Computer Science 105 Introduction to Databases and Data Mining

Boston University, Fall 2021

Unit 1a: Database Fundamentals	
Course Overview	
Fundamental Facts About Data and Databases	. pre-lecture: 13 / in-lecture: 16
The Relational Model: Foundations; Primary and Foreign Keys.	
Constraints and Null Values; Designing a Database	
Unit 1b: The SQL Query Language	

Simple SELECT Commands	53 / 59
Pattern Matching; Comparisons with NULL; Removing Duplicates; Aggregates	
Subqueries; GROUP BY and HAVING	
Data Types; Creating Tables and Inserting Rows	100 / 104
Cartesian Product; Joins	110 / 121
Outer Joins	128 / 131
Other Commands; Practice with Queries	141 / 145
More Practice with Queries	

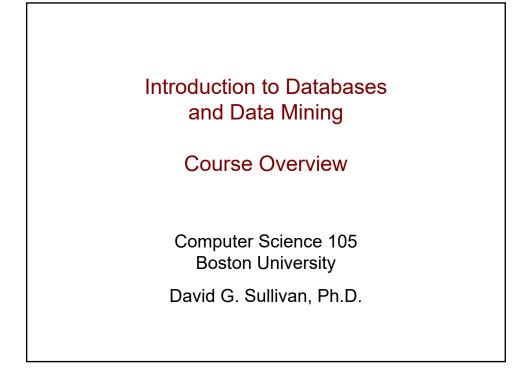
Unit 2: Writing Programs Using Python	
Getting Started; Programming Building Blocks	153 / 165
Built-in Functions; User Input; List Basics; Loops	177 / 185
Writing Functions; Cumulative Computations	
Making Decisions; Working with Numbers	
Working with Strings and Lists	
Using Objects; Splitting and Joining Strings	
Accessing a Database from Python	
Review: Strings and Lists; Accessing a Databases	
Working with Text Files; File Writing	
Reading Text Files	
File-Reading Revisited	

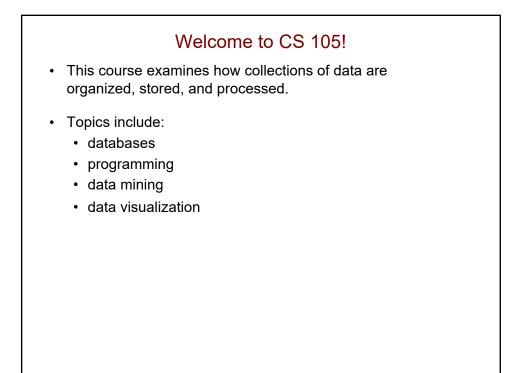
Unit 3: Data Visualization

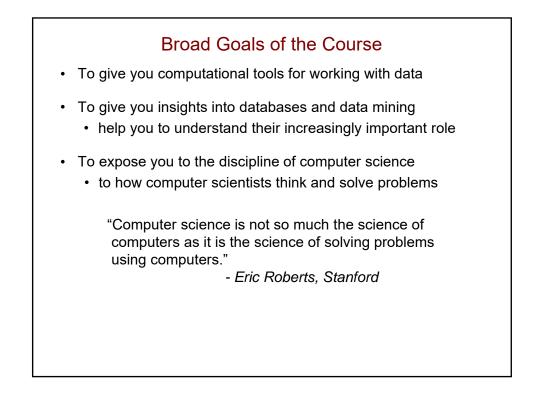
Notes for this unit will be provided separately.

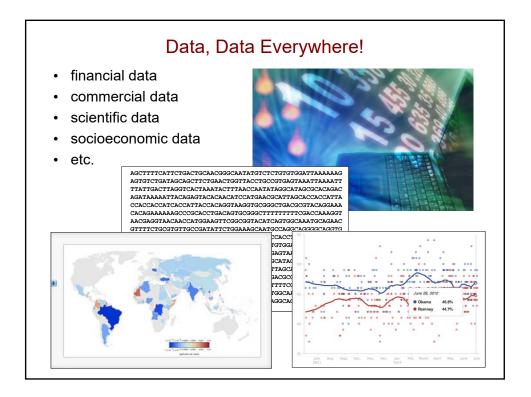
Unit 4: Data Mining

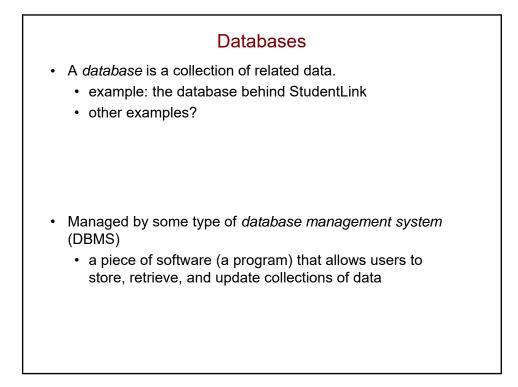
Fundamentals	
Evaluating a Model Learned in Data Mining; More Fundamentals	301 / 307
Classification Learning Using the 1R Algorithm; More on Evaluating Model	s 317 / 325
Classification Learning: Learning a Decision Tree	337 / 343
More Practice with Classification Learning	
Numeric Estimation; Using Weka	
Association Learning	
Discretizing Data	
Preparing Your Data	
Case Study: Predicting Patient Outcomes	

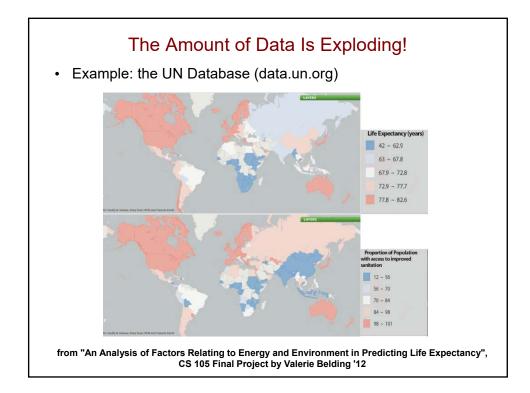






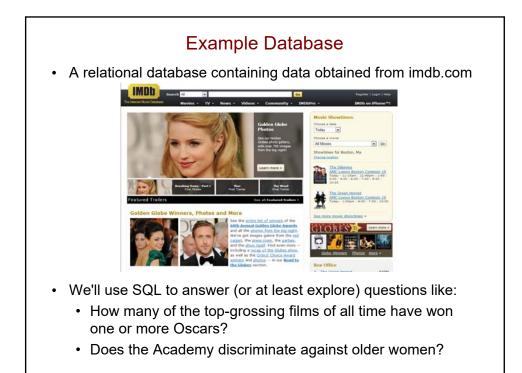


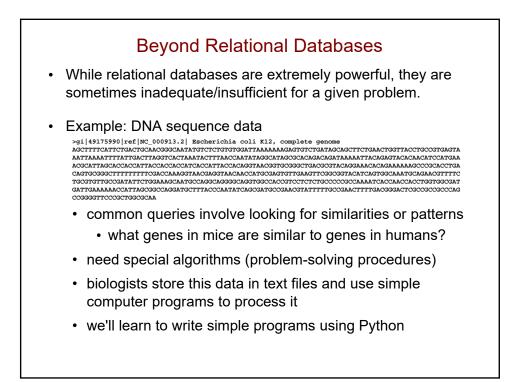


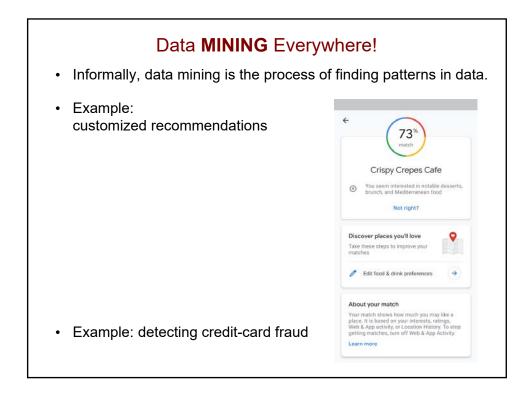


Relational Databases								
a way of • exam	organ ples: I soft A	nizing data kno IBM DB2, Ora ccess	anaged by a DBMS tha own as the <i>relational n</i> cle, Microsoft SQL Se	nodel. rver,				
 In the rel 	ationa	al model, data	is organized into table	es of <i>i</i>	ecords.			
			is organized into <i>table</i> ne or more <i>field</i> s	es of <i>i</i>	records.			
• each	record	d consists of o	•	es of <i>i</i>	records.			
• each	recorc ple: a	d consists of o	ne or more <i>fields</i>	es of <i>i</i>				
eachexam	recorc ple: a	d consists of o table of inforr	ne or more <i>field</i> s nation about students					
• each • exam <i>id</i>	record ple: a 678 J	d consists of o table of inforr name	ne or more <i>fields</i> nation about students <i>address</i>	class 2007	dob			
 each exam id 12345 	record ple: a 6678 J 5225 A	d consists of o table of inforr name Dill Jones	ne or more <i>fields</i> nation about students <i>address</i> warren Towers 100	class 2007	dob 3/10/85			
 each exam id 12345 25252 	record ple: a 6678 J 2525 A 6891 A	d consists of o table of inforr name Jill Jones Alan Turing	ne or more <i>fields</i> nation about students address Warren Towers 100 Student village A210	class 2007 2010 2008	dob 3/10/85 2/7/88			
 each exam id 12345 25252 33566 	record ple: a 678 J 525 A 891 A 3900 J	d consists of o table of inforr name Jill Jones Alan Turing Audrey Chu	ne or more <i>fields</i> nation about students address Warren Towers 100 Student Village A210 300 Main Hall	class 2007 2010 2008	dob 3/10/85 2/7/88 10/2/86			

SQL
 A relational DBMS has an associated <i>query language</i> called SQL that is used to:
define the tables
 add records to a table
 modify or delete existing records
 retrieve data according to some criteria
 example: get the names of all students who live in Warren Towers
 example: get the names of all students in the class of 2024, and the number of courses they are taking
 perform computations on the data
 example: compute the average age of all students who live in Warren Towers









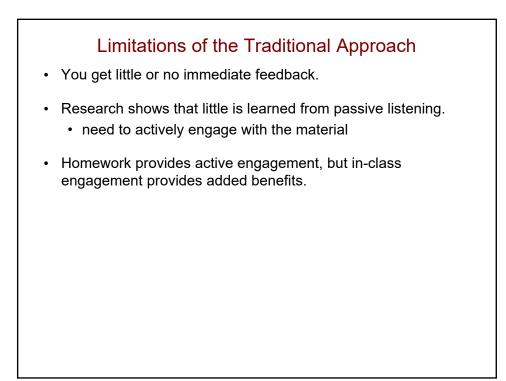
Structure of the Course

- databases (4 weeks)
- programming in Python (4 weeks)
- data graphics/visualization (1 week)
- data mining (4 weeks)

Course Materials

- Required: The CS 105 Coursepack
 - use it during pre-lecture and lecture need to fill in the blanks!
 - PDF version is available on Blackboard
 - recommended: get it printed
 - one option: FedEx Office (Cummington & Comm Ave)
- Required in-class software: Top Hat Pro platform
 - used for pre-lecture quizzes and in-lecture exercises
 - create your account and purchase a subscription ASAP (see Lab 0 for more details)

<section-header><section-header><section-header><list-item><list-item><list-item>

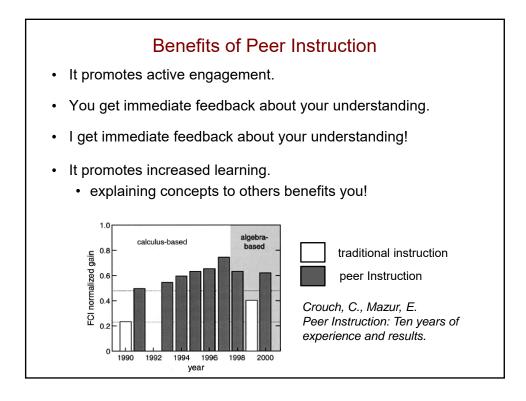


Lectures in this Class

- Based on an approach called *peer instruction*.
 - developed by Eric Mazur at Harvard

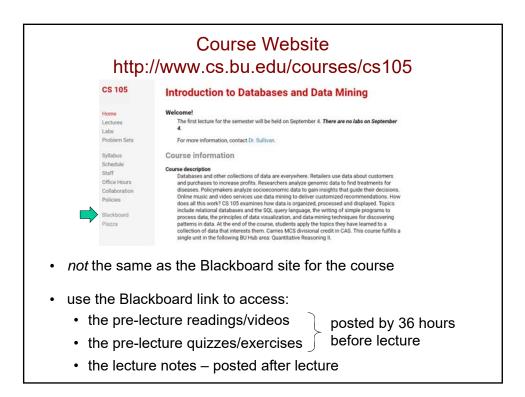
Basic process:

- 1. Question posed (possibly after a short intro)
- 2. Solo vote (no discussion yet)
- 3. Small-group discussions (in teams of 3)
 - explain your thinking to each other
 - · come to a consensus
- 4. Group vote
 - · each person in the group should enter the same answer
- 5. Class-wide discussion



Preparing for Lecture

- Short video(s) and/or readings
 - fill in the blanks as you watch the videos!
- Short online reading quiz or other exercise
 - complete **by 1 p.m.** of the day of lecture (unless noted otherwise)
 - · won't typically be graded for correctness
 - · your work should show that you've prepared for lecture
 - no late submissions accepted
- Preparing for lecture is essential!
 - · gets you ready for the lecture questions and discussions
 - · we won't cover everything in lecture

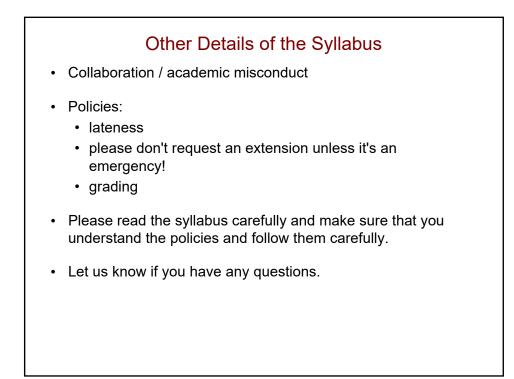


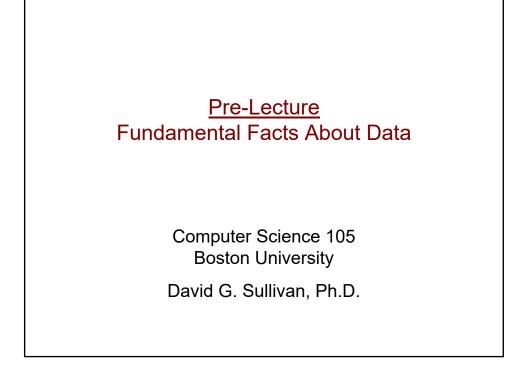
Labs

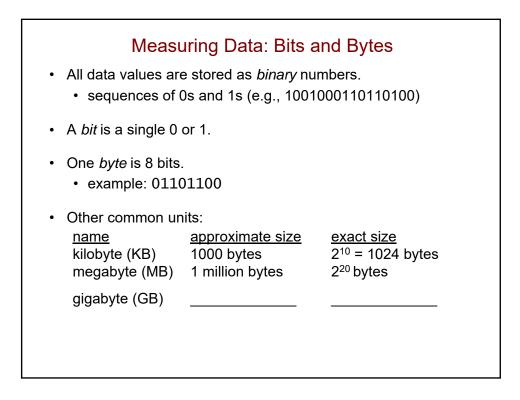
- Attendance is required
 - begin next week
- Will help you prepare for and get started on the assignments
- · Will also reinforce essential skills
- ASAP: Complete Lab 0 (on the course website)
 - setup Top Hat account/subscription
 - setup a CS account before your first lab session
 - · some other tasks to prepare you for the semester

Requirements / Grading Preparation and participation (10%) lecture preparation attendance/participation – full credit if you: make 85% of the votes over the entire semester attend 85% of the labs Nine homework assignments (25%) Final project (10%): done in teams of three Three quizzes (25%) Final exam (30%)

Course Staff • Instructor: Dave Sullivan (dgs@cs.bu.edu) • Teaching fellow • Office hours and contact info. will be available on the main course Web site: http://www.cs.bu.edu/courses/cs105 • For questions on content, homework, etc.: • use Piazza • send e-mail to cs105-staff@cs.bu.edu

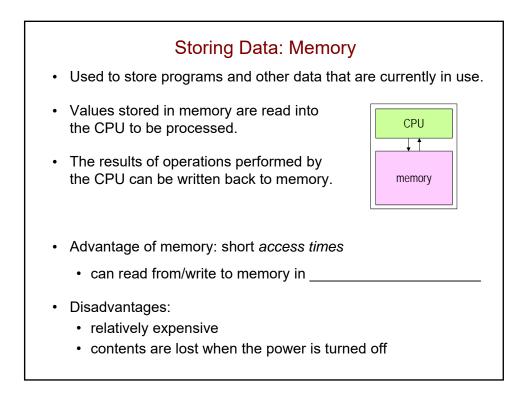


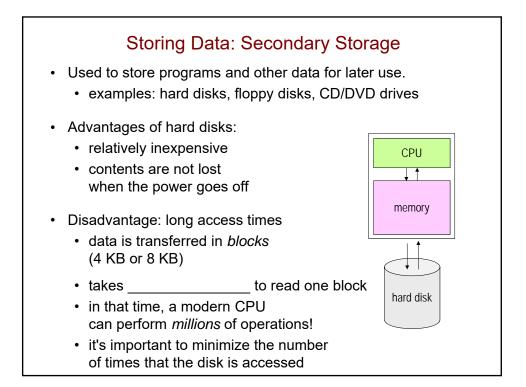


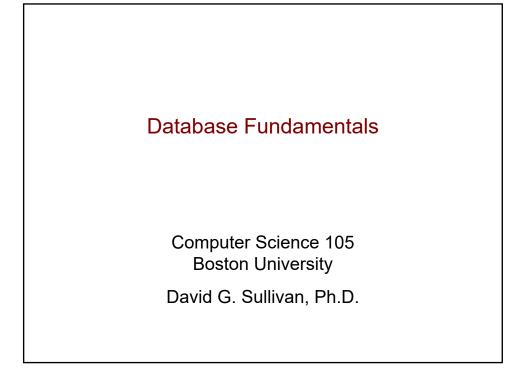


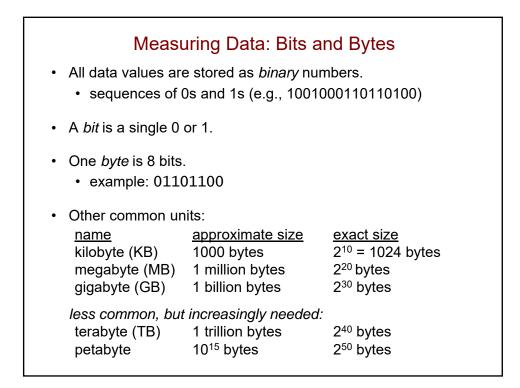
Processing Data: the CPU

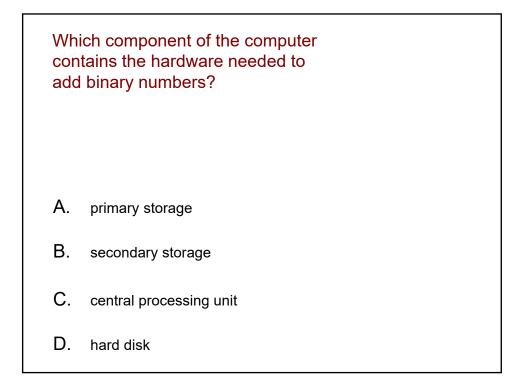
- At the heart of every computer is its CPU.
 - short for central processing unit
- Includes hardware for processing data stored in binary form.
 - · example: a circuit for adding two binary numbers
- The CPU can only store a small amount of data at a time.
 - · the values it is currently processing

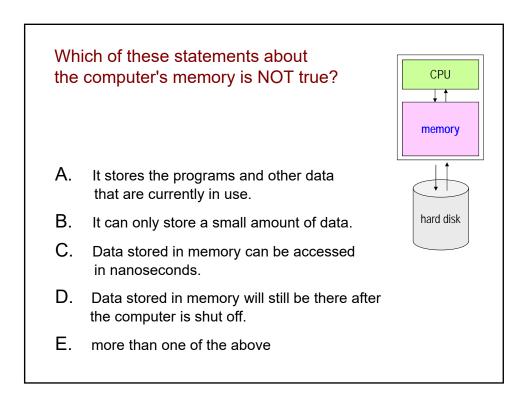


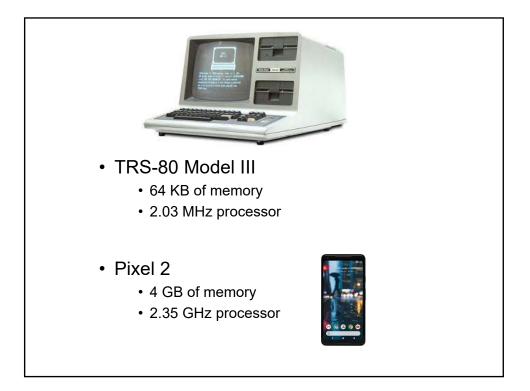


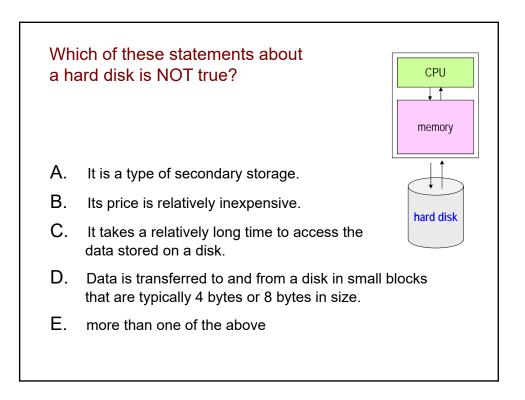


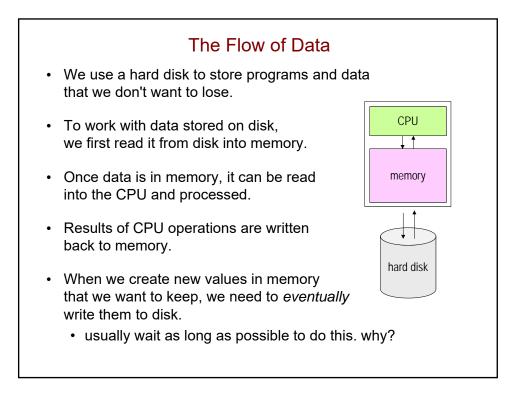


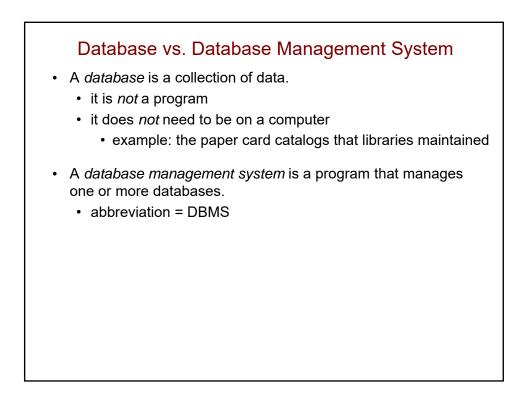


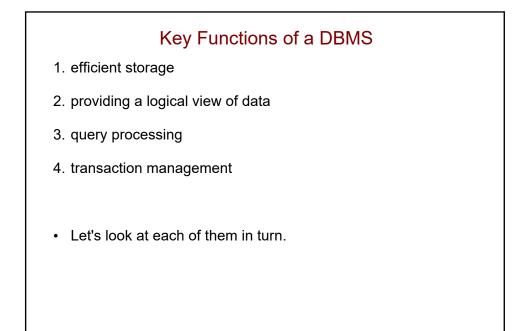


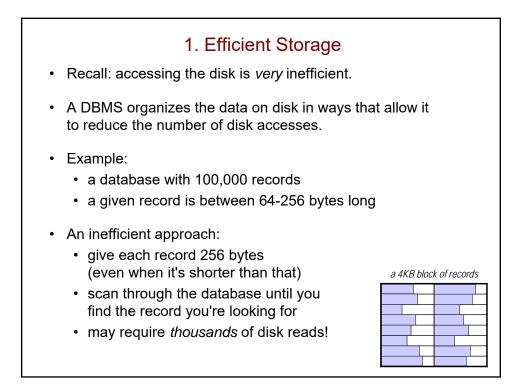


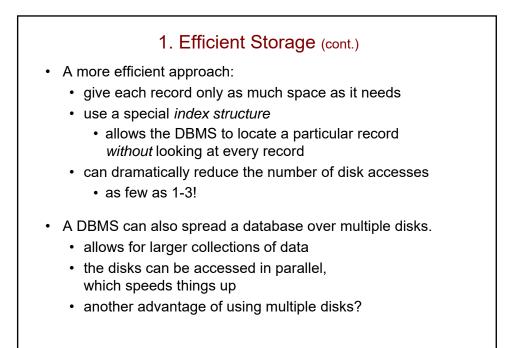


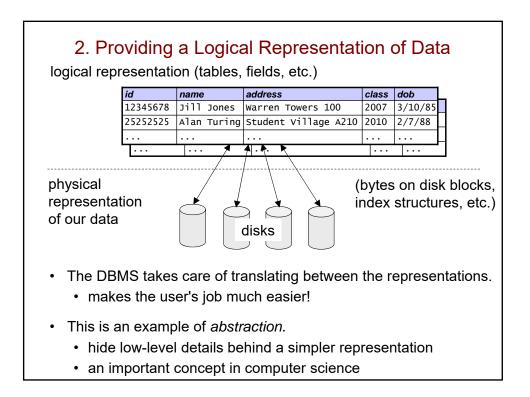






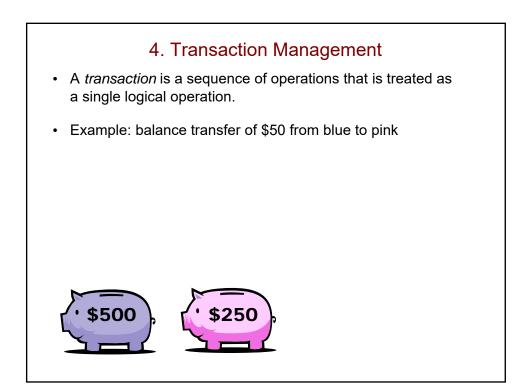


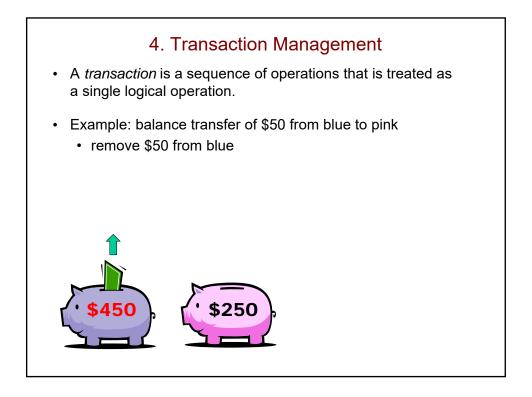


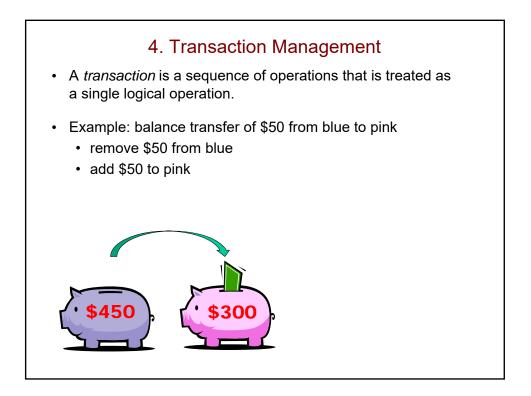


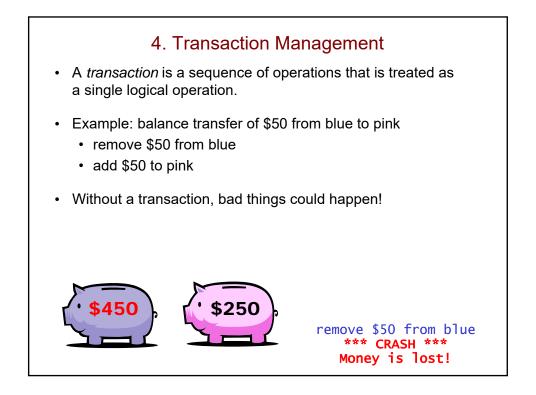
3. Query Processing

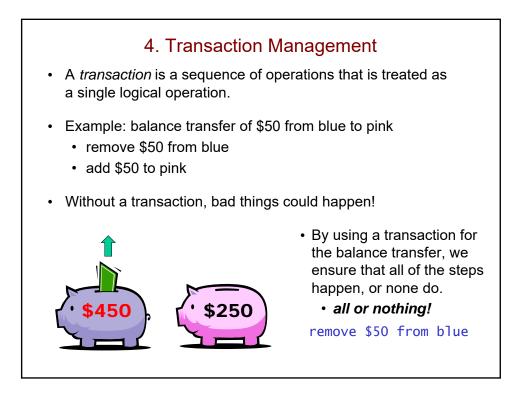
- A DBMS has some type of query language.
 - example: SQL
 - · includes commands for:
 - · adding new records
 - modifying or deleting existing records
 - · retrieving data according to some criteria
- The DBMS performs the low-level steps needed to execute a given command.

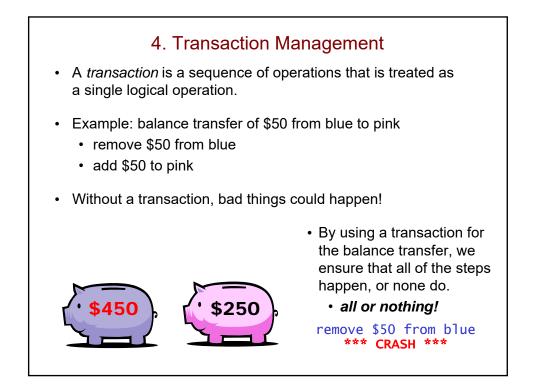


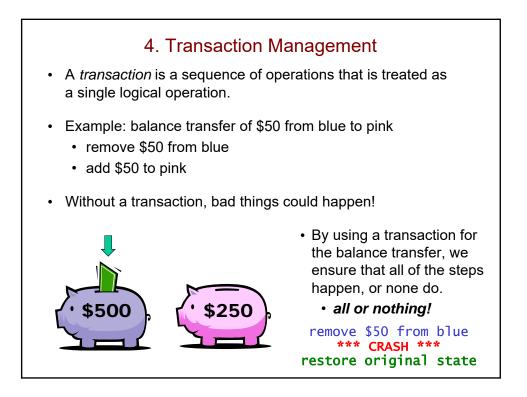


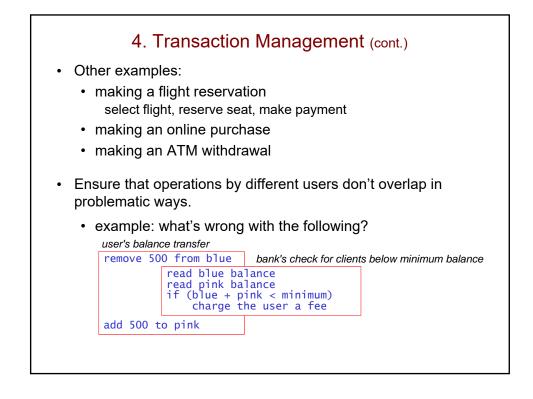


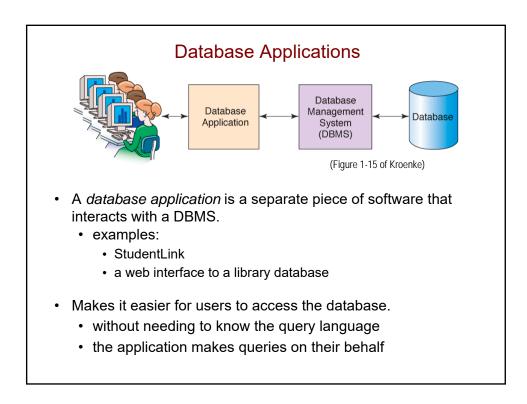


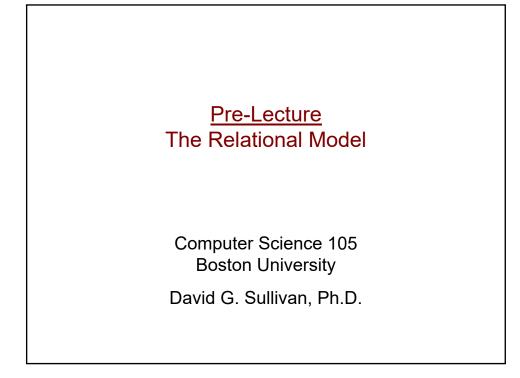












What Is a Data Model?
 A formal way of describing: pieces of data (<i>data items</i>) relationships between data items constraints on the values of data items
 We'll focus on the relational model – the dominant data model in current database systems.

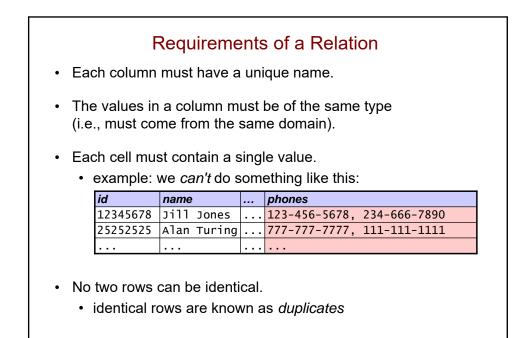
The Relational Model: Basic Concepts

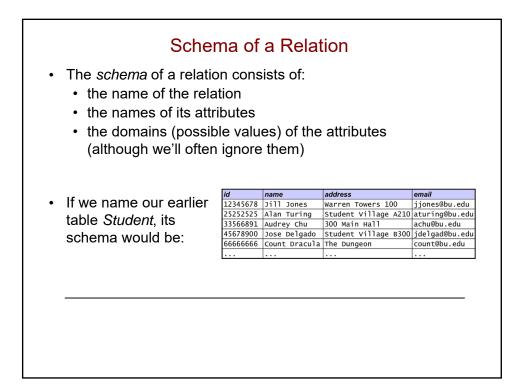
- A database consists of a collection of tables.
- Example of a table:

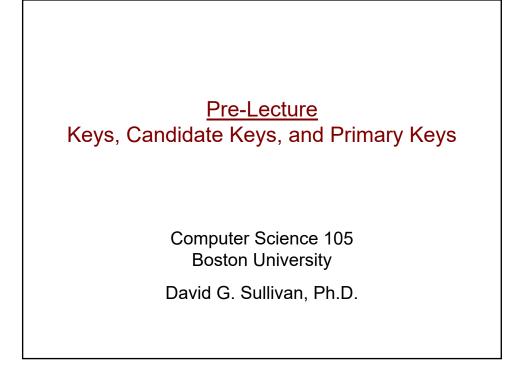
id	name	email	
12345678	Jill Jones	Warren Towers 100	jjones@bu.edu
25252525	Alan Turing	Student Village A210	aturing@bu.edu
33566891	Audrey Chu	300 Main Hall	achu@bu.edu
45678900	Jose Delgado	Student Village B300	jdelgad@bu.edu
66666666	Count Dracula	The Dungeon	count@bu.edu

- Each row in a table holds data that describes either:
 - an *entity* (a person, place, or thing!)
 - a relationship between two or more entities
- Each *column* in a table represents one attribute of an entity.

Re	lational Model: Terminology
Two sets of te	rminology:
table	=
row	=
column	=
 We'll use both 	sets of terms.







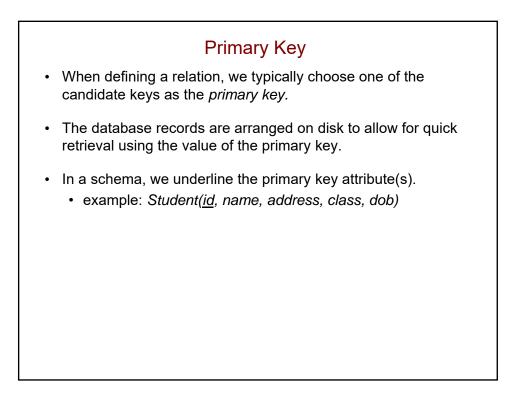
to ur • a	 to uniquely identify a row in a relation. allows us to distinguish one row from another 								
	id	name		email					
	12345678	Jill Jones		jjones@bu.edu					
	25252525	Alan Turing		aturing@bu.edu					

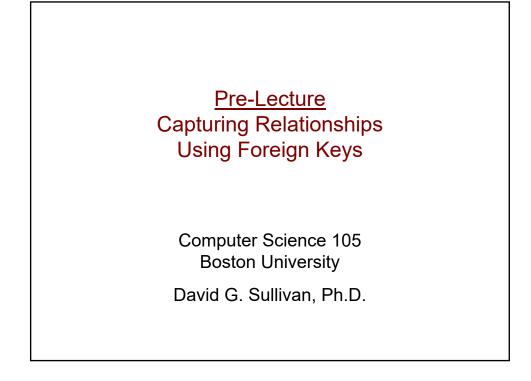
Candidate Key

- A *candidate key* is a *minimal* collection of attributes that is a key.
 - minimal = no unnecessary attributes are included

Candidate Key (cont.) Consider a table describing the courses in which students ٠ are enrolled: course CS 105 student credit_status 12345678 ugrad 25252525 CS 111 ugrad 45678900 CS 460 grad CS 105 33566891 non-credit 45678900 CS 510 grad

 Conside are enr 	er a table o		e Key (cont.) he courses in w	hich students					
	student	course	credit_status						
	12345678	CS 105	ugrad						
	25252525	CS 111	ugrad						
	45678900	CS 460	grad						
	33566891	CS 105	non-credit						
	45678900	CS 510	grad						
	key? candidate key?								
student	student								
student, c	student, course								
student, c	student, course, credit status								



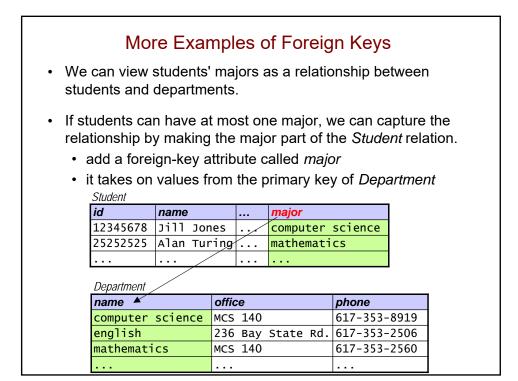


	Relations and Keys									
 Let's satisfies 	 Let's say that we have the following relations: 									
Stu	Student(id, name, address, email)									
id		name		address			email			
123	45678	Jill Jor	nes	Warren T	owe	rs 100	jjone	es@bu.edu		
252	52525	Alan Tu	ring	Student '	vil	lage A210	atur	ing@bu.edu		
Fac	culty(<u>id</u> ,	name, offi	ice, phoi	ne)						
id	nan	ne		office	phe	one				
1111	11 Ted	Codd		MCS 207	617	7-353-1111	1			
555	55 Gra	се Норре	er	MCS 222	617	7-353-5555	í.			
777	77 Edg	ar Dijks	stra	MCS 266	617	7-353-7777	7			
						•				
Dej	partmer	nt <u>(name</u> , or	ffice, ph	one)				_		
nam	e		office			phone				
com	puter	science	MCS 14	10		617-353-8	3919			
eng	lish		236 Ва	ay State I	Rd.	617-353-2	2506			
matl	hemati	cs	MCS 14	10		617-353-2	2560			
								•		

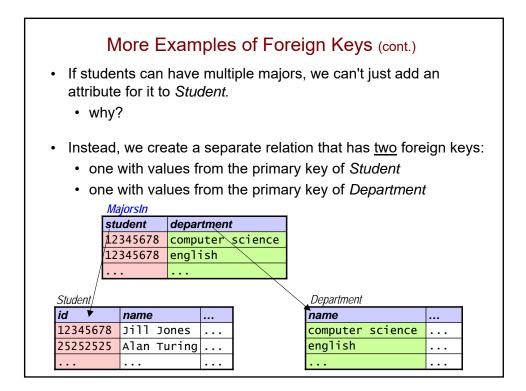
Capturing Relationships • In addition to storing info. about entities, we also use relations to capture relationships between two or more entities.

