

Roll-Up: Summarize Data, dimension reduction, less details: *sum(Sales)*

Drill-Down: Higher level summary to lower level, more dimensions and details: *sum(Sales)*

Slicing: picks a specific value for one of the dimensions: *sum(Sales)*

Dicing: subcube by picking specific values for multiple dimensions: *sum(Sales)*

Pivoting: rotate cube to provide alternative vps. of data: *sum(Sales)*

Decomposing into 3NF: #1 lossless join

- Given FDs F, compute F': minimal cover for F
- Decompose using F' if violating 3NF similar to BCNF - all the way to BCNF
- After each decomp, identify set of dependencies N in F' that are not preserved by the decomp.
- At end, for each X → b in N, create a relation R_n (X U b) & add it to the decomposition

Valid: The patterns hold in general.

Novel: We did not know the pattern before hand.

Useful: We can devise actions from analysis of large quantities of data in order to discover valid, novel, (business intelligence), potentially useful, and ultimately understandable patterns in data.

Understandable: We can interpret and comprehend the patterns in data.

SQL Examples

Find students who've taken all classes (division example)

1) With EXCEPT & NOT EXISTS

```
SELECT sname
FROM student s
WHERE NOT EXISTS (
  (SELECT c.name
   FROM class c)
  EXCEPT
  (SELECT e.cname
   FROM enrolled e
   WHERE e.snum = s.snum))
```

2) Without EXCEPT using NOT EXISTS

```
SELECT sname
FROM student s
WHERE NOT EXISTS (
  (SELECT c.name
   FROM class c)
  WHERE NOT EXISTS (
    (SELECT e.snum
     FROM enrolled e
     WHERE c.name = e.cname
     AND e.snum = s.snum)))
```

Find the names of all movie stars who've been in a movie

```
SELECT DISTINCT Name
FROM StarsIn s, MovieStar m
WHERE s.starsid = m.starsid
```

Find the departments that have more than one faculty member

```
SELECT DISTINCT f1.deptid
FROM Faculty f1, Faculty f2
WHERE f1.fid <> f2.fid AND f1.deptid = f2.deptid
```

Find IDs of MovieStars who've been in a movie in 1944 or 1974

```
SELECT StarId
FROM Movie M, StarsIn S
WHERE M.MovieID = S.MovieID AND Year = 1944
UNION ALL
SELECT StarId
FROM Movie M, StarsIn S
WHERE M.MovieID = S.MovieID AND year=1974
```

Find the name & Age of the oldest student(s)

```
SELECT sname, age
FROM Student s2
WHERE NOT EXISTS (
  (SELECT *
   FROM Student s1
   WHERE s1.age > s2.age))
```

HRU Algorithm

Assuming 'a' is already materialized, what are the best 3 other views that we should materialize?

1. Make table of values. "Gain" is the top (a) subtracted with the current one.

	First Choice	Second Choice	Third Choice
b	50*6=300	-	-
c	25*6=150	25*2=50	25*1=25
d	80*3=240	30*3=90	30*3=90
e	70*4=280	20*4=80	20*4=80
f	60*3=180	60*3=120	-
g	99*2=198	49*2=98	49*2=98
h	90*2=180	40*2=80	30*2=60
i	99	49	39

Remember that node includes itself, and once you've made a choice, you count the subtraction from that point onwards.

Data warehousing: Finding number of tuples

Data warehouse with dimensions D1, D2, D3, D4, D5. They have 999, 99, 24, 4, and 19 values. Sparsity is 10%. What's est. number of tuples? The size of the full cube is (999+1)(99+1)(24+1)(19+1) = 250,000,000 tuples. With a sparsity factor of 10%, the size of the sparse cube is about 250,000,000 x 0.1 = 25,000,000 tuples.

Apriori Algorithm minimum support: 33%.

T1	tree, cup, paper
T2	book, tomato, pen
T3	pen, book, tree
T4	tomato, pen, cup
T5	tree, paper, book
T6	pen, book

1. Find possible itemsets, and support of size C_i denoted as C_i. C₁ = {(tree): 3, {cup}: 2, {paper}: 2, {book}: 4, {tomato}: 2 {pen}: 4}

2. Find frequent itemsets of size k, denoted as F_k. F₁ = {(tree), {cup}, {paper}, {book}, {tomato}, {pen}}

It's these values because they occur more than 33% (2/6 times).

