NORMALISATION

(Relational Database Schema Design Revisited)

Designing an ER Diagram is fairly intuitive, and faithfully following the steps to map an ER diagram to tables may not always result in the 'best' table design.

- This is a 'top-down' approach
- In this lecture we discuss some design criteria, and introduce the concept of **normalisation**, which is the process of transforming a collection of attributes into relations which satisfy a collection of rules.
 - This is a 'bottom-up' approach
- In constructing a relational database there is a considerable freedom which ranges from:
 - putting all of the data in one relation
 - having lots of small relations

How do we decide which makes the best set of relations for our purpose?

MSc/Dip IT - ISD L18 Normalisation (425-456)	425	19/11/2009		MSc/Dip IT - ISD L18 Normalisation (425-456)	426

Ill Effects of Multiple Copies of Data

<u>ni#</u>	name	bDate	dNum	dName	mgrni#
666	Jim Grey	01/12/67	5	Admin	666
777	Jo White	04/06/53	5	Admin	666
888	Sue Low	12/09/88	6	Finance	888
999	Tim Lee	30/07/78	7	HR	999
898	Mary Law	22/02/60	7	HR	999

Insertion. We cannot add a new department which has no employees as yet. Entering a new employee in department 5 means repeating existing departmental data

Deletion. Deleting Sue Low loses details about department 6

Modification. If the manager of department 5 is changed, it must be changed in all tuples

427

19/11/2009

Some Guidelines

Be guided by the semantics or meaning of the data

 put together in one relation data that belongs together (e.g. using the ER approach)

Avoid having too many null values

- if only 10% of employees have offices we would prefer to set up a separate relation pairing employees and offices than add a column to employee

Beware of keeping multiple versions of information

Avoid breaking up relations in such a way that spurious information is created

Preserve known functional dependencies

19/11/2009

It's Easy to Create False Data

	.		. [n No mo	ml a a
	Let's tr	y to split up t	he	<u>ni#</u>		name	рмате	ргос
	data ab	out employee	s	666	J	lim Grey	Jupiter	Perth
	and pro	jects into two)	666	,	lim Grey	Saturn	Dundee
	tables			888		Sue Low	Mars	Perth
1	- 1			ı —			7	
	<u>ni#</u>	name	pName	<u>n</u>	i#	pLoo	Thick	osos tha
	666	Jim Grey	Jupiter	6	66	Perth	info th	nat Jim
	666	Jim Grey	Saturn	6	66	Dundee	works	on Projec
	888	Sue Low	Mars	8	88	Perth	Jupite	r AT Perth
				-				

pLoc

428

Jupiter Perth

Jupiter Dundee

etc etc

Joining them back together, we get NEW TUPLES! Jim works on Jupiter and works at Perth

ni#

666

666

etc

name pName

Jim Grey

Jim Grey

etc

19/11/2009

Preserving Functional Dependencies Some Examples of Functional Dependency Functional Dependencies are known connections between items in the dName, mgrni#, mrgrStartDate all depend on dNum database which can be installed into the database by good schema design dNum, mgrni#, mrgrStartDate all depend on dName - Given ni#, we can determine the employee's name, date of birth, salary and department • i.e. columns name, bDate, salary, department all depend on ni# ni# name bDate dNum dName mgrni# mgrStartDate - Also, from knowing the ni#, we know the department number, and therefore we can find out all about the other department information as well fd1 fd2 fd3 We will proceed to use functional dependencies to guide our database design Another example is that hours depends on ni#, pNum MSc/Dip IT - ISD L18 Normalisation (425-456) 429 19/11/2009 MSc/Dip IT - ISD L18 Normalisation (425-456) 430 19/11/2009 Kev Slide Kev Slide **Defining Functional Dependence Full Functional Dependency** An attribute, X, of a relation is functionally dependent on attributes A, B -X fully depends on A . . . N if it is not dependent on any subset of A . . . N. N if the same values of R_A R_N always lead to the same value of R_X We write $R_A....R_N \rightarrow R_X$. Example for Employees : e.g. ni# \rightarrow name, bDate, dNum bDate is dependent on ni# and name but only fully dependent on ni# ni#, pNum \rightarrow hours However This is a property of the **meaning** of the data, not a property that WorksOn : hours is fully dependent on ni# and pNum emerges from the current values of the data i.e. if you happened to have no two employees with the same name, you may indeed be able to infer birth date from name for this data set, but this would not constitute genuine dependence because next week you might employ someone with the same name as a current staff member