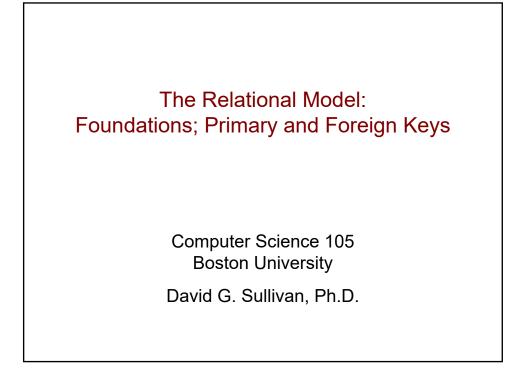
			Capturing	Re	latio	nsł	nips	(cont.))			
•			onship we mig students and th				pture	is the	relationship			
•	to inclu	ude	o so by expand an attribute ca the faculty ID	alled	advis	sor .						
	Student											
	id		name		advis			'				
	123456	578	Jill Jones		11111	L		es' advisor is Ted Codd.				
	252525	525	Alan Turing					Alan Tu	ıring's advisor is			
								Edgar [Dijkstra.			
	Faculty											
	id	nan	ne	offic	office		phone					
	11111	тес	Codd	MCS	207	617	-353-	1111				
	55555	Gra	ce Hopper	MCS	222	617	-353-	5555				
	77777	Edg	ar Dijkstra	MCS	266	617	-353-	7777				

	Fo	reig	n Ke	eys		
Student						
id	name		advise	or		
123456	78 Jill Jones		11111	-		
252525	25 Alan Turing	1	55555	;		
Faculty	name	offic	20	pho	200	1
111111	Ted Codd		207		-353-1111	
55555	Grace Hopper		-		-353-5555	
77777	Edgar Dijkstra		266	-	-353-7777	
		·				
advisor	is an example of	a for	oian l			



If students car attribute for it t	n have m	ultiple m	Ū	<mark>≺eys</mark> (cont.) can't just add a	n
• why?					
id	name		major		
12345678	Jill Jon	es	computer	science,englis	h
25252525	Alan Tur	ing	mathemati	CS	
Department		•			
name		office		phone	
computer :	science	MCS 140		617-353-8919	
<mark>english</mark>		236 Bay	State Rd.	617-353-2506	
mathematic	cs	MCS 140		617-353-2560	
				••••	





		Which of these states	aanta	about	4
		Which of these staten			L
		the Movie table is N		ue?	
	Мо∨	vie(id, <u>name, year</u> , rat	ting,	runti	me)
	id	name	year	rating	runtime
	12345	Star Wars: The Force Awakens	2015	PG-13	138
	78910	Avatar	2009	PG-13	162
	23232	Titanic	1997	PG-13	194
	90210	Finding Dory	2016	PG	97
	55555	Toy Story 3	2010	G	103
	01111	Ocean's Eleven	2001	PG-13	116
А. В.		er name for the Movie table lovie table has five tuples.	e is the	Movie	relatio
C.	The p	rimary key of Movie is the c	combin	ation (I	name, y
D.	more	than one of the above			

id	name	year	rating	runtimes
12345	Star Wars: The Force Awakens	2015	PG-13	138 , 150
78910	Avatar	2009	PG-13	162, 194
23232	Titanic	1997	PG-13	194
90210	Finding Dory	2016	PG	97
55555	Toy Story 3	2010	G	103
01111	Ocean's Eleven	2001	PG-13	116
			he ver	sions.

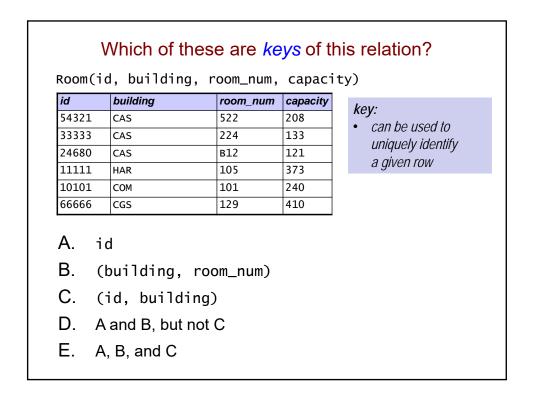
id	name	year	-	runtime1	
12345	Star Wars: The Force Awakens		PG-13		150
78910	Avatar	2009	PG-13	-	194
23232	Titanic	1997	PG-13	194	
90210	Finding Dory	2016	PG	97	
55555	Toy Story 3	2010	G	103	
01111	Ocean's Eleven	2001	PG-13	116	

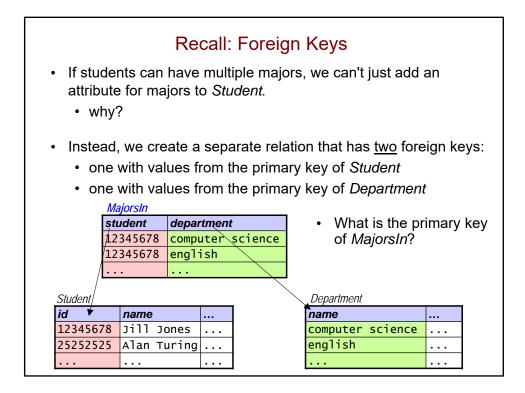
Here's a relation with info about rooms on campus...

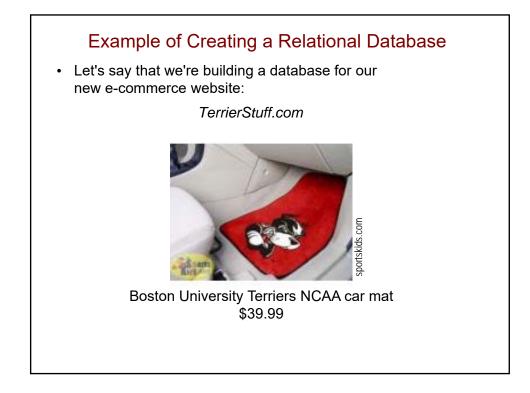
Room(id,	building,	room_num,	capacity)
----------	-----------	-----------	-----------

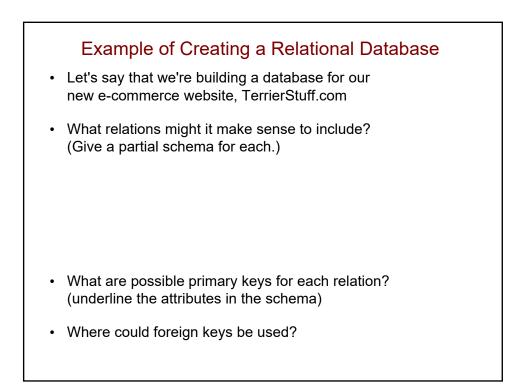
id	building	room_num	capacity
54321	CAS	522	208
33333	CAS	224	133
24680	CAS	в12	121
11111	HAR	105	373
10101	СОМ	101	240
66666	CGS	129	410

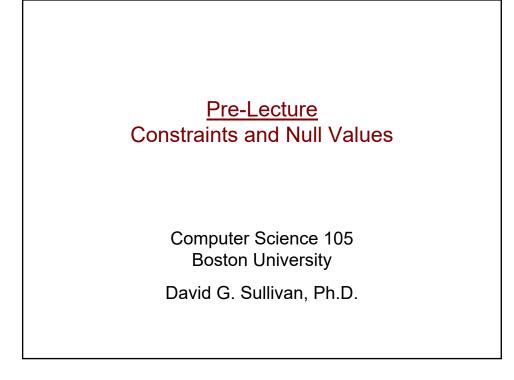
id	d, building, ı <i>building</i>	room_num	capacity	
54321	CAS	522	208	candidate key:
33333	CAS	224	133	can be used to
24680	CAS	в12	121	uniquely identify
11111	HAR	105	373	a given row
10101	СОМ	101	240	• none of the attributes
66666	CGS	129	410	are unnecessary
C. (d building, ro id, building and B, but not)		

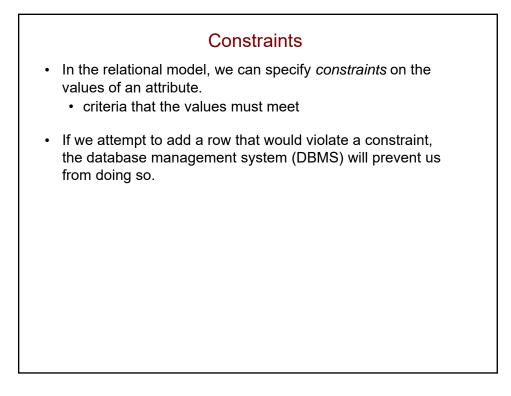












Uniqueness Constraints

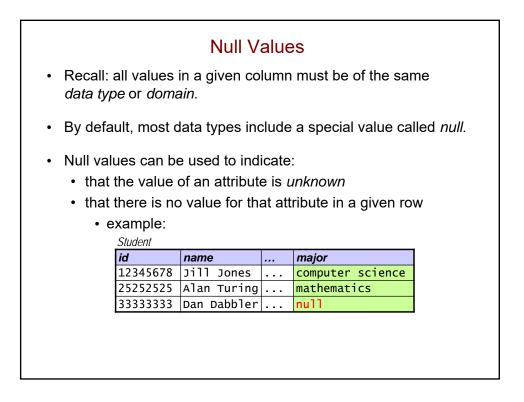
Student(id, name, address, email)

id	name	address	email
12345678	Jill Jones	Warren Towers 100	jjones@bu.edu
25252525	Alan Turing	Student Village A210	aturing@bu.edu
33566891	Audrey Chu	300 Main Hall	achu@bu.edu
45678900	Jose Delgado	Student Village B300	jdelgad@bu.edu
66666666	Count Dracula	The Dungeon	count@bu.edu

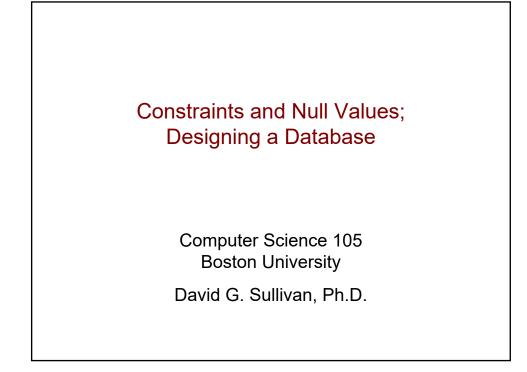
- When we specify a primary key, the DBMS imposes a *uniqueness constraint* on those attribute(s).
 - each row must have unique value(s) for those attribute(s)
 - example: we <u>can't</u> add this row to Student: (25252525, Alex Hamilton, 45B Smith Hall, aham@bu.edu)
 - could we add this row?
 (4444444, Jill Jones, Student Village A100, jill44@bu.edu)

Movie(<u>name</u> , <u>year</u> , rating, runt	ime)		
name	year	rating	runtime
Star Wars: The Force Awakens	2015	PG-13	138
Avatar	2009	PG-13	162
Titanic	1997	PG-13	194
Finding Dory	2016	PG	97
Toy Story 3	2010	G	103
Ocean's Eleven	2001	PG-13	116
he primary key is a combination			

		F	Referential I	nte	grity	Сс	onstr	aints	5
	Student								
i i	id		name		adviso	or			
-	123456	78	Jill Jones	·	11111	-			
1	252525	25	Alan Turing		55555	;			
/	Faculty								
i	id 🔺	nan	ne	offic	ce	pho			
1	11111	тес	Codd	MCS	207	617	-353-	1111	
	55555	Gra	ce Hopper	MCS	222	617	-353-	5555	
7	77777	Edg	ar Dijkstra	MCS	266	617	-353-	7777	
re	eferent	<i>tial</i> ore	pecify a foreig integrity const ign key attribu in the corresp	<i>raint</i> te(s)	on th must	nose tak	e attril ce on	bute(s value	s).
•	(33	333	es: can we ado 3333, Alex Ha 3333, Alex Ha	milto	n,,	222	222)	dent?	?



Null Values (cont.) • We can't put a null value in a primary-key column. • We *can* put a null value in a foreign-key column. • even though null is not in the corresponding primary key · allows us to indicate the absence of a relationship Student id name advisor 11111 12345678 Jill Jones 25252525 Alan Turing ... 77777 48484848 Alex Hamilton ... null Faculty office id name phone 11111 Ted Codd MCS 207 617-353-1111 MCS 222 617-353-5555 55555 Grace Hopper 77777 Edgar Dijkstra MCS 266 617-353-7777 We can also tell the DBMS that we don't want a given column • to include any null values.



id	name	dob	pob
000007	Humphrey Bogart	1899-12-25	New York, NY, USA
0000030	Audrey Hepburn	1929-05-04	Brussels, Belgium
0000133	Geena Davis '79	1956-01-21	Wareham, MA, USA
0000151	Morgan Freeman	1937-06-01	Memphis, TN, USA
0000158	Tom Hanks	1956-07-09	Concord, CA, USA
0000194	Julianne Moore'83	1960-12-03	Fayetteville, NC, USA
	over 2400 peop	le (actors an	d directors)!

Person	(<u>id</u> , name, dob, p	(ac	
id	name	dob	pob
0000007	Humphrey Bogart	1899-12-25	New York, NY, USA
0000030	Audrey Hepburn	1929-05-04	Brussels, Belgium
0000133	Geena Davis '79	1956-01-21	Wareham, MA, USA
0000151	Morgan Freeman		Memphis, TN, USA
0000158	Tom Hanks	1956-07-09	Concord, CA, USA
0000194	Julianne Moore'83	1040 12 02	
		1900-12-03	rayetteviile, NC, USA
A. ((B. (4 C. ((0000007, James Bond, 1444444, Morgan Fre 0000030, Audrey Hep	1920-11-11 eman, 1937-	, London) 06-01, Memphis)
A. ((B. (4 C. (() 0000007, James Bond, 1444444, Morgan Fre	1920-11-11 eman, 1937-	, London) 06-01, Memphis)

movie_id	person_id	type	year
663202	0000138	BEST-ACTOR	2016
8170832	0488953	BEST-ACTRESS	2016
8682448	0753314	BEST-SUPPORTING-ACTOR	2016
810819	2539953	BEST-SUPPORTING-ACTRESS	2016
663202	0327944	BEST-DI RECTOR	2016
895587	NULL	BEST-PI CTURE	2016
—	takes on v	alues from the id column in th rales from the id column in th ole tells us the 2016 Best Act	e Perso

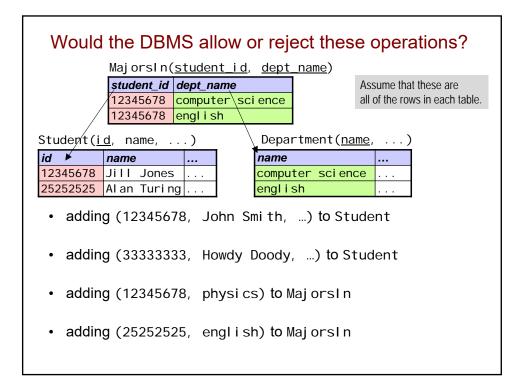
What are movie_id and person_id examples of?

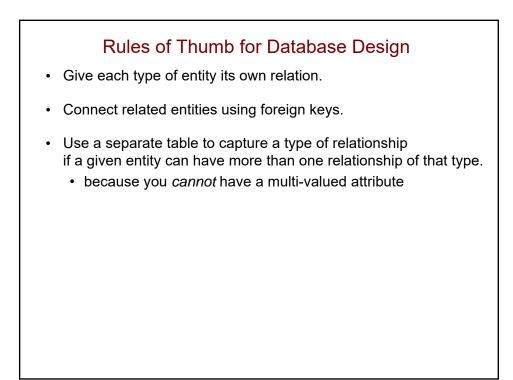
movie_id	person_id	type	year
1663202	0000138	BEST-ACTOR	2016
3170832	0488953	BEST-ACTRESS	2016
3682448	0753314	BEST-SUPPORTING-ACTOR	2016
0810819	2539953	BEST-SUPPORTI NG-ACTRESS	2016
1663202	0327944	BEST-DI RECTOR	2016
1895587	NULL	BEST-PI CTURE	2016
	s NULL me	ean? of a string (a piece of text)!	

movie_id	person_id	type	year
1663202	0000138	BEST-ACTOR	2016
3170832	0488953	BEST-ACTRESS	2016
3682448	0753314	BEST-SUPPORTING-ACTOR	2016
0810819	2539953	BEST-SUPPORTING-ACTRESS	2016
663202	0327944	BEST-DI RECTOR	2016
895587	NULL	BEST-PI CTURE	2016
/hat abou	ut (person_	ork as the primary key? _id, year)? _id, type, year)?	

movie_id	person_id	type	year
1663202	0000138	BEST-ACTOR	2016
3170832	0488953	BEST-ACTRESS	2016
3682448	0753314	BEST-SUPPORTING-ACTOR	2016
0810819	2539953	BEST-SUPPORTING-ACTRESS	2016
1663202	0327944	BEST-DI RECTOR	2016
1895587	1111111	BEST-PI CTURE	2016

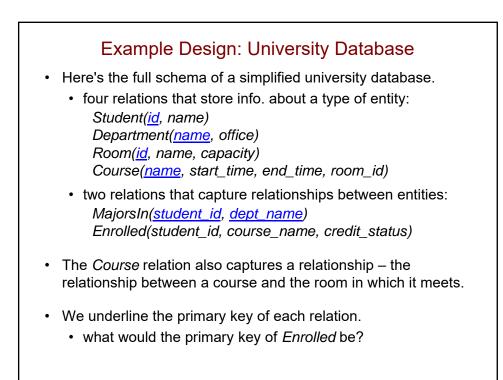
Oscar(n	novie_id,	<u>person_id, type</u> , <u>year</u>)	
movie_id	person_id	type	year
1663202	0000138	BEST-ACTOR	2016
3170832	0488953	BEST-ACTRESS	2016
3682448	0753314	BEST-SUPPORTI NG-ACTOR	2016
0810819	2539953	BEST-SUPPORTI NG-ACTRESS	2016
1663202	0327944	BEST-DI RECTOR	2016
1895587	1111111	BEST-PI CTURE	2016
. (77777	77, 000013	38, BEST-ACTOR, 2017)	
. (22222	22, 11111	11, BEST-ACTRESS, 2016)	
. (44444	44, 04889	53, BEST-ACTRESS, 2016)	
. A and	C, but no	t B	
	,		





Rules of Thumb for Database Design

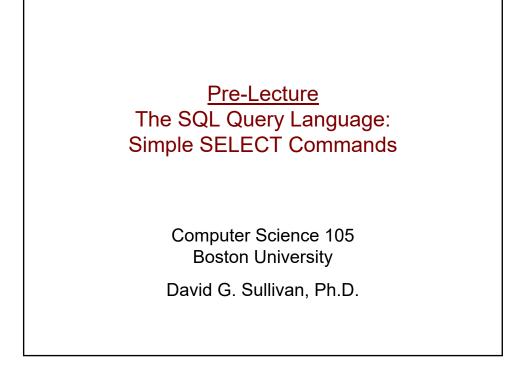
- Give each type of entity its own relation.
- Connect related entities using foreign keys.
- Use a separate table to capture a type of relationship if a given entity can have more than one relationship of that type.
 - because you cannot have a multi-valued attribute



Foreign Keys in the University Database

Student(<u>id</u>, name) Department(<u>name</u>, office) Room(<u>id</u>, name, capacity) Course(<u>name</u>, start_time, end_time, room_id) MajorsIn(<u>student_id</u>, <u>dept_name</u>) Enrolled(<u>student_id</u>, <u>course_name</u>, credit_status)

- Foreign keys we've already discussed:
 - *student_id* in *MajorsIn* (takes on values from *id* in *Student*)
 - *dept_name* in *MajorsIn* (takes on values from *name* in *Department*)
- What other foreign keys make sense?
 - •
 - •
 - •



Student	
id	name
12345678	Jill Jones
25252525	Alan Turing
33566891	Audrey Chu
45678900	Jose Delgado
66666666	Count Dracula

Course

name	start_time	end_time	room_id
CS 105	13:00:00	14:00:00	4000
CS 111	09:30:00	11:00:00	5000
EN 101	11:00:00	12:30:00	1000
CS 460	16:00:00	17:30:00	7000
CS 510	12:00:00	13:30:00	7000
PH 101	14:30:00	16:00:00	NULL

Enrolled

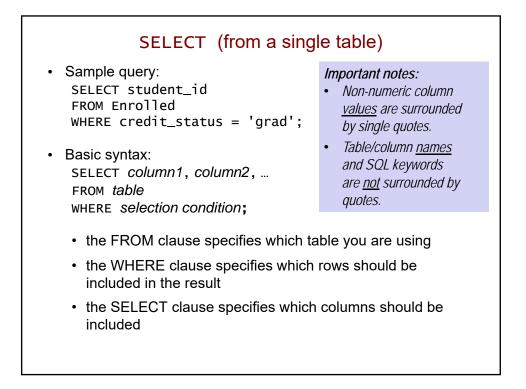
student_id	course_name	credit_status
12345678	CS 105	ugrad
25252525	CS 111	ugrad
45678900	CS 460	grad
33566891	CS 105	non-credit
45678900	CS 510	grad

Room		
id	name	capacity
1000	CAS Tsai	500
2000	CAS BigRoom	100
3000	EDU Lecture Hall	100
4000	CAS 315	40
5000	CAS 314	80
6000	CAS 226	50
7000	MCS 205	30

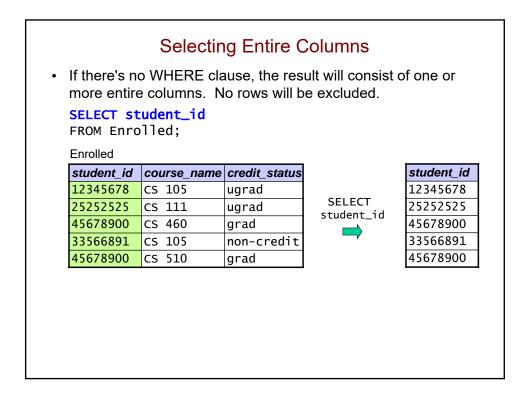
Department

name	office
comp sci	MCS 140
mathematics	MCS 140
the occult	The Dungeon
english	235 Bay State Road

MajorsIn student_id dept_name 12345678 comp sci 45678900 mathematics 25252525 comp sci 45678900 english 66666666 the occult



• Example: SELECT st FROM Enro	udent_id	om a single	e table) (con	t.)
	dit_status	= 'grad';		
Enrolled				
student_id	course_name	credit_status		
12345678	CS 105	ugrad		
25252525	CS 111	ugrad		
45678900	CS 460	grad		
33566891	CS 105	non-credit		
45678900	CS 510	grad		
	credit_stat		SELECT	
student_id	course_name		student_id	student_id
45678900	CS 460	grad		45678900
45678900	CS 510	grad	,	45678900



	Select	ing Entire	Rows
	the result to i we use a * in		e rows (i.e., all of↑ Γ clause:
SELECT * FROM Enro	lled dit_status	- 'arad'.	
Enrolled	uit_status	= yrau ,	
student_id	course_name	credit_status	
12345678	CS 105	ugrad	
25252525	CS 111	ugrad	
45678900	CS 460	grad	
33566891	CS 105	non-credit	
45678900	CS 510	grad	
WHERE	credit_stat	us = 'grad'	;
student_id	course_name	credit_status	
45678900	CS 460	grad	
45678900	CS 510	grad	

The WHERE Clause

SELECT column1, column2, ... FROM table WHERE selection condition;

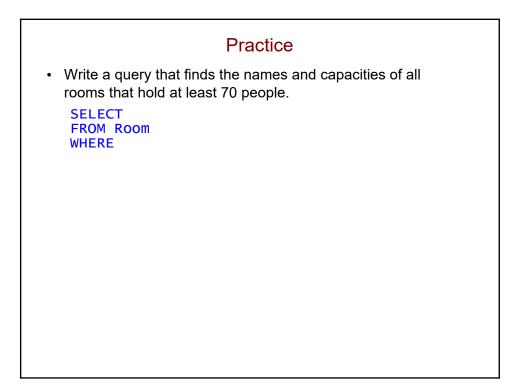
- The selection condition must be an expression that evaluates to either true or false.
 - example: credit_status = 'grad'
 - can include any column from the table(s) in the FROM clause
- The results of the SELECT command will include only those tuples for which the selection condition evaluates to true.

	Simple Comparisons					
	est selection condition is a comparison that uses following <i>comparison operators</i> :					
<u>operator</u> < >	<u>name</u> less than greater than					
<= >=	less than or equal to greater than or equal to					
= !=	equal to not equal to					

Practice

• Write a query that finds the names and capacities of all rooms that hold at least 70 people.

SELECT FROM WHERE



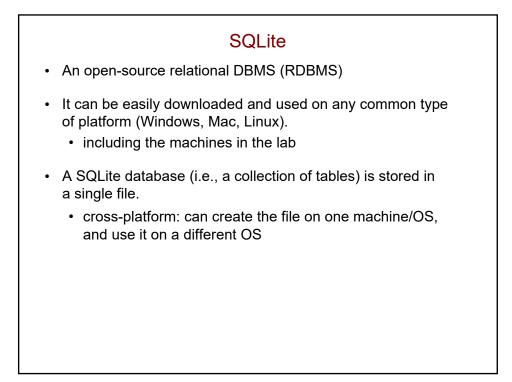
	Practice Write a query that finds the names and capacities of all rooms that hold at least 70 people.					
	SELEC	ст				
	FROM					
	WHERE	Ξ				
Г	id	name	capacity			
	1000	CAS Tsai	500			
	2000	CAS BigRoom	100			
	3000	EDU Lecture Hall	100			
	4000	CAS 315	40			
	5000	CAS 314	80			
	6000	CAS 226	50			
	7000	MCS 205	30			
-	Ţ			•		
Г	id	name	capacity	1	name	capacity
	1000	CAS Tsai	500		CAS Tsai	500
	2000	CAS BigRoom	100		CAS BigRoom	100
-	3000	EDU Lecture Hall	100		EDU Lecture Hall	100
	5000	CAS 314	80	r r	CAS 314	80

The SQL Query Language: Simple SELECT Commands

Computer Science 105 Boston University

David G. Sullivan, Ph.D.

Why Learn SQL? Desktop database systems like Access provide tools for manipulating data in a database. However, these tools don't allow you to perform all possible types of queries. For more flexibility and power, we use SQL. a query language In addition, knowledge of SQL is needed to perform queries from within a program.



De De	Browser for SQL	ite - C:\Users\dgs\Do	cuments\e66\ass	ignments\ps1\mov	ieDB\mov —	o ×	
File E	dit View Help						
6 Ne	w Database 🛛 😹	Open Database	Write Changes	Severt Changes			
Data	abase Structure	Browse Data Edit Pr	agmas Execut	e SQL			
Tabl	e: 🗾 Movie		• 🚳 📓		New Record	Delete Record	
	id	name	year	rating	runtime	ge ^	
	Filter	Filter	Filter	Filter	Filter	Filter	
1	2488496	Star Wars: Th		PG-13	138	A	
2	0499549	Avatar	2009	PG-13	162	AVYS	
3	0120338	Titanic	1997	PG-13	194	DR	
4	0369610	Jurassic World	2015	PG-13	124	A	
5	0848228	The Avengers	2012	PG-13	143	A	
6	0468569	The Dark Kni	2008	PG-13	152	AT	
7	3748528	Rogue One	2016	PG-13	133	A	
8	2771200	Beauty and th	. 2017	PG	129	F	
9	2277860	Finding Dory	2016	PG	97	N ×	
	4 1 - 10 of 68	. 121 121		Go to:	1		
. A.	1 - 10 0F 08			GO to:	4		

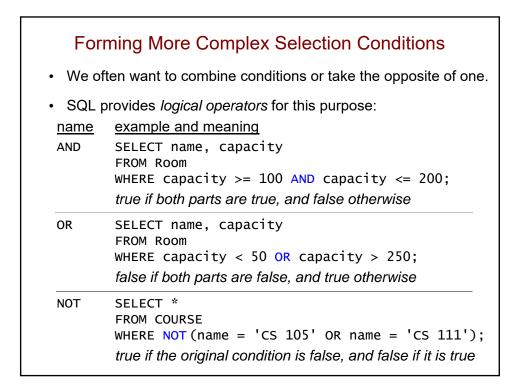
released in 2010?						
_ L		name		year	rating	runtime
	2488496	Star Wars: The Force Awak	ens	2015	PG-13	138
	1228705	Iron Man 2		2010	PG-13	124
	0120338	Titanic		1997	PG-13	194
	0435761	Toy Story 3		2010	G	103
Ī	1323594	Despicable Me		2010	PG	95
Ī	0240772	Ocean's Eleven		2001	PG-13	116
İ						
				OM Mo∨ LECT y		2010;
-	SELECT FROM M		FR	LECT * OM Mov ERE ye	-	010;

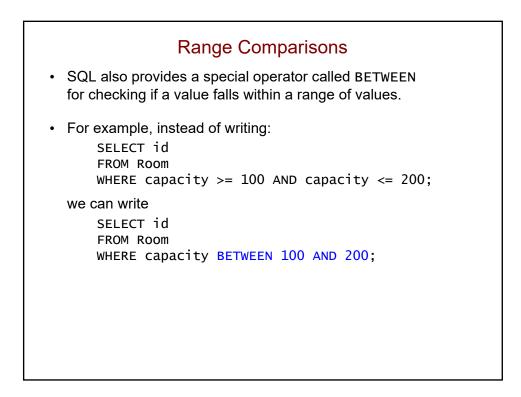
1228705 Iron Man 2 2010 PG-13 124 0120338 Titanic 1997 PG-13 194 0435761 Toy Story 3 2010 G 103 1323594 Despicable Me 2010 PG 95 0240772 Ocean's Eleven 2001 PG-13 116	id 2488496	<i>name</i> Star Wars: The Force Awakens	year 2015	rating	runtime
D120338 Titanic 1997 PG-13 194 D435761 Toy Story 3 2010 G 103 L323594 Despicable Me 2010 PG 95 D240772 Ocean's Eleven 2001 PG-13 116					
L323594 Despicable Me 2010 PG 95 0240772 Ocean's Eleven 2001 PG-13 116					
L323594 Despicable Me 2010 PG 95 0240772 Ocean's Eleven 2001 PG-13 116	0435761	Toy Story 3	2010	G	103
0240772 Ocean's Eleven 2001 PG-13 116			2010	PG	95
			2001	PG-13	116

Movie					
id	name	year	rating	runtime	
2488496	Star Wars: The Force Awakens	2015	PG-13	138	
1228705	Iron Man 2	2010	PG-13	124	
0120338	Titanic	1997	PG-13	194	
0435761	Toy Story 3	2010	G	103	
1323594	Despicable Me	2010	PG	95	
0240772	Ocean's Eleven	2001	PG-13	116	

How could we get the *name and runtime* of movies released *before* 2010?

	name	year	rating	runtime
2488496	Star Wars: The Force Awakens	2015	PG-13	138
1228705	Iron Man 2	2010	PG-13	124
0120338	Titanic	1997	PG-13	194
0435761	Toy Story 3	2010	G	103
1323594	Despicable Me	2010	PG	95
0240772	Ocean's Eleven	2001	PG-13	116





1228705	Star Wars: The Force Awakens	2015	PG-13	120	
			PG-13	138	
0120220	Iron Man 2	2010	PG-13	124	
0120338	Titanic	1997	PG-13	194	
0435761	Toy Story 3	2010	G	103	
<pre>FROM Movie WHERE name = 'Titanic' AND name = 'Toy Story 3'; SELECT name, runtime FROM Movie WHERE name = 'Titanic' OR name = 'Toy Story 3'; SELECT name, runtime</pre>					

id	name		id	name	capacity
12345678	Jill Jone		1000	CAS Tsai	500
			2000	CAS BigRoom	100
25252525	Alan Turi	ng	3000	EDU Lecture Hall	100
33566891	Audrey Ch	าน	4000	CAS 315	40
45678900	Jose Delg	gado	5000	CAS 314	80
66666666	Count Dra	acula	6000	CAS 226	50
			7000	MCS 205	30
Course			-	•	

name

start_time	enu_ume	100111_10
13:00:00	14:00:00	4000
09:30:00	11:00:00	5000
11:00:00	12:30:00	1000
16:00:00	17:30:00	7000
12:00:00	13:30:00	7000
14:30:00	16:00:00	NULL
	13:00:00 09:30:00 11:00:00 16:00:00 12:00:00	13:00:00 14:00:00 09:30:00 11:00:00 11:00:00 12:30:00 16:00:00 17:30:00 12:00:00 13:30:00

Enrolled

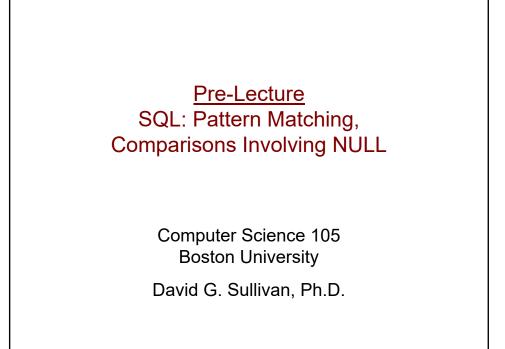
student_id	course_name	credit_status
12345678	CS 105	ugrad
25252525	CS 111	ugrad
45678900	CS 460	grad
33566891	CS 105	non-credit
45678900	CS 510	grad

comp sci	MCS 140				
mathematics	MCS 140				
the occult	The Dungeon				
english	235 Bay State Road				
MaiorsIn					

office

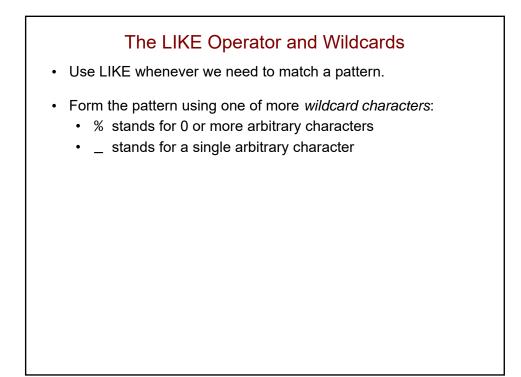
majorani	
student_id	dept_name
12345678	comp sci
45678900	mathematics
25252525	comp sci
45678900	english
66666666	the occult

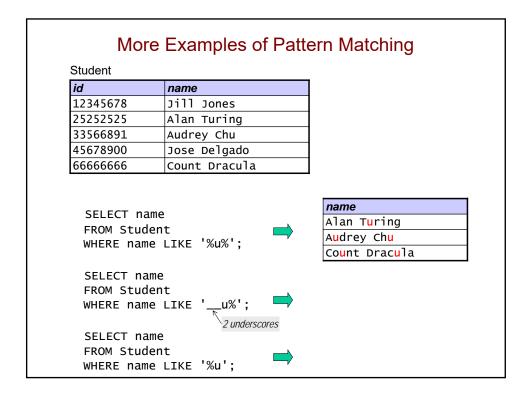
Practice with Simple SQL Queries Write a query that finds all information about CAS 315. Write a query that lists the names and start times of all courses. Write a query that gets the ID numbers of student(s) who are taking CS 105 for undergraduate (ugrad) credit.



