

# Design Theory and Normalization

CSC343

Mark Kazakevich

(with slides from Diane Horton, Jeff Ullman)

# What's the point?

- Most likely, you have some idea of how to make a schema for a data set
- But how do you know it's a good schema?
- How do you make it better?
- Is there an 'optimal' schema for your data?

# Database ‘Design’ Theory

- Improve your (probably bad) schema systematically
- Try to get a schema that is in a “**normal form**” that guarantees **good properties**.
  - “Normal” in the sense of conforming to a standard.
- The process of converting a schema to a normal form is called ***normalization***.

# Database 'Design' Theory

- **General idea:**
  - Express *constraints* on the relationships between attributes
  - Use these constraints to *decompose* the relations

# What are these 'constraints'?

- We've defined constraints before like keys, foreign keys, etc. based on our knowledge of the **problem domain**.
- We can define other domain-specific constraints on the data that can help us design better schemas
  - These constraints are called **functional dependencies**

- A **functional dependency** (FD) on a relation R is a statement of the form:

If two tuples of R **agree** on all of the attributes

$A_1, A_2, \dots, A_n$

then they must also **agree** on all of another list of attributes

$B_1, B_2, \dots, B_m$

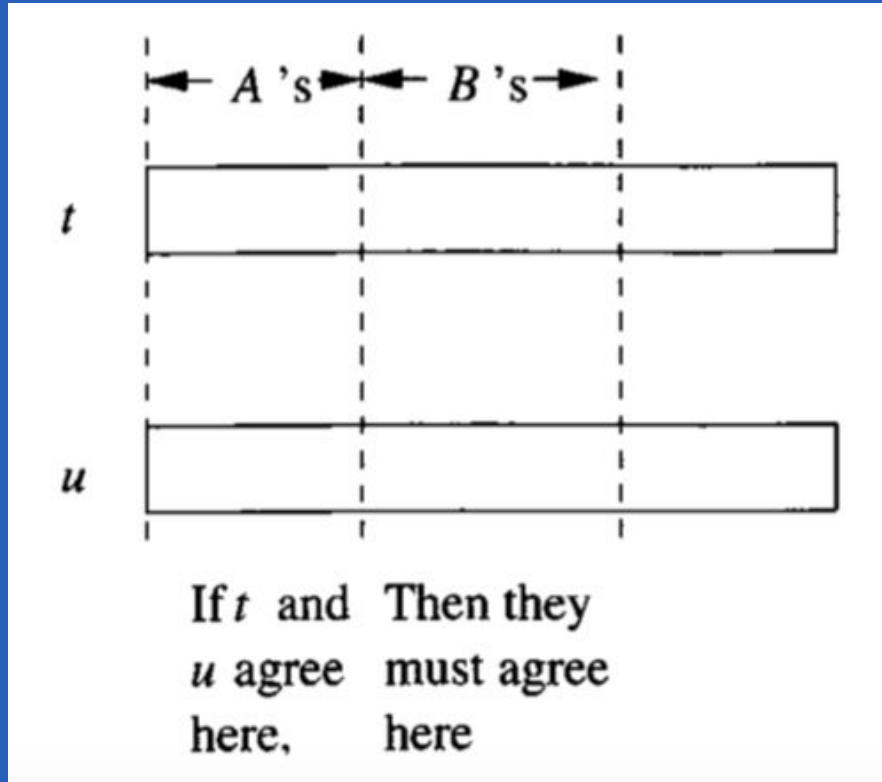
- We write this FD formally as

$A_1, A_2, \dots, A_n \longrightarrow B_1, B_2, \dots, B_m$

and we say:

“ $A_1, A_2, \dots, A_n$  functionally determines  $B_1, B_2, \dots, B_m$ ”

