

# Introduction to Data Management

## Design Theory

Alyssa Pittman

Based on slides by Jonathan Leang, Dan Suciu, et al

Paul G. Allen School of Computer Science and Engineering  
University of Washington, Seattle

# Recap: Other Constraints

- CHECK (condition)
  - Single attribute
  - Single tuples

```
CREATE TABLE User (  
    uid INT PRIMARY KEY,  
    firstName TEXT,  
    lastName TEXT,  
    age INT CHECK (age > 12 AND age < 120),  
    email TEXT,  
    phone TEXT,  
    CHECK (email IS NOT NULL OR phone IS NOT NULL)  
);
```

# Assertions

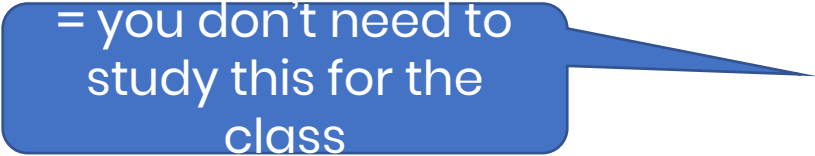
- Hard to support
- Usually impractical
- Usually not supported
  - Simulated with triggers

```
CREATE ASSERTION myAssert CHECK
  (NOT EXISTS (
    SELECT Product.name
      FROM Product, Purchase
     WHERE Product.name = Purchase.prodName
    GROUP BY Product.name
   HAVING count(*) > 200));
```

# Triggers

- Triggers activate on a specified event

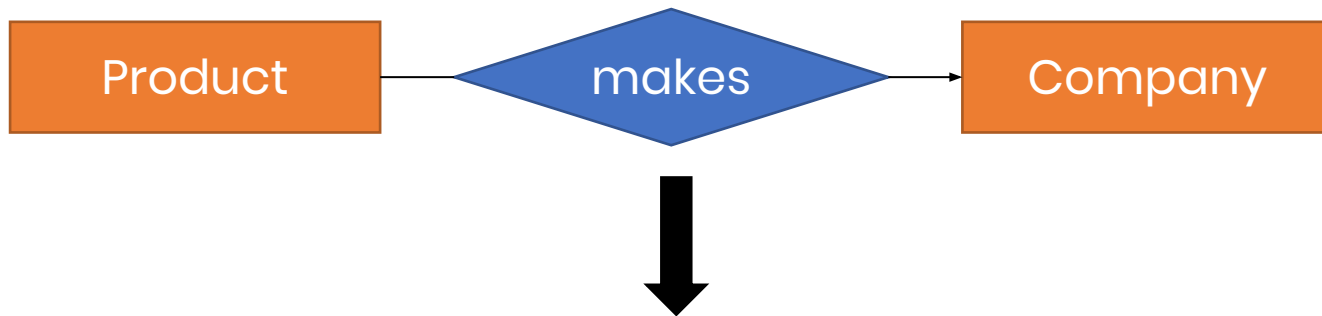
```
CREATE TRIGGER LowCredit ON Purchasing.PurchaseOrderHeader
AFTER INSERT AS
    IF (ROWCOUNT_BIG() = 0) RETURN;
    IF EXISTS (SELECT *
               FROM Purchasing.PurchaseOrderHeader AS p
               JOIN inserted AS i
               ON p.PurchaseOrderID = i.PurchaseOrderID
               JOIN Purchasing.Vendor AS v
               ON v.BusinessEntityID = p.VendorID
               WHERE v.CreditRating = 5
              )
    BEGIN
        RAISERROR ('A vendor''s credit rating is too
                    low to accept new purchase orders.', 16, 1);
        ROLLBACK TRANSACTION;
        RETURN
    END;
GO
```



# Recap

## ▪ ER Diagrams

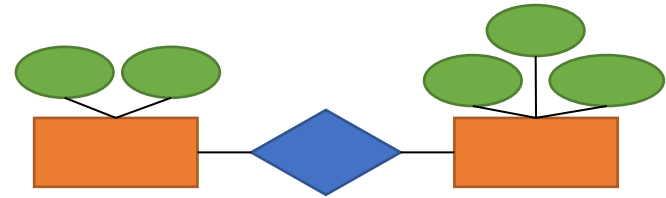
- Conceptual modeling
- Rules of thumb for converting diagram into schema



```
CREATE TABLE Company (  
    name VARCHAR(100) PRIMARY  
    KEY,  
    ...);  
CREATE TABLE Product (  
    name VARCHAR(100) PRIMARY  
    KEY,  
    cname VARCHAR(100)  
    REFERENCES Company  
    ...);
```

# The Database Design Process

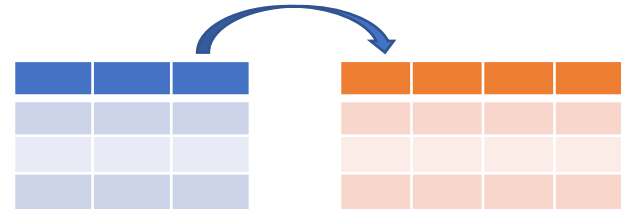
## Conceptual Model



## Relational Model

□ + Schema

□ + Constraints



Today

## Conceptual Schema

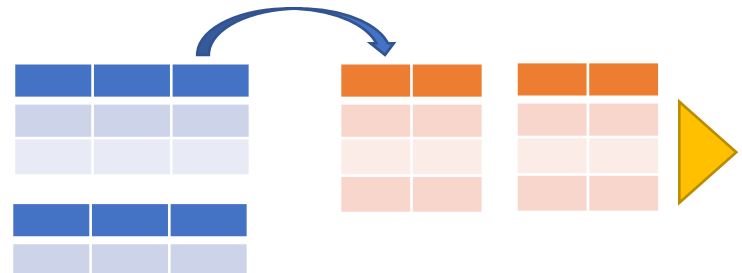
□ + Normalization



## Physical Schema

□ + Partitioning

□ + Indexing



# Goals for Today

- Figure out the fundamentals of what makes a good schema

# Outline

- Background
  - Anomalies, i.e. things we want to avoid
  - Functional Dependencies (FDs)
  - Closures and formal definitions of keys
- Normalization: BCNF Decomposition
- Losslessness