3. Go up one k, and ONLY use possible combinations from F_{k-1} (k-1). C₂ = {(tree, cup): 1, (tree, paper): 1, (tree, book): 2, (tree, tomato): 0, (tree, pen): 1, (cup, paper): 1, (cup, book): 0, (cup, tomato): 1, (cup, pen): 1, (paper, book): 1, (paper, tomato): 0, (paper, pen): 0, (book, tomato): 1, (book, pen): 3, (tomato, pen): 2}

repeat until it's not possible anymore. Apriori will finally return all the items meeting the minimum support.

Apriori speeds up calculating association rules based on the observation that each subset of a frequent itemset must also be a frequent itemset. Ex. rice only appears one time, it can't appear two or more times with anything else.

Relational Algebra Examples

Find name of actors in Jaws and Spongebob "Joining with Relation"

```
Jaws ← π name ((σ title = "Jaws" (Movie)) ⋈ StarsIn) ⋈ MovieStar
Spongebob ← π name ((σ title = "spongebob" (Movie)) ⋈ StarsIn) ⋈ MovieStar
Res ← Jaws ∩ Spongebob
```

Find name of actors in all movies "Has to be connected in some way"

```
π name (π starID, movieID (StarsIn) / π movieID (Movie)) ⋈ MovieStar
```

Find all profs in exactly all the committees prof piper is in

```
R2 ← π commname (σ profname = 'piper' (committee))
R3 ← π commname (committee) - R2
committee / R2 - π profname (committee ⋈ R3)
```

Find all profs who have offices in at least all the buildings that Piper has offices in

```
profname, building (professor ⋈ department) / π building (σ profname = 'piper' (professor) ⋈ department)
```

For each person, find all pizzas the person eats that aren't served by any pizzeria the person frequents

```
Eats - π name, pizza (Frequents ⋈ Serves)
```

Find names of people who frequent only pizzerias serving at least one pizza they eat.

```
π name (Person) - π name (Frequents - π name, pizzeria (Eats ⋈ Serves))
```

Find names of people who frequent all the pizzerias serving at least one pizza they eat.

```
π name (Person) - π name (π name, pizzeria (Eats ⋈ Serves) - Frequents)
```

The most expensive pass cost

```
π cost (Pass) - π cost (P cost → newcost (Pass) ⋈ newcost < cost (Pass))
```

The min price book

```
π isbn (Book) - π isbn (σ b1.price > b2.price (p (b1, Book) x p (b2, Book)))
```

Pizzeria serving cheapest pepperoni pizza. For ties, return all of cheapest pepperoni pizzas.

```
pizzeria (σ pizza = pepperoni (Serves)) - π pizzeria (σ Note (P (Serves, s1) x P (Serves, s2))
```

Aggregate Operations

AVG, MIN, MAX, SUM, COUNT

```
SELECT country, city, count(*)
FROM customers
WHERE country LIKE '%ali'
GROUP BY country, city
HAVING count(*) < 3
ORDER BY count(*) DESC
```

Null Values tuples can have null values. unknown or doesn't exist

- can use IS NULL (IS NOT NULL) in WHERE or others

Null and 3-valued logic

- unknown or true = true
- unknown or false = unknown
- unknown and true = unknown
- unknown and false = false
- not unknown = not unknown
- Any comparison with null returns unknown
- result of WHERE pred. is false if evals to unknown
- All aggregate ops except count(*) ignore nulls

LIKE is used for string matching

- _ stands for any one character
- % is 0 or more arbitrary characters

Relational Algebra Part 2

Find all the professors who are in any of the committee's professor Piper is in

```
SELECT DISTINCT profname
FROM Committee
WHERE commname IN (
  SELECT *
  FROM Committee
  WHERE profname = 'Piper')
```

RA Statement

```
π commname (σ profname = 'Piper' (c1, profname = c2, profname (c1, committee) x c2, committee))
```