						Dec					
Our simp	ole ı	universit	y databa	ase		Roc	om	name		00	pacity
Student						100	0	CAS TS		50	
id		name				200	-	CAS IS		10	-
12345678		Jill Jon	~~			300	-		cture Hal		-
				_		400	-	CAS 31		40	-
25252525	4	Alan Tur	ing				-		-		
33566891		Audrey C	hu			500	0	CAS 31	.4	80	
45678900		Jose Del	gado			600	0	CAS 22	6	50	
66666666		Count Dr	-	-		700	0	MCS 20	5	30	
		counc bi	acura			-					
Course											
name	star	rt time	end time		room id	De	partm	ient			
CS 105		00:00	14:00:00		1000	nam	1e		office		
cs 105		30:00	11:00:00		5000	com	p sci		MCS 140		
						mat	hemat	ics	MCS 140		
CS 460		00:00	17:30:00		7000	the	осси	1+	The Dung	00n	
CS 510	12:	:00:00 13:30:00) 7	7000				3		
CS 999	19:	30:00	21:30:00) (NULL	eng	lish		235 Bay	State Roa	.d
Enrolled								Mai	orsIn		
student_id		course_	name	cred	it_status				ent_id	dept_nam	е
12345678		CS 105		ugra	ad			1234	5678	comp sci	
25252525		CS 111		ugra	ad	1		4567	8900	mathemat	ics
45678900		CS 460		grad		-		2525	2525	comp sci	
33566891		CS 105		5	-credit			4567	8900	english	
						4		6666	6666	the occu	lt
45678900		CS 510		grad	ł			<u> </u>			

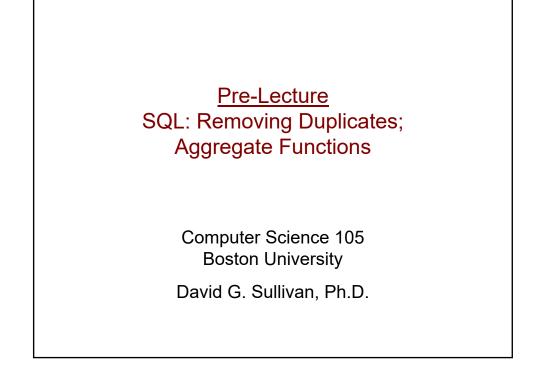
	F	attern M	latching)	
S	nameCAS TsaiCAS BigRoomEDU Lecture HallCAS 315CAS 314CAS 226MCS 205is won't work:ELECT name, capROM Room	40 80 50 30	and ca in CAS • the • nee	ay we want the apacities of all S. e names begin ed to find cours mes matching	rooms with 'CAS' ses with
•	HERE name = 'CA	AS';	Room		
			id	name	capacity
• Th	is will:		1000	CAS Tsai	500
ç	ELECT name, cap	acity	2000	CAS BigRoom	100
5	LELCI Hame, cap	ucicy	4000	CAS 315	
	ROM ROOM	_	4000		40
F	ROM ROOM HERE name LIKE	'CAS%' ·	⇒ 5000	CAS 314	40 80



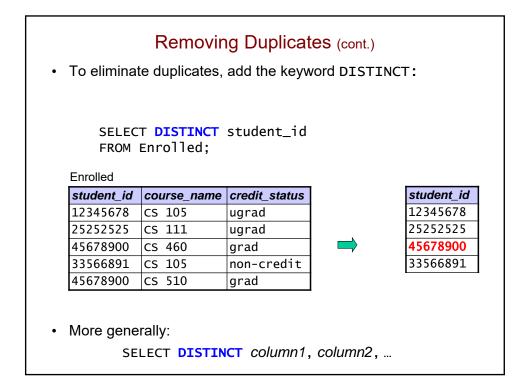


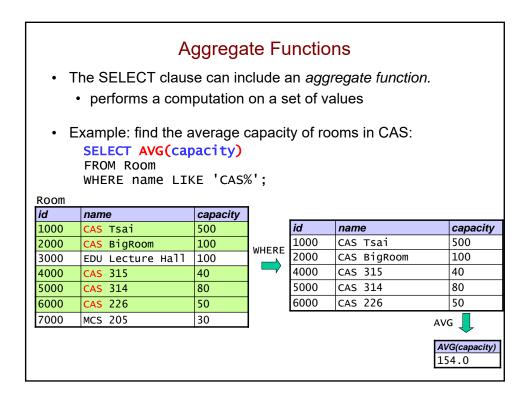
	Com	parison	s Involv	ving NULL
Course				
name	start_time	end_time	room_id	
CS 105	13:00:00	14:00:00	4000	
CS 111	09:30:00	11:00:00	5000	
CS 460	16:00:00	17:30:00	7000	
CS 510	12:00:00	13:30:00	7000	
CS 999	19:30:00	21:30:00	NULL	
This que SELE	uld we find ery produc CT name Course			nly courses?
WHER	E room_i	d = NUL	L;	

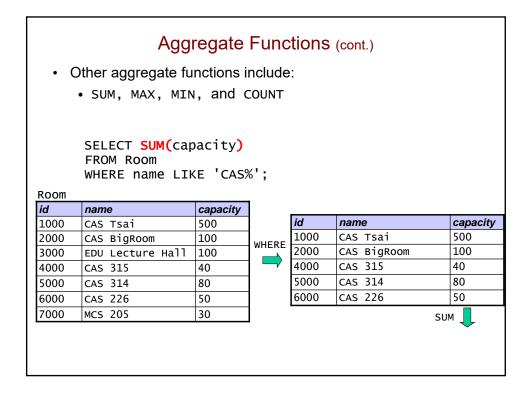
Comparisons Involving NULL Because NULL is a special value, any comparison involving NULL that uses the standard operators is always false. The following will *always* be false: room_id = NULL room_id != NULL NULL = NULL SQL provides special operators: IS NULL IS NOT NULL This query will find the online-only courses: SELECT name FROM Course WHERE room_id IS NULL;



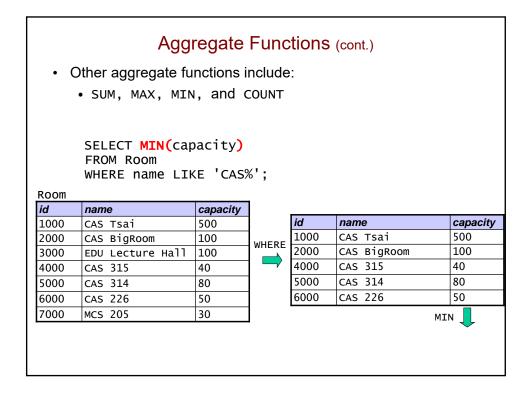
may includ	the relation p le duplicate tu e: find the IDs	ples.		
	T student_i Enrolled;	d		
student id	course name	credit_status	ľ	student_id
	CS 105	ugrad		12345678
25252525	CS 111	ugrad		25252525
45678900	CS 460	grad		45678900
33566891	CS 105	non-credit		33566891
45678900	CS 510	grad		45678900





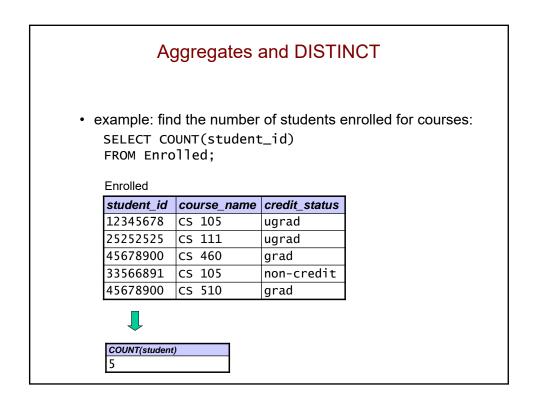


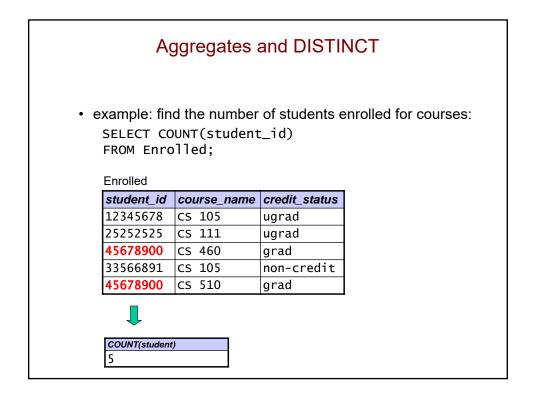
 Aggregate Functions (cont.) Other aggregate functions include: SUM, MAX, MIN, and COUNT 						
SELECT <mark>MAX(</mark> capacity) FROM Room WHERE name LIKE 'CAS%'; Room						
<u> </u>		i	-			
id	name	capacity				
<i>id</i> 1000	CAS Tsai	500		id	name	capacity
<i>id</i> 1000 2000	CAS Tsai CAS BigRoom	500 100	WHERE	1000	CAS Tsai	500
<i>id</i> 1000	CAS Tsai	500	WHERE	1000 2000	CAS Tsai CAS BigRoom	500 100
<i>id</i> 1000 2000	CAS Tsai CAS BigRoom	500 100	WHERE	1000	CAS Tsai	500
<i>id</i> 1000 2000 3000	CAS Tsai CAS BigRoom EDU Lecture Hall	500 100 100	WHERE	1000 2000	CAS Tsai CAS BigRoom	500 100
<i>id</i> 1000 2000 3000 4000	CAS Tsai CAS BigRoom EDU Lecture Hall CAS 315	500 100 100 40	WHERE	1000 2000 4000	CAS Tsai CAS BigRoom CAS 315	500 100 40

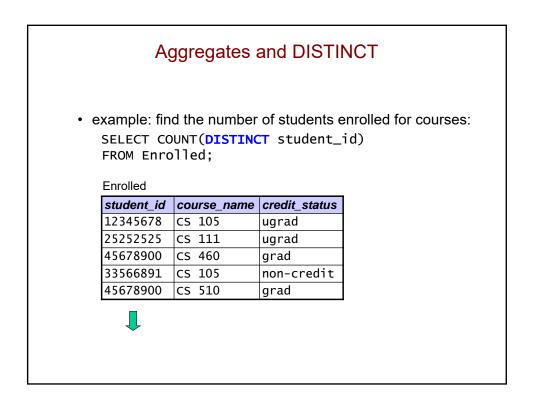


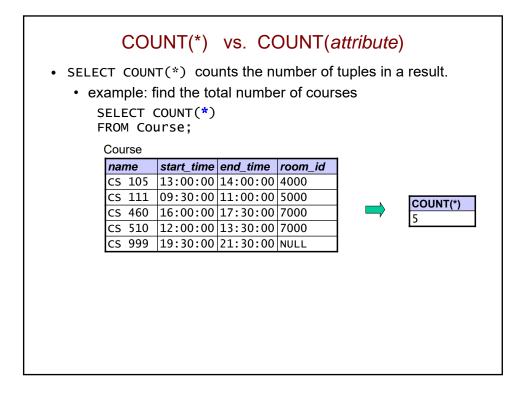
 Aggregate Functions (cont.) Other aggregate functions include: SUM, MAX, MIN, and COUNT 						
SELECT <mark>COUNT(</mark> capacity) FROM Room WHERE name LIKE 'CAS%'; Room						
Room			,			
id	name	capacity	, ,			
<i>id</i> 1000	name CAS Tsai	capacity 500	, ,	id	name	capacity
id	name	capacity	where	1000	CAS Tsai	500
<i>id</i> 1000	name CAS Tsai	capacity 500	-	1000 2000	CAS Tsai CAS BigRoom	500 100
<i>id</i> 1000 2000	name CAS Tsai CAS BigRoom	capacity 500 100	-	1000	CAS Tsai	500
<i>id</i> 1000 2000 3000	nameCAS TsaiCAS BigRoomEDU Lecture Hall	capacity 500 100 100	-	1000 2000	CAS Tsai CAS BigRoom	500 100
<i>id</i> 1000 2000 3000 4000	nameCAS TsaiCAS BigRoomEDU Lecture HallCAS 315	capacity 500 100 100 40	-	1000 2000 4000	CAS Tsai CAS BigRoom CAS 315	500 100 40

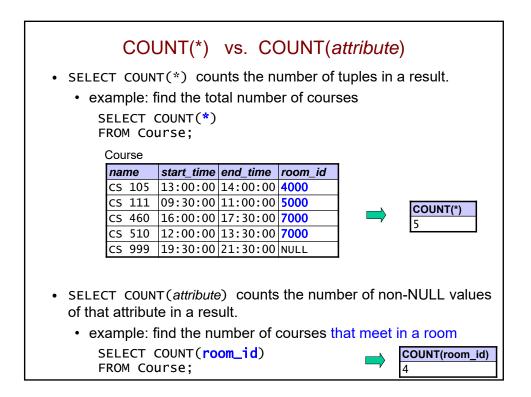
Ą	Aggregates and DISTINCT						
•	UNT(student		enrolled for courses:				
student_id	course name	credit_status					
12345678	CS 105	ugrad					
25252525	CS 111	ugrad					
45678900	CS 460	grad					
33566891	CS 105	non-credit					
45678900	CS 510	grad					

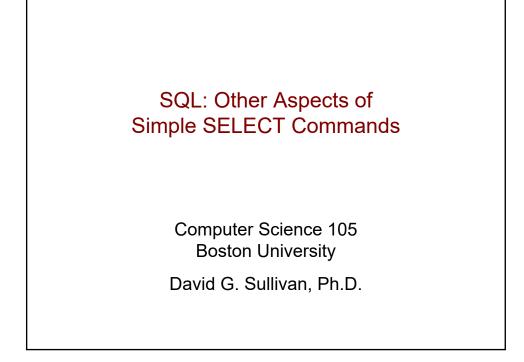




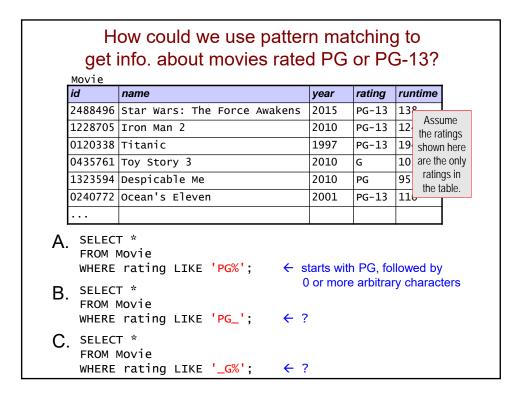








	How could we use pattern matching to get info. about movies rated PG or PG-13?						
	Movie <i>id</i>	name	year	rating	runtime		
	2488496	Star Wars: The Force Awakens	2015	PG-13	138		
	1228705	Iron Man 2	2010	PG-13	12 Assume the ratings		
	0120338	Titanic	1997	PG-13	19 shown here		
	0435761	Toy Story 3	2010	G	10 are the only		
	1323594	Despicable Me	2010	PG	95 ratings in the table.		
	0240772	Ocean's Eleven	2001	PG-13			
A	A. SELECT * FROM Movie WHERE rating LIKE 'PG%'; D. two of the queries at left would work						
B	FROM N			ee of t would	he queries work		
С	FROM N						



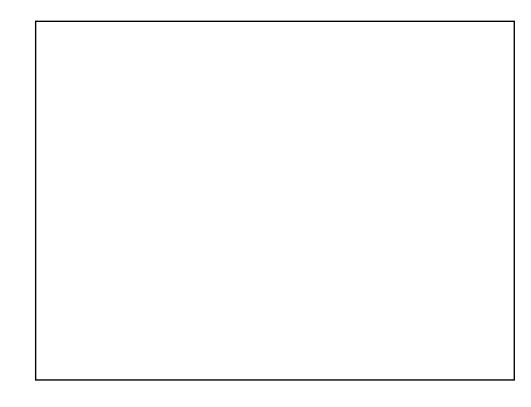
Movie <i>id</i>	name	year	rating	runtime
	Star Wars: The Force Awakens	2015	PG-13	138
	Iron Man 2	2010	PG-13	Assume
	Titanic	1997	PG-13	19 shown her
0435761	Toy Story 3	2010	G	10 are the onl
	Despicable Me	2010	PG	95 ratings in
0240772	Ocean's Eleven	2001	PG-13	the table.
SELEC FROM	Movie rating LIKE '%G%'; T *			

Pattern Matching (cont.)

- DBMSs typically have an operator that performs *case-insensitive* pattern matching.
 - not part of the SQL standard
 - · different implementations use different names for it

• In SQLite:

- the LIKE operator itself is case-insensitive
- there's no easy way to do case-sensitive pattern matching
- the = operator *is* case-sensitive



	How could we find the names of all courses without a room?						
	Course						
	name	start_time	end_time	room_id			
	CS 105	13:00:00	14:00:00	4000			
	CS 111	09:30:00	11:00:00	5000			
	EN 101	11:00:00	12:30:00	1000			
		16:00:00					
	CS 510	12:00:00	13:30:00	7000			
	PH 101	14:30:00	16:00:00	NULL			
A.	SELECT name FROM Course WHERE room_id =	'NULL';	D.		ore of the queries ould work		
В.	SELECT name FROM Course WHERE room_id =	NULL;	E.		the queries ould work		
C.	SELECT name FROM Course WHERE room_id I	S NULL;					

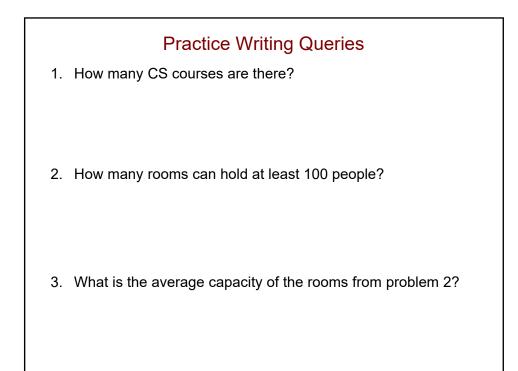
	How could we find the name all courses without a room					
	Course					
		name	start_time	end_time	room_id	
		CS 105	13:00:00	14:00:00	4000	
		CS 111	09:30:00	11:00:00	5000	
		EN 101	11:00:00	12:30:00	1000	
		CS 460	16:00:00	17:30:00	7000	
		CS 510	12:00:00	13:30:00	7000	
		PH 101	14:30:00	16:00:00	NULL	
A.	SELECT n FROM Cou WHERE ro	irse	'NULL';			
В.	SELECT n FROM Cou WHERE ro	irse	NULL;			
C.	SELECT n FROM Cou WHERE ro	irse	S NULL;			



	How could we determine					
	how many people have won Best Actor?					
	Oscar					
	movie_id	person_id	type	year		
	1663202	0000138	BEST-ACTOR	2016		
	3170832	0488953	BEST-ACTRESS	2016		
	3682448	0753314	BEST-SUPPORTING-ACTOR	2016		
	0810819	2539953	BEST-SUPPORTING-ACTRESS	2016		
	1663202	0327944	BEST-DIRECTOR	2016		
	1895587	NULL	BEST-PICTURE	2016		
<pre>A. SELECT COUNT(person_id) FROM Oscar WHERE type = 'BEST-ACTOR';</pre>						
<pre>B. SELECT TOTAL(person_id) FROM Oscar WHERE type = 'BEST-ACTOR';</pre>						
C.	SELECT C FROM OSC WHERE ty	ar	ST-ACTOR';			

movie_id	person_id	type	year
1663202	0000138	BEST-ACTOR	2016
3170832	0488953	BEST-ACTRESS	2016
3682448	0753314	BEST-SUPPORTING-ACTOR	2016
0810819	2539953	BEST-SUPPORTING-ACTRESS	2016
1663202	0327944	BEST-DIRECTOR	2016
1895587	NULL	BEST-PICTURE	2016
FROM OS	car	TINCT person_id) ST-ACTOR';	

movie_id	person_id	type	year
1663202	0000138	BEST-ACTOR	2016
3170832	0488953	BEST-ACTRESS	2016
3682448	0753314	BEST-SUPPORTING-ACTOR	2016
0810819	2539953	BEST-SUPPORTING-ACTRESS	2016
1663202	0327944	BEST-DIRECTOR	2016
1895587	NULL	BEST-PICTURE	2016
FROM OS		FINCT *)	



Student				
id	name			
12345678	Jill Jones			
25252525	Alan Turing			
33566891	Audrey Chu			
45678900	Jose Delgado			
66666666	Count Dracula			

Course

name	start_time	end_time	room_id
CS 105	13:00:00	14:00:00	4000
CS 111	09:30:00	11:00:00	5000
EN 101	11:00:00	12:30:00	1000
CS 460	16:00:00	17:30:00	7000
CS 510	12:00:00	13:30:00	7000
PH 101	14:30:00	16:00:00	NULL

Enrolled

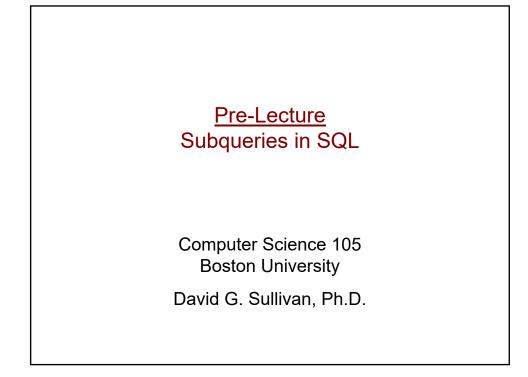
student_id	course_name	credit_status
12345678	CS 105	ugrad
25252525	CS 111	ugrad
45678900	CS 460	grad
33566891	CS 105	non-credit
45678900	CS 510	grad

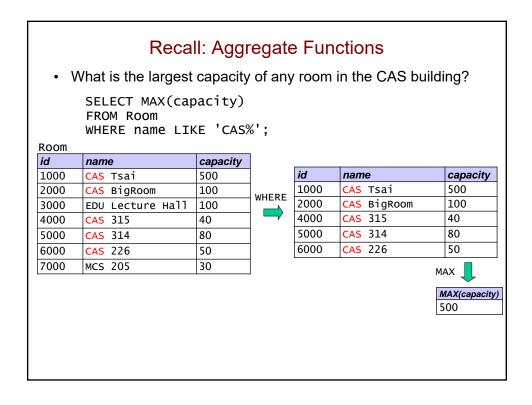
Room		
id	name	capacity
1000	CAS Tsai	500
2000	CAS BigRoom	100
3000	EDU Lecture Hall	100
4000	CAS 315	40
5000	CAS 314	80
6000	CAS 226	50
7000	MCS 205	30

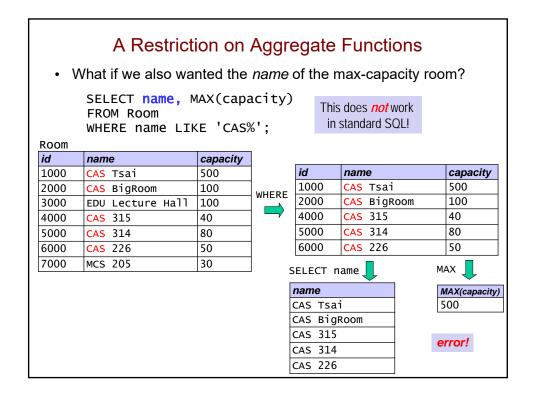
Department

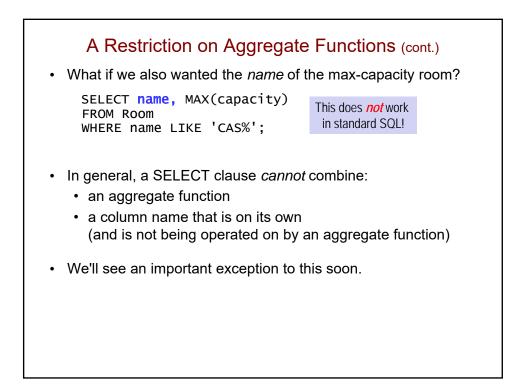
name	office
comp sci	MCS 140
mathematics	MCS 140
the occult	The Dungeon
english	235 Bay State Road

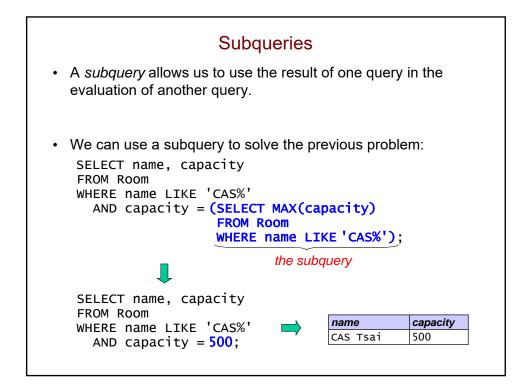
Student_id dept_name 12345678 comp sci 45678900 mathematics 25252525 comp sci 45678900 english 66666666 the occult

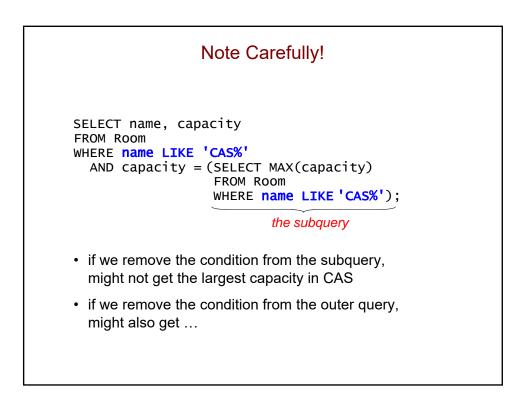


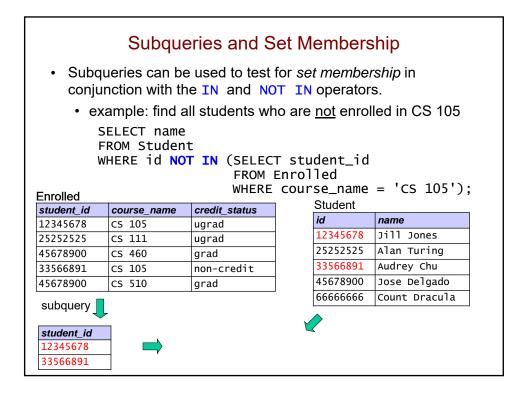


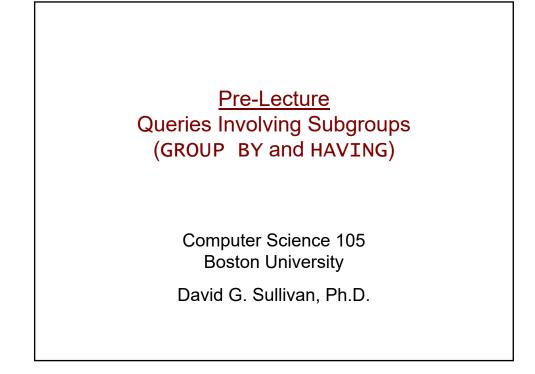


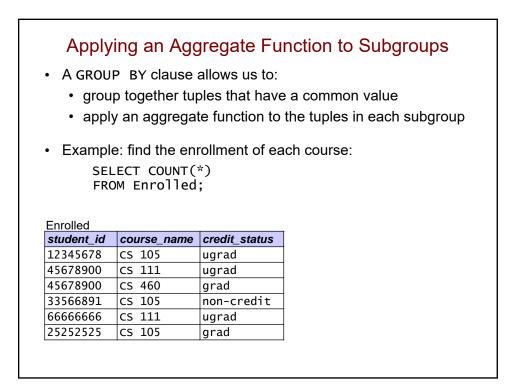


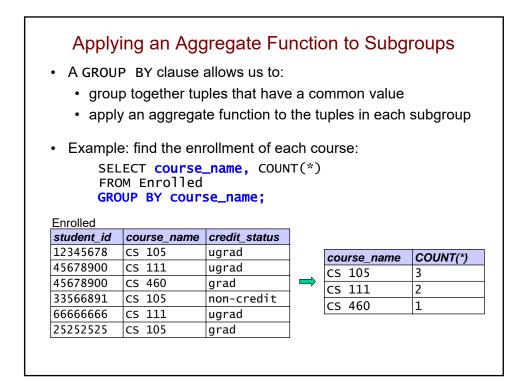


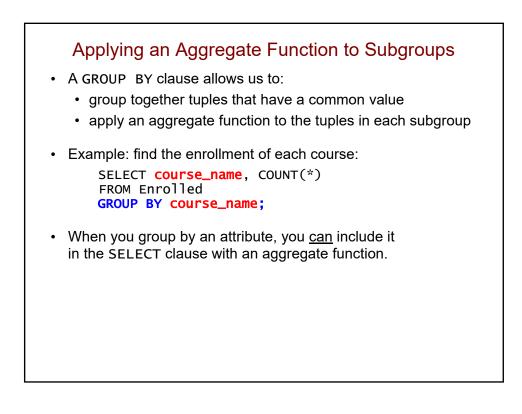








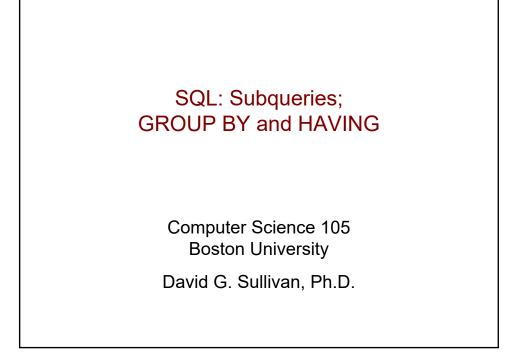




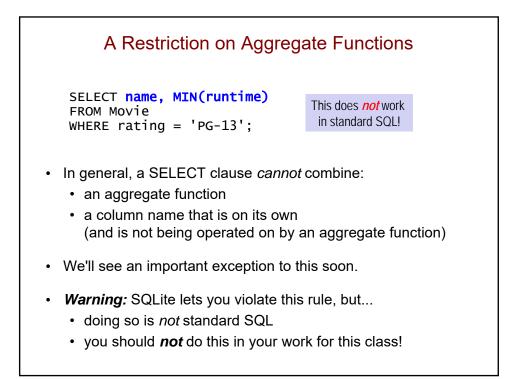
ELECT CO			Evaluating a query with GROUP BY					
	ourse_name,	COUNT(*)						
ROM Enro								
ROUP BY	course_nam	e;						
Inrolled	-		-					
student_id	course_name	credit_status						
12345678	CS 105	ugrad						
45678900	CS 111	ugrad						
45678900	CS 460	grad						
33566891	CS 105	non-credit						
66666666	CS 111	ugrad						
25252525	CS 105	grad]					
	•		-					
student_id	course_name	credit_status		course_name	COUNT(*)			
12345678	CS 105	ugrad		CS 105	3			
33566891	CS 105	non-credit		CS 111	2			
25252525	CS 105	grad	1 1	CS 460	1			
45678900	CS 111	ugrad		ļ	<u> </u>			
66666666	CS 111	ugrad						
45678900	CS 460	grad	11					

FROM Enrol WHERE cred	rse_name, CO		WHERE
student_id	course_name	credit_status	 The WHERE clause
12345678	CS 105	ugrad	is applied <i>before</i>
45678900	CS 111	ugrad	the GROUP BY clause.
45678900	CS 460	grad	the GIVOOF DT clause.
33566891	CS 105	non-credit	
66666666	CS 111	ugrad	
25252525	CS 105	grad	
	WHERE 棏		
student_id	course_name	credit_status	
12345678	CS 105	ugrad	
45678900	CS 111	ugrad	
66666666	CS 111	ugrad	
GR	OUP BY 🜷		
student_id	course_name	credit_status	
12345678	CS 105	ugrad	
45678900	CS 111	ugrad	
66666666	CS 111	ugrad	

Applying a Condition to Subgroups					
What if I only want courses with more than one student? Enrolled					
student_id	course_name	credit_status		course_name	COUNT(*)
12345678	CS 105	ugrad	1	CS 105	3
45678900	CS 111	ugrad		CS 111	2
45678900	CS 460	grad		CS 460	1
33566891	CS 105	non-credit		HAVING	
66666666	CS 111	ugrad			•
25252525	CS 105	grad		course_name	COUNT(*)
			-	CS 105	3
This wo	n't work:			CS 111	2
<pre>SELECT course, COUNT(*) FROM Enrolled WHERE COUNT(*) > 1 GROUP BY course; WHERE is applied before GROUP BY.</pre>					
This will: HAVING is applied after					
FROM EI GROUP I	course, COU nrolled BY course COUNT(*) >		•	OUP BY. used for all co involving agg	

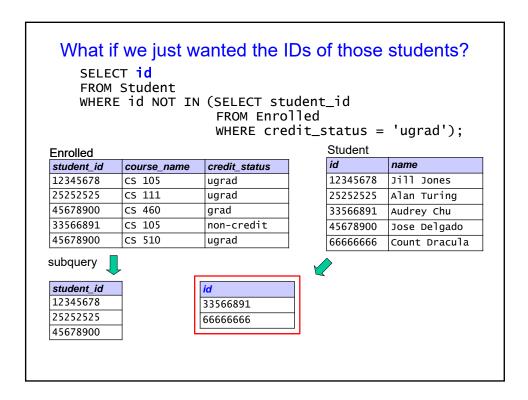


	How could we find the shortest						
	Movie PG-13 movie in the database?						
	id	name	year	rating	runtime		
	2488496	Star Wars: The Force Awakens	2015	PG-13	138		
	1228705	Iron Man 2	2010	PG-13	124		
	0435761	Toy Story 3	2010	G	103		
	1323594	Despicable Me	2010	PG	95		
	0118998	Dr. Dolittle	1998	PG-13	85		
В.	FROM N WHERE SELEC FROM N WHERE	MOVIE rating = 'PG-13'; E. T name, runtime Movie runtime = (SELECT MIN(run WHERE rating =	all thr	FROM M	e would v uld work ovie		
C.	FROM N WHERE	T name, runtime Movie rating = 'PG-13' runtime = (SELECT MIN(run ⁻ WHERE rating =			ovie		



	How could we find the shortest PG-13 movie in the database?						
1	Movie PG-13 MOVIE IN the database ?						
	2488496	Star Wars: The Force Awakens	2015	PG-13	138		
	1228705	Iron Man 2	2010	PG-13	124		
	0435761	Toy Story 3	2010	G	103		
	1323594	Despicable Me	2010	PG	95		
	0118998	Dr. Dolittle	1998	PG-13	85		
	 A. SELECT name, MIN(runtime) FROM Movie WHERE rating = 'PG-13'; B. SELECT name, runtime FROM Movie WHERE runtime = (SELECT MIN(runtime) FROM Movie 						
	WHERE rating = 'PG-13');						
C	FROM N WHERE	T name, runtime Movie rating = 'PG-13' runtime = (SELECT MIN(r WHERE rating			ovie		

WHER		(SELECT st FROM Enro WHERE cre	11ed		'ugrad');
Enrolled			Stu	dent	
student_id	course_name	credit_status	id		name
12345678	CS 105	ugrad	123	45678	Jill Jones
25252525	CS 111	ugrad	252	52525	Alan Turing
45678900	CS 460	grad	335	66891	Audrey Chu
33566891	CS 105	non-credit	456	78900	Jose Delgado
45678900	CS 510	ugrad	666	66666	Count Dracula
					was the query g for?



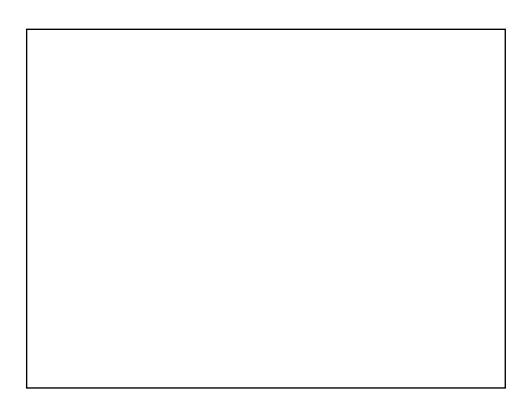
WHERE	Enrolled student_i	FRO	_ECT student_id M Enrolled RE credit_status='ugrad')
Enrolled student_id	course_name	credit_status	
12345678	CS 105	ugrad	
25252525	CS 111	ugrad	omit the Student table!
45678900	CS 460	grad	offic the Student table!
33566891	CS 105	non-credit	
45678900	CS 510	ugrad	

FROM	CT student_ Enrolled	What about id atus != 'ug
student_id	course_name	credit_status
12345678	CS 105	ugrad
25252525	CS 111	ugrad
45678900	CS 460	grad
33566891	CS 105	non-credit
45678900	CS 510	ugrad

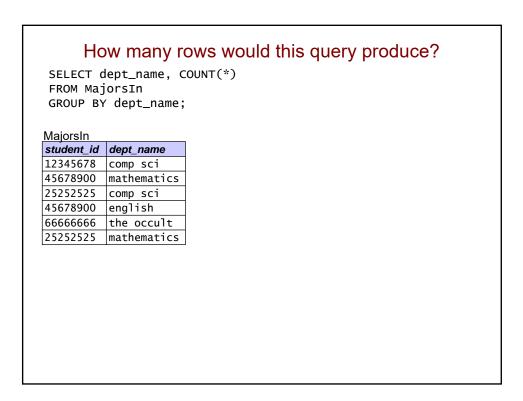
What about this?

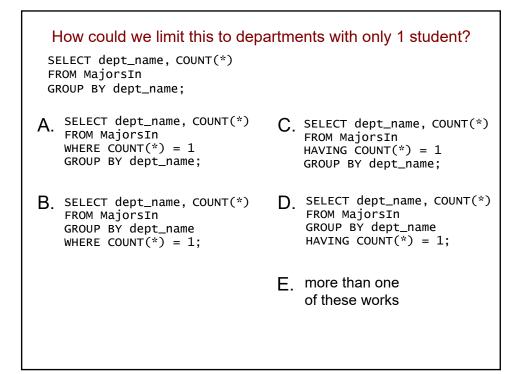
SELECT student_id
FROM Enrolled
WHERE credit_status != 'ugrad';

Enrolled		
student_id	course_name	credit_status
12345678	CS 105	ugrad
25252525	CS 111	ugrad
45678900	CS 460	grad
33566891	CS 105	non-credit
45678900	CS 510	ugrad

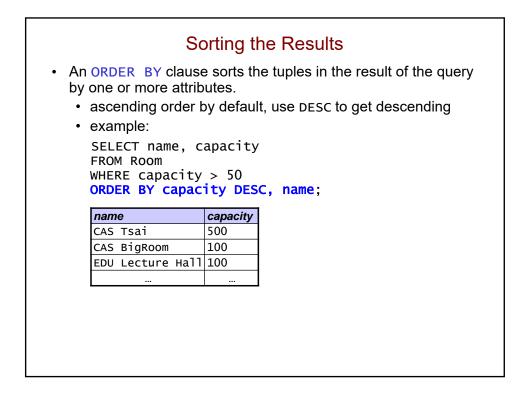


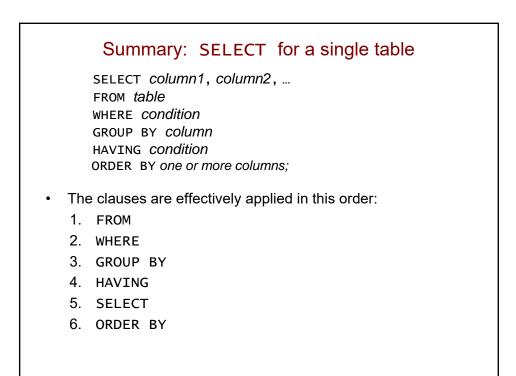
FROM	T student_ Enrolled	Vhat about id atus != 'ug		
student_id	course_name	credit_status		
12345678	CS 105	ugrad		
25252525	CS 111	ugrad		
45678900	CS 460	grad		
33566891	CS 105	non-credit		
45678900	CS 510	ugrad		
student_id	45678900 i	s included		Need to use a subquery and
45678900	even thoug	h he is enrolled i	n	NOT IN for problems
33566891	a course fo	r undergrad cred	it!	like this one!

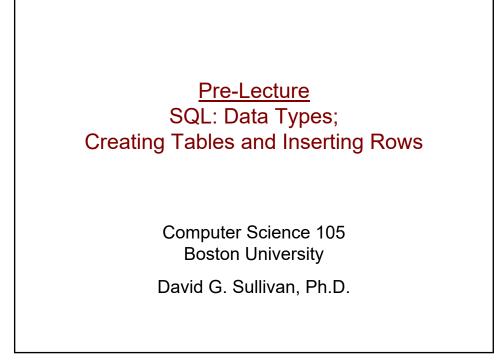


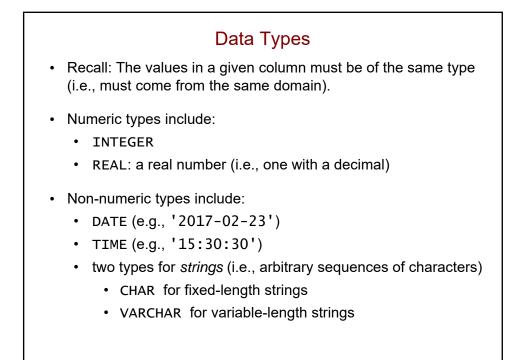


FROM Enrol WHERE cred	rse_name, CO		WH	ERE	
student_id	course_name	credit_status	•	The WHER	E clause
12345678	CS 105	ugrad		is applied be	efore
45678900	CS 111	ugrad	1	the GROUP	
45678900	CS 460	grad	1	THE GROOP	DT clause.
33566891	CS 105	non-credit	1		
66666666	CS 111	ugrad	1		
25252525	CS 105	grad			
	WHERE 棏		-		
student_id	course_name	credit_status]		
12345678	CS 105	ugrad	1		
45678900	CS 111	ugrad	1		
66666666	CS 111	ugrad			
GR	OUP BY 퉞				
student_id	course_name	credit_status		course_name	COUNT(*)
12345678	CS 105	ugrad		CS 105	1
45678900	CS 111	ugrad		CS 111	2
66666666	CS 111	ugrad			





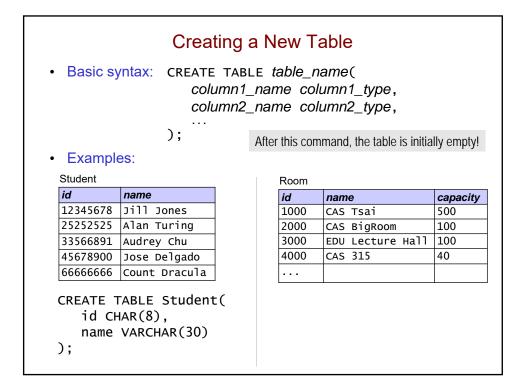


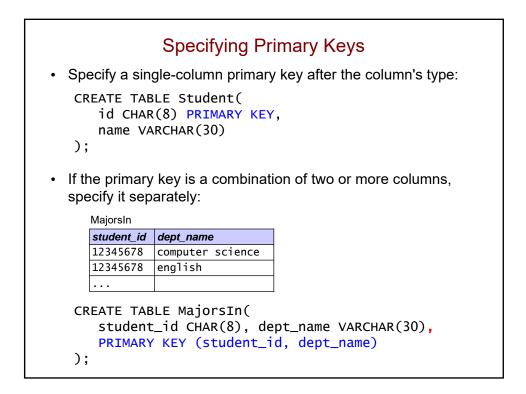


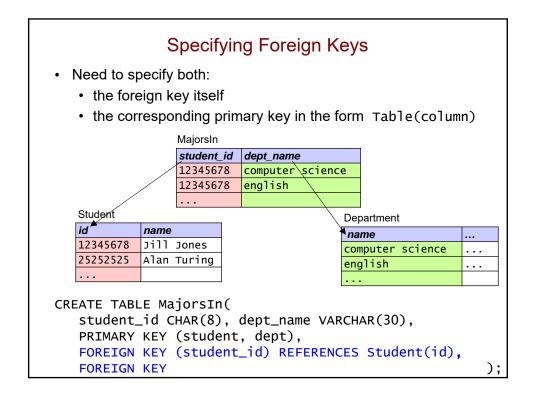
CHAR vs. VARCHAR

- CHAR(*n*) is for *fixed-length* strings of <u>exactly n</u> characters.
- VARCHAR(*n*) is for *variable-length* strings of <u>up to n</u> characters.
 - used for values that can have a wide range of possible lengths
- Example: types for a *Person* table:
 - VARCHAR(64) for the person's name
 - VARCHAR(128) for the street address
 - VARCHAR(32) for the city
 - CHAR(2) for the state abbreviation ('MA', 'NY', etc.)
 - CHAR(5) for the zip code
 - CHAR(8) for the id since every id has the same # of digits
 - example: '00123456'
 - a numeric type would not keep the leading 0s

CHAR	VS. VARCHAR (cont	.)
	and VARCHAR(<i>n</i>), if the <u>more than <i>n</i></u> characters,	•
 examples: 		
type	user-specified value	value stored
CHAR(5)	'123456'	'12345'
VARCHAR(10)	'computer science'	
 if the type is CHAR 	o specify a value of <u>less t</u> (<i>n</i>), the system pads with the system pads with the system doe	ith spaces
type	user-specified value	value stored
CHAR(5)	'123'	'123 '
VARCHAR(10)	'math'	







	Adding a Si	ngle F	Row to	an E	xisting Table	
• S	yntax:					
I	INSERT INTO tak	ole VALI	JES (<i>va</i>	l1, va	al2,);	
• •	vemple:			/ i	d is CHAR(4), so ne	ed quotes
	xample: INSERT INTO RO	0 \/ \	IES (1	2211	'MCS 148' 4	5)
	INSERT INTO RO	UIII VALU		-	MCS 140 , 4	3)
Room				Room	1	
id	name	capacity		id	name	capacity
1000	CAS Tsai	500		1000	CAS Tsai	500
2000	CAS BigRoom	100		2000	CAS BigRoom	100
3000	EDU Lecture Hall	100	1 1	3000	EDU Lecture Hall	100
4000	CAS 315	40		4000	CAS 315	40
			1	1234	MCS 148	45
• N	otes:					
	need to specify	the valu	ies in th	e ann	ropriate order	
	• •				CREATE TABLE)
•	non-numeric va	lues are	e surrou	nded b	by single quotes	
•	the DBMS won	-				
	Unique e antelouv	IDDDCC (nr rotoro	ntial_ir	ntegrity constrain	t .