# Visual representation



**Note:** it is not a requirement that the attributes B have to follow A directly - they can be anywhere

# Example

<i>title</i>	<i>year</i>	<i>length</i>	<i>genre</i>	<i>studioName</i>	<i>starName</i>
Star Wars	1977	124	SciFi	Fox	Carrie Fisher
Star Wars	1977	124	SciFi	Fox	Mark Hamill
Star Wars	1977	124	SciFi	Fox	Harrison Ford
Gone With the Wind	1939	231	drama	MGM	Vivien Leigh
Wayne's World	1992	95	comedy	Paramount	Dana Carvey
Wayne's World	1992	95	comedy	Paramount	Mike Meyers

What is a functional dependency here?  
(think about the domain)



# Example

<i>title</i>	<i>year</i>	<i>length</i>	<i>genre</i>	<i>studioName</i>	<i>starName</i>
Star Wars	1977	124	SciFi	Fox	Carrie Fisher
Star Wars	1977	124	SciFi	Fox	Mark Hamill
Star Wars	1977	124	SciFi	Fox	Harrison Ford
Gone With the Wind	1939	231	drama	MGM	Vivien Leigh
Wayne's World	1992	95	comedy	Paramount	Dana Carvey
Wayne's World	1992	95	comedy	Paramount	Mike Meyers

title year -> length genre studioName

*What assumptions did we make?*

# Example

<i>title</i>	<i>year</i>	<i>length</i>	<i>genre</i>	<i>studioName</i>	<i>starName</i>
Star Wars	1977	124	SciFi	Fox	Carrie Fisher
Star Wars	1977	124	SciFi	Fox	Mark Hamill
Star Wars	1977	124	SciFi	Fox	Harrison Ford
Gone With the Wind	1939	231	drama	MGM	Vivien Leigh
Wayne's World	1992	95	comedy	Paramount	Dana Carvey
Wayne's World	1992	95	comedy	Paramount	Mike Meyers

What is **not** a functional dependency here?

# Coincidence or FD?

- An FD is an assertion about **every** instance of the relation.
- You can't know it holds just by looking at one instance.
- You must use **knowledge of the domain** to determine whether an FD holds.

# FDs are closely related to **keys**

- Suppose  $K$  is a set of attributes for relation  $R$ .
- Our old definition of **superkey**:
  - *a set of attributes for which no two rows can have the same values.*
- Using FDs
  - $K$  is a **superkey** for  $R$  iff  
 $K$  functionally determines **all of  $R$** .

# FDs are a generalization of keys



key:

$X \rightarrow R$  ← Every attribute



Functional dependency:

$X \rightarrow Y$  ← Not necessarily every attribute



An FD can be more subtle.

# FD Properties

# Trivial FDs

- Suppose  $A_1, A_2 \dots A_n \rightarrow B_1, B_2 \dots B_n$  holds
- Trivial FDs are those where the right side is a subset of the left side, that is,

$$\{B_1, B_2 \dots B_n\} \subseteq \{A_1, A_2 \dots A_n\}$$

Examples:

- $\text{title} \rightarrow \text{title}$
- $\text{title year} \rightarrow \text{title}$
- $\text{title year} \rightarrow \{\}$

# Splitting Rules

- An FD such as  $A \rightarrow BCD$  can be split up into a set of **equivalent** FDs:
  - $A \rightarrow B$
  - $A \rightarrow C$
  - $A \rightarrow D$
- We are 'splitting' the Right Hand Side of the FD (can also combine the three into the original)
  - Can we split the left hand side?
    - No.



# Inferring FDs

- Given a set of FDs, we can often **infer** further FDs.
- This will be handy when we apply FDs to the problem of database design.
- Big task: given a set of FDs,
  - infer every other FD that must also hold.
- Simpler task: given a set of FDs,
  - check whether a given FD must also hold.

