - \{\alpha \rightarrow \beta\}) \cup \{\alpha \rightarrow (\beta - A)\}$ and compute α^+ (the closure of α) under F' ; if α^+ includes A , then A is extraneous in β .

Example:

$F = \{AB \rightarrow CD, A \rightarrow E, E \rightarrow C\}$, check if C is extraneous in $AB \rightarrow CD$ or not?

formula, if $F = \{P \rightarrow QR, Q \rightarrow R\}$ then R is extraneous in $P \rightarrow QR$

- If $A \in \alpha$, to check if A is extraneous, let $\gamma = \alpha - \{A\}$, and compute γ^+ (the closure of γ) under F ; if γ^+ includes all attributes in β , then A is extraneous in α .

Example:

$F = \{P \rightarrow Q, PQ \rightarrow R\}$, check if Q is extraneous in $PQ \rightarrow R$?

1.5) Minimal Cover(No redundancy):

Given a set F of FDs, we say another set E of FDs is a minimal cover of F if

- ▶ Every FDs in E has a single attribute on the RHS.
- ▶ F and E are equivalent, that is, every FD in E can be inferred from the FDs in F , and every FD in F can be inferred from the FDs in E .
- ▶ Every FD $A \rightarrow b$ in E is minimal in its LHS, that is, there is no proper subset C of A such that $C \rightarrow b$
- ▶ There is no redundant FD in E . That is removing any FD from E will result in a set of FD that is not equivalent to F .

Algorithm:

Initially $E=F$

Step 1: rewrite each FD that has m attributes on the RHS into m FDs where the RHS is a single attribute.

Step 2: remove trivial FDs.

Step 3: minimize LHS of each FD. For each FD $X \rightarrow y$ in E , and for each attribute x in X , if $X-\{x\} \rightarrow y$ is implied by E , then replace $X \rightarrow y$ with $X-\{x\} \rightarrow y$.

Step 4: remove redundant FDs. For each FD in E , if it is implied by other FDs in E , then remove it from E .

Example:

If $R=(A, B, C, D, E, F)$ and $F=\{ABC \rightarrow CDEF, C \rightarrow E, A \rightarrow B, D \rightarrow F\}$

Step 1	Step 2	Step 3	Step 4
$ABC \rightarrow C$ $ABC \rightarrow D$ $ABC \rightarrow E$ $ABC \rightarrow F$ $C \rightarrow E$ $A \rightarrow B$ $D \rightarrow F$	$ABC \rightarrow C$ (cancel) $ABC \rightarrow D$ $ABC \rightarrow E$ $ABC \rightarrow F$ $C \rightarrow E$ $A \rightarrow B$ $D \rightarrow F$	$C \rightarrow E$ $A \rightarrow B$ $D \rightarrow F$ $AC \rightarrow D$ $AC \rightarrow F$ $ABC \rightarrow E$ (cancel)	$C \rightarrow E$ $A \rightarrow B$ $D \rightarrow F$ $AC \rightarrow D$ $AC \rightarrow F$ (cancel)

Final minimal cover, $F = \{AC \rightarrow D, C \rightarrow E, A \rightarrow B, D \rightarrow F\}$

Practices:

1. $F = \{AB \rightarrow CD, B \rightarrow C, BC \rightarrow D, CD \rightarrow EF, E \rightarrow F\}$. Find minimal cover for this FD set.

Solution: $F = \{B \rightarrow C, CD \rightarrow E, E \rightarrow F\}$ is minimal cover.

2. $F = \{A \rightarrow BC, CD \rightarrow E, E \rightarrow C, D \rightarrow AEH, ABH \rightarrow BD, DH \rightarrow BC\}$. Find minimal cover for this FD set.

Solution: $F = \{A \rightarrow BC, D \rightarrow AEH, AH \rightarrow D, E \rightarrow C\}$ is minimal cover

3. $F = \{AB \rightarrow C, C \rightarrow A, BC \rightarrow D, ACD \rightarrow B, D \rightarrow E, D \rightarrow G, BE \rightarrow C, CG \rightarrow B, CG \rightarrow D, CE \rightarrow A, CE \rightarrow G\}$

Solution 1: $\{AB \rightarrow C, C \rightarrow A, BC \rightarrow D, CD \rightarrow B, D \rightarrow E, D \rightarrow G, BE \rightarrow C, CG \rightarrow D, CE \rightarrow G\}$

Solution 2: $\{AB \rightarrow C, C \rightarrow A, BC \rightarrow D, D \rightarrow E, D \rightarrow G, BE \rightarrow C, CG \rightarrow B, CE \rightarrow G\}$