Find all professors who are in at least all those committees that professor Piper is in

```
SELECT DISTINCT c.profname
FROM Committee c
WHERE NOT EXISTS (
  (SELECT *
   FROM Committee a
   WHERE profname = 'Piper')
  EXCEPT
  (SELECT commname
   FROM committee a
   WHERE a.profname = c.profname))
```

Find all enclosures that were never visited by the visitors born before 2000-01-01

```
SELECT E.id
FROM Employee E
EXCEPT
(SELECT EZ.id
 FROM Visitor v, Visits vs, Enclosure EZ
 WHERE v.visitorID = vs.visitorID AND vs.enclosureID = EZ.id
 AND v.DateOfBirth < Date('2001-01-01'))
```

Find all the animals of species cat that are taken care of by 'Sam Smith'

```
SELECT employeeID
FROM TakesCareOf NATURAL INNER JOIN Employee
WHERE animalID IN (SELECT id from Animal WHERE Species='Cat') AND
firstName = 'Same' AND lastName = 'Smith'
```

CREATE VIEW Temp(major, average) AS

```
SELECT s.major, AVG(s.age) AS average
FROM Student S
GROUP BY major
```

grouping by major, getting average age selecting minimum over all majors

Get all students using natural join

```
SELECT *
FROM Student s, enrolled e
WHERE s.snum = e.snum
```

Find the names of all movie stars who've been in a movie

```
SELECT DISTINCT Name
FROM StarsIn s, MovieStar m
WHERE s.starsid = m.starsid
```

Find the departments that have more than one faculty member

```
SELECT DISTINCT f1.deptid
FROM Faculty f1, Faculty f2
WHERE f1.fid <> f2.fid AND f1.deptid = f2.deptid
```

Find IDs of MovieStars who've been in a movie in 1944 or 1974

```
SELECT StarId
FROM Movie M, StarsIn S
WHERE M.MovieID = S.MovieID AND Year = 1944
UNION ALL
SELECT StarId
FROM Movie M, StarsIn S
WHERE M.MovieID = S.MovieID AND year=1974
```

Find the name & Age of the oldest student(s)

```
SELECT sname, age
FROM Student s2
WHERE NOT EXISTS (
  (SELECT *
   FROM Student s1
   WHERE s1.age > s2.age))
```

Or with aggregate operators

```
SELECT sname, age
FROM Student s2
WHERE s2.age = (SELECT MAX(s2.age) FROM Student S2)
```

Set operations

- UNION** usually used for or
- INTERSECT** usually used for and
- EXCEPT** usually used for A but didn't B

Natural Join

Default is INNER JOIN - only include matches.

FROM Student S NATURAL LEFT OUTER JOIN Enrolled E

Views

Relations defined w/ a create table, statement existing in the physical layer

- Hides data from users, make queries easier, modularity

CREATE VIEW CourseWithFalls(dept, course #, mark) AS

```
SELECT c.dept, c.course#, mark
FROM Course C, Enrolled E
WHERE c.course# = e.course# and mark < 50
```

View updates must occur at base tables.

- DBMS restrict view updates only to some simple views on single tables (updatable views)

DROP VIEW <view name> does NOT affect any tuples in underlying relation

We can assign commands to DROP table for view.

DROP TABLE student RESTRICT/CASCADE

- drop table unless there's a view on it
- drop table & recursively drops any view referencing it

Updating with condition

```
update SuppProd
set Price = Price * 0.9
where PID in
(select PID
 from Product
 where Name = 'Coronary Stent');
```

Find the names of sailors who have reserved at least 2 boats

```
SELECT s.sname
FROM Sailors s, Reserves r1, Reserves r2
WHERE r1.sid=s.sid AND r2.sid=s.sid AND r1.bid <> r2.bid
```

Find the sailor id of the sailors with the highest rating

```
SELECT s.sid
FROM sailors s
WHERE s.rating = (SELECT MAX(s2.rating) FROM Sailors s2)
```