# What we'd like to be able to find out:

- If  $A \rightarrow B$  and  $B \rightarrow C$  hold,  
must  $A \rightarrow C$  hold?
- If  $A \rightarrow H$ ,  $C \rightarrow F$ , and  $FG \rightarrow AD$  hold,  
must  $FA \rightarrow D$  hold?  
must  $CG \rightarrow FH$  hold?
- If  $H \rightarrow GD$ ,  $HD \rightarrow CE$ , and  $BD \rightarrow A$  hold,  
must  $EH \rightarrow C$  hold?
- Note: we are not generating new FDs,  
but **testing** a specific possible one.

# How do we figure them out?

- **Method 1** (first principles):
  - Refer to FDs you know and use definition of functional dependency
  - This can get complicated
- **Method 2**
  - The **Closure Test**
  - Systematic, and much easier!

# Closure $\{ \}^+$

- Suppose  $\{A_1, A_2, \dots, A_n\}$  is a set of attributes and  $S$  is a set of FDs
- The **closure** of  $\{A_1, A_2, \dots, A_n\}$  under the FD's in  $S$  is the set of attributes  $B$  such that:
  - Every relation that satisfies all the FD's in set  $S$  also satisfies  $A_1A_2 \dots A_n \rightarrow B$ .
- We denote the closure of a set of attributes  $A_1A_2 \dots A_n$  by  $\{A_1, A_2, \dots, A_n\}^+$ .

# Closure Algorithm

- How do we find  $Y^+$ , the closure of a set of attributes  $Y$  under a set of FDs  $S$ .

Attribute\_closure( $Y, S$ ):

Initialize  $Y^+$  to  $Y$ ; split RHS's of FDs if necessary

Repeat until no more changes occur:

    If there is an FD  $LHS \rightarrow RHS$  in  $S$   
    such that  $LHS$  is in  $Y^+$ :

        Add  $RHS$  to  $Y^+$

Return  $Y^+$

# Example

- Relation **R**(A, B, C, D, E, F)
- FDs:  $AB \rightarrow C$ ,  $BC \rightarrow AD$ ,  $D \rightarrow E$ ,  $CF \rightarrow B$
- What is  **$\{A, B\}^+$**  (the closure of  $\{A, B\}$ )

$X = \{A, B\}$ ; split  $BC \rightarrow AD$  to  $BC \rightarrow A$ ,  $BC \rightarrow D$

LHS of  $AB \rightarrow C$  in  $X$ , so  $X = \{A, B, C\}$

LHS of  $BC \rightarrow A$  and  $BC \rightarrow D$  in  $X$ , so  $X = \{A, B, C, D\}$

LHS of  $D \rightarrow E$  in  $X$ , so  $X = \{A, B, C, D, E\}$

(can not use  $CF \rightarrow B$ , since  $F$  not in  $X$ )

So we're done,  $\{A, B\}^+ = \{A, B, C, D, E\}$

# Closure Test

- $S$  is a set of FDs;  $LHS \rightarrow RHS$  is a single FD.
  - Return true iff  $LHS \rightarrow RHS$  follows from  $S$ .

```
FD_follows(S, LHS → RHS):  
    Y+ = Attribute_closure(LHS, S)  
    return (RHS is in Y+)
```

# Property that follows:

- Transitive Rule
  - $(A \rightarrow B \text{ and } B \rightarrow C) \text{ implies } A \rightarrow C$
  - Check this using closure test
    - C would appear in  $A^+$



