

Can you have too much data?

Or: Even a blind squirrel may occasionally find a nut

Or: How do you know that needle you found in the haystack is the one you were looking for?



Potential Problems on the Analysis Side

Bonferroni's Principle

In a nutshell: *if you look harder than the data supports, you will find a pattern that "fits" the pattern you are looking for.*

If you look for events of a certain type in a data set you can expect events of that type to occur -- **even if the data is completely random** -- and the number of occurrences will grow as the size of the data grows.



Example

- Consider a country with a population of 300,000,000 people. Any individual could be a "bad-guy".
- On average each person stays in a hotel once every 100 days, or equivalently, 1% of the population stays in a hotel on any given day.
- Hotels hold 100 people on average. There are just enough hotels for everyone
 $\Rightarrow (300,000,000/100)/100 = 30,000$ hotels
- An unnamed government agency collects data on hotel stays for 1000 days, looking for pairs of people who stayed at the same hotel on two different days.
 Hypothesis: "bad-guys" like to meet at hotels to plan evil-doing.



What does Bonferroni tell us?

- What if there are actually no "bad-guys" and everyone just behaves normally (randomly).
- What if everyone stays in a hotel 1% of the time and chooses one of the 30,000 hotels randomly? What will the data indicate?
- $P(\text{Alice and Bob both choose the same day to stay in a hotel}) = .01 * .01 = .0001.$
- $P(\text{Alice and Bob both stay in a given hotel on the same day}) = .0001 / 30000 = 3.3 * 10^{-9}$
- $P(\text{Alice and Bob stay in same hotel on 2 days}) = 3.3 * 10^{-9} * 3.3 * 10^{-9} = 10^{-17}$



How many events will look "suspicious"?

- The number of pairs of individuals is $(300,000,000 \text{ choose } 2) = 4.5 * 10^{16}$
- The number of pairs of days is $(1000 \text{ choose } 2) = 5 * 10^5$
- The number of events that look "suspicious" will be $4.5 * 10^{16} * 5 * 10^5 * 10^{-17} = 225,000$
- That is, the data will identify 225,000 suspicious events, even if there are no "bad guys".
- Hard Question: If there are 100 "bad guys" in the population, is it okay to investigate 225,000 others to find them?



Potential Problems on the Analysis Side

Bonferroni's Principle

In a nutshell: *if you look harder than the data supports, you will find a pattern that "fits" the pattern you are looking for.*

If you look for events of a certain type in a data set you can expect events of that type to occur -- **even if the data is completely random** -- and the number of occurrences will grow as the size of the data grows.

WAKE FOREST UNIVERSITY

Why do we need to know about relational databases in 2014?

WAKE FOREST UNIVERSITY

Relational Databases

- Continue to be the most common way to organize large datasets, although not the very largest data sets
- Work well for structured data
- Provide layers of separation between physical storage and a user's perspective
- Work well for very high speed storage and retrieval of moderately large size datasets (up to a few TB)
- Handle concurrent access very well
- Ensure that data is consistent and valid:
 - ACID properties

WAKE FOREST UNIVERSITY

Relational Databases


- Come with well-developed query languages such as SQL, QBE, etc.
- Associated with a simple, well-defined, and well understood *data model* including constraints and operations
 - Data model – relations over domains; tables with rows and columns
 - Constraints – each value in a column is of a specific data type; all rows have same number of columns, etc
 - Operations – union, intersection, difference, Cartesian product, etc.

WAKE FOREST UNIVERSITY

Relational Databases

- The set of operations we perform on relations is referred to as the *relational algebra*
- The relational database model allows developers to separate structure from content in the same way that *classes* are used in O-O design.
 - Database schema – the design, identifies the domains (data types), the number of columns, the relationships between tables, constraints, etc.
 - Database instance – like objects in O-O, refers to a snapshot of the data in a DB at any point in time; might be none!