# Think About This



Make a simple directory that can:

- Hold information about name, SSN, phone, and city
- Associate **people** with the **city** they live in
- Associate **people** with any **phone numbers** they have

Name	SSN	Phone	City
Fred	123-45-6789	206-555-9999	Seattle
Fred	123-45-6789	206-555-8888	Seattle
Joe	987-65-4321	415-555-7777	San Francisco

The above instance does the job, but are there issues?

# Think About This



Make a simple directory that can:

- Hold information about name, SSN, phone, and city
- Associate **people** with the **city** they live in
- Associate **people** with any **phone numbers** they have

Name	SSN	Phone	City
Fred	123-45-6789	206-555-9999	Seattle
Fred	123-45-6789	206-555-8888	Seattle
Joe	987-65-4321	415-555-7777	San Francisco

Anomalies:

- **Redundancy** □ **Slow Update**
  - Change Fred's city to Bellevue (two rows!)
- **Deletion Anomalies**
  - How to delete Joe's phone without deleting Joe?

# Think About This



Make a simple directory that can:

- Hold information about name, SSN, phone, and city
- Associate **people** with the **city** they live in
- Associate **people** with any **phone numbers** they have

Name	SSN	Phone	City
Fred	123-45-6789	206-555-9999	Seattle
Fred	123-45-6789	206-555-8888	Seattle
Joe	987-65-4321	415-555-7777	San Francisco

Anomalies:

- **Redundancy** □ **Slow Update**
  - Change Fred's city to Bellevue (two rows!)
- **Deletion Anomalies**
  - How to delete Joe's phone without deleting Joe?

# Think About This



Make a simple directory that can:

- Hold information about name, SSN, phone, and city
- Associate **people** with the **city** they live in
- Associate **people** with any **phone numbers** they have

Name	SSN	Phone	City
Fred	123-45-6789	206-555-9999	Seattle
Fred	123-45-6789	206-555-8888	Seattle
Joe	987-65-4321	415-555-7777	San Francisco

Anomalies:

- **Redundancy** □ **Slow Update**
  - Change Fred's city to Bellevue (two rows!)
- **Deletion Anomalies**
  - How to delete Joe's phone without deleting Joe?

# Think About This



We can solve the anomalies by converting this

Name	SSN	Phone	City
Fred	123-45-6789	206-555-9999	Seattle
Fred	123-45-6789	206-555-8888	Seattle
Joe	987-65-4321	415-555-7777	San Francisco

into this

Name	SSN	City
Fred	123-45-6789	Seattle
Joe	987-65-4321	San Francisco

SSN	Phone
123-45-6789	206-555-9999
123-45-6789	206-555-8888
987-65-4321	415-555-7777

**How can we systematically avoid anomalies?**

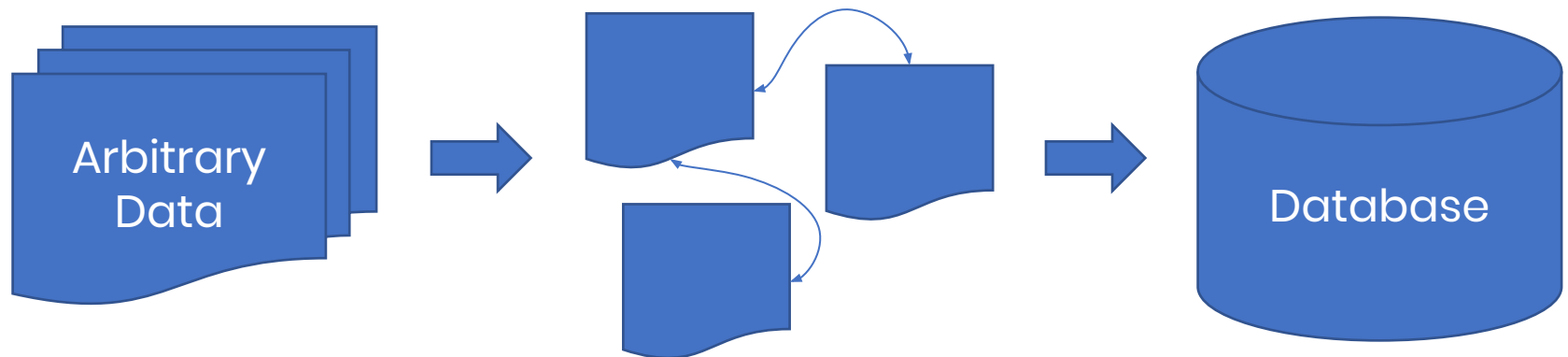
# Informal Design Guidelines

- Semantics of attributes should be self-evident
- Avoid redundant information in tuples
- Avoid NULL values in tuples
- Disallow the generation of “spurious” tuples
  - If certain tuples shouldn’t exist, don’t allow them

# Database Design

## Database Design

**Database Design** or **Logical Design** or **Relational Schema Design** is the process of organizing data into a database model. This is done by considering **what data needs to be stored** and the **interrelationship of the data**.