# Projecting FDs

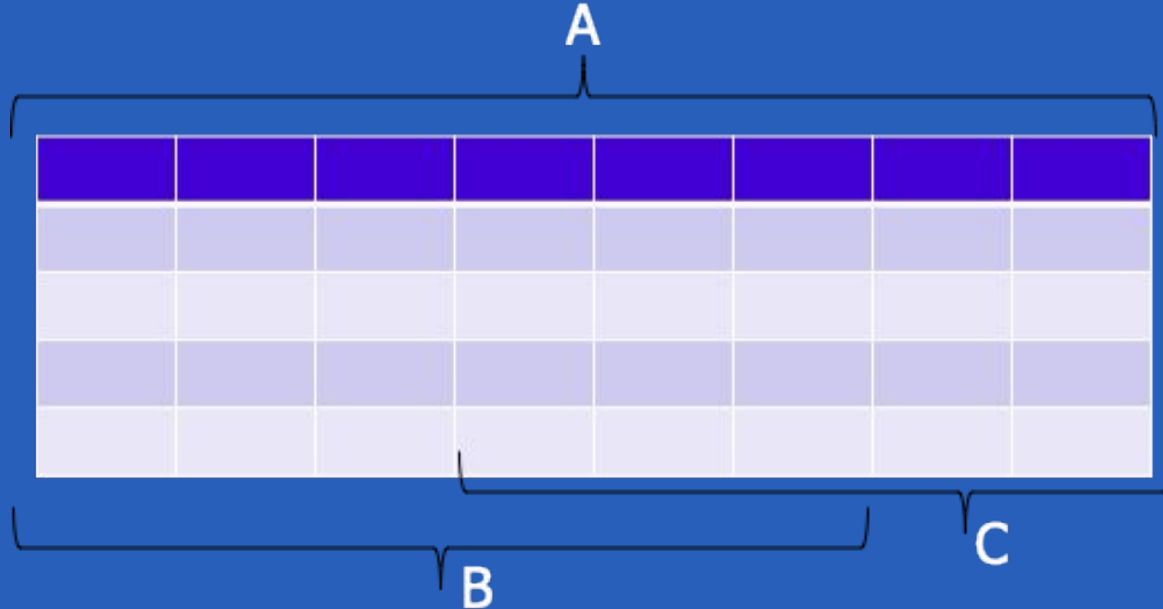
# Projecting FDs

- Later, we will learn how to ***normalize*** a schema by ***decomposing*** relations into smaller ones
  - This is the whole point of this theory.
- We will need to know what FDs hold in the new, smaller relations.
  - We must **project** our FDs onto the attributes of our new relations.