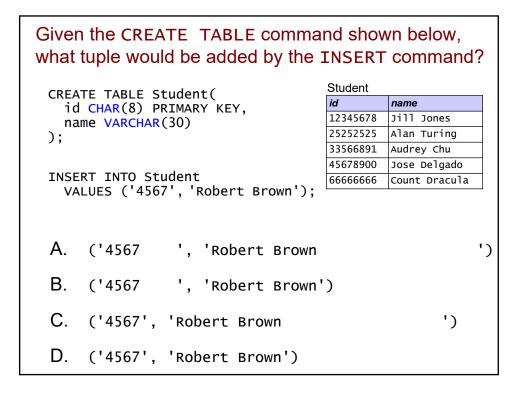
SQL: Data Types; Creating Tables and Inserting Rows

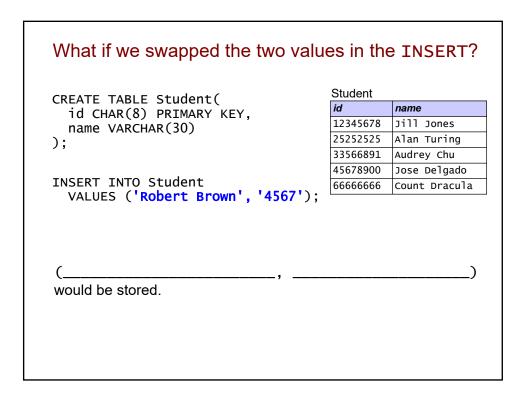
Computer Science 105 Boston University

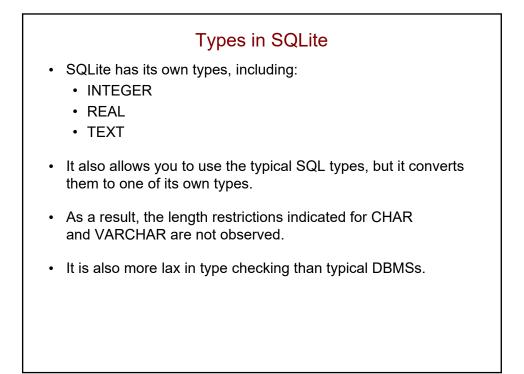
David G. Sullivan, Ph.D.

SQL Data Types

- Numeric types include:
 - INTEGER
 - REAL: a real number (i.e., one that may have a fractional part)
- Non-numeric types include:
 - DATE (e.g., '2017-02-23')
 - TIME (e.g., '15:30:30')
 - two types for *strings* (i.e., arbitrary sequences of characters)
 - CHAR
 - VARCHAR

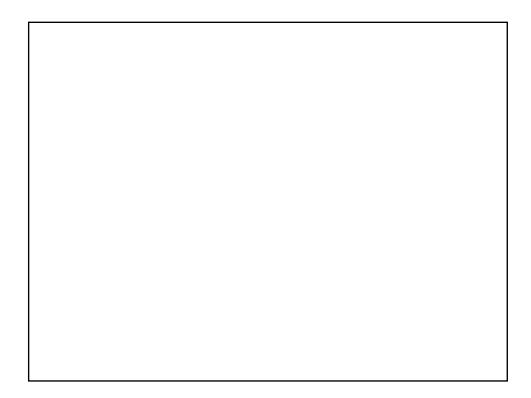






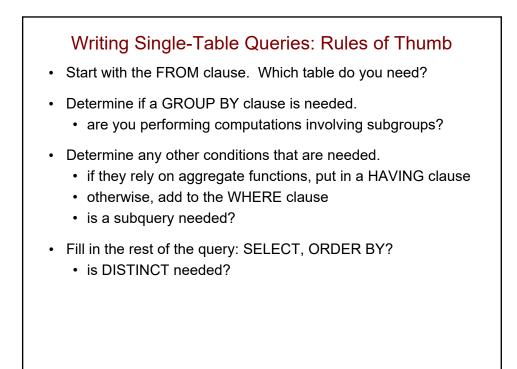
Enrolled				Student	1
student_id	course_name	credit_status		id	name
12345678	CS 105	ugrad		12345678	Jill Jones
25252525	CS 111	ugrad		25252525	Alan Turing
45678900	CS 460	grad		33566891	Audrey Chu
33566891	CS 105	non-credit		45678900	Jose Delgado
45678900	CS 510	ugrad			
studen	ABLE Enrol t_id CHAR(a _status VA	led(8), course_	_name	VARCHAR	Count Dracula

Enrolled				Student	-
student_ia	course_name	credit_status		id	name
12345678	CS 105	ugrad		12345678	Jill Jones
25252525	CS 111	ugrad		25252525	Alan Turing
45678900	CS 460	grad		33566891	Audrey Chu
33566891	CS 105	non-credit		45678900	Jose Delgado
45678900	CS 510	ugrad		66666666	Count Dracula
stud cred	TABLE Enro ent_id CHAR it_status V ARY KEY (stu	(8), course_ ARCHAR(10),			(10),
stud cred	ent_id CHAR(it_status V/	(8), course_ ARCHAR(10),			(10),);
stud cred PRIM	ent_id CHAR(it_status V/	(8), course_ ARCHAR(10), Ident_id, co	ourse_	name),);
Stud cred PRIM A. F	ent_id CHAR it_status V/ ARY KEY (stu	(8), course_ ARCHAR(10), udent_id, co (student_id)	ourse_) REFE	name),); Student(id)
Stud cred PRIM A. F B. F	ent_id CHAR it_status VA ARY KEY (stu OREIGN KEY	(8), course_ ARCHAR(10), udent_id, co (student_id) (student_id)) REFE	name), RENCES S d IN Stu); Student(id) udent

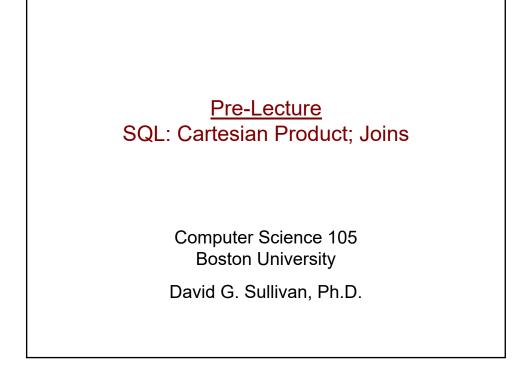


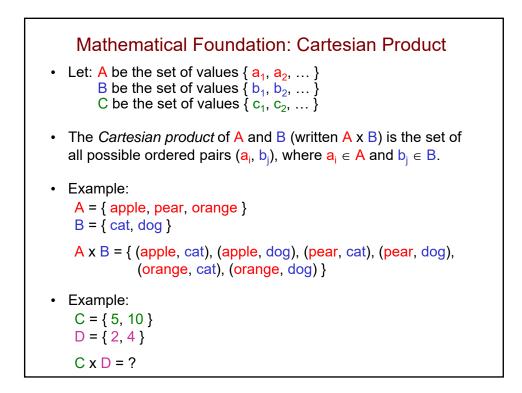
Enrolled				Student	
student_id	course_name	credit_stat	us	id	name
12345678	CS 105	ugrad		12345678	Jill Jones
25252525	CS 111	ugrad		25252525	Alan Turing
45678900	CS 460	grad		33566891	Audrey Chu
33566891	CS 105	non-credi	t	45678900	Jose Delgado
45678900	CS 510	ugrad		66666666	Count Dracula
student credit_ PRIMARY	BLE Enrol _id CHAR(status VA KEY (stu KEY (stu	8), cour RCHAR(10 dent_id,), course_	name),	
student credit_ PRIMARY	_id CHAR(_status VA	8), cour RCHAR(10 dent_id,), course_	name),	
student credit_ PRIMARY FOREIGN	_id CHAR(status VA KEY (stu	8), cour RCHAR(10 dent_id,), course_	name),	
student credit_ PRIMARY FOREIGN	:_id CHAR(status VA KEY (stu KEY (stu	8), cour RCHAR(10 dent_id, dent_id)), course_ REFEREN	name), CES Stud	
student credit_ PRIMARY FOREIGN	:_id CHAR(status VA KEY (stu KEY (stu se name	8), cour RCHAR(10 dent_id, dent_id) start_time), COUTSE_ REFEREN end_time	name), ICES Stud	
student credit_ PRIMARY FOREIGN	:_id CHAR(status VA KEY (stu KEY (stu KEY (stu se <u>name</u> cs 105	8), cour RCHAR(10 dent_id, dent_id) start_time 13:00:00), COURSE_ REFEREN end_time 14:00:00	name), ICES Stud room_id 4000	
student credit_ PRIMARY FOREIGN	s_id CHAR(status VA KEY (stu KEY (stu KEY (stu se <u>name</u> cs 105 cs 111 cs 460	8), cour RCHAR(10 dent_id, dent_id) <u>start_time</u> 13:00:00 09:30:00), COURSE_ REFEREN end_time 14:00:00 11:00:00	name), ICES Stud room_id 4000 5000	

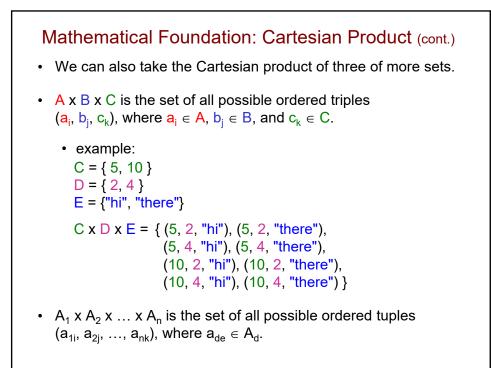
Enrolled			Student	
student_id	course_name	credit_status	id	name
12345678	CS 105	ugrad	12345678	Jill Jones
25252525	CS 111	ugrad	25252525	Alan Turing
45678900	CS 460	grad	33566891	Audrey Chu
33566891	CS 105	non-credit	45678900	Jose Delgado
45678900	CS 510	ugrad	66666666	Count Dracula
, ,		led VALUES(' nt VALUES ('	', 'CS 1	L05', 'grad'
) INSERT	INTO Studer	nt VALUES ('	', 'CS 1	L05', 'grad'
	INTO Studer	-	', 'CS 1	L05', 'grad'
) insert A. (INTO Studer	nt VALUES ('	', 'CS 1	L05', 'grad'



Practice Writing Queriesname1) Find the start times of CS 105CSand CS 111.CScscs		start timo		
 1) Find the start times of CS 105 and CS 111. 2) Find the course(s) that end latest in the course 		start timo		
 1) Find the start times of CS 105 and CS 111. 2) Find the course(s) that end latest in the course 		start_time	end_time	room_id
and CS 111. 2) Find the course(s) that end latest in the	105	13:00:00	14:00:00	4000
and CS 111. (CS (CS (CS (CS (CS (CS (CS (CS	111	09:30:00	11:00:00	5000
 2) Find the course(s) that end latest in the course (s) the cour	-	11:00:00		1000
PH2) Find the course(s) that end latest in the		16:00:00		7000
2) Find the course(s) that end latest in the		12:00:00		7000
, , , , , , , , , , , , , , , , , , , ,	101	14:30:00	16:00:00	NULL
 Find the ids of all rooms that have two The result should be tuples of the forr 	query o or	y!) more co		







		Cartesian F						
	 The Cartesian product of two or more relations forms all possible combinations of <i>rows</i> from the relations. 							
•	The result i	s itself a rela	tion.					
	• its rows	contain all of	the columns	fro	m the con	nbined relations		
•	Example:							
	Enrolled				MajorsIn			
	student_id	course_name	credit_status		student_id	dept_name		
	12345678	CS 105	ugrad		12345678	comp sci		
	25252525	CS 111	ugrad		45678900	mathematics		
	Enrolled x I	VlajorsIn						
	Enrolled. student_id	course_name	credit_status		<mark>jorsIn.</mark> dent_id	dept_name		
	12345678	CS 105	ugrad	12	845678	comp sci		
	12345678	CS 105	ugrad	456	578900	mathematics		
	1	1	1	1				

. . .

Example:					
Enrolled			_	MajorsIn	
student_id	course_name	credit_status		student_id	l dept_name
12345678	CS 105	ugrad		12345678	comp sci
25252525	CS 111	ugrad		45678900	mathemat
45678900	CS 460	grad	1 \\\`	25252525	comp sci
	CS 105	non-credit		45678900	english
33566891	CS 105 CS 510	non-credit grad		45678900 66666666	
33566891 45678900 Enrolled x Enrolled.	CS 510 MajorsIn	grad	Ma	66666666 iorsIn.	the occu
33566891 45678900 Enrolled x <i>Enrolled.</i> <i>student_id</i>	CS 510 MajorsIn course_name		stu	66666666 iorsIn. dent_id	the occu dept_name
33566891 45678900 Enrolled x <i>Enrolled.</i> <i>student_id</i>	CS 510 MajorsIn	grad	stu	66666666 iorsIn.	the occu
33566891 45678900 Enrolled x <i>Enrolled.</i> <i>student_id</i> 12345678	CS 510 MajorsIn course_name	grad	stu 123	66666666 iorsIn. dent_id	the occu dept_name comp sci
33566891 45678900 Enrolled x <i>Enrolled.</i> <i>student_id</i> 12345678 12345678	CS 510 MajorsIn course_name	grad credit_status ugrad	stu 123 456	66666666 iorsIn. dent_id 845678	the occu dept_name comp sci
33566891 45678900	CS 510 MajorsIn course_name CS 105 CS 105	grad credit_status ugrad ugrad	stu 123 456 252	66666666 iorsIn. dent_id 845678 578900	the occu dept_name comp sci mathematic

Example:					
Enrolled				MajorsIn	
student_id	course_name	credit_status		student_id	dept_name
12345678	CS 105	ugrad		12345678	comp sci
25252525	CS 111	ugrad	\vdash	45678900	mathematic
45678900	CS 460	grad	\mathbb{N}	25252525	comp sci
33566891	CS 105	non-credit	$1 \setminus$	45678900	english
45678900	CS 510	grad	1	66666666	the occult
Enrolled x Enrolled. student_id	MajorsIn course_name	credit_status	Maj stu		dept_name
Enrolled.		e credit_status ugrad	stu	iorsIn. dent_id	
Enrolled. student_id	course_name	_	stu 123	iorsIn. dent_id 45678	dept_name comp_sci
Enrolled. student_id 12345678	course_name	ugrad	stu 123 456	iorsIn. dent_id 45678 578900	dept_name comp_sci
<i>Enrolled.</i> <i>student_id</i> 12345678 12345678	course_name cs 105 cs 105	ugrad ugrad	stu 123 456 252	iorsIn. dent_id 845678 678900 852525	dept_name comp sci mathematics
Enrolled. student_id 12345678 12345678 12345678	course_name CS 105 CS 105 CS 105	ugrad ugrad ugrad	stu 123 456 252 456	orsIn. dent_id 45678 578900 252525 578900	dept_name comp sci mathematics comp sci
Enrolled. student_id 12345678 12345678 12345678 12345678	course_name cs 105 cs 105 cs 105 cs 105 cs 105 cs 105	ugrad ugrad ugrad ugrad	stu 123 456 252 456 666	iorsIn. dent_id 445678 578900 52525 578900 566666	dept_name comp sci mathematics comp sci english
Enrolled. student_id 12345678 12345678 12345678 12345678 12345678	course_name cs 105 cs 105	ugrad ugrad ugrad ugrad ugrad	stu 123 456 252 456 666 123	iorsIn. dent_id 445678 78900 252525 578900 666666 845678	dept_name comp sci mathematics comp sci english the occult

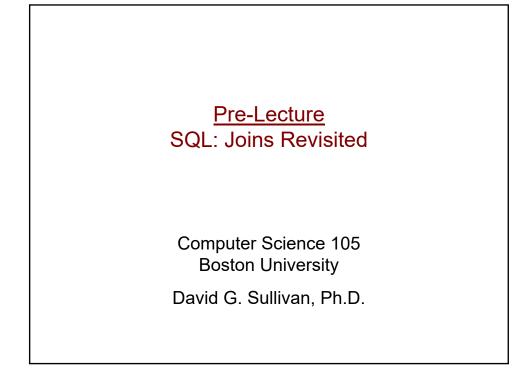
Joining Multiple Tables SELECT column1, column2, ... FROM table1, table2, When the FROM clause specifies multiple tables, the ٠ resulting operation is known as a join. • The result is equivalent to: · forming the Cartesian product of the tables in the FROM clause table1 x table2 x ... · applying the remaining clauses to the Cartesian product, in the same order as for a single-table command: WHERE **GROUP BY** HAVING SELECT ORDER BY

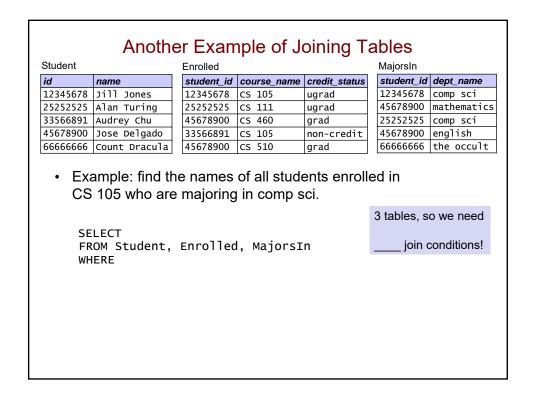
	Joining Multiple Tables (cont.)									
• Ex	ample: fir	nd Alan Turing'	s majo	r.						
Student	id	name		student_id	dept_name	MajorsIn				
	12345678	Jill Jones		12345678	comp sci					
	25252525	Alan Turing		45678900	mathematics					
	33566891	Audrey Chu		25252525	comp sci					
	45678900	Jose Delgado		45678900	english					
	66666666	Count Dracula		66666666	the occult					
	SELECT d FROM Stu WHERE na	ery that works: ept_name dent, Majors me = 'Alan T = student_i	uring	,						
•	used to r selects th	ent_id is a <i>joi</i> natch up "relat ne tuples in the les, you typical	ed" tup Carte	les from the	t that "make	sense"				

					student_id;
Student			MajorsIn		
id	name		student_i	d	dept_name
	Jill Jones		12345678		comp sci
25252525	Alan Turing		45678900		mathematic
33566891	Audrey Chu		25252525		comp sci
45678900	Jose Delgado		45678900		english
66666666	Count Dracula		66666666		the occult
Student x N	lajorsIn				
id	name	stu	udent_id	de	ept_name
12345678	Jill Jones	12	345678	со	omp sci
12345678	Jill Jones	45	678900	ma	thematics
12345678	Jill Jones	25	252525	со	omp sci
12345678	Jill Jones	45	678900	en	nglish
12345678	Jill Jones	66	666666	th	ne occult
25252525	Alan Turing	12	345678	со	omp sci
25252525	Alan Turing	45	678900	ma	thematics
25252525	Alan Turing	25	252525	сс	omp sci
25252525	Alan Turing	45	678900	en	nglish

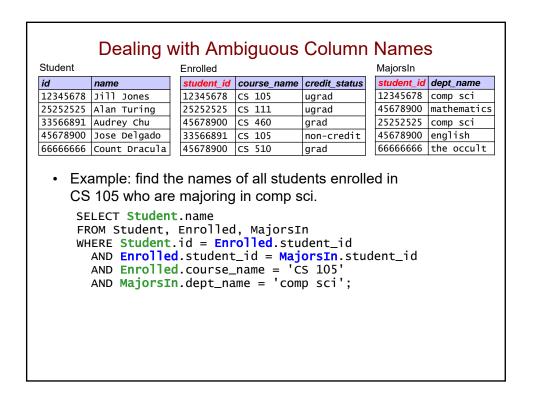
Student <i>id</i>	name	MajorsIn student	d	dept_name	7
	Jill Jones	12345678	_	comp sci	-
	Alan Turing	45678900		mathematics	
	Audrey Chu	25252525		comp sci	-
	Jose Delgado	45678900	_	english	-
	Count Dracula	66666666	_	the occult	-
	lajorsIn WHERE ic		_		
id	name	student id		ept_name	
12345678	Jill Jones	12345678	-	mp sci	
25252525	Alan Turing	25252525	-	mp sci	
	1 -	45678900	-	thematics	
45678900	Jose Delgado	15070500	1		
45678900	Jose Delgado Jose Delgado	45678900		glish	

id	name		student_i	d	dept_name	
12345678	Jill Jones		12345678		comp sci	
	Alan Turing		45678900		mathematics	s
33566891	Audrey Chu		25252525		comp sci	
45678900	Jose Delgado		45678900		english	
66666666	Count Dracula		66666666		the occult	
After select	ing only tuples tha	t sat	tisfv the W	HE	RE clause:	
	ing only tuples tha <i>nam</i> e	-	-		RE clause: pt_name	
id	U U U	Sti	-	de		





	Dealing w	vith Am	biguous	Column	Name	6
Student	_	Enrolled	-		MajorsIn	
id	name	student_id	course_name	credit_status	student_id	dept_name
12345678	Jill Jones	12345678	CS 105	ugrad	12345678	comp sci
25252525	Alan Turing	25252525	CS 111	ugrad	45678900	mathematics
33566891	Audrey Chu	45678900	CS 460	grad	25252525	comp sci
45678900	Jose Delgado	33566891	CS 105	non-credit	45678900	english
66666666	Count Dracula	45678900	CS 510	grad	66666666	the occult
SE FR WH	105 who are LECT name OM Student, IERE id = Enr AND Enrolled AND course_r AND dept_nam	Enrolled colled.st student ame = 'C	l, MajorsI udent_id _id = <mark>Maj</mark> S 105'	n	dent_id	



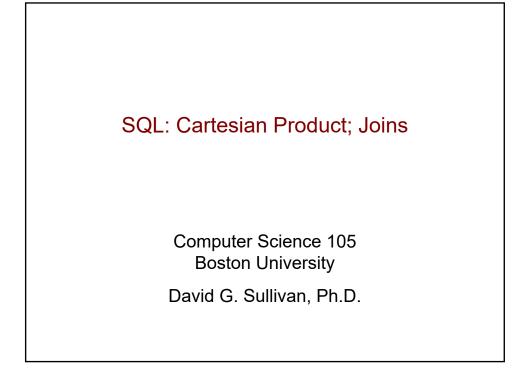
	A	liases f	or Table	Names		
Student		Enrolled			MajorsIn	
id	name	student_id	course_name	credit_status	student_id	dept_name
12345678	Jill Jones	12345678	CS 105	ugrad	12345678	comp sci
25252525	Alan Turing	25252525	CS 111	ugrad	45678900	mathematics
33566891	Audrey Chu	45678900	CS 460	grad	25252525	comp sci
45678900	Jose Delgado	33566891	CS 105	non-credit	45678900	english
66666666	Count Dracula	45678900	CS 510	grad	66666666	the occult
SE FR WH	105 who are LECT S.name OM Student A IERE S.id = E AND E.studen AND E.course AND M.dept_n	S S , Enr E.student ht_id = M e_name =	olled <mark>AS</mark> _id .student_ 'CS 105'	E, Majors:	In <mark>AS M</mark>	

WHERE S.i AND E.s AND E.c	name dent S, Enrol d = E.studen student_id = I course_name = dept_name = '	t_id M.studen 'CS 105	t_io ';				Meiorolo	
id	name	1			e credit statu	•	MajorsIn student i	d dept name
	Jill Jones	12345		course_name	ugrad	5	12345678	= 1
	Alan Turing	25252		CS 111	ugrad	\neg	45678900	
	Audrey Chu	456789		CS 460	grad	$\exists $	252525252	
	Jose Delgado	335668	391	CS 105	non-credit	: \)	45678900) english
66666666	Count Dracula	456789	900	CS 510	grad		66666666	5 the occult
Student x <i>id</i>	Enrolled x Majors name		id	course_name	credit_status	M.s	student_id	dept_name
12345678	Jill Jones	12345678	(CS 105	ugrad	123	345678	comp sci
12345678	Jill Jones	12345678	(CS 105	ugrad	456	578900	mathematics
12345678	Jill Jones	12345678	(CS 105	ugrad	252	252525	comp sci
12345678	Jill Jones	12345678	(CS 105	ugrad	456	578900	english
12345678	Jill Jones	12345678	(CS 105	ugrad	666	566666	the occult

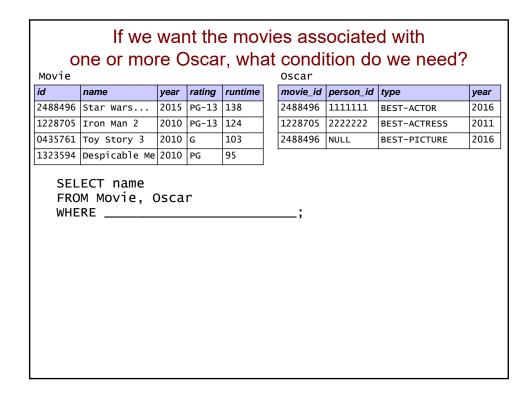
WHERE S. AND E.S AND E.G	name dent S, Enrol d = E.studen student_id = course_name = dept_name = '	t_id M.student_i 'CS 105'			Mainala	
id	name	student id	l course name	e credit statu	MajorsIn	id dept_name
12345678	Jill Jones	12345678		ugrad	1234567	-
25252525	Alan Turing	25252525	CS 111	ugrad	4567890	
33566891	Audrey Chu	45678900		grad	2525252	5 comp sci
45678900	Jose Delgado	33566891	CS 105	non-credit	4567890	0 english
66666666	Count Dracula	45678900	CS 510	grad	6666666	6 the occult
	Enrolled x Majors		s in all!			
id	name	_	course_name	credit_status		
12345678		12345678	CS 105	ugrad	12345678	comp sci
12345678		12345678	CS 105	ugrad	45678900	mathematics
12345678		12345678	CS 105	ugrad	25252525	comp sci
12345678	Jill Jones	12345678	CS 105	ugrad	45678900	english
12345678	Jill Jones	12345678	CS 105	ugrad	66666666	the occult
12345678	Jill Jones	25252525	CS 111	ugrad	12345678	comp sci
12345678	Jill Jones	25252525	CS 111	ugrad	45678900	mathematics
12345678	Jill Jones	25252525	CS 111	ugrad	25252525	comp sci
12345678	Jill Jones	25252525	CS 111	ugrad	45678900	english
12345678	Jill Jones	25252525	CS 111	ugrad	66666666	the occult

WHERE S. ⁻ AND E AND E	.name dent S, Enrol id = E.studer student_id = course_name = dept_name = '	nt_id M.student_i 'CS 105'					
Student		Enrolled				MajorsIn	
id	name	student_i	d course_name	e credit_statu	s	student_i	d dept_name
12345678	Jill Jones	12345678	CS 105	ugrad		12345678	comp sci
25252525	Alan Turing	25252525	CS 111	ugrad		45678900	mathematics
33566891	Audrey Chu	45678900	CS 460	grad		25252525	comp sci
45678900	Jose Delgado	33566891	CS 105	non-credit		45678900	english
66666666	Count Dracula	45678900	CS 510	grad		66666666	the occult
Student x	Enrolled x Major <i>name</i>		, , , , , , , , , , , , , , , , , , ,	ions <i>credit_status</i>	M.s	tudent_id	dept_name
12345678	Jill Jones	12345678	CS 105	ugrad	123	45678	comp sci
25252525	Alan Turing	25252525	CS 111	ugrad	252	52525	comp sci
45678900	Jose Delgado	45678900	CS 460	grad	456	78900	mathematics
45678900	Jose Delgado	45678900	CS 460	grad	456	78900	english
45678900	Jose Delgado	45678900	CS 510	grad	456	78900	mathematics
45678900	Jose Delgado	45678900	CS 510	grad	456	78900	english

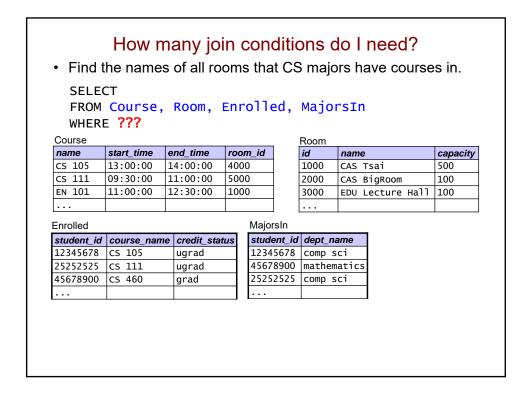
WHERE S. AND E.S AND E.G	.name dent S, Enroli id = E.studen student_id = M course_name = dept_name = 'd	t_id 4.student_i 'CS 105'				
Student		Enrolled			MajorsIr	1
id	name	student_ic	l course_name	e credit_statu	s student	_id dept_name
12345678	Jill Jones	12345678	CS 105	ugrad	123456	78 comp sci
25252525	Alan Turing	25252525	CS 111	ugrad	4567890	00 mathematics
	Audrey Chu	45678900		grad	2525252	
45678900	Jose Delgado	33566891	CS 105	non-credit	4567890	· · · · ·
66666666	Count Dracula	45678900	CS 510	grad	6666666	66 the occult
	Enrolled x Majors		,			
id	name		course_name	_	-	
12345678	Jill Jones	12345678	CS 105	ugrad	12345678	comp sci
after name Jill J	ones					

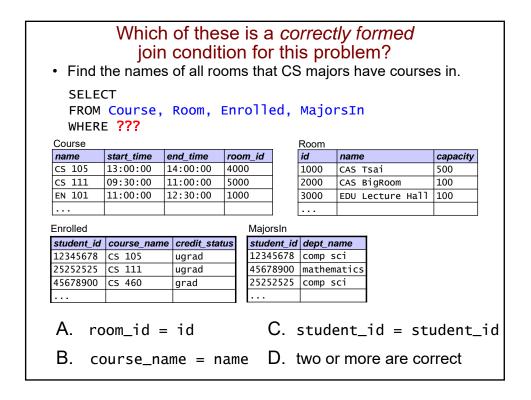


Movie						Oscar			
id	name	year	rating	runtime		movie_id	person_id	type	yea
2488496	Star Wars	2015	PG-13	138] [2488496	1111111	BEST-ACTOR	201
1228705	Iron Man 2	2010	PG-13	124	1 [1228705	2222222	BEST-ACTRESS	201
0425761	Toy Story 3	2010	C	100					
0433761	,, .	2010	G	103		2488496	NULL	BEST-PICTURE	2016
1323594 SEL	Despicable Me ECT name M Movie, (2010	PG	95		2488496	NULL	BEST-PICTURE	2016
1323594 SEL	Despicable Me ECT name	2010	PG			2488496	NULL	BEST-PICTURE	2016
1323594 SEL	Despicable Me ECT name	2010	PG			2488496	NULL	BEST-PICTURE	2016



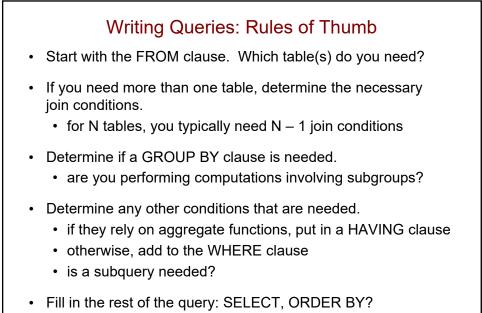
	• Fin	d the na		Which s of all r						ave cours	es in.
		_ECT									
	WHE										
	Course							Room	-		
	name	start_t	ime	end_time	1	room_id		id	name		capacity
	CS 105	13:00	:00	14:00:00	4	1000		1000	CAS T	sai	500
	CS 111			11:00:00		5000		2000	CAS E	igRoom	100
	EN 101	11:00	:00	12:30:00	1	L000		3000	EDU L	ecture Hal	1 100
	CS 460	16:00	:00	17:30:00	7	7000		4000	GCB 2	04	40
	CS 510	12:00	:00	13:30:00	7	7000		5000	CAS 3	14	80
	РН 101	14:30	:00	16:00:00	Ν	NULL		6000	CAS 2	26	50
								7000	MCS 2	05	30
Stu	dent			Enrolled	I					MajorsIn	
id		name		studen	t_id	course_r	name	credit_s	status	student_id	dept_name
123	45678	Jill Jon	es	123456	78	CS 105		ugrad		12345678	comp sci
252	252525	Alan Tur	ing	252525	25	CS 111		ugrad		45678900	mathematic
335	66891	Audrey C	hu	456789	00	CS 460		grad		25252525	comp sci
456	578900	Jose Del	gado	335668	91	CS 105		non-cr	edit	45678900	english
666	666666	Count Dr	acula	456789	00	CS 510		grad		66666666	the occult





• Find 1	the names							ave course	es in
SELE FROM WHER	Course,	Room,	Enr	rolle	ed,	Major	rsIn		
Course						Room			
Course name	start_time	end_time	roo	om_id		Room <i>id</i>	name		capacity
	_	<i>end_time</i> 14:00:00	roo	_			name CAS T	sai	capacity 500
name	13:00:00	-	_	00		id	CAS T	rsai BigRoom	
<i>name</i> CS 105	13:00:00 09:30:00	14:00:00	400	00 00		<i>id</i> 1000	CAS T CAS E		500
name CS 105 CS 111	13:00:00 09:30:00	14:00:00 11:00:00	400	00 00		<i>id</i> 1000 2000	CAS T CAS E	BigRoom	500 100
name CS 105 CS 111 EN 101	13:00:00 09:30:00	14:00:00 11:00:00	400	00 00		<i>id</i> 1000 2000 3000	CAS T CAS E	BigRoom	500 100
name CS 105 CS 111 EN 101 Enrolled	13:00:00 09:30:00	14:00:00 11:00:00 12:30:00	400	00 00 00 00 Majoi	rsIn	<i>id</i> 1000 2000 3000	CAS T CAS E EDU L	BigRoom	500 100
name CS 105 CS 111 EN 101 Enrolled	13:00:00 09:30:00 11:00:00	14:00:00 11:00:00 12:30:00	400	00 00 00 00 Major <i>stude</i>	rsIn	<i>id</i> 1000 2000 3000 	CAS T CAS E EDU L	BigRoom	500 100
name CS 105 CS 111 EN 101 Enrolled student_id	13:00:00 09:30:00 11:00:00	14:00:00 11:00:00 12:30:00	400	00 00 00 00 Major <i>stude</i> 1234	rsIn ent_id	id 1000 2000 3000 dept_na	CAS T CAS E EDU L ame ci	aigRoom Lecture Hall	500 100
name CS 105 CS 111 EN 101 Enrolled student_id 12345678	13:00:00 09:30:00 11:00:00 course_name CS 105	14:00:00 11:00:00 12:30:00 e credit_sta ugrad	400	Major 50 00 00 50 50 1234 4567	rsIn ent_id 5678	<i>id</i> 1000 2000 3000 <i>dept_na</i> comp_s	CAS I CAS E EDU L EDU L CI atics	aigRoom Lecture Hall	500 100

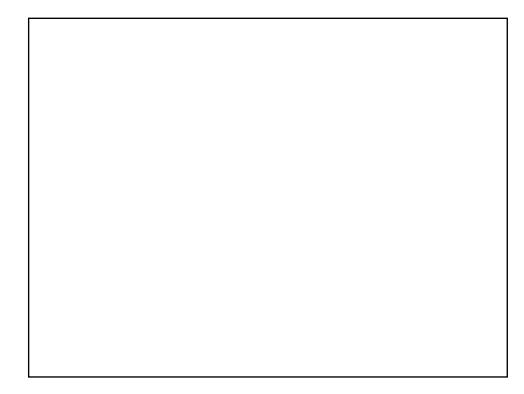




• is DISTINCT needed?

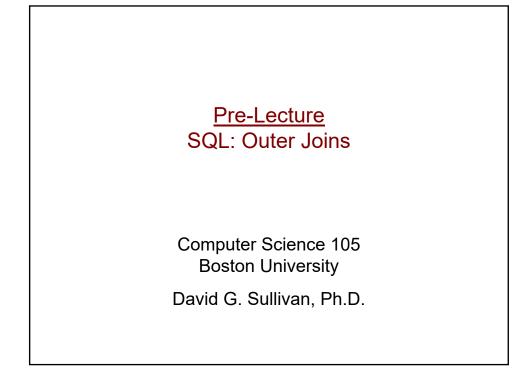
	Practice Writing C	lueries
Course(name, start	Department(name, office) _time, end_time, room_id) , course_name, credit_status)	Room(id, name, capacity) MajorsIn(student_id, dept_name))
1) Find the nan	nes of all courses taken	by comp sci majors.
	nber of students majorin hould be tuples of the fo	g in each department. orm (dept name, # students).)