# Database Design

Database Design is about  
(1) characterizing data and (2) organizing data



# Database Design

Database Design is about  
(1) characterizing data and (2) organizing data

# Database Design

Database Design is about  
(1) characterizing data and (2) organizing data

How to talk about properties  
we know or see in the data

# Data Interrelationships

How do we start talking about data interrelationships?

- What rules govern our data?
  - Domain knowledge
    - Dimension vs measure
  - Pattern analysis

# Data Interrelationships

How do we start talking about data interrelationships?

- What rules govern our data?
  - Domain knowledge
    - Dimension vs measure
  - Pattern analysis

The rules that are known to us since we **made them up** or they correlate to **things in the real world**



# Data Interrelationships

How do we start talking about data interrelationships?

- What rules govern our data?

- Domain knowledge
  - Dimension vs measure
- Pattern analysis

The rules that are known to us since we **made them up** or they correlate to **things in the real world**



[ex] An engineer knows that a plane model determines the plane's wingspan

# Data Interrelationships

How do we start talking about data interrelationships?

- What rules govern our data?
  - Domain knowledge
    - Dimension vs measure
  - Pattern analysis

# Data Interrelationships

How do we start talking about data interrelationships?

- What rules govern our data?
  - Domain knowledge
    - Dimension vs measure
  - Pattern analysis

Rules that are found  
by finding  
correlations within  
the given data



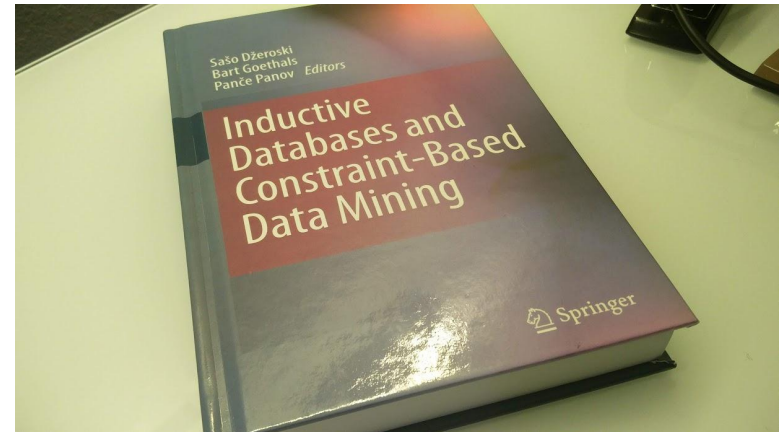
# Data Interrelationships

How do we start talking about data interrelationships?

- What rules govern our data?

- Domain knowledge
  - Dimension vs measure
- Pattern analysis

Rules that are found  
by finding  
correlations within  
the given data



- Data mining
- Knowledge Discovery in Databases (KDD)



# Data Interrelationships

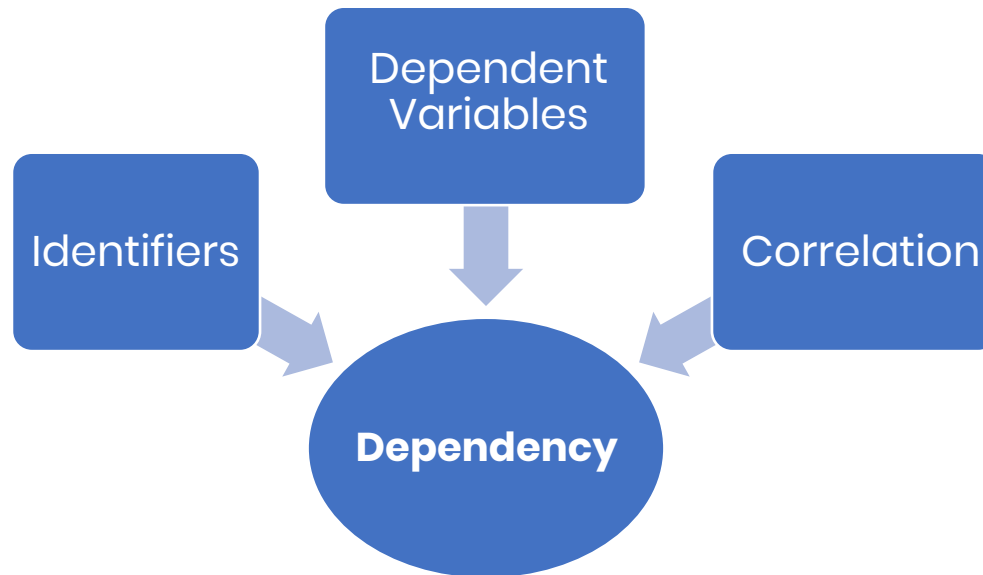
How do we start talking about data interrelationships?

- What **rules** govern our data?
  - Domain knowledge
    - Dimension vs measure
  - Pattern analysis

# Data Interrelationships

How do we start talking about data interrelationships?