Find names of sailors who have reserved all boats whose name starts w/ "typhoon"

```
SELECT s.sname
FROM Sailors s
WHERE NOT EXISTS (
  (SELECT b.bid FROM Boats b WHERE name LIKE '%typhoon%')
  EXCEPT
  (SELECT r.bid FROM Reserves r WHERE r.sid = s.sid))
```

Find the name & age of the oldest sailor

```
SELECT s.sname, s.age
FROM Sailor s
WHERE s.age = (SELECT MAX(s2.age) FROM Sailor s2)
```


Entity Set: Collection of entities (ex. cats) **Entity Set:** Collection of entities (ex. cats)

- All entities in entity set have same set of attributes
- All entity sets have keys

Domain: value type (ex. float, date, int)

Key: minimal set of attributes which can identify an entity in entity set.

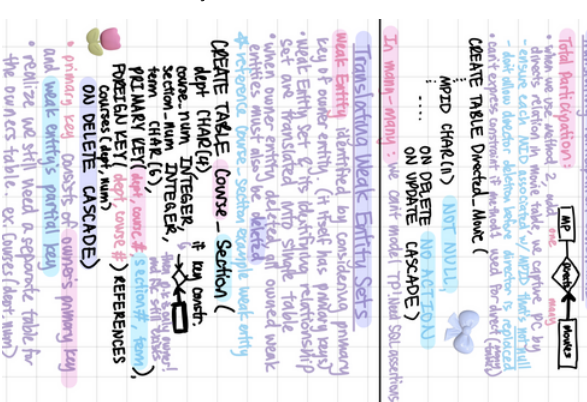
Primary Key: main key to identify entity in entity set. has to be as minimal as possible. (ex. cat_id). MUST BE UNIQUE, and NON-NULL.

Candidate key: One more key or keys in a relation

Super key: A key + zero or more other additional attributes. (Ex. {sid, name}, or {CWL, major}, or {name, major, age})

Minimal Key: Smallest set of keys to identify entity in entity set

Cardinality ratio: the number of relationships in the set that an entity can participate in. These are called cardinality constraints:



Key Constraints

Key constraints are shown with **arrow** in ER diagram.

one-to-one: entity in A is associated at most once with one entity in B (ex. A: Student, B: student ID #)

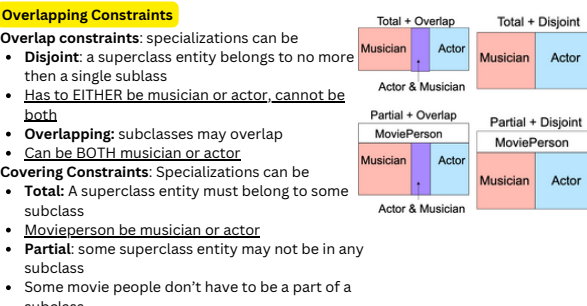
one-to-many: entity in A is associated with any number of entities in B (ex. A: Mom, B: children)

many-to-many: entity in A is associated with any number of entities in B (ex. A: Mom, B: children)

one-to-one: opposite of one-to-many

one-to-many: (B is one, A is many) is primary key of A

many-to-many: entities in A can be associated with or many to many is Primary Key of AND Primary Key of OR



Participation Constraints Whether a relationship has to be populated or not. Important for updates.

Aggregation

Having a relationship between relationships is forbidden.

Aggregation allows us to treat relationship set as entity set, letting us participate in other relationships.

Exercise

The (minimal) key of Evaluates is iid + course# + term.

TP vs. OLAP

Insert

DELETE FROM Student WHERE name = 'Smith'

Update

Update Student SET age = age + 2 WHERE age < 98

Translating Entity Set: mapped to a table. The attributes of table must include:

- Keys for each participating entity set
- Foreign keys for each relationship
- Primary key for the participating entity set
- Can't really reduce the 3 tables each

One-to-many: combine primary side of the relationship, takes place

many-to-many: combine primary side of the relationship, takes place

ON DELETE CASCADE: if one entity is deleted, all related entities are also deleted.

REFERENCES: FOREIGN KEY (movie_id) REFERENCES movie (id)

Exercise: Choose one side to combine w/ relationship

One-to-one: choose one side to combine w/ relationship

Exercise: Choose one side to combine w/ relationship

Exercise: Choose one side to combine w/ relationship

Normalization Removing redundancy from data.