Practice Writing Que	ries (cont.)
Student(id, name) Department(name, office) Course(name, start_time, end_time, room_id) Enrolled(student_id, course_name, credit_status)	Room(id, name, capacity) MajorsIn(student_id, dept_name)
 Find the names and ids of all student in GCB 204. 	s who have a course
4) Find the names of all rooms in which CS courses meet.	one or more



										_			
									F	From e	earlier	in the lec	ture:
				н			ma	n	w row		uld ha	in this re	cult?
-	SELEC					vv	IIIC		1910	103 000			suit
			vie, Oscar d = movie										
		10	J = 100010	īu,						0.0.0.0			
id	ovie	na	mo	yea	ar	rati	na	****	ntime	Oscar movie id	person_id	type	year
	88196		ar Wars	201		PG-	5	138		2488496	11111111	BEST-ACTOR	2016
			on Man 2	201		PG-	-	124	-	1228705	2222222	BEST-ACTOR BEST-ACTRESS	2010
			y Story 3	201		G G		103	-	2488496	NULL	BEST-PICTURE	2011
-		_	spicable Me			PG		<u>10</u> 95	>	2400490	NULL	DEST-PICTORE	2010
тэ,	23394	De	spicable Me	201		PG		93					
	Movie	х	Oscar										
	id		name		Мо	vie.	ratir	ŋg	runtime	movie_id	person_id	type	Oscar.
				_	yea								year
	248849		Star Wars	-	201		PG-3			2488496	1111111	BEST-ACTOR	2016
	248849	96	Star Wars		201	.5	PG-3	-		1228705	2222222	BEST-ACTRESS	2011
	248849	96	Star Wars		201	.5	PG-3	13	138	2488496	NULL	BEST-PICTURE	2016
	122870)5	Iron Man 2		201	.0	PG-2	13	124	2488496	1111111	BEST-ACTOR	2016
	122870)5	Iron Man 2		201	0.	PG-3	13	124	1228705	2222222	BEST-ACTRESS	2011
	122870)5	Iron Man 2		201	.0	PG-3	13	124	2488496	NULL	BEST-PICTURE	2016
	043576	51	Toy Story 3		201	.0	G		103	2488496	1111111	BEST-ACTOR	2016
	043576	51	Toy Story 3		201	.0	G		103	1228705	2222222	BEST-ACTRESS	2011
ľ	043576	51	Toy Story 3		201	.0	G		103	2488496	NULL	BEST-PICTURE	2016
ľ	132359	94	Despicable	Ме	201	.0	PG		95	2488496	1111111	BEST-ACTOR	2016
Ì	132359	94	Despicable	Ме	201	.0	PG		95	1228705	2222222	BEST-ACTRESS	2011
Ì	132359	94	Despicable	Ме	201	.0	PG		95	2488496	NULL	BEST-PICTURE	2016

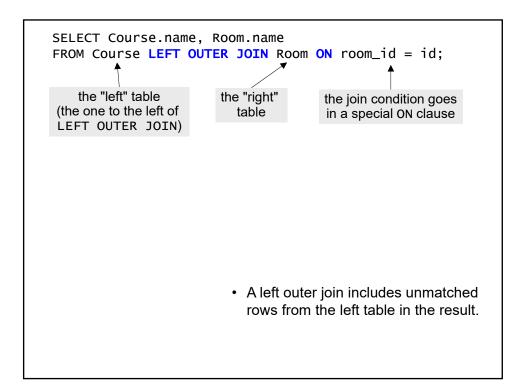
	٩o	name /ie, Oscar 1 = movie_		SW	mar				in the lec in this re	
Movie							Oscar			
id	nai	me	year	rati	ng rui	ntime	movie_id	person_id	type	year
2488496	Sta	ar Wars	2015	i PG-	13 13	8	2488496	1111111	BEST-ACTOR	2016
1228705	Ire	on Man 2	2010) PG-	13 12	4	1228705	2222222	BEST-ACTRESS	2011
0435761	То	y Story 3	2010) G	10	3	2488496	NULL	BEST-PICTURE	2016
		spicable Me Oscar, fol			95 join	condit	ion			
id		name		lovie. ear	rating	runtime	movie_id	person_id	type	Oscar. year
248849	96	Star Wars	. 2	015	PG-13	138	2488496	1111111	BEST-ACTOR	2016
248849	96	Star Wars	. 2	015	PG-13	138	2488496	NULL	BEST-PICTURE	2016
122870)5	Iron Man 2	2	010	PG-13	124	1228705	2222222	BEST-ACTRESS	2011
	S Wa Wa	after ELECT ars ar 2								



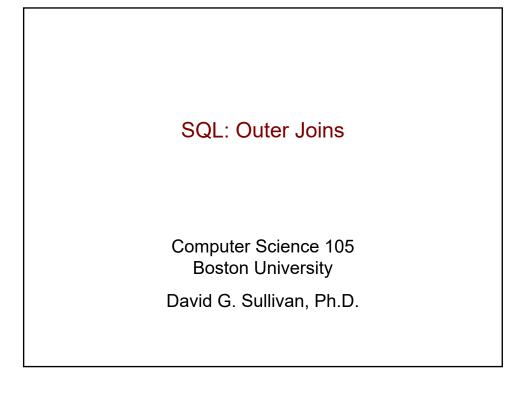
 Need 	l a query	that form	s (cours	e n	ame,	room name) pa	irs.
Course					Room		
name	start_time	end_time	room_id	1	id	name	capacity
CS 105	13:00:00	14:00:00	4000	1	1000	CAS Tsai	500
CS 111	09:30:00	11:00:00	5000	1	2000	CAS BigRoom	100
EN 101	11:00:00	12:30:00	1000	1	3000	EDU Lecture Hall	100
CS 460	16:00:00	17:30:00	7000	1	4000	GCB 204	40
CS 510	12:00:00	13:30:00	7000	1	5000	CAS 314	80
PH 101	14:30:00	16:00:00	NULL	1	6000	CAS 226	50
				-	7000	MCS 205	30
	Mult of the que Room.na GCB 204 CAS 314 CAS TSa	me	SE		Cours	irse.name, Room se, Room	n.name
CS 460	MCS 205		WH	IERE	room	u_id = id;	
CS 510	MCS 205						
PH 101	NULL						

Course					Room			
name	start_time	end_time	room_id	·	id	name		capacity
CS 105	13:00:00	14:00:00	4000		1000	CAS Tsai		500
CS 111	09:30:00	11:00:00	5000		2000	CAS BigRoom		100
EN 101	11:00:00	12:30:00	1000		3000	EDU Lecture	Hall	100
CS 460	16:00:00	17:30:00	7000		4000	GCB 204		40
CS 510	12:00:00	13:30:00	7000		5000	CAS 314		80
PH 101	14:30:00	16:00:00	NULL		6000	CAS 226		50
Course x F	Room 42	2 rows in all	!		7000	MCS 205		30
Course.na	me start_time	end_time	room_id	id	Room.	name	capad	ity
CS 105	13:00:00	14:00:00	4000	1000	CAS T	sai	500	
CS 105	13:00:00	14:00:00	4000	2000	CAS B	igRoom	100	
CS 105	13:00:00	14:00:00	4000	3000	EDU L	ecture Hall	100	
CS 105	13:00:00	14:00:00	4000	4000	GCB 2	04	40	
CS 105	13:00:00	14:00:00	4000	5000	CAS 3	14	80	
CS 105	13:00:00	14:00:00	4000	6000	CAS 2	26	50	
CS 105	13:00:00	14:00:00	4000	7000	MCS 2	05	30	
CS 111	09:30:00	11:00:00	5000	1000	CAS T		500	
CS 111	09:30:00	11:00:00	5000	2000	CAS B	igRoom	100	
CS 111	09:30:00	11:00:00	5000	3000	EDU L	ecture Hall	100	
CS 111	09:30:00	11:00:00	5000	4000	GCB 2	04	40	
CS 111	09:30:00	11:00:00	5000	5000	CAS 3	14	80	
				10000		± 1	100	

Course					Room			
name	start_time	end_time	room_id		id	name		capacity
CS 105	13:00:00	14:00:00	4000		1000	CAS Tsai		500
CS 111	09:30:00	11:00:00	5000		2000	CAS BigRoom		100
EN 101	11:00:00	12:30:00	1000		3000	EDU Lecture	Hall	100
CS 460	16:00:00	17:30:00	7000		4000	GCB 204		40
CS 510	12:00:00	13:30:00	7000		5000	CAS 314		80
PH 101	14:30:00	16:00:00	NULL		6000	CAS 226		50
Course x F	Room, followed	by the join o	condition		7000	MCS 205		30
	me start_time			id	Room.	.name	capac	ity
CS 105	13:00:00	14:00:00	4000	4000	GCB 2	04	40	
CS 111	09:30:00	11:00:00	5000	5000	CAS 3	14	80	
EN 101	11:00:00	12:30:00	1000	1000	CAS T	sai	500	
CS 460	16:00:00	17:30:00	7000	7000	MCS 2	05	30	
CS 510	12:00:00	13:30:00	7000	7000	MCS 2	05	30	
Course.na	me Room.nan	1e	• Th	e las	t row o	of Course		
CS 105	GCB 204		do	osn't	have	a match in	Roo	m
CS 111	CAS 314							
EN 101	CAS Tsai		•	it is	an "ui	nmatched r	°ow"	
	MCS 205		 thus it's not in the result of the 					



Course					Room			
name	start_time	end_time	room_id		id	name		capacity
CS 105	13:00:00	14:00:00	4000		1000	CAS Tsai		500
CS 111	09:30:00	11:00:00	5000		2000	CAS BigRoom		100
EN 101	11:00:00	12:30:00	1000		3000	EDU Lecture	Hall	100
CS 460	16:00:00	17:30:00	7000		4000	GCB 204		40
CS 510	12:00:00	13:30:00	7000		5000	CAS 314		80
PH 101	14:30:00	16:00:00	NULL		6000	CAS 226		50
result of th	e LEFT OUTE	r Join			7000	MCS 205		30
	me start_time		room_id	id	Room	.name	capad	city
CS 105	13:00:00	14:00:00	4000	4000	GCB 2	204	40	
CS 111	09:30:00	11:00:00	5000	5000	CAS 3	314	80	
EN 101	11:00:00	12:30:00	1000	1000	CAS T	ſsai	500	
CS 460	16:00:00	17:30:00	7000	7000	MCS 2	205	30	
CS 510	12:00:00	13:30:00	7000	7000	MCS 2	205	30	
PH 101	14:30:00	16:00:00	NULL	NULL	NULL		NULL	
	me Room.nam	e	• A	left ol	iter jo	<i>in</i> adds an	extra	a row to
CS 105	GCB 204				•	any row fror		
CS 111	CAS 314					•		
EN 101	CAS Tsai		tha	at doe	esnith	nave a mato	cn in	the rigi
CS 460	MCS 205		•	uses	S NULL	s for the rigl	ht-tał	ble
CS 510	MCS 205					in the extra r		
PH 101	NULL			aun	Juicos		0,43	



Movie				Oscar			
id	name	rating	runtime		person_id		year
		 		2488496		BEST-ACTOR	2016
		 PG-13	124	1228705	2222222	BEST-ACTRESS	2011
	Toy Story 3 Despicable Me	 G PG	103 95	2488496	NULL	BEST-PICTURE	2016

How can we get just the movies that won Oscars? SELECT name FROM Movie, Oscar;

2

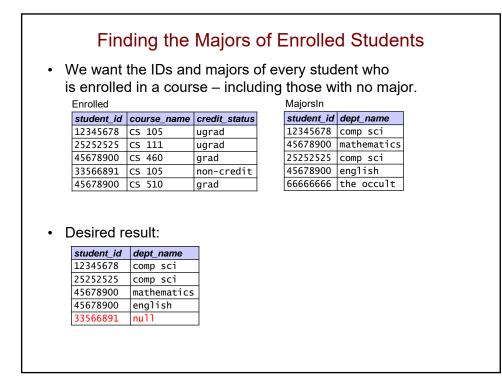
id		year	rating	runtime
2488496	Star Wars	2015	PG-13	138
1228705	Iron Man 2	2010	PG-13	124
0435761	Toy Story 3	2010	G	103
1323594	Despicable Me	2010	PG	95

Oscar		-	
movie_id	person_id	type	year
2488496	1111111	BEST-ACTOR	2016
1228705	2222222	BEST-ACTRESS	2011
2488496	NULL	BEST-PICTURE	2016

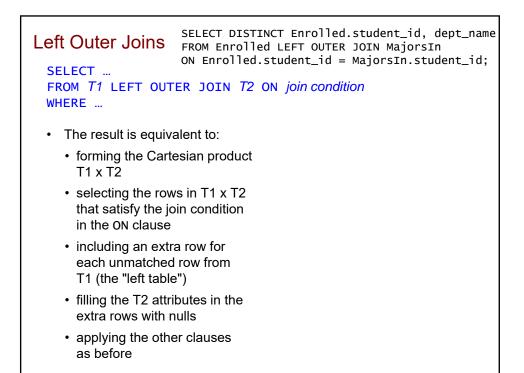
	Movie, Oscar id = movie_ BY name;				Oscar			
id	name	vear	rating	runtime		person_id	type	year
			PG-13		2488496	-	BEST-ACTOR	2016
1228705	Iron Man 2	2010	PG-13	124	1228705	2222222	BEST-ACTRESS	2011
0435761	Toy Story 3	2010	G	103	2488496	NULL	BEST-PICTURE	2016
	Despicable Me	2010	PG	95		•		·

F V C	What if we wanted a count for each movie SELECT name, COUNT(*) FROM Movie, Oscar WHERE id = movie_id GROUP BY name; Movie													
id	ovie	na	mo	vear	rati	na ri	Intime		id	person_id	tuno		year	
	88496			2015	_	13 1		248849		11111111	BEST-ACTO)R	2016	
				2010	-	-	24	122870	-	22222222	BEST-ACTR		2011	
043	35761	TO		2010	-	1	03	248849	6	NULL	BEST-PICT	URE	2016	
	Movie		spicable Me Oscar, fol	lowe	d by									1
	id		name		ovie. ear	rating	runtime	movie_	id	person_id	type		Oscar. year	
	248849	96	Star Wars	. 20)15	PG-1	3 138	248849	6	1111111	BEST-ACTO	R	2016	1
	248849	96	Star Wars	. 20)15	PG-1	3 138	248849	6	NULL	BEST-PICT	URE	2016	1
	122870	05	Iron Man 2	20	010	PG-1	3 124	122870	5	2222222	BEST-ACTR	ESS	2011	Γ
		_	after SELE	ст										
		[name	C	OUN	T(*)		nan	ıе		COUNT			
			Star Wars	. 2				Sta	r	wars	2			
			Iron Man 2	1				Iro	n		1			
								-		tory 3 cable Me	0			

	l.	Nhi	ch of	fthe	se	W	วมได	d work	?		
Movie					_	Osc	ar				
id	name		-	runtime	•			person_id	type		year
	Star Wars							1111111	BEST-ACTO	R	2016
1228705	Iron Man 2	2010	PG-13	124		122	8705	2222222	BEST-ACTR	ESS	2011
0435761	Toy Story 3	2010	G	103		248	8496	NULL	BEST-PICT	URE	2016
1323594	Despicable Me	2010	PG	95					•		
Λ	SELECT name,	COU	NT(*)								
	FROM Movie,						name	•	COUNT]	
	WHERE id = movie_id Star wars 2							2]		
	GROUP BY nam	e;					Iron	Man 2	1]	
D.	SELECT name, FROM Movie, WHERE id = m GROUP BY nam	Osca ovie	r	be)				Story 3 icable Me	0]	
	SELECT name, FROM Movie L ON id = movi GROUP BY nam	e_id	NT(typ OUTER	De) JOIN	0s(car					
D.	SELECT name, FROM Movie L ON id = movi GROUP BY nam	EFT (e_id		JOIN	0s	car					



	Whic	ch of th	iese wou	ıld work'	?	
		Enrolled			MajorsIn	
		student_id	course_name	credit_status	student_id	dept_name
	ne IDs and majors	12345678	CS 105	ugrad	12345678	comp sci
of every st	udent who is	25252525	CS 111	ugrad	45678900	mathematics
enrolled in	a course –	45678900	CS 460	grad	25252525	
including t	hose with no major.	33566891	CS 105	non-credit	45678900	english
	nose with no major.	45678900	CS 510	grad	66666666	the occult
В.	FROM Enrolled, WHERE Enrolled SELECT DISTINC FROM MajorsIn WHERE Enrolled	l.student CT Enroll LEFT OUT	t_id = Majo led.studen1 FER JOIN Er	t_id, dept <u>.</u> nrolled	_name	
C.	SELECT DISTING FROM Enrolled ON Enrolled.st	LEFT OUT	FER JOIN Ma	ajorsIn		
D.	SELECT DISTING FROM MajorsIn ON Enrolled.st	LEFT OUT	FER JOIN E	nrolled		



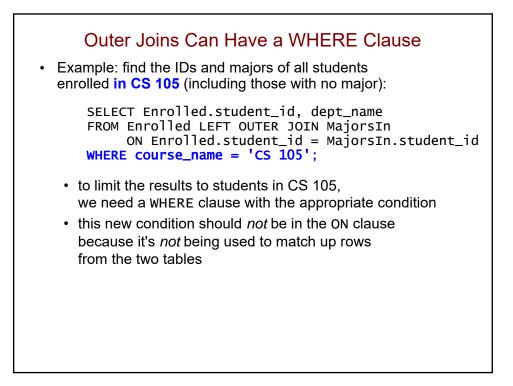
Left Outer Joins SELECT DIS FROM Enrol ON Enrolle FROM <i>T1</i> LEFT OUTER JOIN <i>T</i>	led LEFT d.student	OUTER _id =	JOIN Major	MajorsI	n
WHERE	Enrolled. student_id	course_ name	credit_ status	MajorsIn. student_id	dept_name
 The result is equivalent to: 	12345678	CS 105	ugrad	12345678	comp sci
The result is equivalent to.	12345678	CS 105	ugrad	45678900	math
 forming the Cartesian product 	t 12345678	CS 105	ugrad	25252525	comp sci
T1 x T2	12345678	CS 105	ugrad	45678900	english
11 × 12	12345678	CS 105	ugrad	66666666	the occult
 selecting the rows in T1 x T2 	25252525	CS 111	ugrad	12345678	comp sci
that satisfy the join condition	25252525	CS 111	ugrad	45678900	math
in the ON clause	25252525	CS 111	ugrad	25252525	comp sci
In the on clause	25252525	CS 111	ugrad	45678900	english
 including an extra row for 	25252525	CS 111	ugrad	66666666	the occult
each unmatched row from	45678900	CS 460	grad	12345678	comp sci
	45678900	CS 460	grad	45678900	math
T1 (the "left table")	45678900	CS 460	grad	25252525	comp sci
 filling the T2 attributes in the 	45678900	CS 460	grad	45678900	english
extra rows with nulls	45678900	CS 460	grad	66666666	the occult
	33566891	CS 105	non-cr	12345678	comp sci
 applying the other clauses 	33566891	CS 105	non-cr	45678900	math
as before	33566891	CS 105	non-cr	25252525	comp sci
	33566891	CS 105	non-cr	45678900	english
	22566001	CC 105		CCCCCCCC	

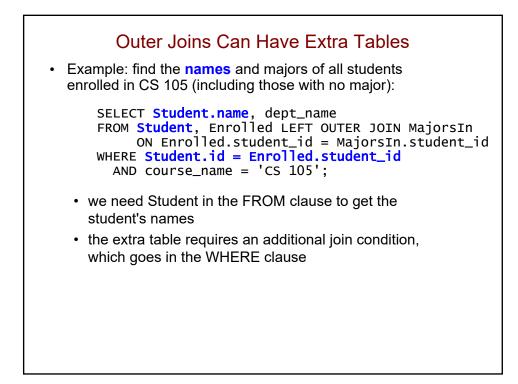
Left Outer Joins SELECT DIST FROM Enroll ON Enrolled SELECT FROM T1 LEFT OUTER JOIN T2	ed LEFT .student	OUTER _id =	JOIN Major	MajorsI	n
WHERE	Enrolled. student_id	course_ name	credit_ status	MajorsIn. student_id	dept_name
 The result is equivalent to: 	12345678		ugrad	12345678	
·	25252525	CS 111	ugrad	25252525	comp sci
 forming the Cartesian product 	45678900	CS 460	grad		math
T1 x T2	45678900	CS 460	grad	45678900	english
 selecting the rows in T1 x T2 	45678900	CS 510	grad	45678900	math
that satisfy the join condition in the ON clause	45678900	CS 510	grad	45678900	english
 including an extra row for each unmatched row from T1 (the "left table") 					
 filling the T2 attributes in the extra rows with nulls 					
 applying the other clauses as before 					

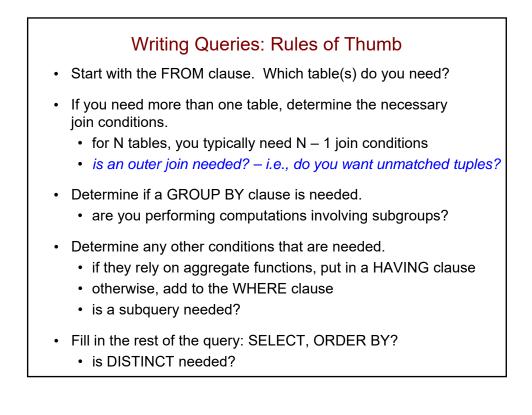
Left Outer Joins SELECT DISTINCT Enrolled.student_id, dept_name FROM Enrolled LEFT OUTER JOIN MajorsIn ON Enrolled.student_id = MajorsIn.student_id; FROM T1 LEFT OUTER JOIN T2 ON join condition									
WHERE	Enrol		course name	e_ credit_ status		orsIn. lent id	dept_na	ame	
The manufactor and include the		5678	CS 10		-	45678	comp s	ci	
 The result is equivalent to: 		2525				52525	comp s		
 forming the Cartesian product 			CS 46	. 5		78900	math		
T1 x T2	4567		CS 46	5		78900	englis		
	4567		CS 51	5		78900	math		
 selecting the rows in T1 x T2 		8900	CS 51	3		78900	englis	-	
		6891	CS 10	3	-	/ 8900	engris		
that satisfy the join condition	3330	0091	102 10						
in the ON clause									
, including on outro row for		Enroll							
 including an extra row for 				course_n	ame	credit	_status		
each unmatched row from		1234	5678	CS 105		ugrad			
T1 (the "left table")		2525	2525	CS 111		ugrad			
, , ,		4567	8900	CS 460		grad			
 filling the T2 attributes in the 		3356	6891	CS 105		non-c	redit		
extra rows with nulls		4567	8900	CS 510		grad			
 applying the other clauses as before 									

Left Outer Joins SELECT FROM <i>T1</i> LEFT OUTI	SELECT DIST FROM Enrolle ON Enrolled ER JOIN 72	ed LEFT .student	OUTER _id =	JOIN Major	MajorsIn	n	
WHERE		Enrolled. student_id	course_ name	credit_ status	MajorsIn. student_id	dept_name	
The result is equiva	alent to:	12345678		ugrad	12345678	comp sci	
		25252525	CS 111	ugrad	25252525	comp sci	
 forming the Carte 	45678900	CS 460	grad		math		
T1 x T2	45678900	CS 460	grad	45678900	english		
		CS 510	grad		math		
 selecting the row 	45678900	CS 510	grad	45678900	english		
that satisfy the jo in the ON clause	33566891	CS 105	non-cr	null	null		
 including an extra each unmatched T1 (the "left table 							
filling the T2 attributes in the extra rows with nulls							
 applying the othe as before 	r clauses						

Left Outer Joins SELECT DIST FROM Enroll ON Enrolled SELECT FROM <i>T1</i> LEFT OUTER JOIN <i>T2</i> WHERE	ed LEFT .student	OUTER _id =	JOIN Major:	MajorsI	n
WITERE	Enrolled. student id	course_ name	credit_ status	MajorsIn. student id	dept_name
The result is equivalent to:	12345678	CS 105	ugrad	12345678	comp sci
·	25252525	CS 111	ugrad	25252525	comp sci
 forming the Cartesian product 	45678900	CS 460	grad	45678900	math
T1 x T2	45678900	CS 460	grad	45678900	english
	45678900	CS 510	grad	45678900	math
 selecting the rows in T1 x T2 	45678900	CS 510	grad	45678900	english
that satisfy the join condition	33566891	CS 105	non-cr	null	null
in the ON clause					
 including an extra row for 					
each unmatched row from			nrolled. tudent_id	dept_na	me
T1 (the "left table")		1	2345678	comp so	:i
filling the TO attributes in the		2	5252525	comp so	
 filling the T2 attributes in the 		4	5678900	mathema	atics
extra rows with nulls		4	5678900	english	1 I
 applying the other clauses 		3	3566891	null	
as before					



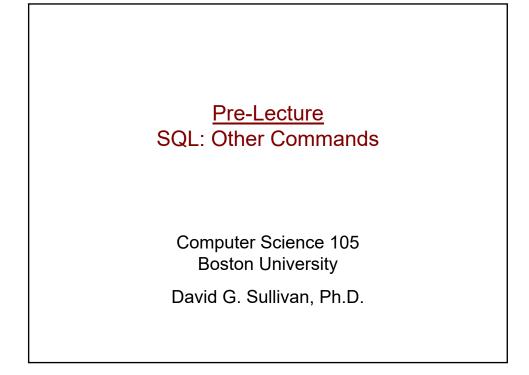


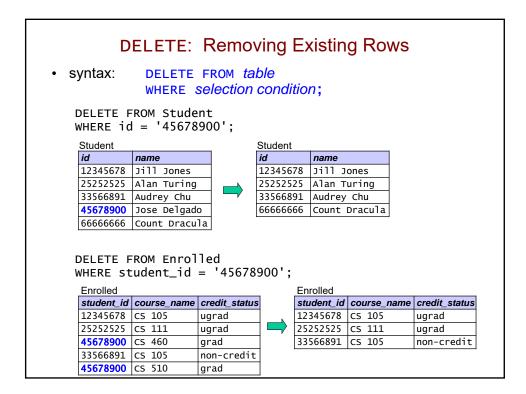


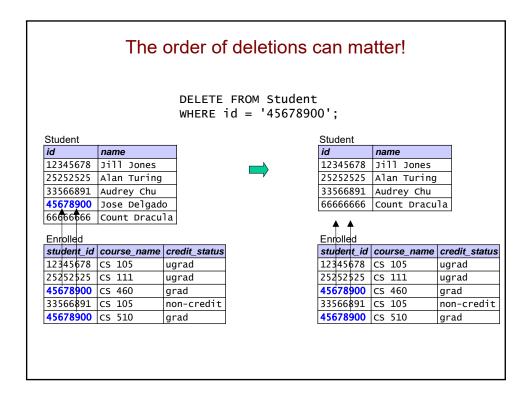


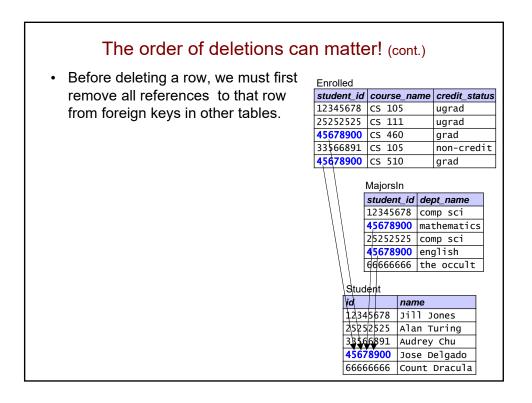
Movie						Oscar				
d	name	year	ratin	g rur	ntime	movie_id	person_id	type		year
2488496	Star Wars	2015	PG-1	.3 13	8	2488496	1111111	BEST-AC	TOR	2016
228705	Iron Man 2	2010	PG-1	.3 12	4	1228705	2222222	BEST-AC	TRESS	2011
0435761	Toy Story 3	2010	G	10	3	2488496	NULL	BEST-PI	CTURE	2016
.323594	Despicable Me	2010	PG	95						
Movie LEFT OUTER JOIN Oscar ON id = movie_id GROUP BY name										
id	name	Мс	vie. I	rating	runtime	movie_id	person_id	type		Oscar.
		yea								year
24884	96 Star Wars	. 20	15 F	PG-13	138	2488496	1111111	BEST-ACTOR		2016
24884	96 Star Wars	. 20	15 F	PG-13	138	2488496	NULL	BEST-PICTURE		2016
12287	05 Iron Man 2	20	10 F	PG-13	124	1228705	2222222	BEST-ACTRESS		2011
04357	61 Toy Story 3	20	10 0	3	103	NULL	NULL	IULL NULL		NULL
13235	94 Despicable	Me 20	10 I	۶G	95	NULL	NULL	IULL NULL		NULL
name COUNT(type)										
C S	ELECT name,	COU	NT(t	:ype)			Star Wa	rs	2	
U. _г	ROM Movie L	EFT (OUTE	R JC	IN OSC	ar 🛋	Iron Ma	n 2	1	
C	N id = movi	e_id					Toy Sto	ry 3	0	
Ģ	ROUP BY nam	e;					Despica	ble Me	0	
		COLU		.			name		COUNT	Γ(*)
	ELECT name,						Star Wa	rs	2	
	ROM Movie L N id = movi			K JU	USC NT		Iron Ma	n 2	1	
		_					Toy Sto	ry 3	1	
GROUP BY name; Despicable Me 1										

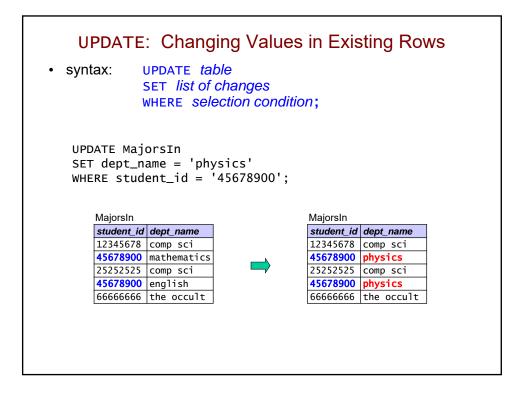
	MajorsIn								
	student_id	course_name	credit_status	student_id	dept_name				
We want the IDs and majors	12345678	CS 105	ugrad	12345678	comp sci				
of every student who is	25252525	CS 111	ugrad	45678900	mathematics				
enrolled in a course –	45678900	CS 460	grad	25252525	comp sci				
including those with no major.	33566891	CS 105	non-credit	45678900	english				
including those with no major.	45678900	CS 510	grad	66666666	the occult				
				student id	dant nome				
SELECT DISTINCT Enrolle		_id, dept_	name		dept_name				
FROM Enrolled, MajorsIn					comp sci				
WHERE Enrolled.student_		comp sci							
	45678900	mathematics							
SELECT DISTINCT Enrolle	45678900	english							
FROM MajorsIn LEFT OUTE	FROM MajorsIn LEFT OUTER JOIN Enrolled								
WHERE Enrolled.student_	student_id	dept_name							
	12345678	comp sci							
SELECT DISTINCT Enrolle	SELECT DISTINCT Enrolled.student_id, dept_name								
FROM Enrolled LEFT OUTE	45678900	mathematics							
ON Enrolled.student_id =	= MajorsIr	1.student_i	d;	45678900	english				
get unmatched rows from the	get unmatched rows from the "left" table Enrolled								
got dimitatorio di forto il orin di o									
				student_id	dept_name				
SELECT DISTINCT Enrolle	SELECT DISTINCT Enrolled.student_id, dept_name								
FROM MajorsIn LEFT OUTE	R JOIN Er	nrolled		25252525	comp sci				
ON Enrolled.student_id =	= MajorsIr	n.student_i	d;	45678900	mathematics				
get unmatched rows from the		45678900	english						
get annatched tows norr the	ion tubic	najor 511		null	the occult				

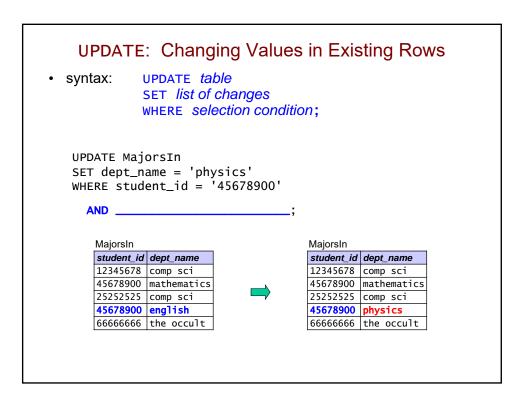


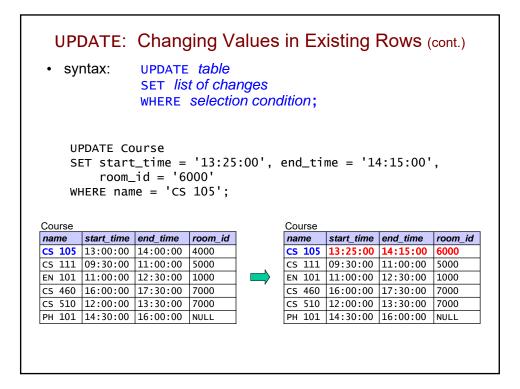


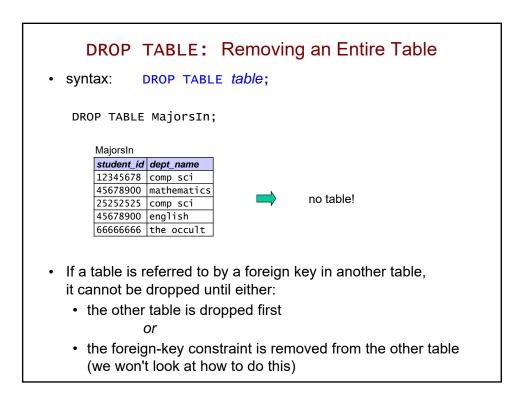


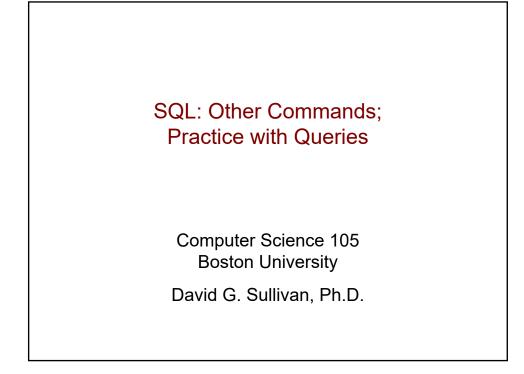








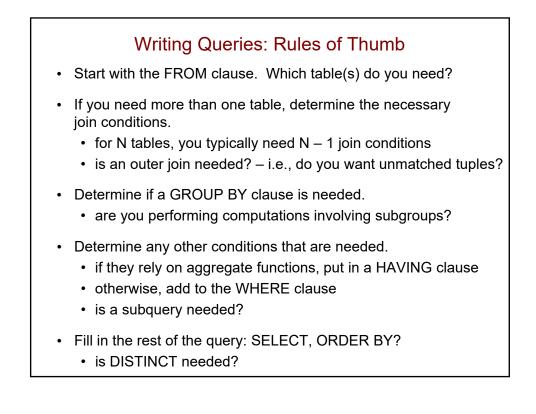


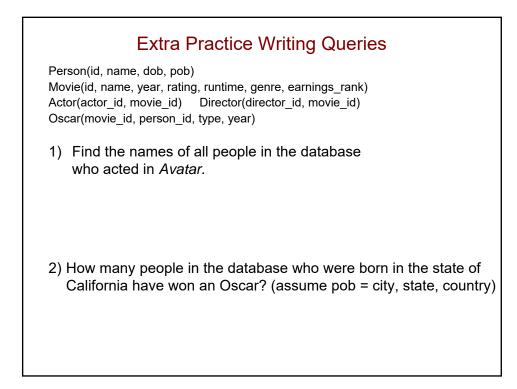


Course Enrolled							
name	start_time	end_time	room_id		student_id	course_name	credit_statu
CS 105	13:00:00	14:00:00	4000		12345678	CS 105	ugrad
CS 111	09:30:00	11:00:00	5000		25252525	CS 111	ugrad
EN 101	11:00:00	12:30:00	1000		45678900	CS 460	grad
CS 460	16:00:00	17:30:00	7000		33566891	CS 105	non-credi
CS 510	12:00:00	13:30:00	7000		45678900	CS 510	grad
PH 101	14:30:00	16:00:00	NULL	nam			
<u>рн 101</u> 1 DE		16:00:00	NULL		ie = 'C	s 111';	
<u>рн 101</u> 1 DE	LETE FRO	16:00:00	e WHERE	RE C	ie = 'C	s 111';	

Course		.			Enrolled			
name	start_time	end_time	room_id		student_i	d course_name	credit_statu	
CS 105	13:00:00	14:00:00	4000		12345678	CS 105	ugrad	
CS 111	09:30:00	11:00:00	5000		25252525	CS 111	ugrad	
EN 101	11:00:00	12:30:00	1000 \		45678900	CS 460	grad	
CS 460	16:00:00	17:30:00	7000		33566891	. CS 105	non-credit	
CS 510	12:00:00	13:30:00	7000		45678900	CS 510	grad	
PH 101	14:30:00	16:00:00	NULL		Room	I I		
					id i	name	capacity	
					1000	CAS Tsai	500	
					2000	CAS BigRoom	100	
					3000	EDU Lecture Ha	all 100	
					4000	GCB 204	40	
					5000	CAS 314	80	
					6000	CAS 226	50	
					7000 I	ACS 205	30	
ELETE FR	ROM ROOM	led WHER WHERE id	= '500	0';	1	CS 111'; 10! when dele hat includes a	•	

name				Room		
	start_time	end_time	room_id	id	name	capacity
CS 105	13:00:00	14:00:00	4000	1000	CAS Tsai	500
CS 111	09:30:00	11:00:00	5000	2000	CAS BigRoom	100
EN 101	11:00:00	12:30:00	1000	3000	EDU Lecture Hall	100
CS 460	16:00:00	17:30:00	7000	4000	GCB 204	40
CS 510	12:00:00	13:30:00	7000	5000	CAS 314	80
PH 101	14:30:00	16:00:00	NULL	6000	CAS 226	50
				7000	MCS 205	30
n	DELETE FRO UPDATE Cou				; HERE room_id =	'7000'

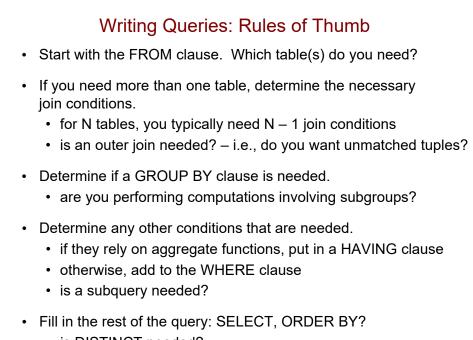


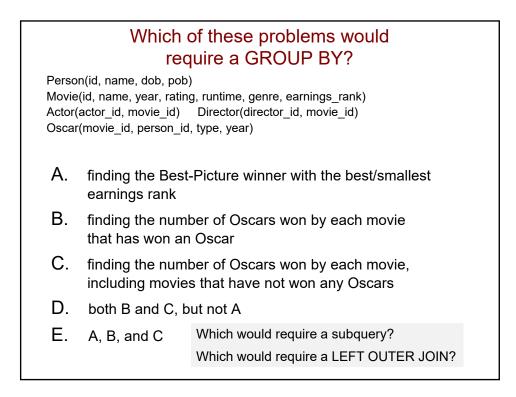


Extra Practice Writing Queries (cont.)

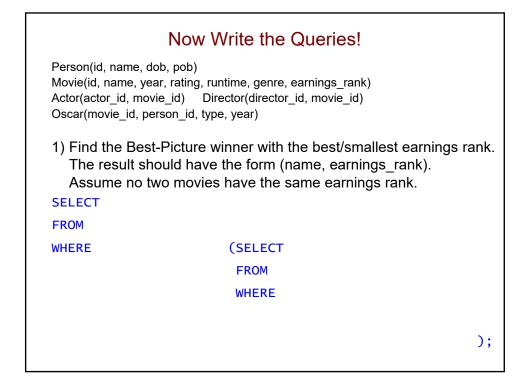
```
Person(id, name, dob, pob)
Movie(id, name, year, rating, runtime, genre, earnings_rank)
Actor(actor_id, movie_id) Director(director_id, movie_id)
Oscar(movie_id, person_id, type, year)
3) How many people in the database did <u>not</u> act in Avatar?
Why won't this work?
SELECT COUNT(*)
FROM Person P, Actor A, Movie M
WHERE P.id = A.actor_id AND M.id = A.movie_id
AND M.name != 'Avatar';
What will?
```

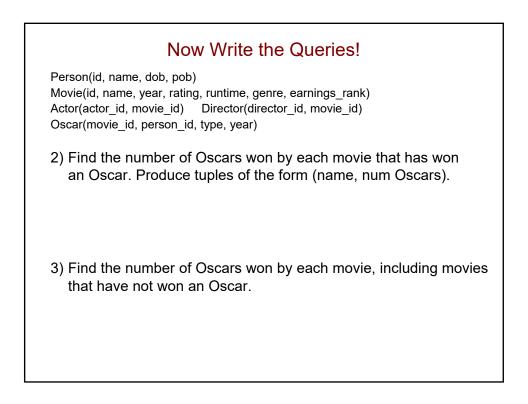
SQL: More Practice with Queries Computer Science 105 Boston University David G. Sullivan, Ph.D.