- What **rules** govern our data?
  - Domain knowledge
    - Dimension vs measure
  - Pattern analysis



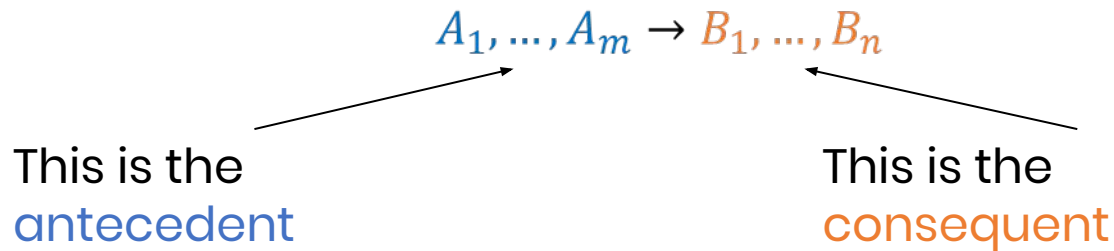
# Data Interrelationships

## Functional Dependency

A **Functional Dependency**  $A_1, \dots, A_m \rightarrow B_1, \dots, B_n$  holds in the relation  $R$  if:

$$\forall t, t' \in R, (t.A_1 = t'.A_1 \wedge \dots \wedge t.A_m = t'.A_m \rightarrow t.B_1 = t'.B_1 \wedge \dots \wedge t.B_n = t'.B_n)$$

Informally, **some attributes determine other attributes**.



**Warning! Dependency does not imply causation!**

# Fundamentals of FDs

## Armstrong's Axioms

- Axiom of **Reflexivity** (**Trivial FD**)
- Axiom of **Augmentation**
- Axiom of **Transitivity**

# Fundamentals of FDs

## Armstrong's Axioms

- Axiom of **Reflexivity** (**Trivial FD**)

If  $B \subseteq A$  then  $A \rightarrow B$

- Axiom of **Augmentation**

- Axiom of **Transitivity**

# Fundamentals of FDs

## Armstrong's Axioms

- Axiom of **Reflexivity** (**Trivial FD**)

If  $B \subseteq A$  then  $A \rightarrow B$

[ex]  $\{name\} \subseteq \{name, job\}$  so  $\{name, job\} \rightarrow \{name\}$

- Axiom of **Augmentation**

- Axiom of **Transitivity**

# Fundamentals of FDs

## Armstrong's Axioms

- Axiom of **Reflexivity** (**Trivial FD**)

If  $B \subseteq A$  then  $A \rightarrow B$

[ex]  $\{name\} \subseteq \{name, job\}$  so  $\{name, job\} \rightarrow \{name\}$

- Axiom of **Augmentation**

If  $A \rightarrow B$  then  $\forall C, AC \rightarrow BC$

- Axiom of **Transitivity**

# Fundamentals of FDs

## Armstrong's Axioms

- Axiom of **Reflexivity (Trivial FD)**

If  $B \subseteq A$  then  $A \rightarrow B$

[ex]  $\{name\} \subseteq \{name, job\}$  so  $\{name, job\} \rightarrow \{name\}$

- Axiom of **Augmentation**

If  $A \rightarrow B$  then  $\forall C, AC \rightarrow BC$

[ex]  $\{ID\} \rightarrow \{name\}$  so  $\{ID, job\} \rightarrow \{name, job\}$

- Axiom of **Transitivity**



# Fundamentals of FDs

## Armstrong's Axioms

- Axiom of **Reflexivity (Trivial FD)**

If  $B \subseteq A$  then  $A \rightarrow B$

[ex]  $\{name\} \subseteq \{name, job\}$  so  $\{name, job\} \rightarrow \{name\}$

- Axiom of **Augmentation**

If  $A \rightarrow B$  then  $\forall C, AC \rightarrow BC$

[ex]  $\{ID\} \rightarrow \{name\}$  so  $\{ID, job\} \rightarrow \{name, job\}$

- Axiom of **Transitivity**

If  $A \rightarrow B$  and  $B \rightarrow C$  then  $A \rightarrow C$

# Fundamentals of FDs

## Armstrong's Axioms

- Axiom of **Reflexivity** (**Trivial FD**)

If  $B \subseteq A$  then  $A \rightarrow B$

[ex]  $\{name\} \subseteq \{name, job\}$  so  $\{name, job\} \rightarrow \{name\}$

- Axiom of **Augmentation**

If  $A \rightarrow B$  then  $\forall C, AC \rightarrow BC$

[ex]  $\{ID\} \rightarrow \{name\}$  so  $\{ID, job\} \rightarrow \{name, job\}$

- Axiom of **Transitivity**

If  $A \rightarrow B$  and  $B \rightarrow C$  then  $A \rightarrow C$

[ex]  $\{ID\} \rightarrow \{name\}$  and  $\{name\} \rightarrow \{initials\}$   
so  $\{ID\} \rightarrow \{initials\}$

# Fundamentals of FDs

## Interesting Secondary Rules

- ***Pseudo Transitivity***

If  $A \rightarrow BC$  and  $C \rightarrow D$  then  $A \rightarrow BD$

- ***Extensivity***

If  $A \rightarrow B$  then  $A \rightarrow AB$

# Fundamentals of FDs

Can I do this to FDs?

I only know  $\{ID\} \rightarrow \{name\}$

So  $\{ID, \textit{hair color}\} \rightarrow \{name\}$

# Fundamentals of FDs

Can I do this to FDs?

I only know  $\{ID\} \rightarrow \{name\}$

So  $\{ID, \textit{hair color}\} \rightarrow \{name\}$

Yes!