# Projecting FDs

Relation  $\mathbf{R}(A1 \dots An)$ , Set of attributes:  $A$

Sets of attributes  $B$  and  $C$  such that  $B \cup C = A$



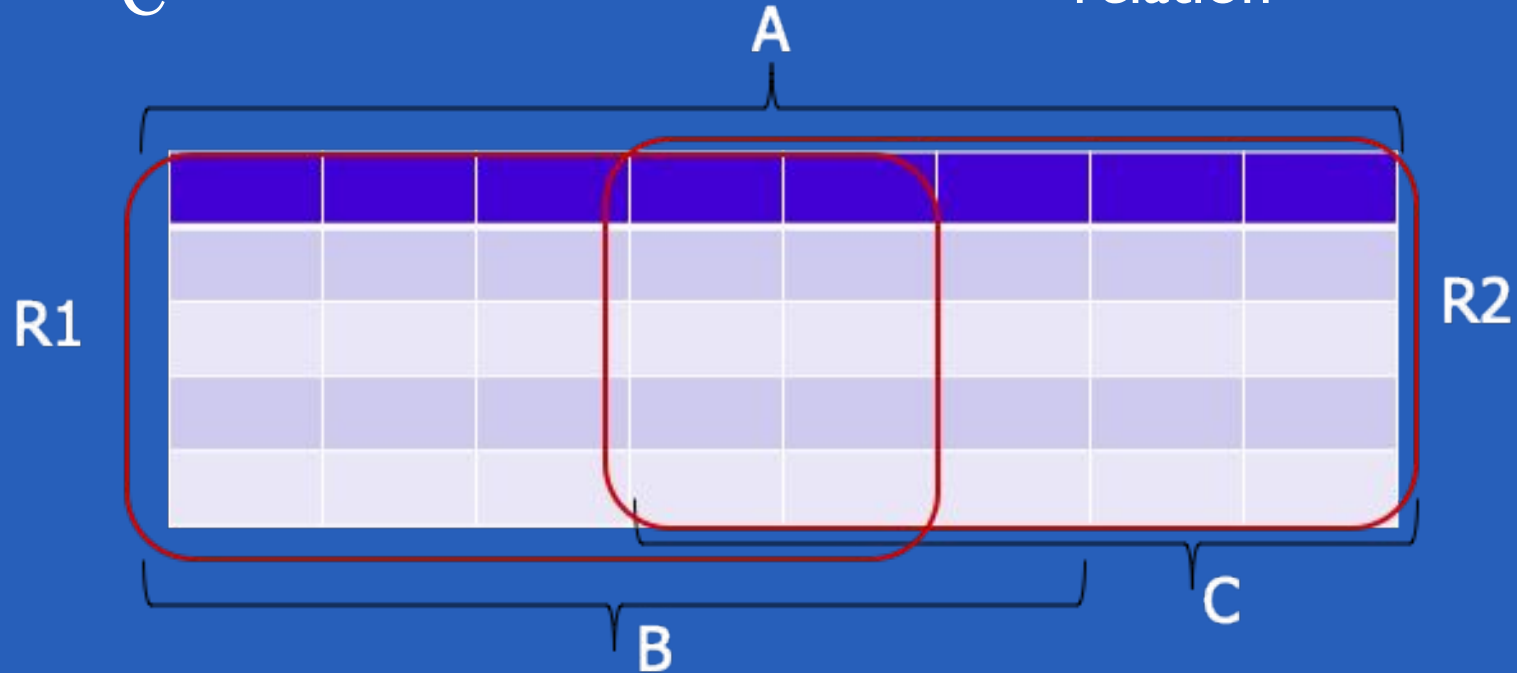
# Projecting FDs

$$R_1 = \pi_B(R)$$

$$R_2 = \pi_C(R)$$

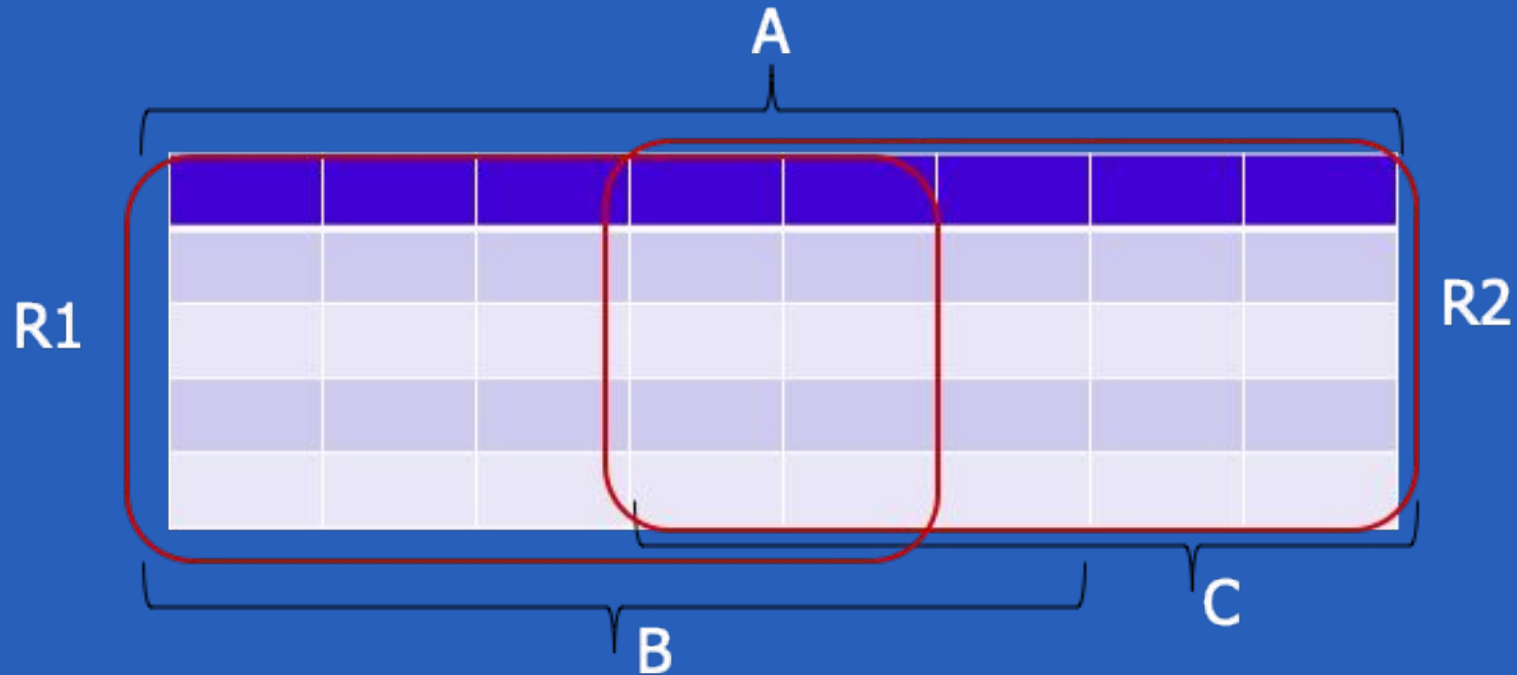
$$R_1 \bowtie R_2 = R$$

**Question:** What FDs hold in the new relation



# Projecting FDs

**Question:** If a set of FDs  $S$  hold in  $R$ ,  
what FDs hold in the new relations  $R_1$  and  $R_2$ ?



# Projection Algorithm

- **S** is a set of FDs; **L** is a set of attributes.
- Return the projection of **S** onto **L**:  
all FDs that follow from **S** and involve only attributes from **L**.

**Project(S, L):**

Initialize T to {}.

For each subset X of L:

    Compute  $X^+$       # Close X and see what we get.

For every attribute A in  $X^+$ :

    If A is in L: #  $X \rightarrow A$  is only relevant if A is in  
                    # L (we know X is).

        add  $X \rightarrow A$  to T.

Return T.

# Some algorithm speed-ups

- No need to add  $X \rightarrow A$  if  $A$  is in  $X$  itself. It's a trivial FD.
- These subsets of  $X$  won't yield anything, so no need to compute their closures:
  - the empty set
  - the set of all attributes in  $X$

Neither are big savings, but ...

# A big speed-up

- If we find  $X^+ = \text{all attributes}$ , we can ignore any superset of  $X$ .