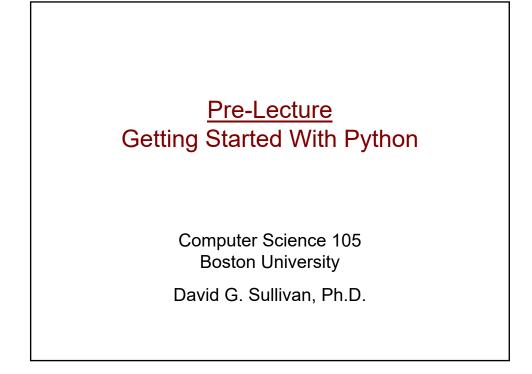


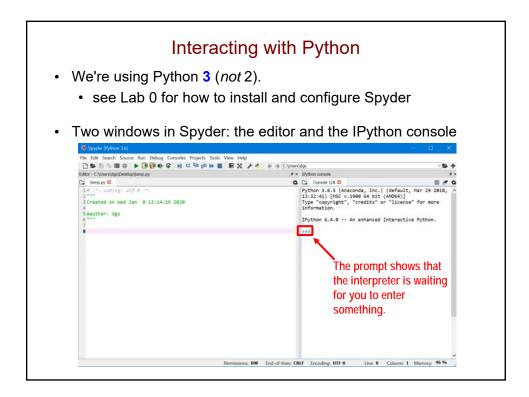
Even More Practice!

Person(id, name, dob, pob) Movie(id, name, year, rating, runtime, genre, earnings_rank) Actor(actor_id, movie_id) Director(director_id, movie_id) Oscar(movie_id, person_id, type, year)

4) Which movie ratings have an avg runtime greater than 120 min?

Even More Practice! (cont.) Person(id, name, dob, pob) Movie(id, name, year, rating, runtime, genre, earnings_rank) Actor(actor_id, movie_id) Director(director_id, movie_id) Oscar(movie_id, person_id, type, year) 5) For each person in the database born in Boston, find the number of movies in the database (possibly 0) in which the person has acted.





Arithmetic in Python

- Numeric operators include:
 - + addition
 - subtraction
 - * multiplication
 - / division
 - ** exponentiation
 - % modulus: gives the remainder of a division

Arithmetic in Python (cont.)

- The operators follow the standard order of operations.
 - example: multiplication before addition
- You can use parentheses to force a different order.

Data Types

• Different kinds of values are stored and manipulated differently.

- Python data types include:
 - integers
 - example: 451
 - floating-point numbers
 - numbers that include a decimal
 - example: 3.1416

Data Types and Operators

- There are really two sets of numeric operators:
 - one for integers (ints)
 - one for floating-point numbers (floats)
- In most cases, the following rules apply:
 - if at least one of the operands is a float, the result is a float
 - if both of the operands are ints, the result is an int
- One exception: division!

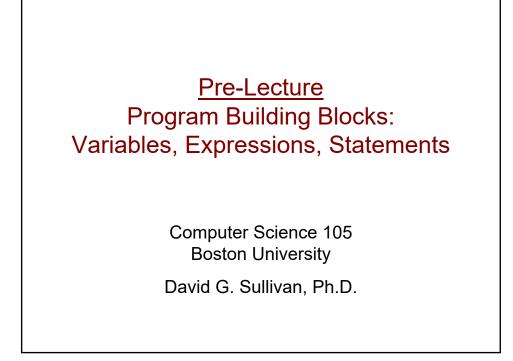
Two Types of Division

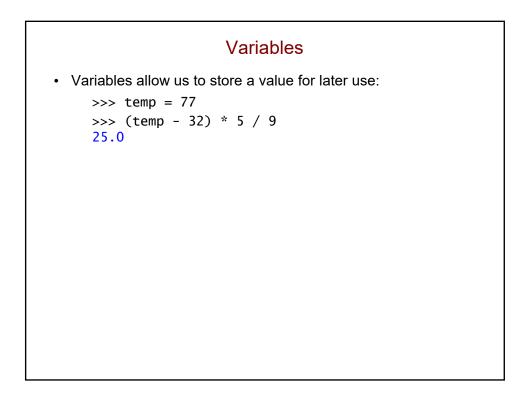
• The / operator always produces a float result.

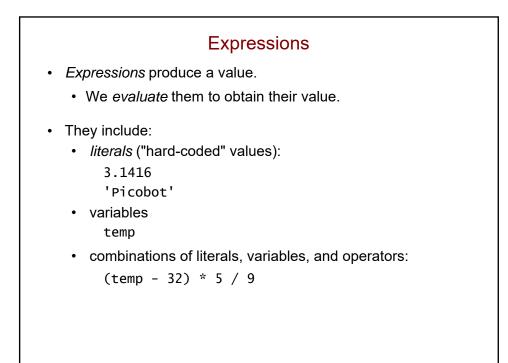
Two Types of Division (cont.) There is a separate // operator for *integer* division. >>> 6 // 3 Integer division *discards* any fractional part of the result: >>> 11 // 5 >>> 5 // 3 Note that it does *not* round!

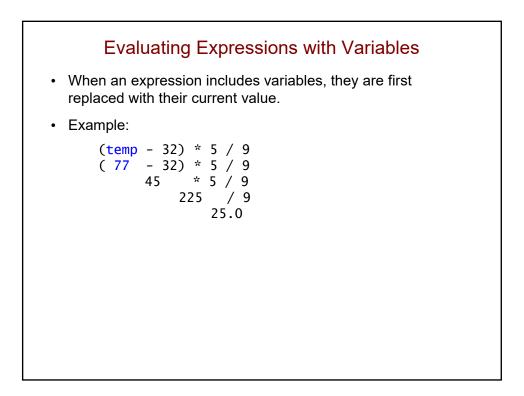
Another Data Type

- A string is a sequence of characters/symbols
 - surrounded by single or double quotes
 - examples: "hello" 'Picobot'





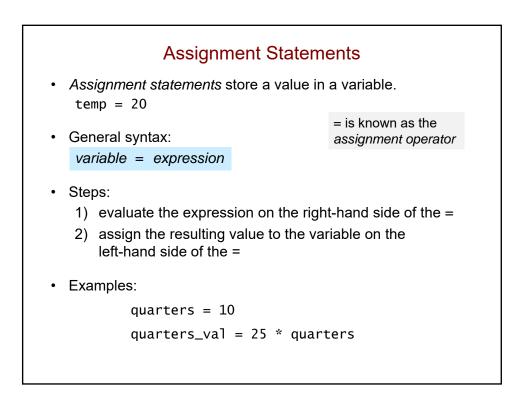


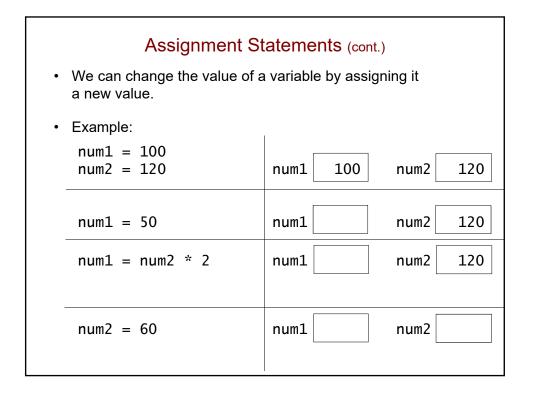


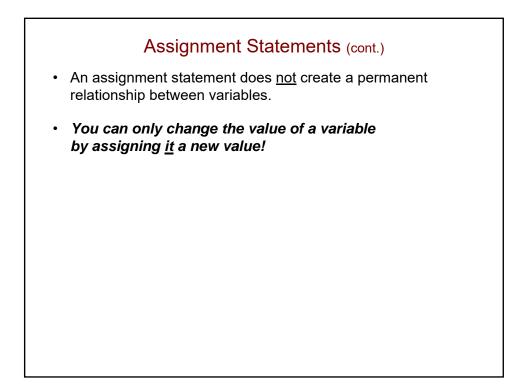
Statements

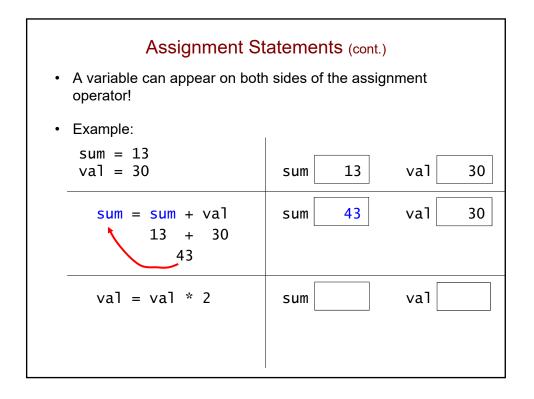
- A statement is a command that carries out an action.
- A program is a sequence of statements.

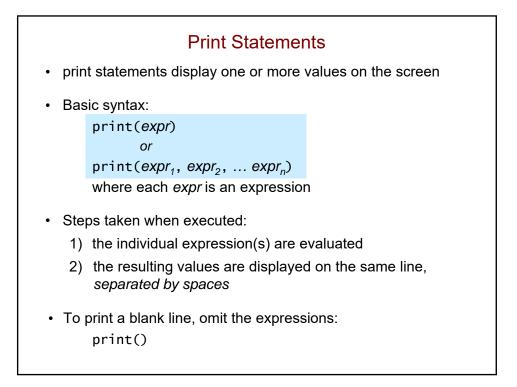
```
quarters = 2
dimes = 3
nickels = 1
pennies = 4
cents = quarters*25 + dimes*10 + nickels*5 + pennies
print('you have', cents, 'cents')
```

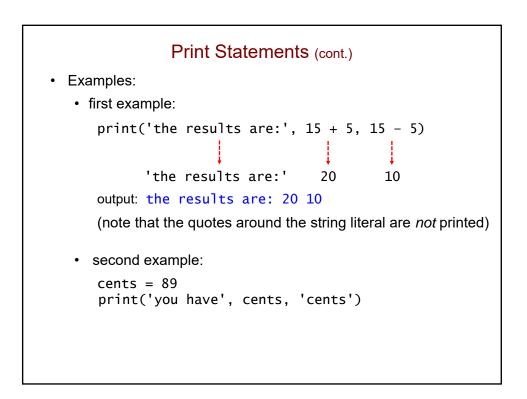


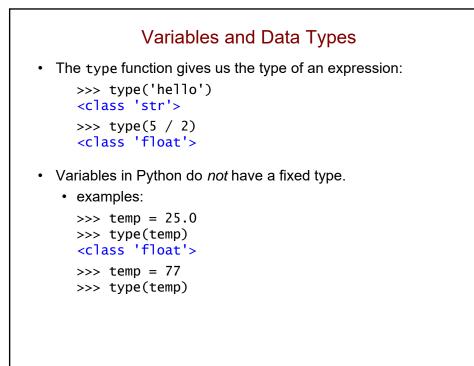


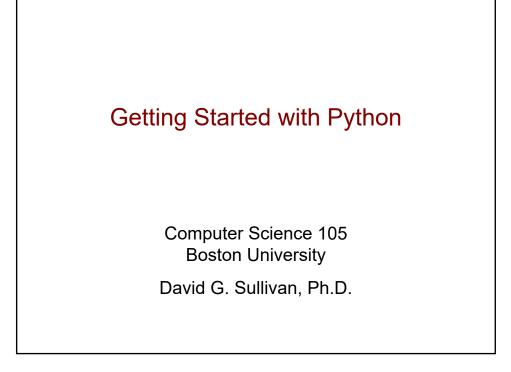


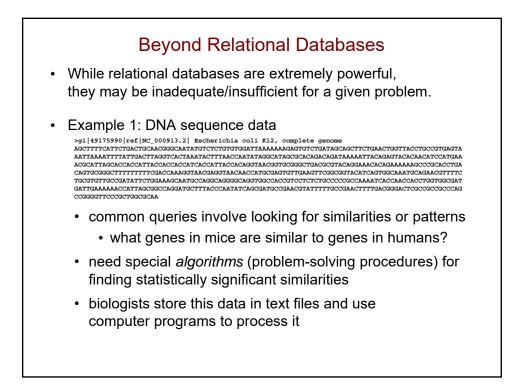


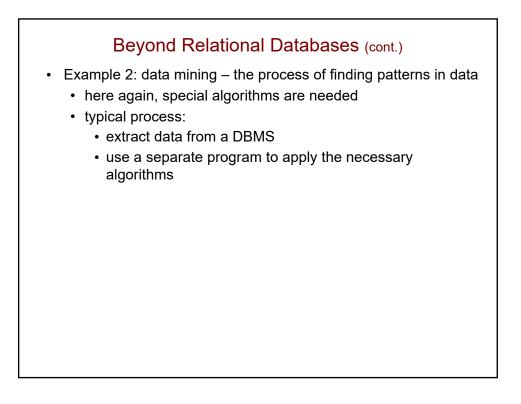


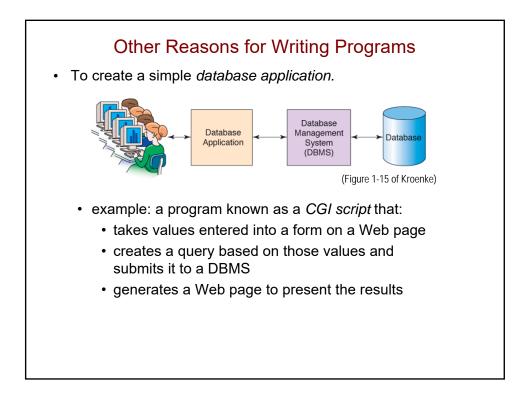


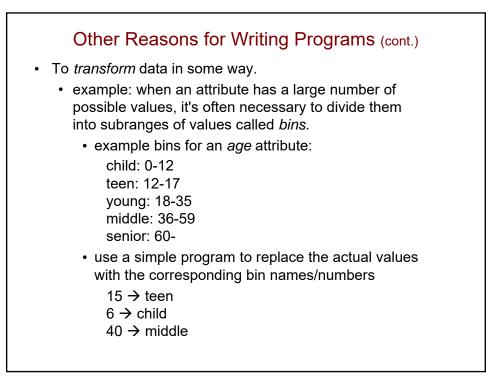


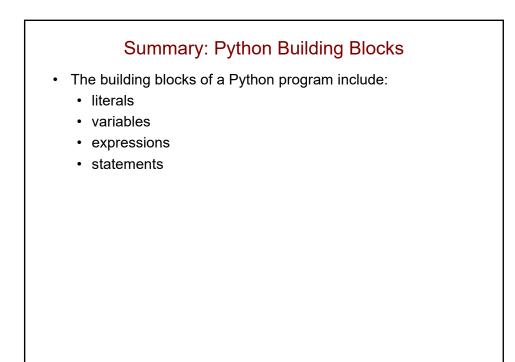








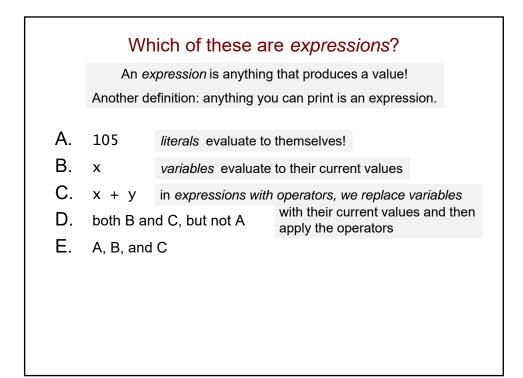


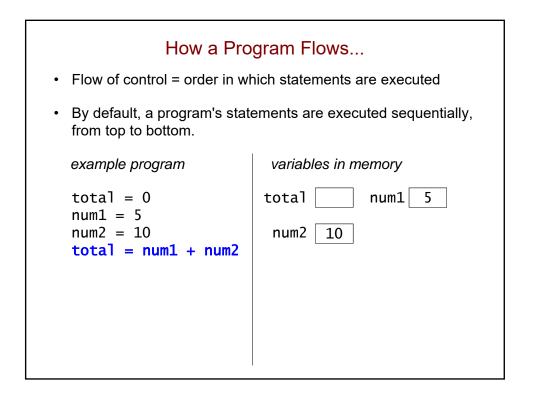




- A. 105
- B. x
- C. x + y
- D. both B and C, but not A
- E. A, B, and C





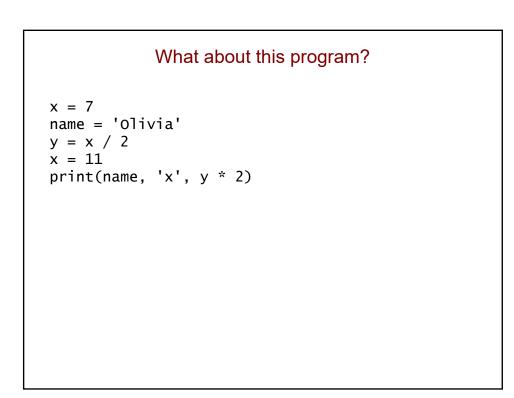


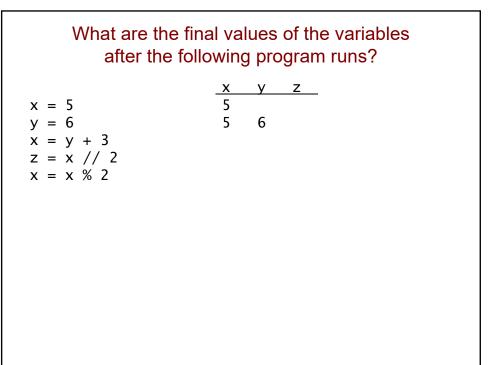
What is the output of the following program?

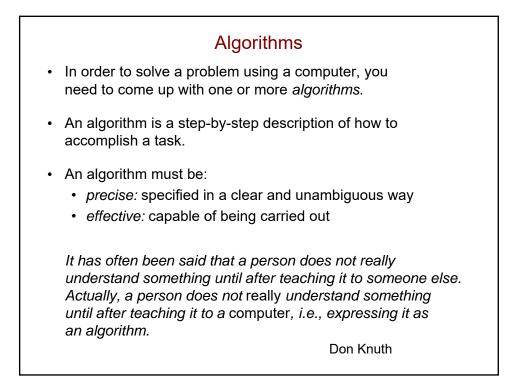
x = 7 name = 'Olivia' y = x / 2 x = 11 print('name', x, y)

note: we do *not* print:

- commas between expressions
- quotes around string literals





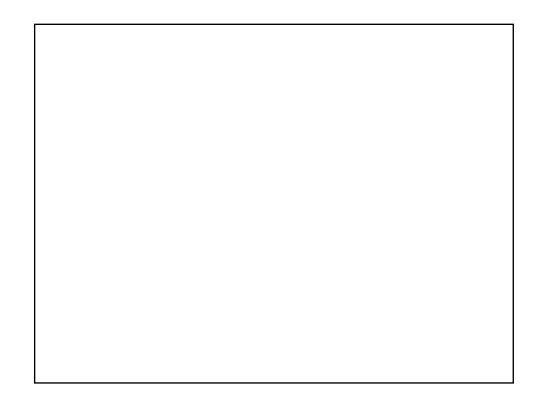


Is This An Algorithm?

• Recipe for preparing a meat roast:

Sprinkle the roast with salt and pepper. Insert a meat thermometer and place in oven preheated to 150 degrees C. Cook until the thermometer registers 80-85 degrees C. Serve roast with gravy prepared from either meat stock or from pan drippings if there is sufficient amount.

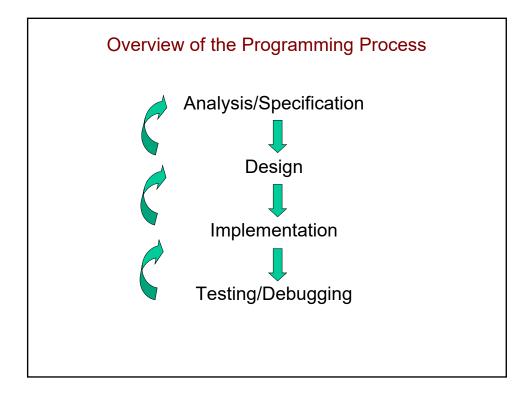
(taken from a book on programming by Pohl and McDowell)

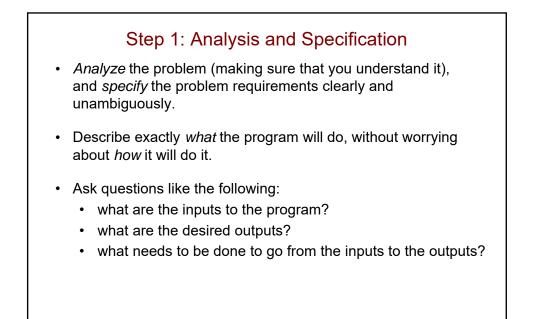


Here's the Algorithm...

- Recipe for preparing a meat roast:
 - 1. Sprinkle roast with 1/8 teaspoon salt and pepper.
 - 2. Turn oven on to 150 degrees C.
 - 3. Insert meat thermometer into center of roast.
 - 4. Wait a few minutes.
 - 5. If oven does not yet register 150 degrees, return to step 4.
 - 6. Place roast in oven.
 - 7. Wait a few minutes.
 - 8. Check meat thermometer. If temperature is less than 80 degrees C, go back to step 7.
 - 9. Remove roast from oven.
 - 10. If there is at least $\frac{1}{2}$ cup of pan drippings, go to step 12.
 - 11. Prepare gravy from meat stock and go to step 13.
 - 12. Prepare gravy from pan drippings.
 - 13. Serve roast with gravy.

(also from Pohl and McDowell)

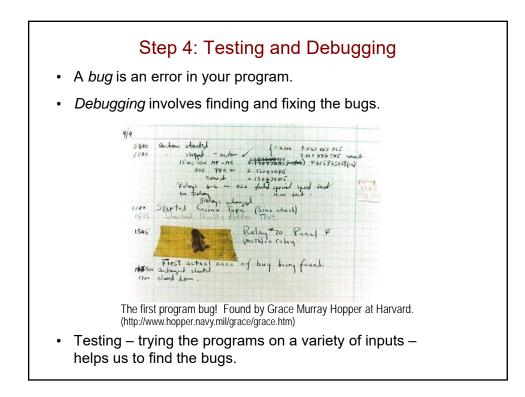


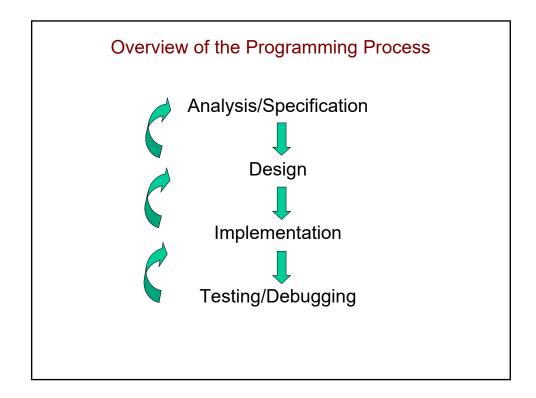


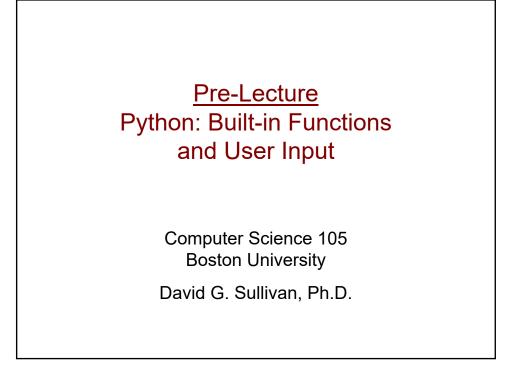
Step 2: Design
 Determine the necessary algorithms (and possibly other aspects of the program) and sketch out a design for them.
 This is where we figure out <i>how</i> the program will solve the problem.
 Algorithms are often designed using <i>pseudocode</i>. more informal than an actual programming language allows us to avoid worrying about the <i>syntax</i> of the language example for our change-adder problem from the video: get the number of quarters get the number of dimes get the number of nickels get the number of pennies compute the total value of the coins

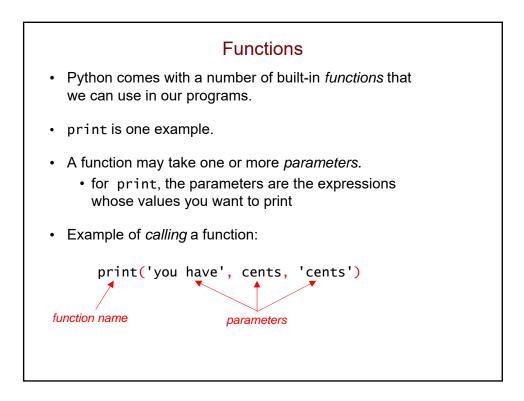
Step 3: Implementation

- Translate your design into the programming language.
 pseudocode → code
- We need to learn more Python before we can do this!

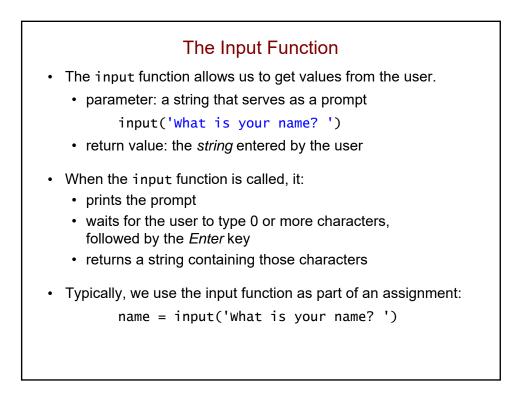


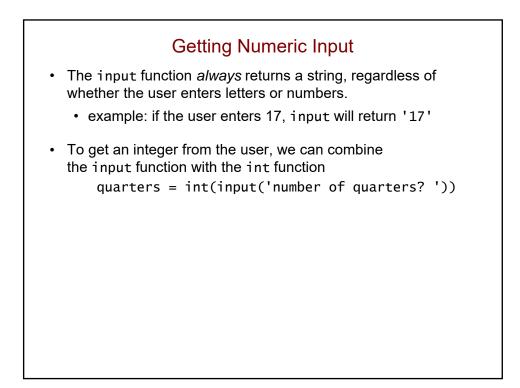


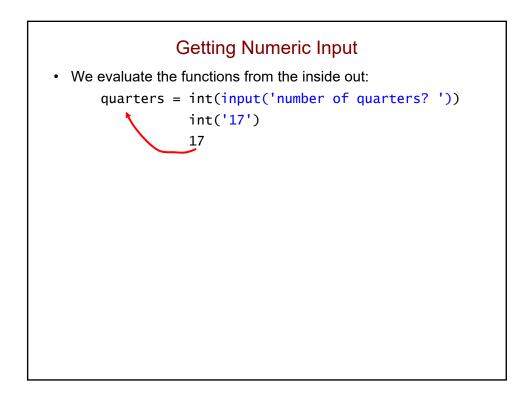


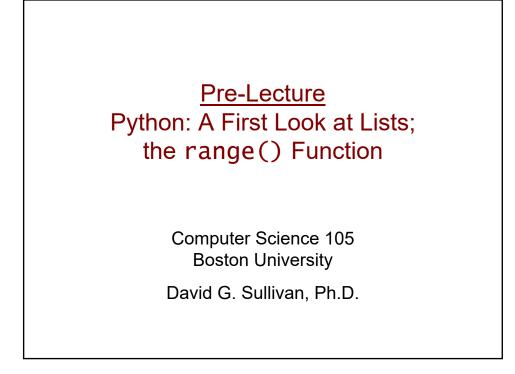


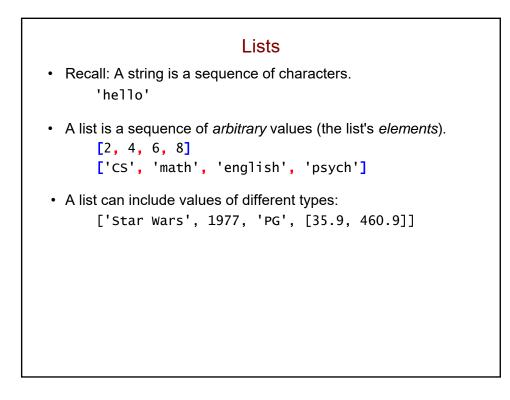
Functions (cont.) Some functions return (i.e., output) a value. Example: the abs function parameter: a number n return value (output): the absolute value of n Example: the int function parameter: a string representing a number return value (output): the number as a value of type int examples: int('15') returns 15 int('3.75') returns _____

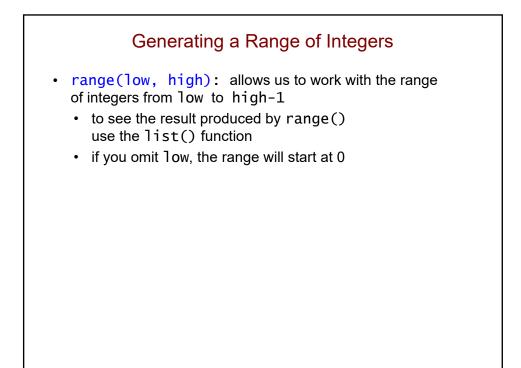


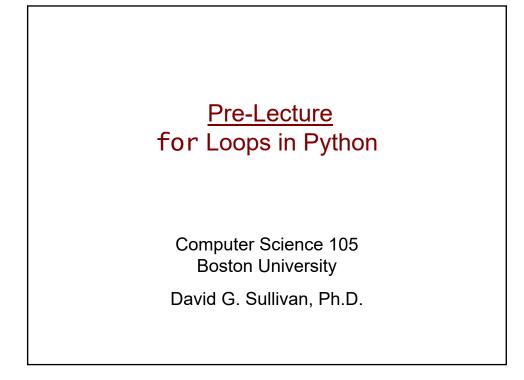


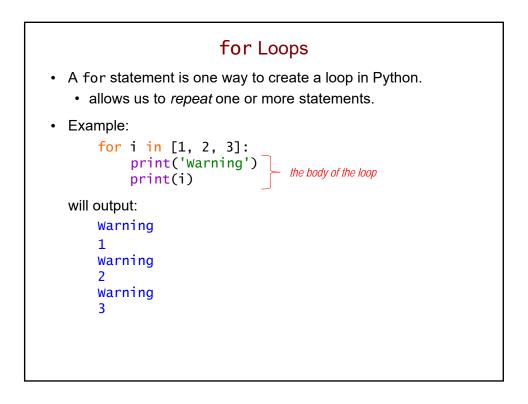


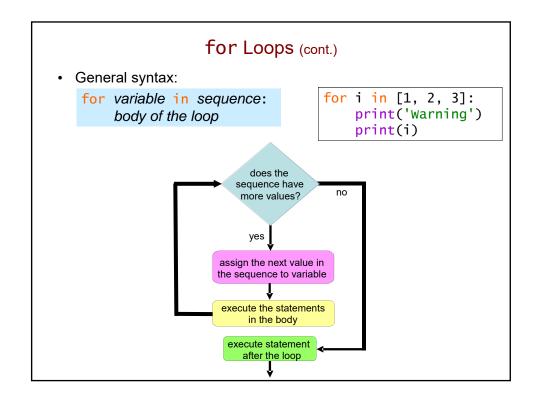


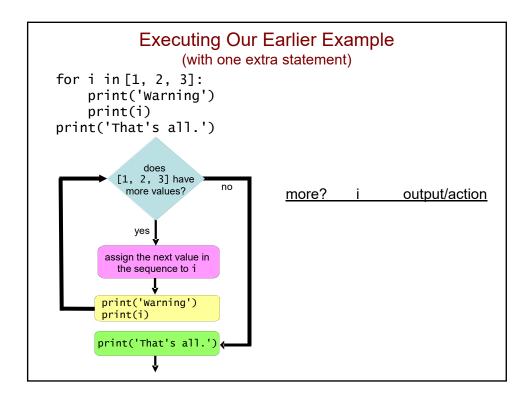


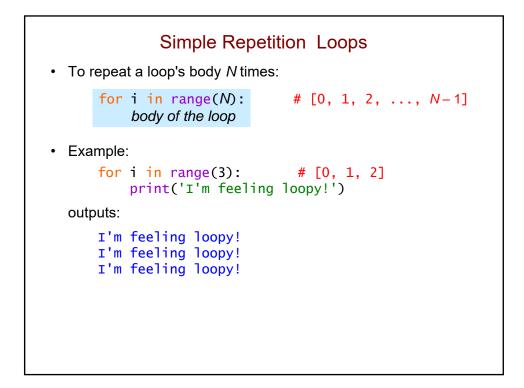


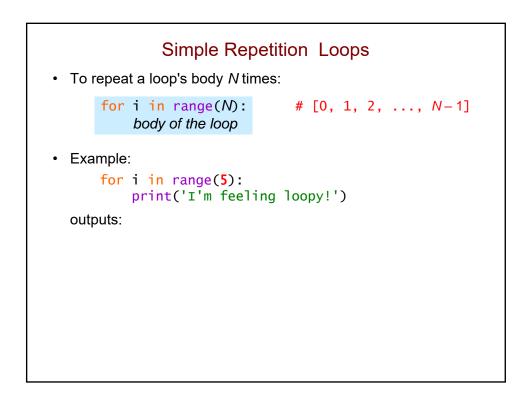


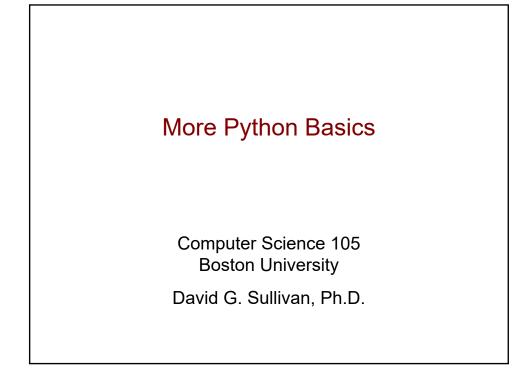


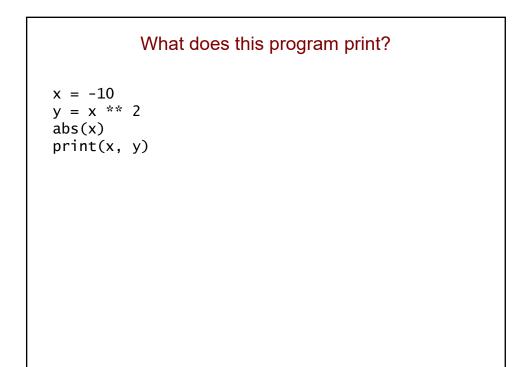






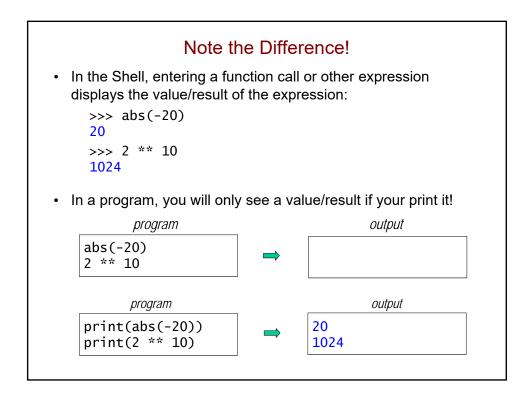






How could we make the program print 10 100

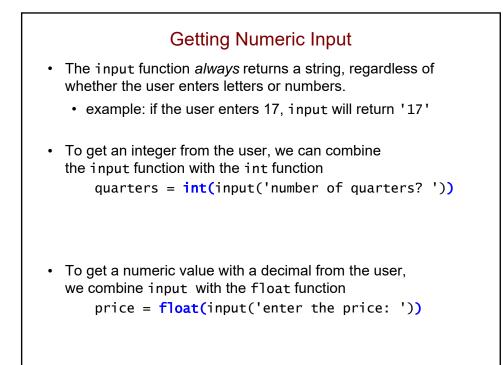
x = -10 y = x ** 2 abs(x) print(x, y)



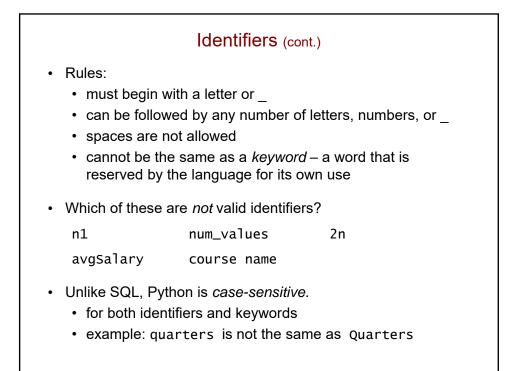
Recall Our Earlier Example Program...

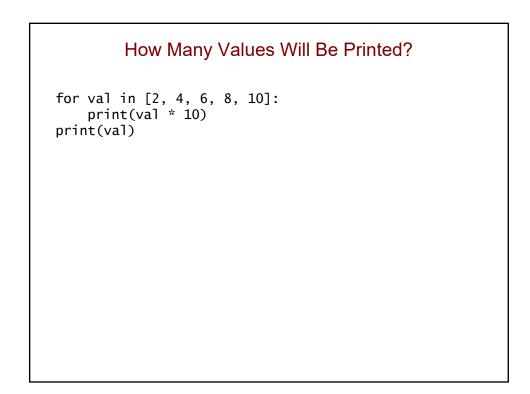
```
quarters = 2
dimes = 3
nickels = 1
pennies = 4
cents = quarters*25 + dimes*10 + nickels*5 + pennies
print('you have', cents, 'cents')
```

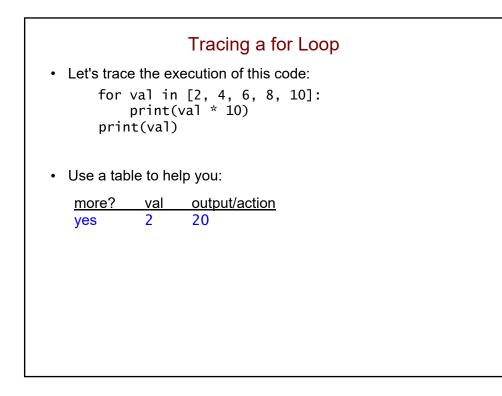
An Improved Version with User Input! quarters = int(input('number of quarters? ')) dimes = int(input('number of nickels? ')) nickels = int(input('number of pennies? ')) pennies = int(input('number of pennies? ')) cents = quarters*25 + dimes*10 + nickels*5 + pennies print('you have', cents, 'cents') • Note the use of the int() function to convert the user's inputs to integers.