# Fundamentals of FDs

Can I do this to FDs?

I only know  $\{ID\} \rightarrow \{name\}$

So  $\{ID, \textit{hair color}\} \rightarrow \{name\}$

Yes!

Adding more attributes to the antecedent can never remove attributes in the consequent.

# Fundamentals of FDs

What about this?

I only know  $\{ID\} \rightarrow \{name\}$

So  $\{ID\} \rightarrow \{name, \textit{hair color}\}$

# Fundamentals of FDs

What about this?

I only know  $\{ID\} \rightarrow \{name\}$

So  $\{ID\} \rightarrow \{name, \textit{hair color}\}$

No!



# Fundamentals of FDs

What about this?

I only know  $\{ID\} \rightarrow \{name\}$

So  $\{ID\} \rightarrow \{name, \textit{hair color}\}$

No!

No way to use the axioms to introduce hair color to the consequent without also introducing it to the antecedent.

# Finding Keys

All this talk about FDs sounds awfully similar to keys...

# Closure

## Closure

The **Closure** of the set  $\{A_1, \dots, A_m\}$ , written as  $\{A_1, \dots, A_m\}^+$ , is the set of attributes  $B$  is such that  $A_1, \dots, A_m \rightarrow B$ .

A closure finds **everything a set of attributes determines**.

## Closure (example)

Given the functional dependencies:

- $SSN \rightarrow Name$
- $Name \rightarrow Initials$

We can derive some closures:

- $Name^+ =$
- $SSN^+ =$
- $Initials^+ =$
- $\{SSN, Initials\}^+ =$

# Closure

## Closure

The **Closure** of the set  $\{A_1, \dots, A_m\}$ , written as  $\{A_1, \dots, A_m\}^+$ , is the set of attributes  $B$  is such that  $A_1, \dots, A_m \rightarrow B$ .

A closure finds **everything a set of attributes determines**.

## Closure (example)

Given the functional dependencies:

- $SSN \rightarrow Name$
- $Name \rightarrow Initials$

We can derive some closures:

- $Name^+ = \{Name, Initials\}$
- $SSN^+ =$
- $Initials^+ =$
- $\{SSN, Initials\}^+ =$

# Closure

## Closure

The **Closure** of the set  $\{A_1, \dots, A_m\}$ , written as  $\{A_1, \dots, A_m\}^+$ , is the set of attributes  $B$  is such that  $A_1, \dots, A_m \rightarrow B$ .

A closure finds **everything a set of attributes determines**.

## Closure (example)

Given the functional dependencies:

- $SSN \rightarrow Name$
- $Name \rightarrow Initials$

We can derive some closures:

- $Name^+ = \{Name, Initials\}$
- $SSN^+ = \{SSN, Name, Initials\}$
- $Initials^+ =$
- $\{SSN, Initials\}^+ =$

# Closure

## Closure

The **Closure** of the set  $\{A_1, \dots, A_m\}$ , written as  $\{A_1, \dots, A_m\}^+$ , is the set of attributes  $B$  is such that  $A_1, \dots, A_m \rightarrow B$ .

A closure finds **everything a set of attributes determines**.

## Closure (example)

Given the functional dependencies:

- $SSN \rightarrow Name$
- $Name \rightarrow Initials$

We can derive some closures:

- $Name^+ = \{Name, Initials\}$
- $SSN^+ = \{SSN, Name, Initials\}$
- $Initials^+ = \{Initials\}$
- $\{SSN, Initials\}^+ =$

# Closure

## Closure

The **Closure** of the set  $\{A_1, \dots, A_m\}$ , written as  $\{A_1, \dots, A_m\}^+$ , is the set of attributes  $B$  is such that  $A_1, \dots, A_m \rightarrow B$ .

A closure finds **everything a set of attributes determines**.

## Closure (example)

Given the functional dependencies:

- $SSN \rightarrow Name$
- $Name \rightarrow Initials$

We can derive some closures:

- $Name^+ = \{Name, Initials\}$
- $SSN^+ = \{SSN, Name, Initials\}$
- $Initials^+ = \{Initials\}$
- $\{SSN, Initials\}^+ = \{SSN, Name, Initials\}$

# Closure

## Closure Algorithm

Find the closure of  
 $\{A_1, \dots, A_m\}$

$X = \{A_1, \dots, A_m\}$

**Repeat until  $X$  does not change:**

**if**  $B_1, \dots, B_n \rightarrow C$  is a FD **and**  $B_1, \dots, B_n \in X$

**then**  $X \leftarrow X \cup C$

In practice:

Repeated use of transitivity



# Closure

## Closure Algorithm

Find the closure of  
 $\{A_1, \dots, A_m\}$

$X = \{A_1, \dots, A_m\}$

**Repeat until  $X$  does not change:**

**if**  $B_1, \dots, B_n \rightarrow C$  is a FD **and**  $B_1, \dots, B_n \in X$

**then**  $X \leftarrow X \cup C$

In practice:

Repeated use of transitivity

If a FD applies, add the  
consequent to the  
answer

# Closure Example

Let's say we have the following relations and FDs:

Restaurants(*rid*, name, rating, popularity)

$rid \rightarrow name$

$rid \rightarrow rating$

$rating \rightarrow popularity$

Compute  $\{rid\}^+$

# Closure Example

Let's say we have the following relations and FDs:

Restaurants(*rid*, name, rating, popularity)

$rid \rightarrow name$

$rid \rightarrow rating$

$rating \rightarrow popularity$

Compute  $\{rid\}^+$

$\{rid\}^+ = \{rid, name, rating, popularity\}$

it's a key!