First Normal Form (1NF)

Each attribute in a tuple has only one value (can't be an array). E.g., for "postal code" you can't have both V6T 1Z4 and V6S 1W6

Codd's original version allowed multi-valued attributes.

Client ID	Postal Code
1	INCORRECT. V6T 1Z4, V6S 1W6
Client ID	Postal Code
1	In 1NF. V6T 1Z4
1	V6S 1W6

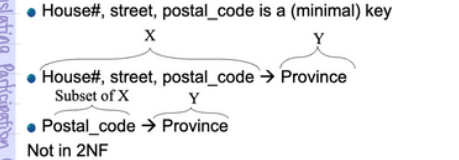
Second Normal Form (2NF)

No partial key dependencies.

A relation is in 2NF if it is in 1NF and for every FD, $X \rightarrow Y$ where X is a (minimal) key and Y is a non-key attribute, then no proper subset of X determines Y.

e.g., the address relation is not in 2NF:

- House#, street, postal_code is a (minimal) key



Boyce-Codd Normal Form (BCNF)

For all non-trivial functional dependencies $X \rightarrow Y$, X must be a superkey for a relation to be in BCNF.

Ex. Whenever a set of attributes R determine another attribute, it should determine all the attributes of R.

Check if all the left parts of the FD are a superkey. If they are, then it is in BCNF, and if not, they need to be decomposed.

Third Normal Form (3NF)

A Relation R is in 3NF if:

If $X \rightarrow B$ is a non-trivial dependency in R, then X is a superkey for R, or B is part of a (minimal) key.

Note: B must be part of a key not part of a superkey (if a key exists, all attributes are part of a superkey)

Example:

Each CK has 3 attributes but each FD has 2, so it's not in BCNF.

3NF?

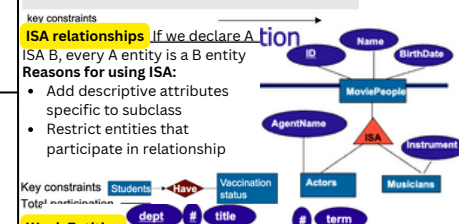
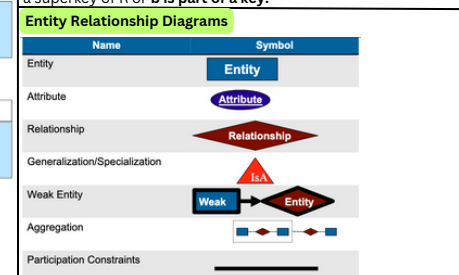
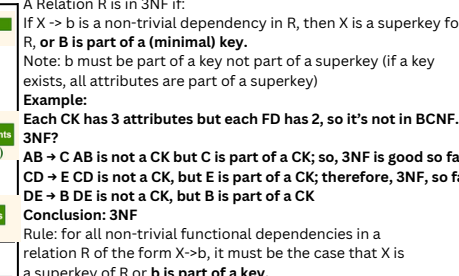
AB + C AB is not a CK but C is part of a CK; so, 3NF is good so far.

CD + E CD is not a CK, but E is part of a CK; therefore, 3NF, so far

DE + B DE is not a CK, but B is part of a CK

Conclusion: 3NF

Rule: for all non-trivial functional dependencies in a relation R of the form $X \rightarrow B$, it must be the case that X is a superkey for R or B is part of a key.



Entity Relationship Diagrams

Entity

Attribute

Relationship

Generalization/Specialization

Weak Entity

Aggregation

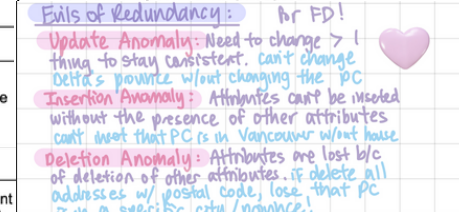
Participation Constraints

key constraints

ISA relationships If we declare A ISA B, every A entity is a B entity

Reasons for using ISA:

- Add descriptive attributes specific to subclass
- Restrict entities that participate in relationship



Weak Entities

Exercise: Choose one side to combine w/ relationship

Exercise: Choose one side to combine w/ relationship

Exercise: Choose one side to combine w/ relationship

Exercise: Choose one side to combine w/ relationship

Decomposing into BCNF

Relation: R(ABCD) FD: B -> C, D -> A

Keys?