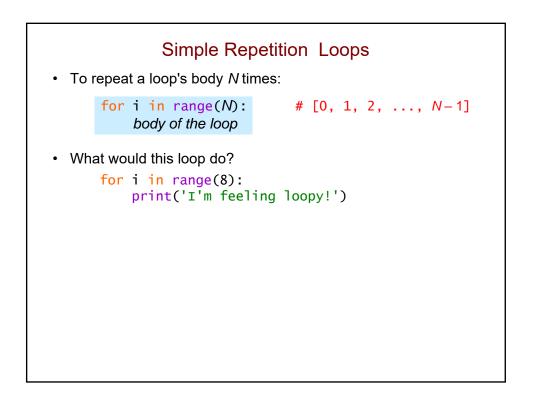


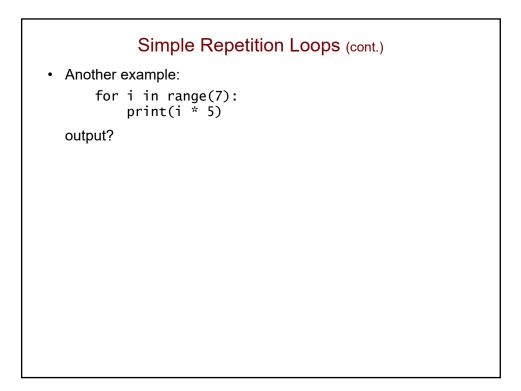
Identifiers quarters = int(input('number of guarters? ')) dimes = int(input('number of dimes? ')) nickels = int(input('number of nickels? ')) pennies = int(input('number of pennies? ')) cents = quarters*25 + dimes*10 + nickels*5 + pennies print('you have', cents, 'cents') Identifiers are words that are used to name components of a Python program. • They include: · variables, which give a name to a value dimes nickels pennies quarters cents function names like int, input and print

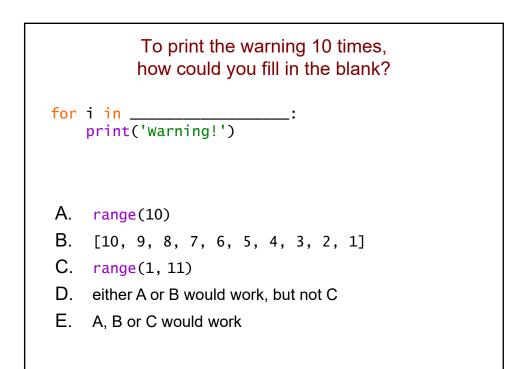


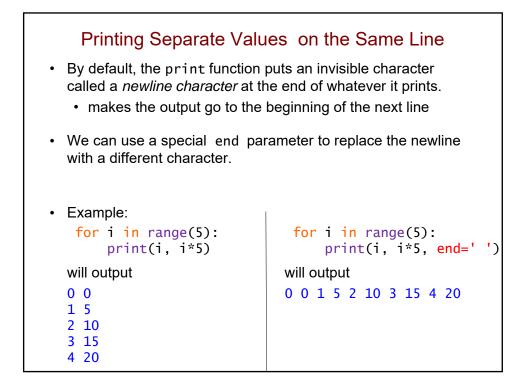


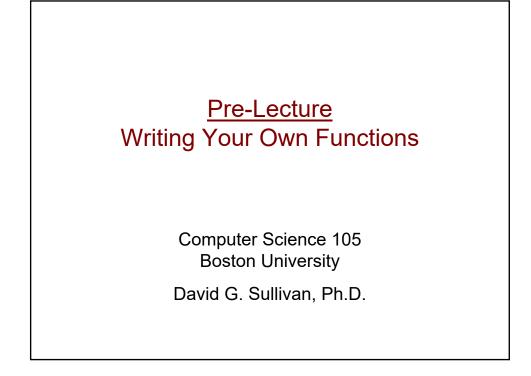


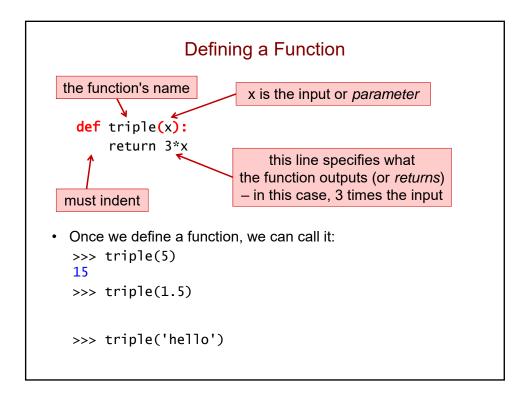




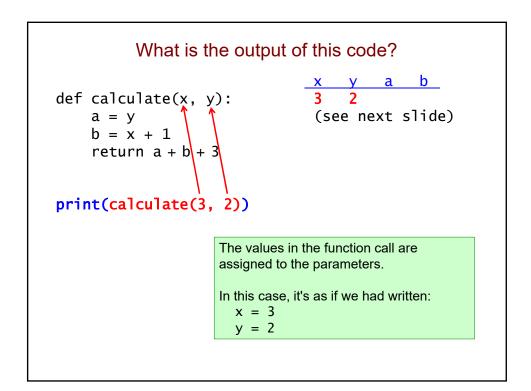


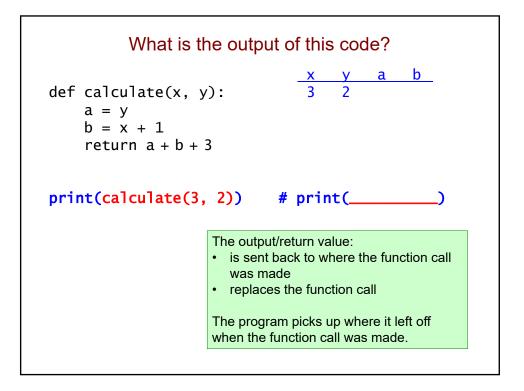


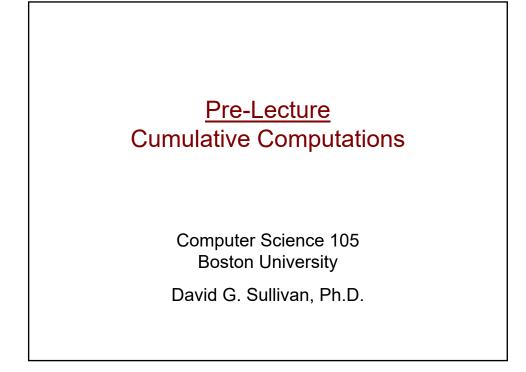


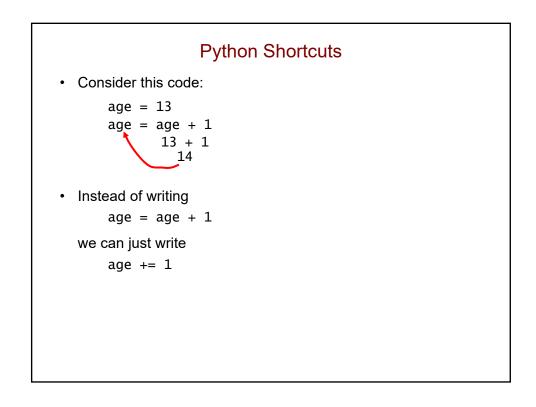


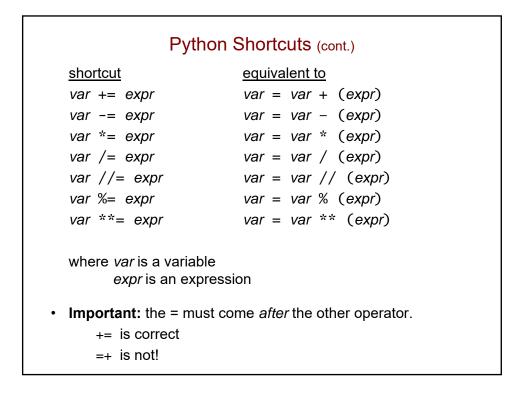
Multiple Lines, Multiple Parameters def circle_area(diam): """ Computes the area of a circle with a diameter diam. radius = diam / 2area = 3.14159 * (radius**2) return area def rect_perim(1, w): """ Computes the perimeter of a rectangle with length 1 and width w. return 2*1 + 2*w • Examples: >>> circle_area(20) 314.159 >>> rect_perim(5, 7)

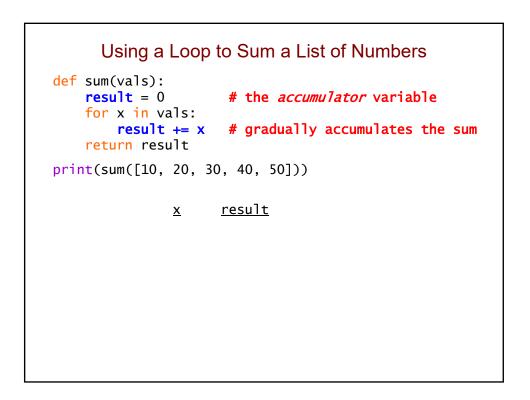


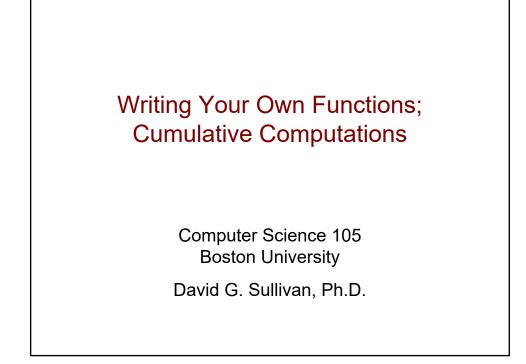


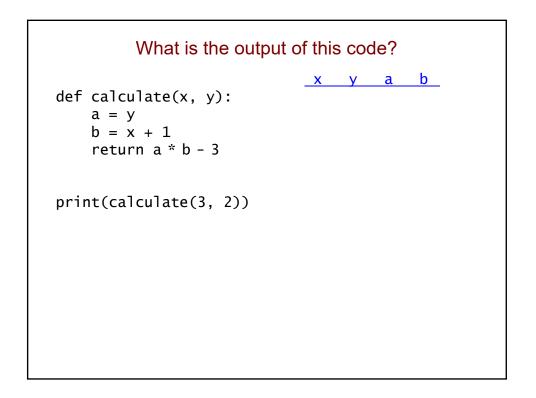


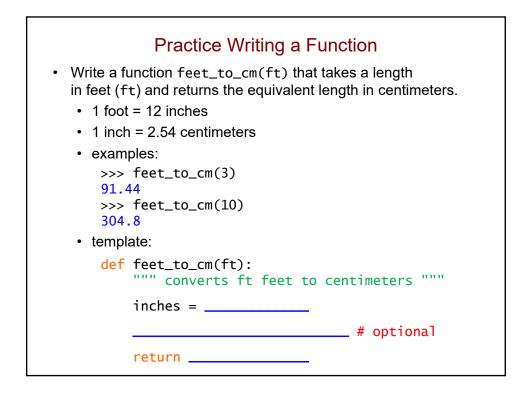














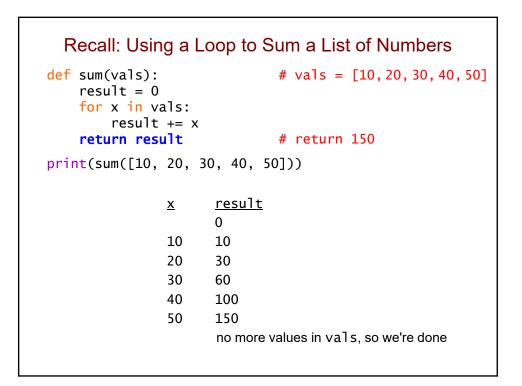
All of these work!

```
def feet_to_cm(ft):
    """ converts ft feet to centimeters """
    inches = ft * 12
    cm = inches * 2.54
    return cm

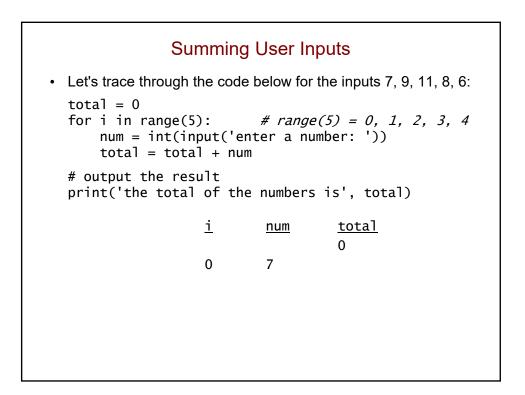
def feet_to_cm(ft):
    """ converts ft feet to centimeters """
    inches = ft * 12
    return inches * 2.54

def feet_to_cm(ft):
    """ converts ft feet to centimeters """
    return ft * 12 * 2.54
```

```
These are not the same!
def feet_to_cm(ft):
    """ converts ft feet to centimeters """
    inches = ft * 12
    cm = inches * 2.54
    return cm
def feet_to_cm(ft):
    """ converts ft feet to centimeters """
    inches = ft * 12
    cm = inches * 2.54
    print(cm)
```



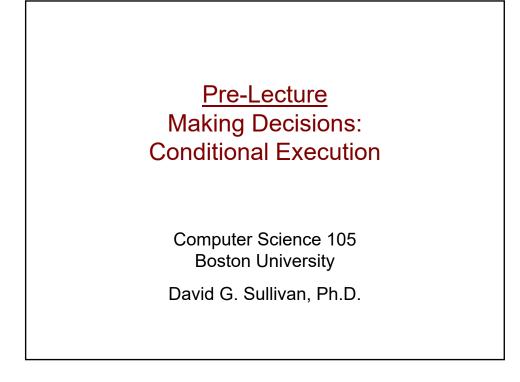
```
Cumulative Computations
def sum(vals):
                        # the accumulator variable
    result = 0
    for x in vals:
         result += x
                        # gradually accumulates the sum
    return result
print(sum([10, 20, 30, 40, 50]))
                       <u>result</u>
                <u>X</u>
                       0
                10
                       10
                20
                       30
                30
                       60
                40
                       100
                50
                       150
                       no more values in vals, so we're done
                       output: 150
```

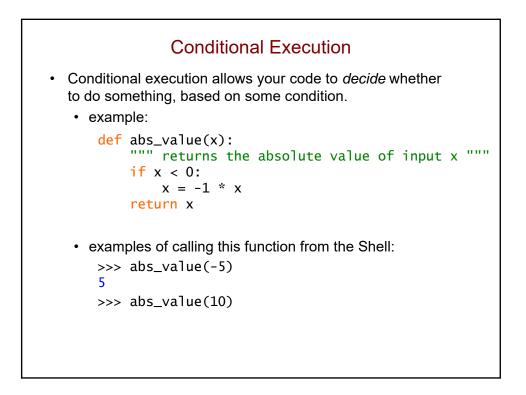


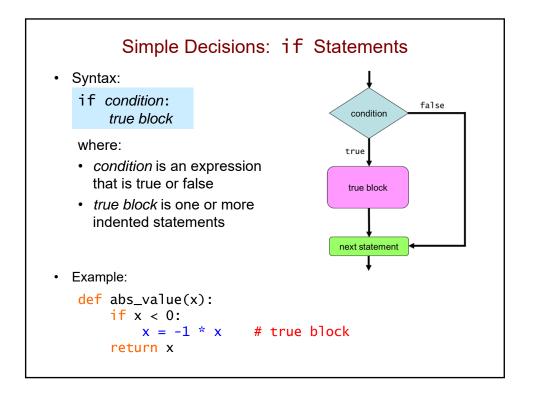
Making the Program More Flexible

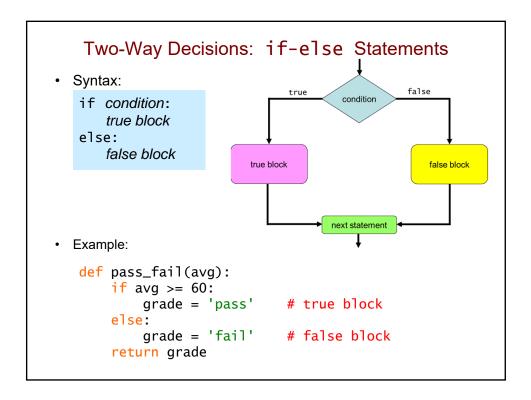
• How could we change the program to allow the user to specify the number of values to be summed?

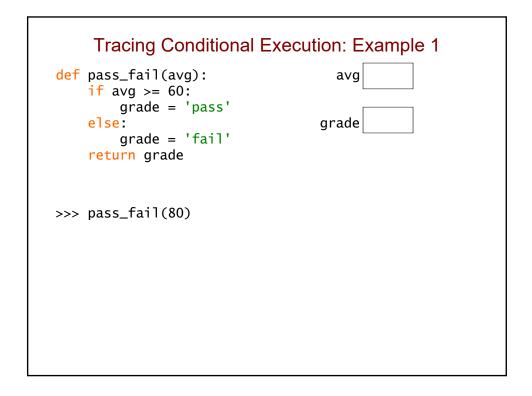
```
total = 0
for i in range(5):
    num = int(input('enter a number: '))
    total = total + num
# output the result
print('the total of the numbers is', total)
```

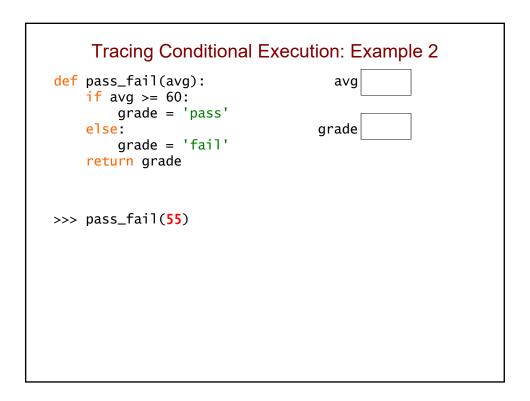




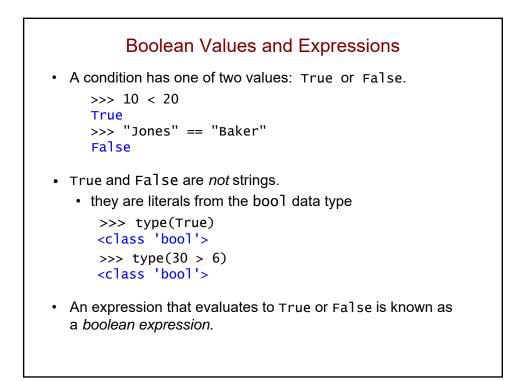


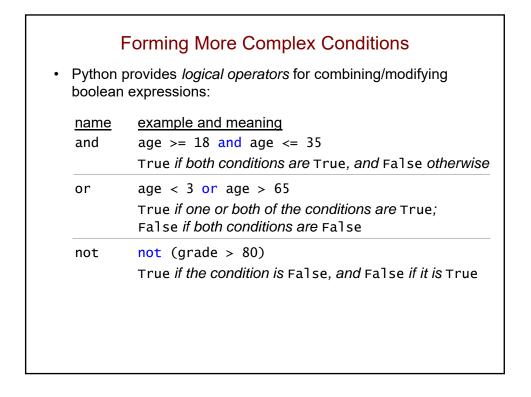


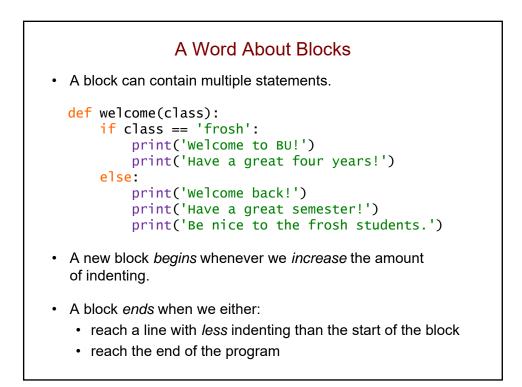




	Expressing Simple C	Conditions
• •	ovides a set of <i>relational op</i> g comparisons:	perators
<u>operator</u> <	<u>name</u> less than	<u>examples</u> val < 10 price < 10.99
>	greater than	num > 60 state > 'Ohio'
<=	less than or equal to	average <= 85.8
>=	greater than or equal to	name >= 'Jones'
== (don't confi	equal to use with =)	total == 10 letter == 'P'
!=	not equal to	age != my_age



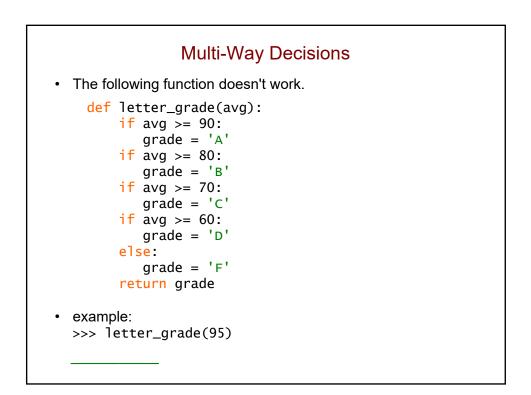


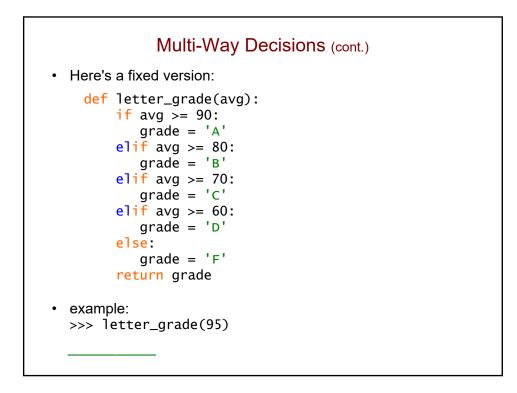


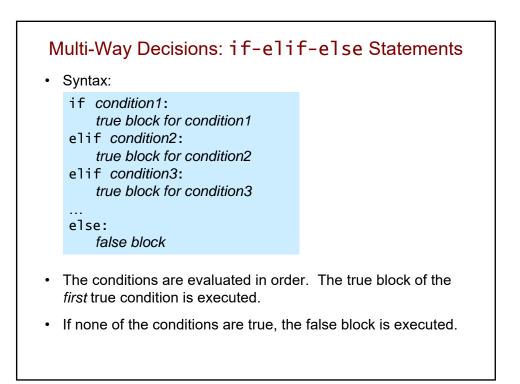
Nesting

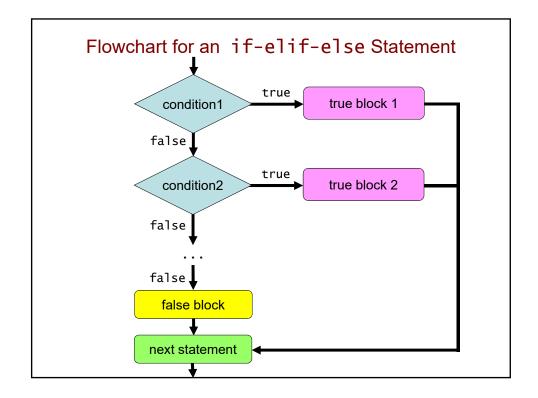
• We can "nest" one conditional statement in the true block or false block of another conditional statement.

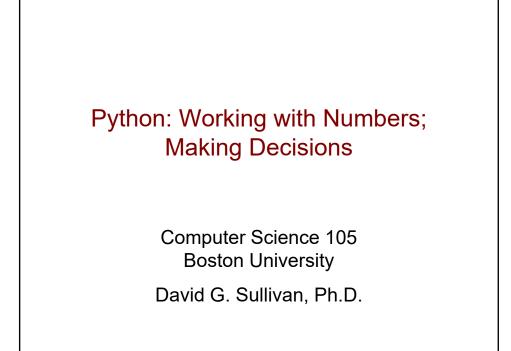
```
def welcome(class):
    if class == 'frosh':
        print('Welcome to BU!')
        print('Have a great four years!')
    else:
        print('Welcome back!')
        if class == 'senior':
            print('Have a great last year!')
        else:
            print('Have a great semester!')
        print('Be nice to the frosh students.')
```

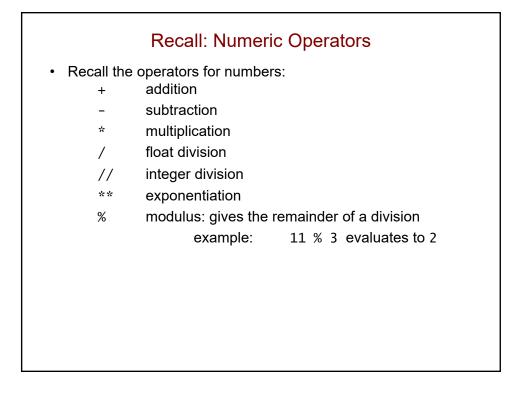


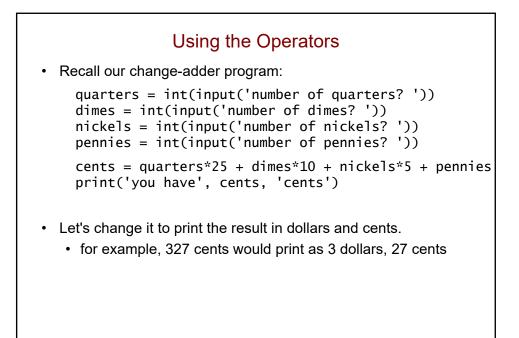




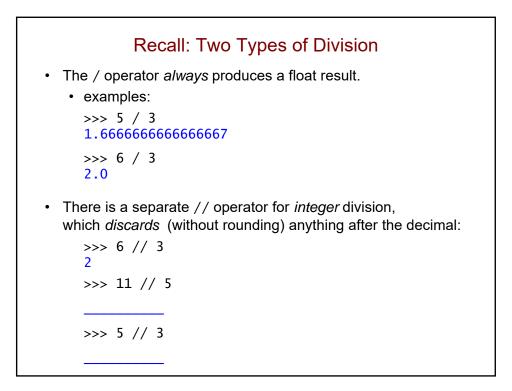


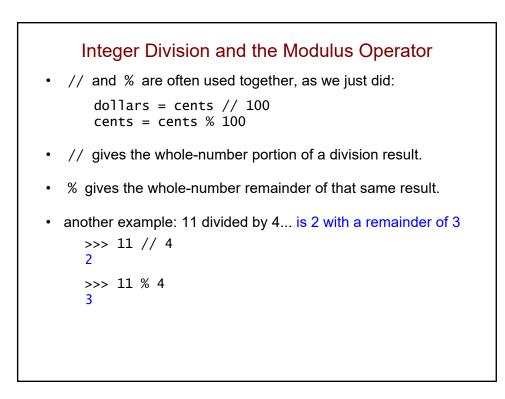






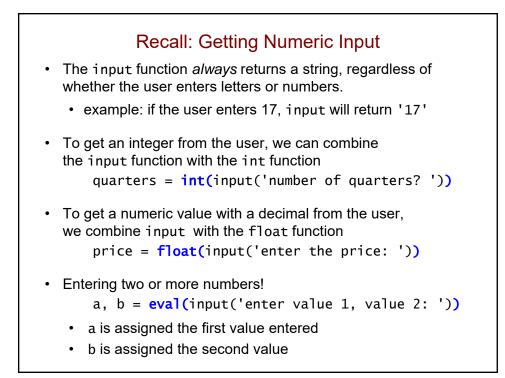
```
How Would Your Complete This Program?
quarters = int(input('number of guarters? '))
dimes = int(input('number of dimes? '))
nickels = int(input('number of nickels? '))
pennies = int(input('number of pennies? '))
cents = quarters*25 + dimes*10 + nickels*5 + pennies
dollars = _____
cents = _
print('you have', dollars, 'dollars,', cents, 'cents')
     <u>first blank</u>
                        second blank
A. cents / 100
                       cents % 100
B. cents // 100
                       cents % 100
C. cents / 100
                       cents % dollars
D. cents // 100
                       cents % dollars
```

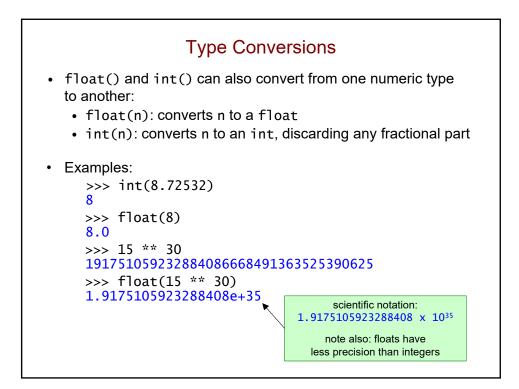


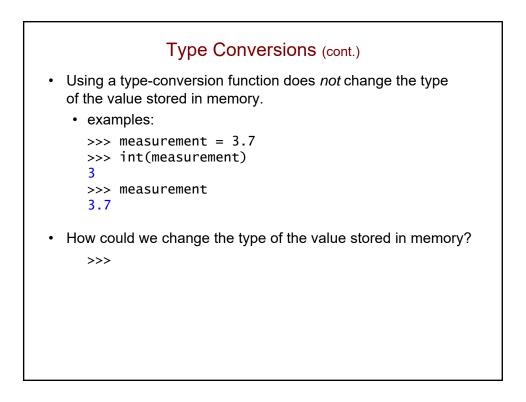


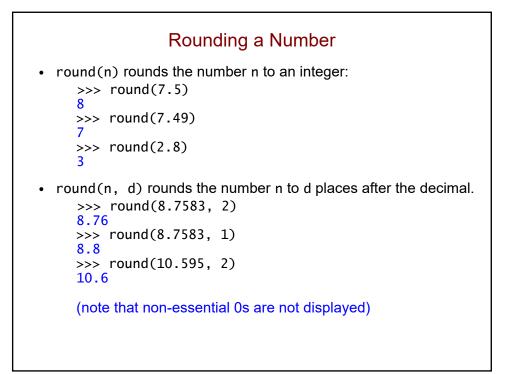
Other Uses of the Modulus Operator

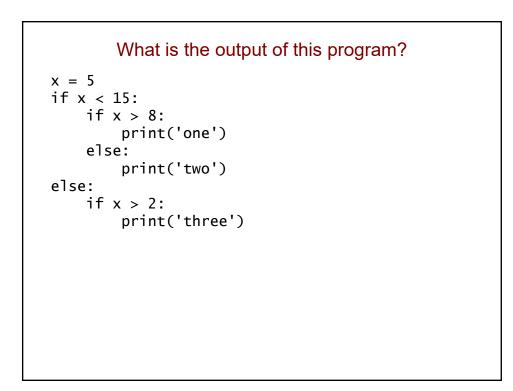
- Determining if an integer n is even or odd:
 - n % 2 == 0 if n is even
 - n % 2 == 1 if n is odd
- Determining if an integer n is a multiple of another integer m:
 - n % m == 0 if n is a multiple of m
 - n % m != 0 if n is not a multiple of m











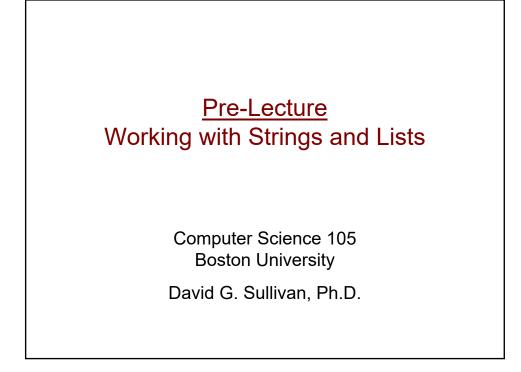
```
What does this print? (note the changes!)
x = 5
if x < 15:
    if x > 8:
        print('one')
else:
        print('two')
if x > 2:
    print('three')
```

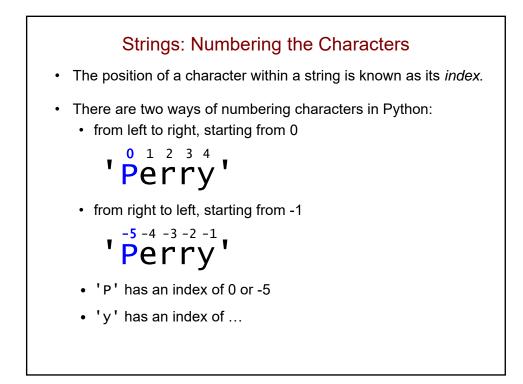
```
What does this print? (note the new changes!)
x = 5
if x < 15:
    if x > 8:
        print('one')
else:
    print('two')
if x > 2:
    print('three')
```

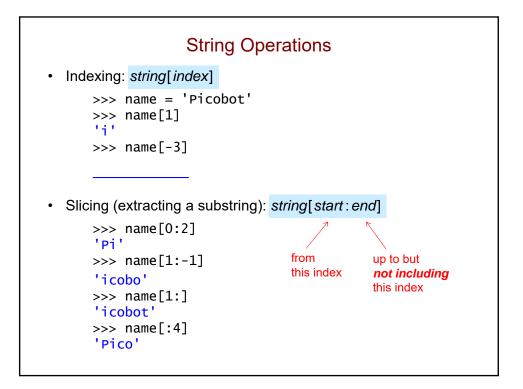
```
How many lines does this print?
```

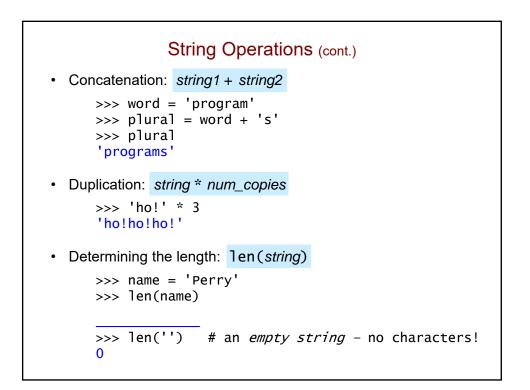
```
x = 5
if x == 8:
    print('how')
elif x > 1:
    print('now')
elif x < 20:
    print('wow')
print('cow')</pre>
```

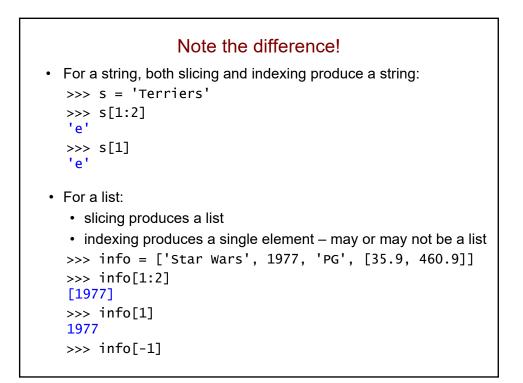
How many lines does this print?
x = 5
if x == 8:
 print('how')
if x > 1:
 print('now')
if x < 20:
 print('wow')
print('cow')</pre>

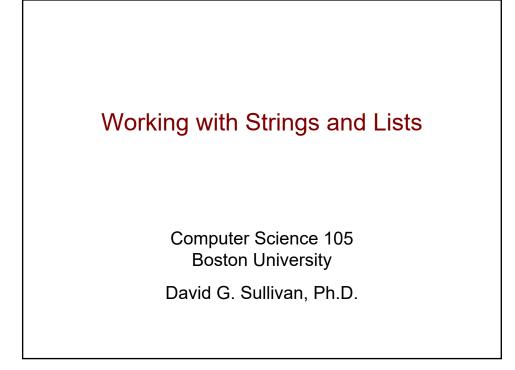


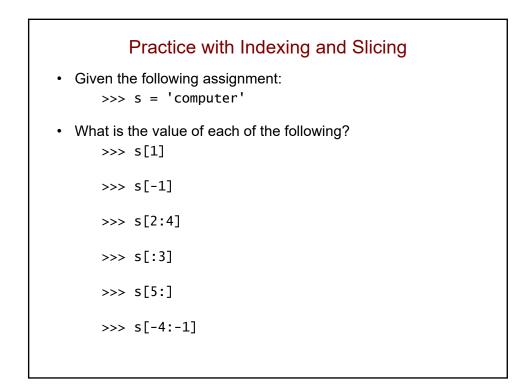




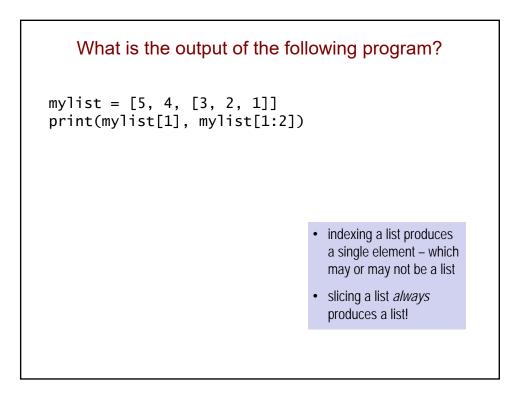






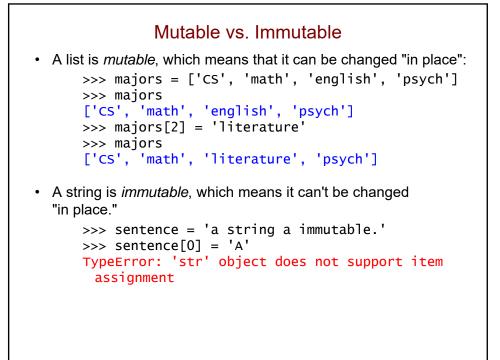


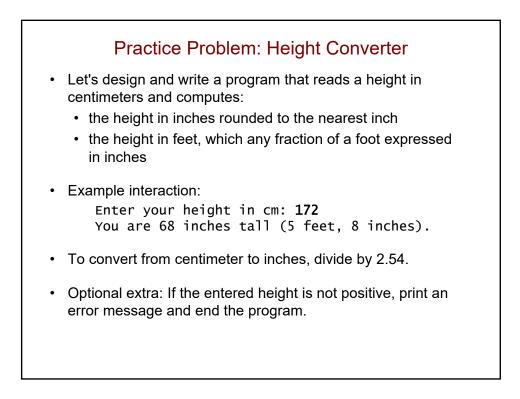
What is the value of s after the following code runs?
s = 'def'
s = ('a' * 3) + s
s = s[2:-2]



```
How could you fill in the blank
to produce [105, 111]?
intro_cs = [101, 103, 105, 108, 109, 111]
dgs_courses = ______
A. intro_cs[2:3] + intro_cs[-1:]
B. intro_cs[-4] + intro_cs[5]
C. intro_cs[-4:3] + intro_cs[5:6]
D. more than one of the above
E. none of the above
```

How could you fill in the blank to produce [105,111]?	
intro_cs = [101, 103, 105, 108, 109, 111]	
dgs_courses =	
<pre>What about this? intro_cs[-4] + intro_cs[-1:]</pre>	





One Possible Solution

```
cm = int(input('Enter your height in cm: '))
inches = cm / 2.54
inches = round(inches)
feet = inches // 12
remaining = inches % 12
print('You are', inches, 'inches tall (' + str(feet),
    'feet,', remaining, 'inches).')
```

```
A Solution That Handles Inputs Less Than 0
cm = int(input('Enter your height in cm: '))
if cm < 0:
    print('Heights must be positive')
else:
    inches = cm / 2.54
    inches = round(inches)
    feet = inches // 12
    remaining = inches % 12
    print('You are', inches, 'inches tall (' + str(feet),
        'feet,', remaining, 'inches).')</pre>
```

```
Extra Practice: Fill in the blank to
make the code print compute!

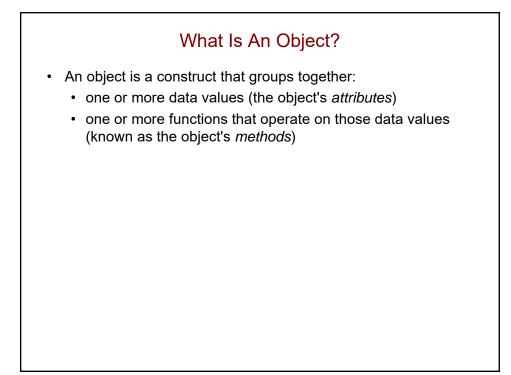
subject = 'computer science!'
verb = _____
print(verb)

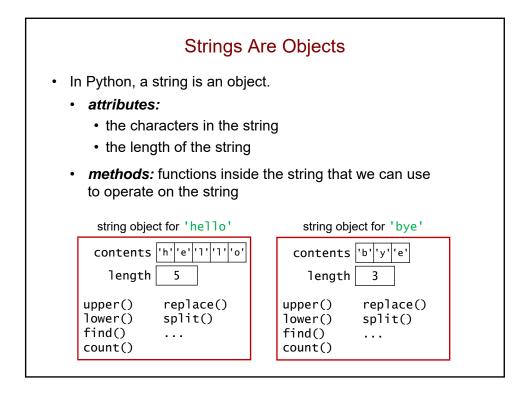
A. subject[:7] + subject[-1]
B. subject[:7] + subject[-1]
C. subject[:8] + subject[-1]
D. subject[:8] + subject[:-1]
E. none of these
```

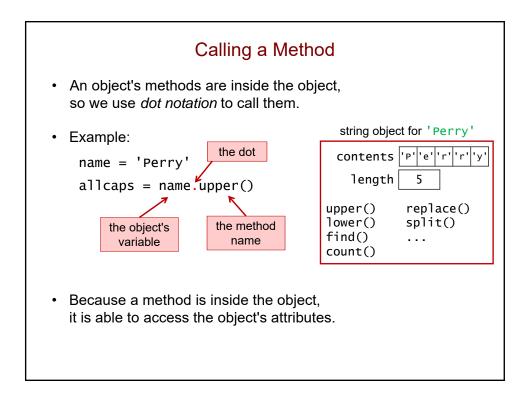
<u>Pre-Lecture</u> Using Objects; Splitting and Joining Strings

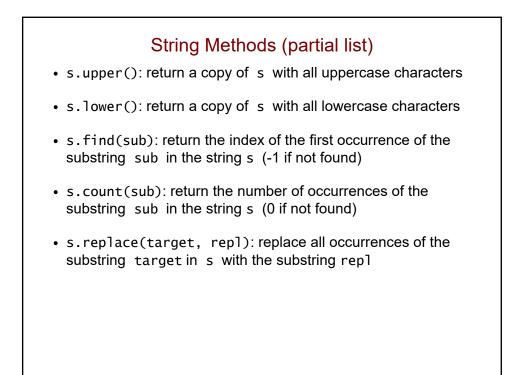
> Computer Science 105 Boston University

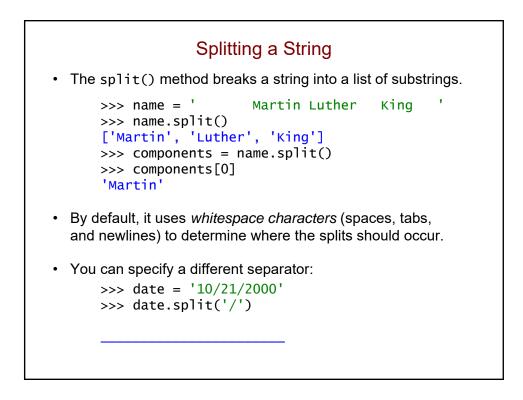
David G. Sullivan, Ph.D.