- It can only give use “weaker” FDs (with more on the LHS).

This is a big time saver!



# Minimal basis

- We saw earlier that we can very often rewrite sets of FDs in equivalent ways.
  - Example:  $S1 = \{A \rightarrow BC\}$  is equivalent to  $S2 = \{A \rightarrow B, A \rightarrow C\}$ .
- Given a set of FDs  $S$ , we may want to find a **minimal basis**: A set of FDs that is equivalent, but has
  - no redundant FDs, and
  - no FDs with unnecessary attributes on the LHS.
  - All RHS will be single attributes

# Minimal Basis

- $S$  is a set of FDs. Return a **minimal basis** for  $S$ .

**Minimal\_basis( $S$ ):**

1. Split the RHS of each FD
2. For each FD  $X \rightarrow Y$  where  $|X| \geq 2$ :  
If you can remove an attribute from  $X$   
(use closures!) and get an FD that  
follows from  $S$ :  
Do so! (It's a stronger FD.)
3. For each FD  $f$  :  
If  $S - \{f\}$  implies  $f$  (use closures!) :  
Remove  $f$  from  $S$ .

# Some comments on Minimal Basis

- Often there are multiple possible results.
  - Depends on the order in which you consider the possible simplifications.
- After you identify a redundant FD, you must not use it when computing subsequent closures.

and less intuitive..

When you are computing closures to decide whether the LHS of an FD

$X \rightarrow Y$

can be simplified, continue to use that FD.