# Finding Keys

What do FDs and Closures do for us?

- Characterize the interrelationships of data
- Able to find keys

# Finding Keys

## Superkey

A **Superkey** is a set of attributes  $A_1, \dots, A_n$  s.t. for any single attribute  $B$ :

$$A_1, \dots, A_n \rightarrow B$$

In other words, for the set of all attributes  $C$  in the relation  $R$ , the set  $\{A_1, \dots, A_n\}$  is a superkey iff  $\{A_1, \dots, A_n\}^+ = C$

## Key

A **Key** is a minimal superkey, i.e. no subset of a key is a superkey.

## Candidate Key

When a relation has multiple keys, each key is a **Candidate Key**.

# Usefulness of Keys in Design

What intuitions do we get from data interrelationships?

- FDs that are not superkeys hint at redundancy
  - If a FD antecedent is **not** a superkey, we can remove redundant information, i.e. the FD consequent

# Usefulness of Keys in Design

What intuitions do we get from data interrelationships?

- FDs that are not superkeys hint at redundancy
  - If a FD antecedent is **not** a superkey, we can remove redundant information, i.e. the FD consequent
- Rephrased
  - $A \rightarrow B$  is fine if  $A$  is a superkey
  - Otherwise, we can extract  $B$


# Usefulness of Keys in Design

Restaurants(rid, name, rating, popularity)

rid  $\square$  name

rid  $\square$  rating

rating  $\square$  popularity




rid	name	rating	popularity
1	Mee Sum Pastry	3	OK
2	Café on the Ave	4	High
3	Guanaco's Tacos	4	High
4	Aladdin Gyro-Cery	5	High



# Usefulness of Keys in Design

Restaurants(rid, name, rating, popularity)

rid □ name  
rid □ rating } Fine because rid is a  
superkey  
rating □ popularity



rid	name	rating	popularity
1	Mee Sum Pastry	3	OK
2	Café on the Ave	4	High
3	Guanaco's Tacos	4	High
4	Aladdin Gyro-Cery	5	High


# Usefulness of Keys in Design

Restaurants(rid, name, rating, popularity)

rid □ name  
rid □ rating

} Fine because rid is a  
superkey

rating □ popularity



rid	name	rating	popularity
1	Mee Sum Pastry	3	OK
2	Café on the Ave	4	High
3	Guanaco's Tacos	4	High
4	Aladdin Gyro-Cery	5	High


# Usefulness of Keys in Design

Restaurants(rid, name, rating, popularity)

rid  $\square$  name  
rid  $\square$  rating

Fine because rid is a superkey

rating  $\square$  popularity



rid	name	rating	popularity
1	Mee Sum Pastry	3	OK
2	Café on the Ave	4	High
3	Guanaco's Tacos	4	High
4	Aladdin Gyro-Cery	5	High

Redundancy!

# Database Design

Database Design is about  
(1) characterizing data and (2) organizing data

How to talk about properties  
we know or see in the data

# Database Design

Database Design is about  
(1) characterizing data and (2) organizing data

How to organize data to promote  
ease of use and efficiency

# Normal Forms

## Normal Forms

- 1NF □ Flat
- 2NF □ No partial FDs (obsolete)
- 3NF □ Preserve all FDs, but allow anomalies
- BCNF □ No transitive FDs, but can lose FDs
- 4NF □ Considers multi-valued dependencies
- 5NF □ Considers join dependencies (hard to do)

# Normal Forms

## Normal Forms

- 1NF □ Flat
- 2NF □ No partial FDs (obsolete)
- 3NF □ Preserve all FDs, but allow anomalies
- BCNF □ No transitive FDs, but can lose FDs
- 4NF □ Considers multi-valued dependencies
- 5NF □ Considers join dependencies (hard to do)

# Normal Forms

## 1NF

A relation  $R$  is in **First Normal Form** if all attribute values are atomic. Attribute values cannot be multivalued. Nested relations are not allowed.

We call data in 1NF “flat.”



## BCNF

A relation  $R$  is in **Boyce-Codd Normal Form (BCNF)** if for every non-trivial dependency,  $X \rightarrow A$ ,  $X$  is a superkey.

Equivalently, a relation  $R$  is in BCNF if  $\forall X$  either  $X^+ = X$  or  $X^+ = C$  where  $C$  is the set of all attributes in  $R$

# Decomposition

- “Extracting” attributes can be done with **decomposition** (split the schema into smaller parts)
- For this class, decomposition means the following:

$$R(A_1, \dots, A_n, B_1, \dots, B_m, C_1, \dots, C_k) \begin{matrix} \swarrow \\ \searrow \end{matrix} \begin{matrix} R_1(A_1, \dots, A_n, B_1, \dots, B_m) \\ R_2(A_1, \dots, A_n, C_1, \dots, C_k) \end{matrix}$$

# Decomposition

- “Extracting” attributes can be done with **decomposition** (split the schema into smaller parts)
- For this class, decomposition means the following:

$$R(A_1, \dots, A_n, B_1, \dots, B_m, C_1, \dots, C_k) \begin{cases} R_1(A_1, \dots, A_n, B_1, \dots, B_m) \\ R_2(A_1, \dots, A_n, C_1, \dots, C_k) \end{cases}$$