431

Kev Slide Normalisation **First Normal Form** Normalisation is the process of taking a set of relations and decomposing them A relation is in First Normal Form (1NF) if all values are **atomic** into more relations which satisfy our criteria better - i.e. all the cells hold single values. Specifically the following are not permitted The decomposition is essentially a series of projections so that the original data can be reconstituted using joins Multi-valued attributes / sets e.g. An attribute locations which is intended to hold more than one location per department cannot be held in a single cell The form of the relations which satisfy our conditions is called a normal form Composite attributes / nested relations e.g. An attribute name which is divided up into hold a person's title, surname and There are a number of these of increasing stringency forenames cannot be held in a single cell MSc/Dip IT - ISD L18 Normalisation (425-456) 433 19/11/2009 MSc/Dip IT - ISD L18 Normalisation (425-456) 434 19/11/2009 Kev Slide **Universal Relation Candidate Keys** We start by creating the Universal Relation which is one giant relation schema A candidate key is a set of attributes which uniquely identifies each entry in the including every database attribute (using unique names) relation, and therefore can be the primary key of that relation It must be in first normal form (or it can't go in the database at all!) A relation can have more than one candidate key A (usually composite) primary key can be chosen which uniquely identifies each (e.g. dNum and dName), row. in which case one (usually the simplest) is chosen to be the primary key The following universal relation is for a subset of the Employee database consisting of Employee, Department and Project only. Composite attributes have been broken down, multi-valued attributes are shown as one attribute. What is the primary key? ni#, lastname, firstnames, street, city, postcode, salary, bDate, gender, dNum, dName, mgrni#, mgrStartDate, pNum, pName, pLocation, hours MSc/Dip IT - ISD L18 Normalisation (425-456) 435 19/11/2009 MSc/Dip IT - ISD L18 Normalisation (425-456) 436 19/11/2009





PropertyID is chosen to be the primary key

County and Lot Number are a composite candidate key, but they were not chosen to be the primary key

Consider each combination of the composite primary key attributes.

Which other attributes depend on just a subset of them?

Repeat for any candidate keys



A Final Example : Medical Records

Patients have a unique NI#, each having a name and a GP;

Each GP has an id number, name & address;

- A hospital appointment connects a patient, a date, a hospital and a consultant;
- A consultant has a phone number and visits one hospital on any given day;

A hospital has an address.

MSc/Dip IT - ISD L18 Normalisation (425-456)

The following example assumes a patient will only have one appointment on a particular day, and that consultant names are unique.

449

- The Universal Relation (choose primary key)
- U (**ni#, appDate**, appTime, patName, gp#, gpName, gpAddress, cName, cPhone, hosp,hospAdd)

F	rom a particular patient on a particular day we can determine the consultant
	information, the time of the appointment and which was the hospital:
	ni#, appDate \rightarrow cName, cPhone, appTime, hosp, hospAdd

Given the consultant, appointment date and time, we can determine the patient: $cName, appDate, appTime \rightarrow ni#$

Each consultant is at one hospital on any day: $cName,appDate \rightarrow hosp, hospAddress$

Each consultant has one phone number

 $cName \rightarrow cPhone$

Func	ctional Dependencies
From a patient's national insura and GP	ance number, we can determine the patient's name
ni# → patName, gp	<mark># , gpAddress, gpName</mark>
From a GP's ID we determine t	the GP's name and address:
gp# → gpAddress, g	gpName
Each hospital has only one add	ress:
hosp → hospAddres	SS Contraction of the second
MSc/Dip IT – ISD L18 Normalisation (425-456)	450 19/11/2009
Nori	malising the Example

Starting with the Universal Relation we find all kinds of redundancy. U (**ni#, appDate**, appTime, patName, gp#, gpName, gpAddress, cName, cPhone, hosp,hospAdd)

So moving to 2NF, split off those attributes only dependent on part of the primary key:

Patient (**ni#**, patName, gp#, gpAddress, gpName) Appt (**ni#**, **appDate**, appTime, cName, cPhone, hosp, hospAddress)

451

19/11/2009

Patient is 2NF but not 3NF since there is a transitive dependency ni# → gp# → gpAdd, gpName So split off the GP information Pat (ni#, patName, gp#) GP (gp#, gpAdd, gpName)	ni#, appDate → cName → cPhone ni#, appDate → hosp → hospAdd Appt (ni# , appDate , appTime, cName, cPhone, hosp, hospAddress) Appt becomes App (ni# , appDate , apptime, cName, hosp) Con (cName , cPhone) Hospital (hosp , hospAdd)
Sc/Dip IT – ISD L18 Normalisation (425-456) 453	19/11/2009 MSc/Dip IT – ISD L18 Normalisation (425-456) 454 19/11/2009
	Final Tables

455

456