WAKE FOREST UNIVERSITY



Relational Databases

Relational Algebra

Select, project, join

Jennifer Widom

Duplicates

Relational Algebra

List of application majors and decisions

Col l ege			Student				Appl y			
cNa me	state	enr	sID	sNam e	GP A	H S	sID	cNam e	maj or	de c

Jennifer Widom

Cross-product: combine two relations (a.k.a. Cartesian product)

Relational Algebra

Col l ege			Student				Appl y			
cNa me	state	enr	sID	sNam e	GP A	H S	sID	cNam e	maj or	de c

Jennifer Widom

Cross-product: combine two relations (a.k.a. Cartesian product)

Relational Algebra

Names and GPAs of students with HS>1000 who applied to CS and were rejected

Col l ege			Student				Appl y			
cNa me	state	enr	sID	sNam e	GP A	H S	sID	cNam e	maj or	de c

Jennifer Widom

Natural Join

Relational Algebra

- Enforce equality on all attributes with same name
- Eliminate one copy of duplicate attributes

Col l ege			Student				Appl y			
cNa me	state	enr	sID	sNam e	GP A	H S	sID	cNam e	maj or	de c

Jennifer Widom

Natural Join

Relational Algebra

Names and GPAs of students with HS>1000 who applied to CS and were rejected

Names and GPAs of students with HS>1000 who applied to CS at college with enr>20,000 and were rejected

Col l ege			Student				Appl y			
cNa me	state	enr	sID	sNam e	GP A	H S	sID	cNam e	maj or	de c

Jennifer Widom

Natural Join

Relational Algebra

Col l ege			Student				Appl y			
cNa me	state	enr	sID	sNam e	GP A	H S	sID	cNam e	maj or	de c

Jennifer Widom

Theta Join

Relational Algebra

- Basic operation implemented in DBMS
- Term “join” often means theta join

Col l ege			Student				Appl y			
cNa me	state	enr	sID	sNam e	GP A	H S	sID	cNam e	maj or	de c

Jennifer Widom

Relational Algebra

Query (expression) on set of relations produces relation as a result

- Simplest query: relation name
- Use operators to filter, slice, combine
- Operators so far: select, project, cross-product, natural join, theta join

Jennifer Widom



Relational Databases

Relational Algebra
Set operators, renaming,
notation

Jennifer Widom

Relational Algebra

Relational algebra query (expression) on set of relations produces relation as a result

Col l ege(cName, state, enrol l ment)

Student(sID, sName, GPA, sl zeHS)

Appl y(sID, cName, maj or, decl si on)

Col l ege			Student				Appl y			
cNa me	state	enr	sID	sNam e	GP A	H S	sID	cNam e	maj or	de c

Jennifer Widom

Union operator

Relational Algebra

List of college and student names

Col l ege			Student				Appl y			
cNa me	state	enr	sID	sNam e	GP A	H S	sID	cNam e	maj or	de c

Jennifer Widom

Difference operator

Relational Algebra

IDs of students who didn't apply anywhere

IDs and names of students who didn't apply anywhere

Col l ege			Student				Appl y			
cNa me	state	enr	sID	sNam e	GP A	H S	sID	cNam e	maj or	de c

Jennifer Widom

Relational Algebra

Intersection operator

Names that are both a college name and a student name

Col l ege

cNa me	state	enr

Student

sID	sNam e	GP A	H S

Appl y

sID	cNam e	maj or	de c

Jennifer Widom

Relational Algebra

Intersection doesn't add expressive power (1)

Col l ege

cNa me	state	enr

Student

sID	sNam e	GP A	H S

Appl y

sID	cNam e	maj or	de c

Jennifer Widom

Relational Algebra

Intersection doesn't add expressive power (2)

Col l ege

cNa me	state	enr

Student

sID	sNam e	GP A	H S

Appl y

sID	cNam e	maj or	de c

Jennifer Widom

Relational Algebra

Rename operator

- 1.
- 2.
- 3.

Col l ege

cNa me	state	enr

Student

sID	sNam e	GP A	H S

Appl y

sID	cNam e	maj or	de c

Jennifer Widom

Relational Algebra

Rename operator

To unify schemas for set operators

List of college and student names

Col l ege

cNa me	state	enr

Student

sID	sNam e	GP A	H S

Appl y

sID	cNam e	maj or	de c

Jennifer Widom

Relational Algebra

Rename operator

For disambiguation in "self-joins"

Pairs of colleges in same state

Col l ege

cNa me	state	enr

Student

sID	sNam e	GP A	H S

Appl y

sID	cNam e	maj or	de c

Jennifer Widom

Relational Algebra summary	Relational Algebra

Jennifer Widom