Some common attributes are present so we can rejoin data

## BCNF Decomposition Algorithm

<p> <i>Normalize(R)</i>  <math>C \leftarrow</math> the set of all attributes in <math>R</math>  <b>find</b> <math>X</math> <b>s.t.</b> <math>X^+ \neq X</math> <b>and</b> <math>X^+ \neq C</math>  <b>if</b> <math>X</math> is not found  <b>then</b> "<math>R</math> is in BCNF"  <b>else</b>              decompose <math>R</math> into <math>R_1(X^+)</math> and <math>R_2((C - X^+) \cup X)</math>              <i>Normalize(R<sub>1</sub>)</i>              <i>Normalize(R<sub>2</sub>)</i> </p>	
---	--

## BCNF Decomposition Algorithm

*Normalize(R)*

$C \leftarrow$  the set of all attributes in  $R$

**find**  $X$  **s.t.**  $X^+ \neq X$  **and**  $X^+ \neq C$

**if**  $X$  is not found

**then** "R is in BCNF"

**else**

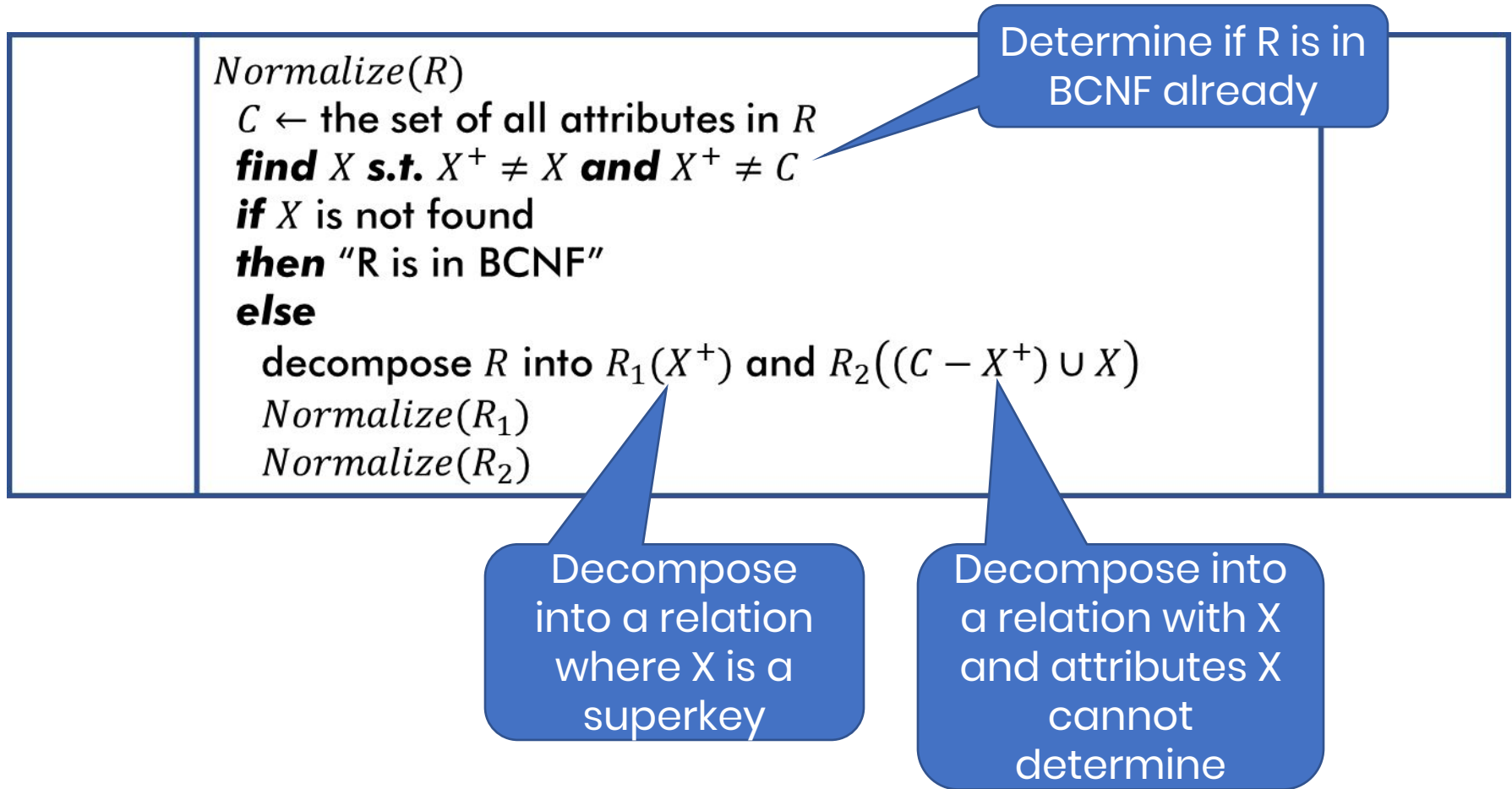
decompose  $R$  into  $R_1(X^+)$  and  $R_2((C - X^+) \cup X)$