A* = {A}

B* = {B, C}

C* = {C}

D* = {A, D}

BD* = {B, D, C, A}

BD is the only key

Look at FD B -> C. Is B a superkey?

No. Decompose R1(B, C), R2(A, B, D)

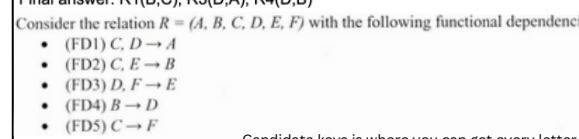
Look at FD D -> A. Is D a superkey for R2?

No. Decompose R3(D, A), R4(D, B)

Final answer: R1(B, C), R3(D, A), R4(D, B)

Consider the relation R = (A, B, C, D, E, F) with the following functional dependencies:

- (FD1) C, D -> A
- (FD2) C, E -> B
- (FD3) D, F -> E
- (FD4) B -> D
- (FD5) C -> F



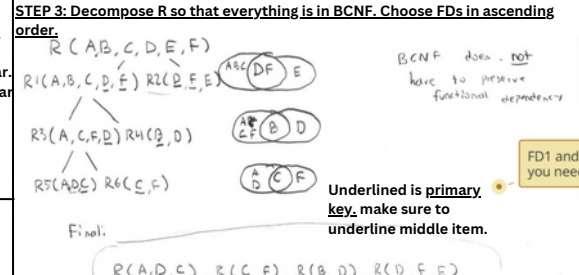
STEP 1: Find all candidate keys.

Candidate keys are where you can get every letter.

Candidate keys: {C, F}, {C, D}, {C, B}

STEP 2: Identify all the FDs that violate the BCNF condition.

BCNF is that for $X \rightarrow Y$, it has to be a superkey. Since {C, E}, {C, D}, {C, B} are keys, FD1 and FD2 are ruled out, meaning only 3 options left.



Decompose into minimal cover

Given the relation A(P, Q, R, S, T, U, V, W) and set of FDs:

P, Q -> T

Q -> U

U -> V

V -> W

W -> S

S -> T

U, V -> W

U, V -> S

U, V -> V

STEP 1: Put FDs into standard form.

Standard form: split output into individual FDs.

V, U -> A, B gets turned into V, U -> A and V, U -> B.

P, Q -> T gets turned into P -> T, Q -> T.

STEP 2: Minimize the LHS (left-hand-side) of FDs.

Since U -> V, V, U -> W can become U -> W

Since U -> W, V, U -> V can become U -> V, since already exists on LHS.

Look for areas that can be changed to make LHS small as possible. The "V" in V, U -> W is useless.

STEP 3: Delete Redundant FDs.

Since V, U -> V is trivial, can remove it

Since Q -> U and U -> V, Q -> V is redundant and can be removed

Final: P, Q -> T Q -> U U -> V V, U -> W R -> S

Statements which are proven by other statements can be removed.

Reflexivity: If $X \subseteq Y$, then $X \rightarrow Y$

Augmentation: If $X \rightarrow Y$, then $XZ \rightarrow YZ$ for any Z

Transitivity: If $X \rightarrow Y$ and $Y \rightarrow Z$, then $X \rightarrow Z$

Union: If $X \rightarrow Y$ and $X \rightarrow Z$, then $X \rightarrow YZ$

Decomposition: If $X \rightarrow YZ$, then $X \rightarrow Y$ and $X \rightarrow Z$

Normal Form, and proving decompositions

1. name, pf -> name, pf

2. name -> city

3. name -> status

4. name, pf -> city, pf

5. name, pf -> status, pf

6. name, pf -> name, pf, status

7. name, pf -> name, pf, status, city

8. name, pf -> name, pf, status, city, qty

9. name, pf -> name, pf, status, city, qty, pname