*Normalize(R<sub>1</sub>)*

*Normalize(R<sub>2</sub>)*

Determine if R is in  
BCNF already

## BCNF Decomposition Algorithm



# BCNF Decomposition Example

*Normalize(R)*

$C \leftarrow$  the set of all attributes in  $R$

**find**  $X$  s.t.  $X^+ \neq X$  **and**  $X^+ \neq C$

**if**  $X$  is not found

**then** "R is in BCNF"

**else**

decompose  $R$  into  $R_1(X^+)$  and  $R_2((C - X^+) \cup X)$

*Normalize(R<sub>1</sub>)*

*Normalize(R<sub>2</sub>)*

Restaurants(rid, name,  
rating, popularity,  
recommended)

rid  $\square$  name, rating

rating  $\square$  popularity

popularity  $\square$  recommended

# BCNF Decomposition Example

*Normalize(R)*

$C \leftarrow$  the set of all attributes in  $R$

**find**  $X$  s.t.  $X^+ \neq X$  **and**  $X^+ \neq C$

**if**  $X$  is not found

**then** "R is in BCNF"

**else**

decompose  $R$  into  $R_1(X^+)$  and  $R_2((C - X^+) \cup X)$

*Normalize(R<sub>1</sub>)*

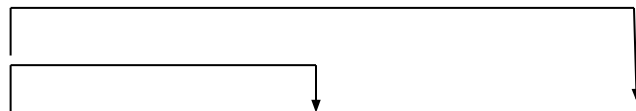
*Normalize(R<sub>2</sub>)*

Restaurants(rid, name,  
rating, popularity,  
recommended)

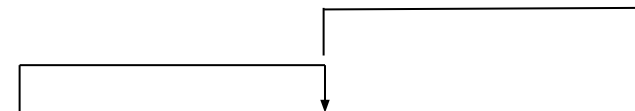
rid  $\square$  name, rating

rating  $\square$  popularity

popularity  $\square$  recommended



rid	name	rating
1	Mee Sum Pastry	3
2	Café on the Ave	4
3	Guanaco's Tacos	4
4	Aladdin Gyro-Cery	5



rating	popularity	recommended
3	OK	no
4	High	yes
5	High	yes



# BCNF Decomposition Example

*Normalize(R)*

$C \leftarrow$  the set of all attributes in  $R$

**find**  $X$  s.t.  $X^+ \neq X$  **and**  $X^+ \neq C$

**if**  $X$  is not found

**then** "R is in BCNF"

**else**

decompose  $R$  into  $R_1(X^+)$  and  $R_2((C - X^+) \cup X)$

*Normalize(R<sub>1</sub>)*

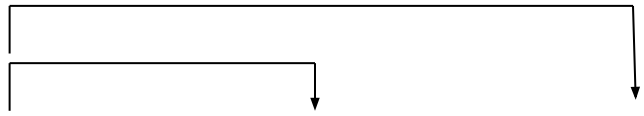
*Normalize(R<sub>2</sub>)*

Restaurants(rid, name,  
rating, popularity,  
recommended)

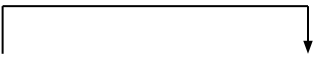
rid  $\square$  name, rating

rating  $\square$  popularity

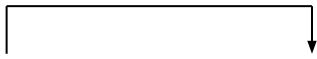
popularity  $\square$  recommended



rid	name	rating
1	Mee Sum Pastry	3
2	Café on the Ave	4
3	Guanaco's Tacos	4
4	Aladdin Gyro-Cery	5



rating	popularity
3	OK
4	High
5	High



popularity	rec
OK	no
High	yes

# Losslessness

## Definition

**Lossless Decomposition** is a reversible decomposition, i.e. rejoining all decomposed relations will always result exactly with the original data.

This is the opposite of a **Lossy Decomposition**, an irreversible decomposition, where rejoining all decomposed relations may result something other than the original data, specifically with extra tuples.

This concept might be familiar if you have ever encountered lossless data compression (e.g. Huffman encoding or PNG) or lossy data compression (e.g. JPEG).

# Losslessness

Is BCNF decomposition lossless?

# Losslessness

Is BCNF decomposition lossless?

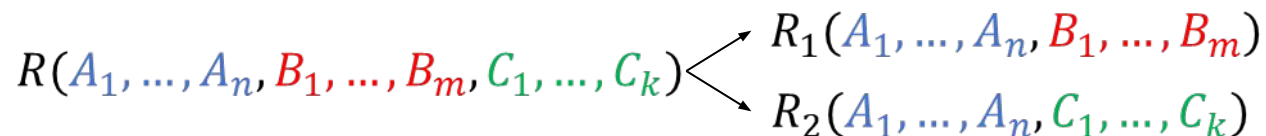
Yes!

# Losslessness

## Definition – Heath's Theorem

Suppose we have the relation  $R$  and three disjoint subsets of the attributes of  $R$  we will write as  $A_1, \dots, A_n$ ,  $B_1, \dots, B_m$ , and  $C_1, \dots, C_k$ . Suppose we also have a FD that is  $A_1, \dots, A_n \rightarrow B_1, \dots, B_m$ .

**Heath's Theorem** states that the decomposition of  $R$  into  $R_1(A_1, \dots, A_n, B_1, \dots, B_m)$  and  $R_2(A_1, \dots, A_n, C_1, \dots, C_k)$  is lossless where  $R_1$  and  $R_2$  are the projections of  $R$  on their respective attributes.



By reflection, the same decomposition of  $R$  under the alternate FD  $A_1, \dots, A_n \rightarrow C_1, \dots, C_k$  is also lossless.



- You may inherit a database that could be lossy. Before you use it, it may be worth your time to check if it is lossy.
- Full normalization is nice but can be inefficient
  - Denormalization □ don't normalize all the way

# Takeaways

- We can characterize the relationships in our data using domain knowledge or pattern analysis.
- Functional dependencies give us ways of normalizing our data to avoid anomalies.