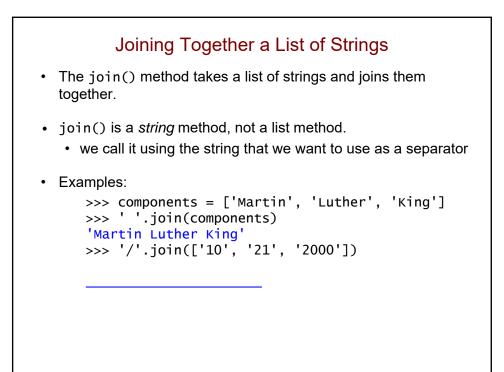


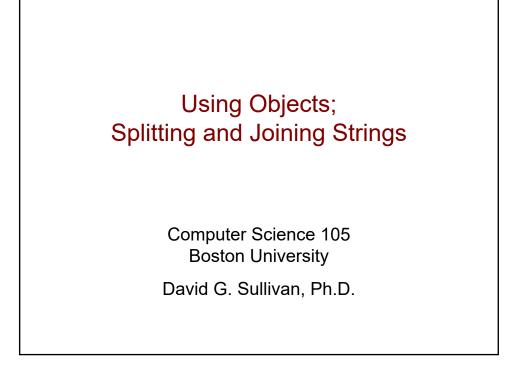


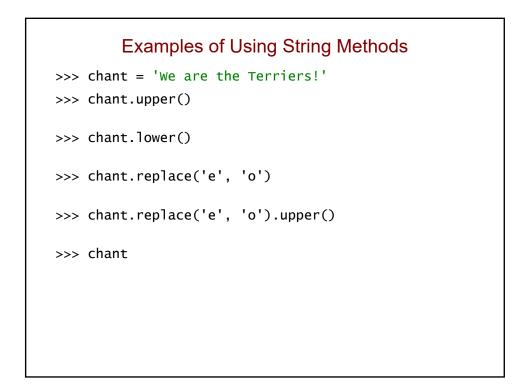








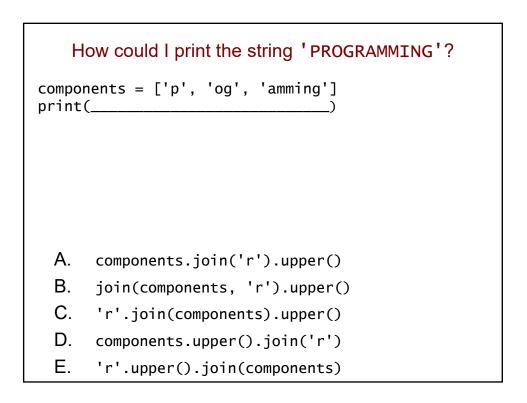




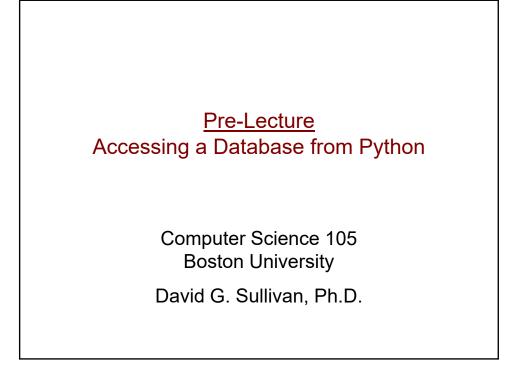
What is the output of this program?

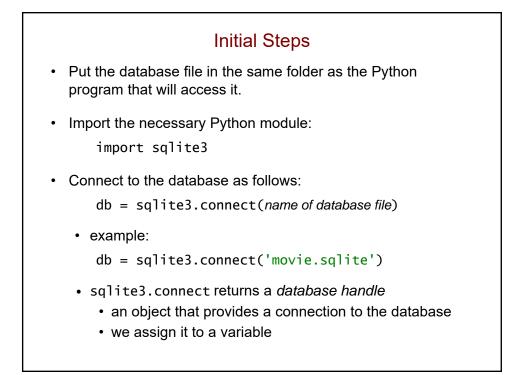
```
s = 'Programming'
s = s.lower()
s.upper()
print(s.split('r'))
```

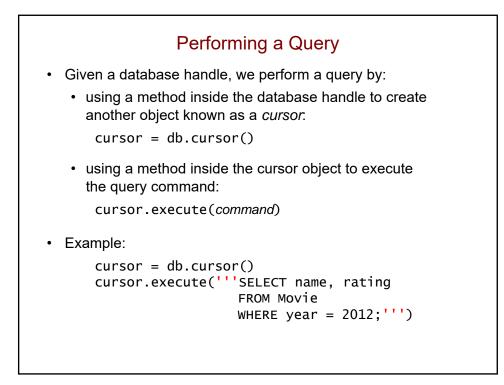
```
A. ['P', 'og', 'amming']
B. ['p', 'og', 'amming']
C. ['P', 'OG', 'AMMING']
D. ['PR', 'OGR', 'AMMING']
E. ['pr', 'ogr', 'amming']
```

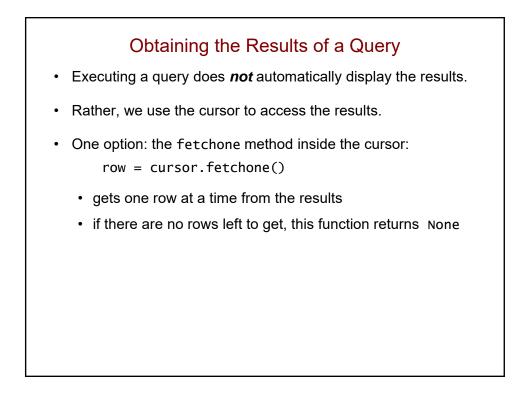


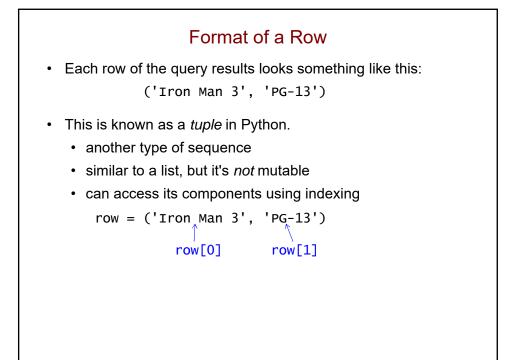
Practice: Analyzing a Name . Write a program that analyzes a person's name. . Here's a sample run of the program: Inter your full name: George Alexander Louis Wales Your name has 2 characters (including spaces). Your name has 4 components. first name: George last name: Wales other names: Alexander Louis Enter a letter: r That letter occurs 2 times in your name. the first occurrence is at position 3 in the name.

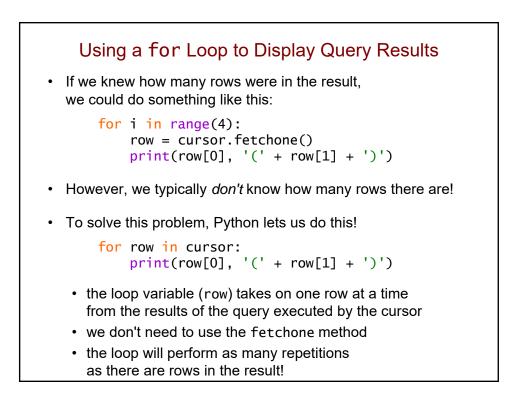


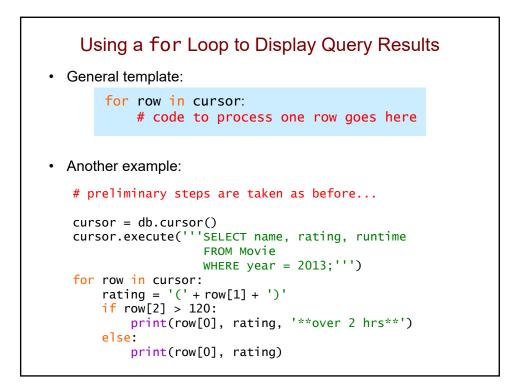


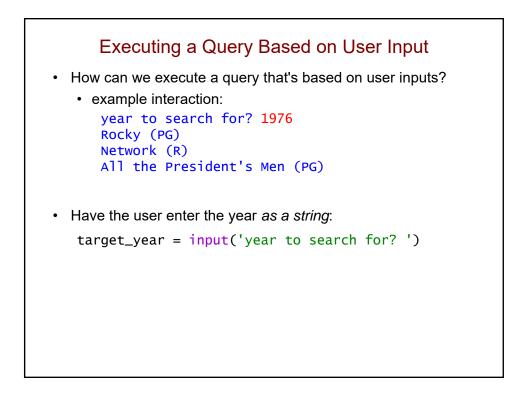


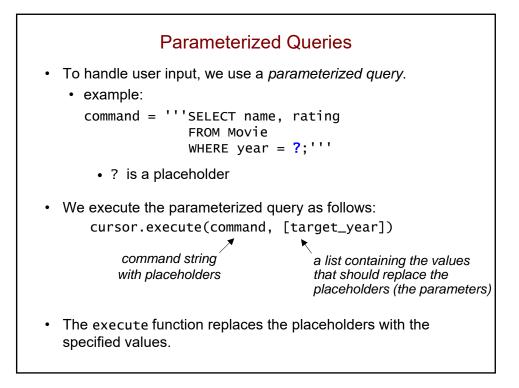


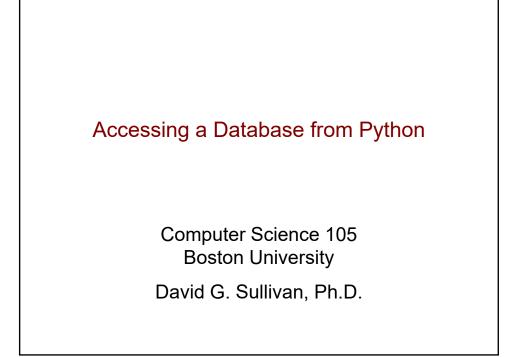


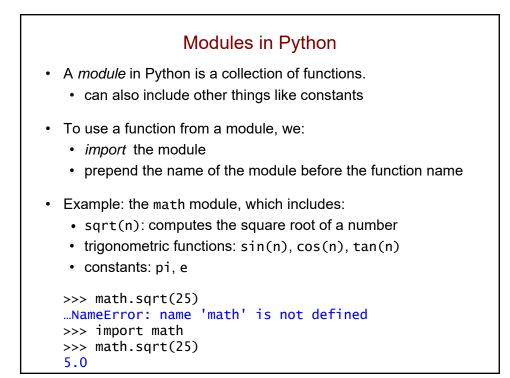


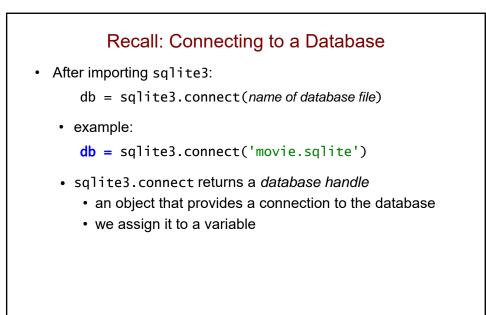


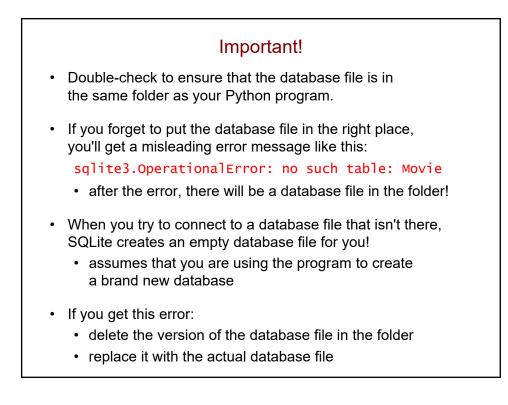


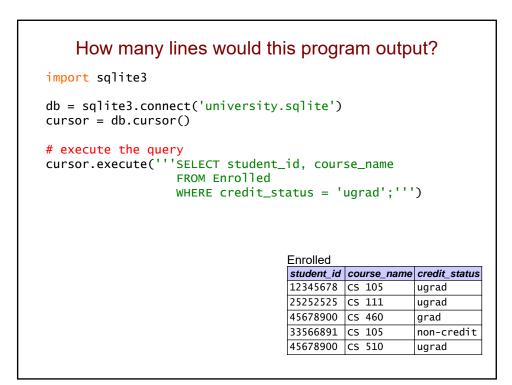




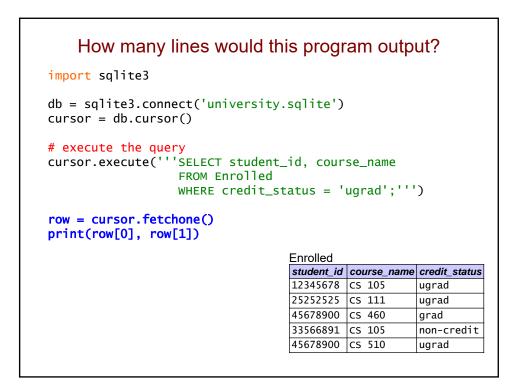


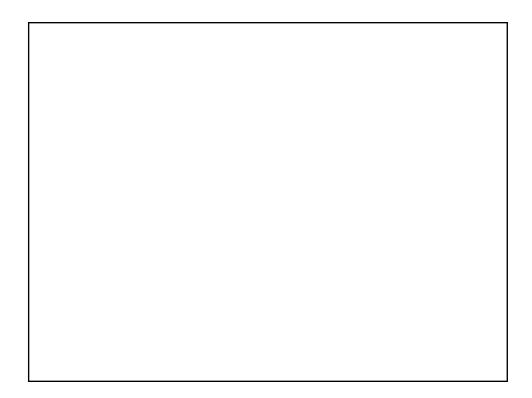


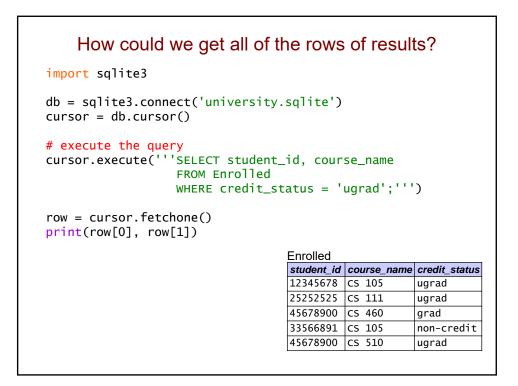




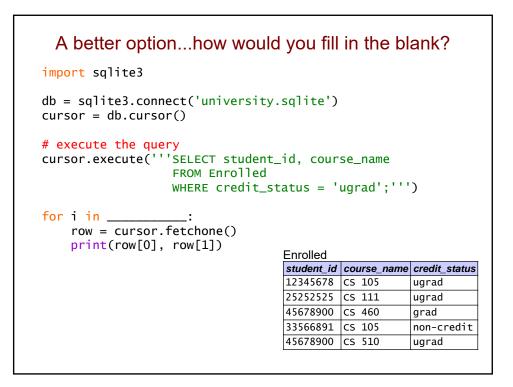


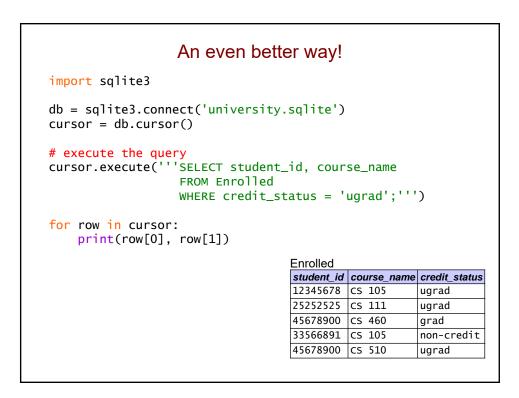


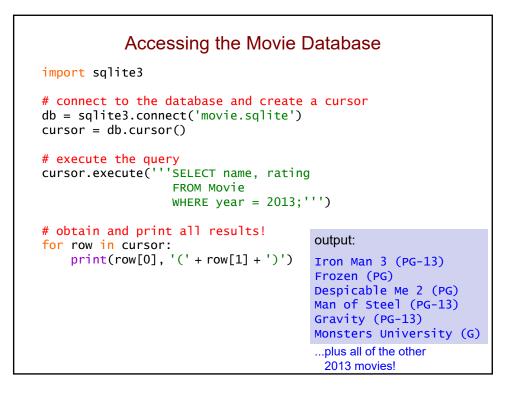


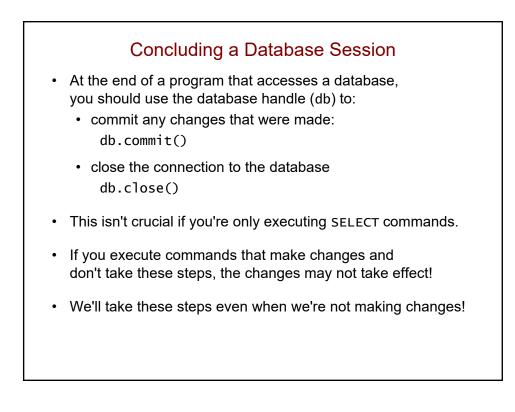


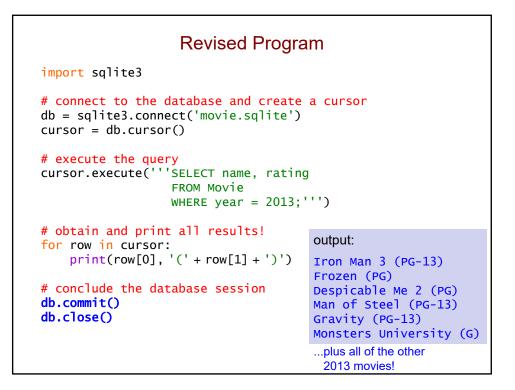
One option				
<pre>import sqlite3</pre>				
<pre>db = sqlite3.connect('university.sqlite') cursor = db.cursor() # execute the query cursor.execute('''SELECT student_id, course_name</pre>				
row = cursor.fetchone()	student_id	course_name	credit_status	
<pre>print(row[0], row[1])</pre>	12345678	CS 105	ugrad	
princ(iow[0], iow[1])	25252525	CS 111	ugrad	
	45678900	CS 460	grad	
	33566891	CS 105	non-credit	
	45678900	CS 510	ugrad	

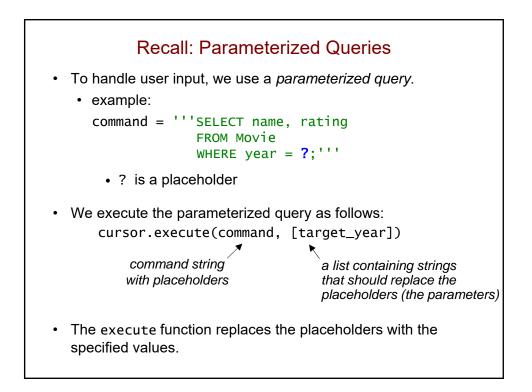




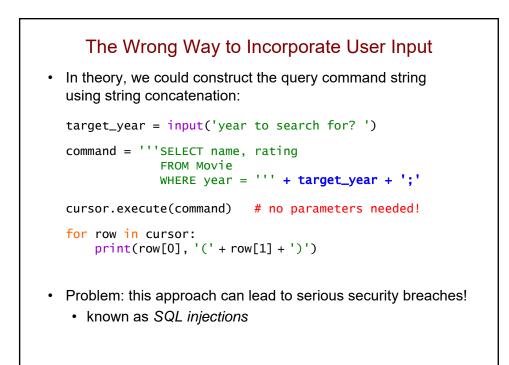


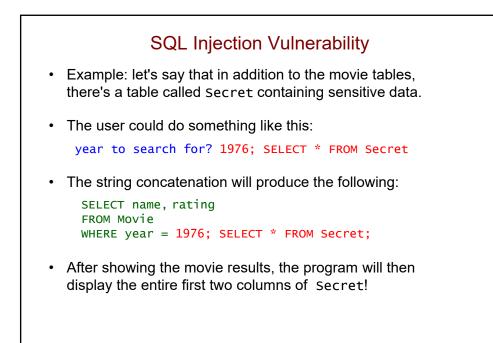


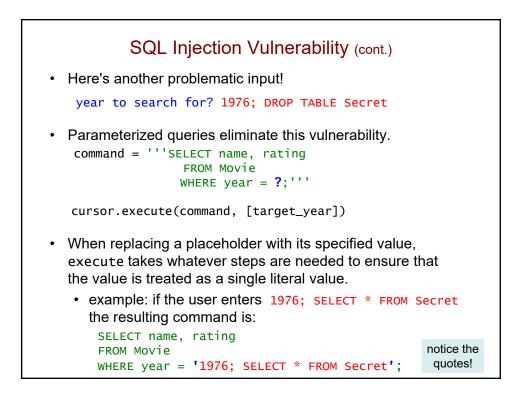


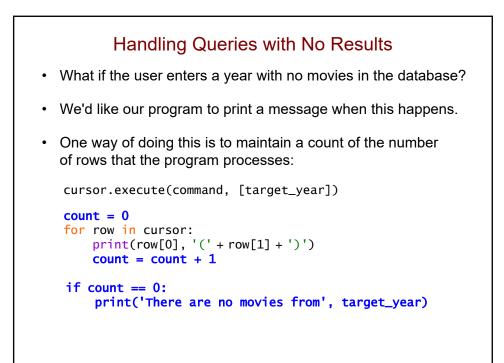


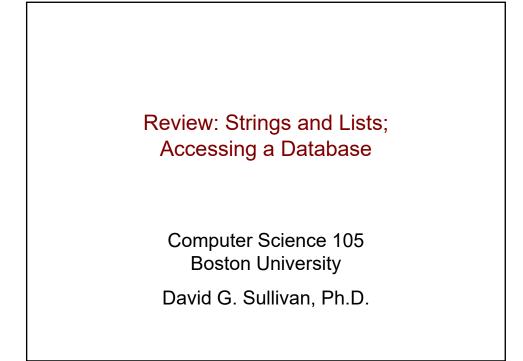
Recall: Example Program: Final Version import sqlite3 # connect to the database and create a cursor db = sqlite3.connect('movie.sqlite') cursor = db.cursor() # get the year from the user as a string target_year = input('year to search for? ') # execute the parameterized query command = '''SELECT name, rating FROM Movie WHERE year = **?**;''' cursor.execute(command, [target_year]) # obtain and print all results! for row in cursor: print(row[0], '(' + row[1] + ')') # conclude the database session db.commit() db.close()

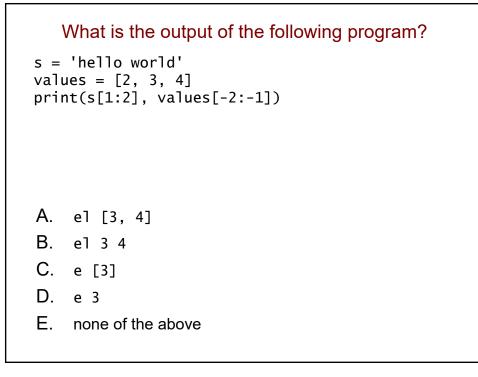


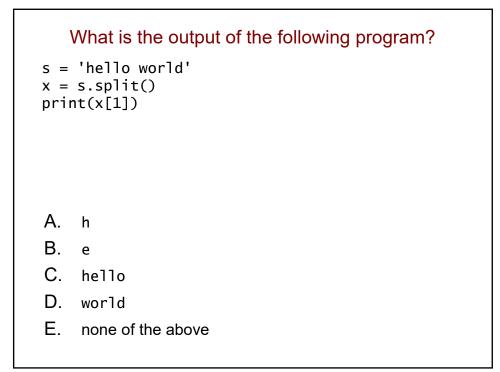


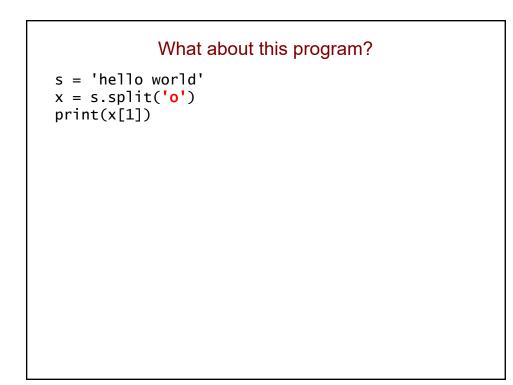


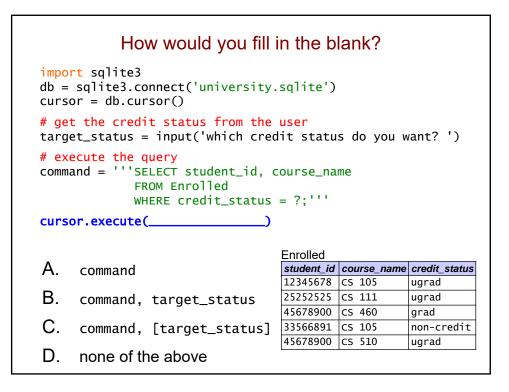




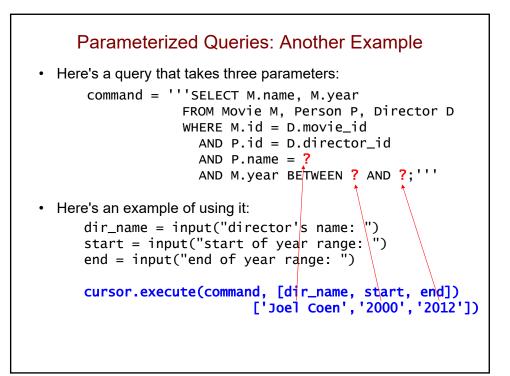




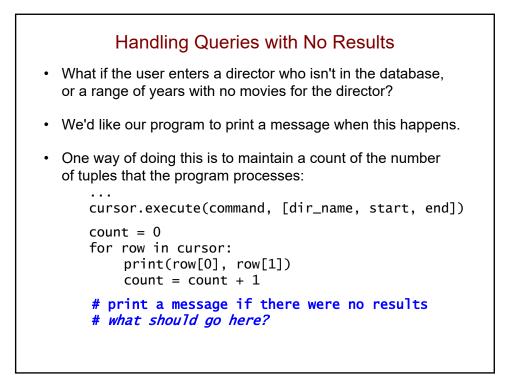


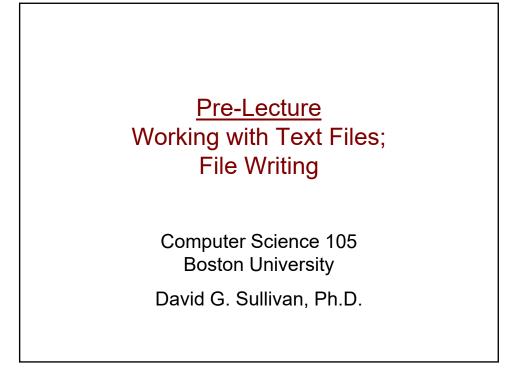


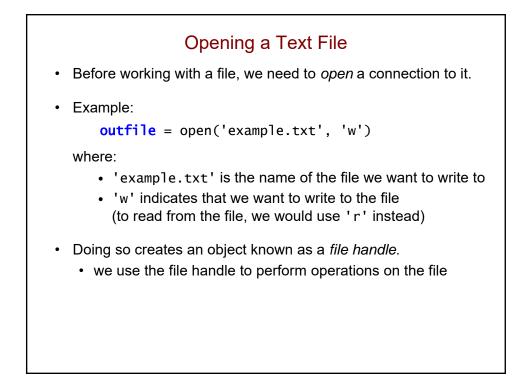
What about this blank?								
<pre>import sqlite3 db = sqlite3.connect('university.sqlite') cursor = db.cursor()</pre>								
<pre># get the credit status from the target_status = input('which cred</pre>		s do you w	ant? ')					
FROM Enrolled	<pre>command = '''SELECT student_id, course_name</pre>							
<pre>cursor.execute() # filled in from previous slide for row in cursor: print(</pre>								
A. row[0], row[1]	12345678		ugrad					
B. row.fetchone()								
	43678900 CS 460 grad							
C. cursor[0], cursor[1]	33566891		non-credit					
_	45678900	CS 510	ugrad					
D. none of the above								



The Full Program (getMoviesByDirector.py)
import sqlite3
filename = input("name of database file: ") db = sqlite3.connect(filename) cursor = db.cursor()
<pre>dir_name = input("director's name: ") start = input("start of year range: ") end = input("end of year range: ")</pre>
<pre>command = '''SELECT M.name, M.year FROM Movie M, Person P, Director D WHERE M.id = D.movie_id AND P.id = D.director_id AND P.name = ? AND M.year BETWEEN ? AND ?;''' cursor.execute(command, [dir_name, start, end])</pre>
<pre>for row in cursor: print(row[0], row[1])</pre>
db.commit() db.close()

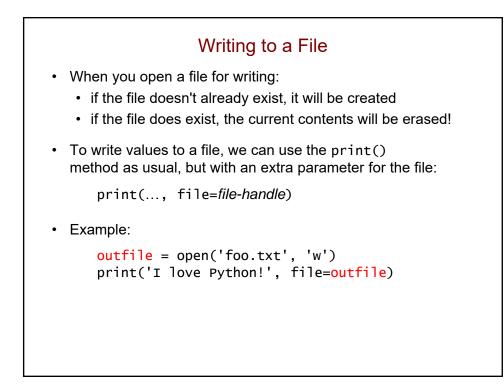


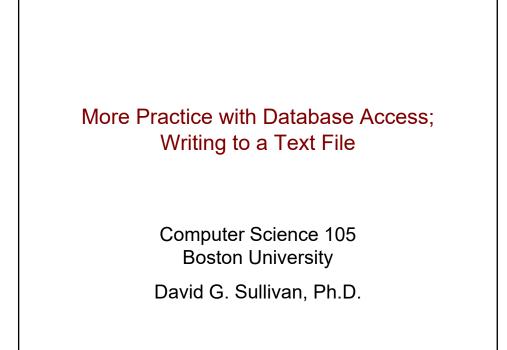


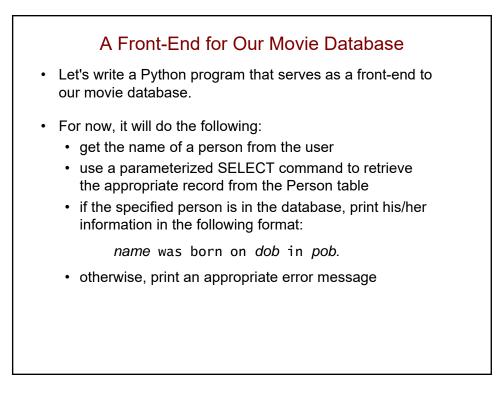


Closing a File

- Here's our previous example: outfile = open('example.txt', 'w')
- When we're done working with the file, we close its handle: outfile.close()
- **Important:** Text that you write to a file may not make it to disk until you close the file handle!





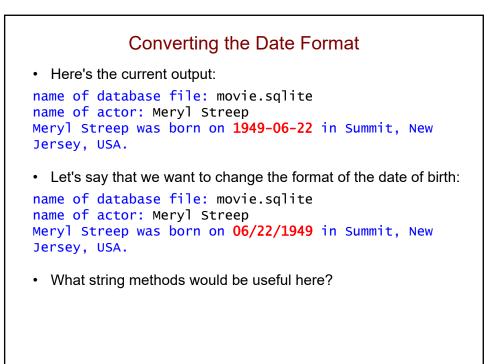


Sample Runs of the Program

```
name of database file: movie.sqlite
name of actor: Dave Sullivan
Dave Sullivan is not in the database.
```

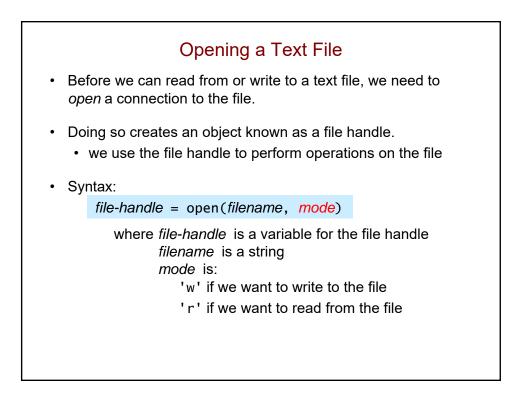
name of database file: movie.sqlite name of actor: Meryl Streep Meryl Streep was born on 1949-06-22 in Summit, New Jersey, USA.

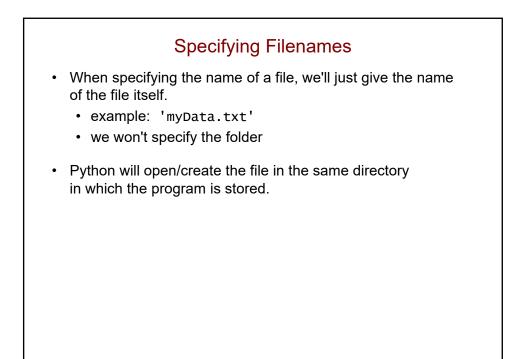
file db =	ort sqlite3 name = input('name of database file: ') sqlite3.connect(filename) or = db.cursor()
# Ge	t the actor's name.
	ecute the command, get the result, and print it. and = '''
curs	or.execute()
for	: print()

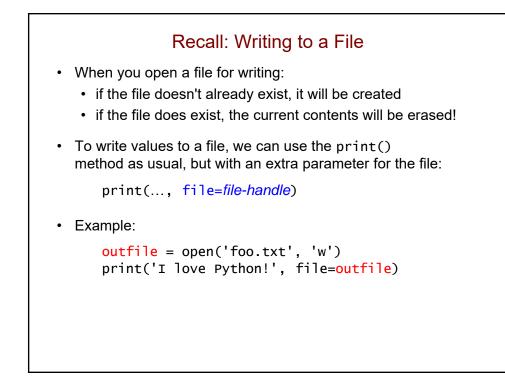


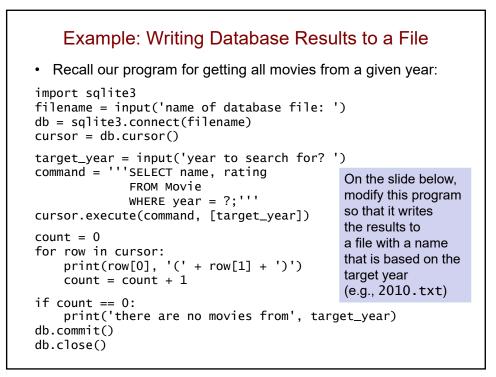
How would you fill in the blanks?							
cı	ommand = '''SELECT do FROM Pers WHERE nar WSor.execute(command ount = 0	<i>Example:</i> We want to go from '1949-06-22' to '06/22/1949'					
fo	<pre>count = 0 for row in cursor: comps = dob = print(name, 'was born on', dob, 'in', row[1] + '.') count = count + 1</pre>						
	first blank second blank						
A.	row[0].split()	'/'.join(comp	s)				
В.	row[0].split('-')	'/'.join(comps)					
C.	row[0].split()	'/'.join([comps[1], comps[2], comps[0]])					
D.	<pre>row[0].split('-')</pre>	<pre>'/'.join([comps[1], comps[2], comps[0]])</pre>					

Revised Front-End for Our Movie Database import sqlite3 filename = input('name of database file: ') db = sqlite3.connect(filename) cursor = db.cursor() # Get the actor's name. name = input('name of actor: ') # Execute the command, get the result, and print it. command = 'SELECT dob, pob FROM Person WHERE name = ?;' cursor.execute(command, [name]) count = 0for row in cursor: comps = row[0].split('-') dob = '/'.join([comps[1], comps[2], comps[0]]) print(name, 'was born on', dob, 'in', row[1] + '.') count = count + 1if count == 0: print(name, 'is not in the database.') db.commit() db.close()

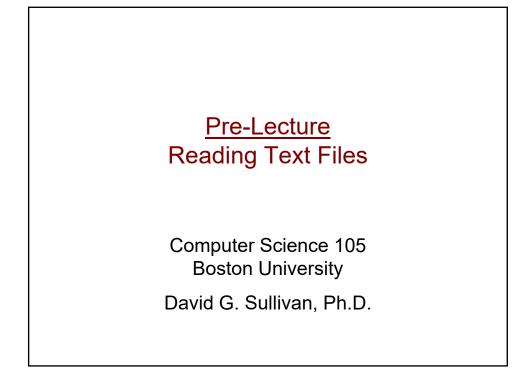


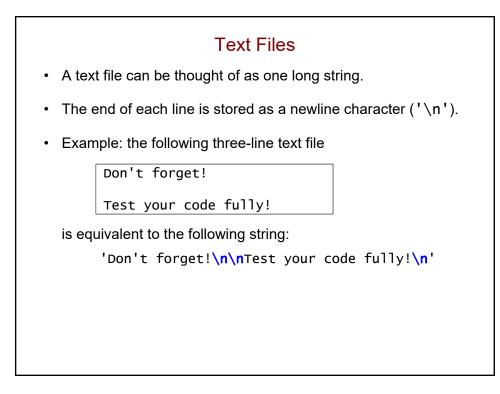


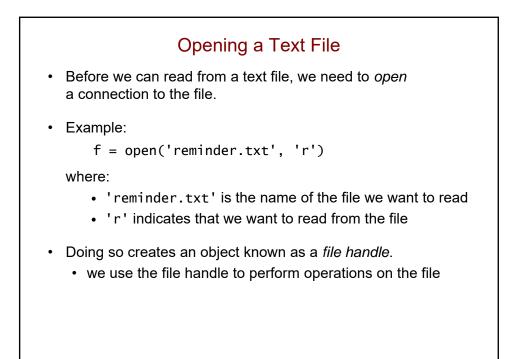


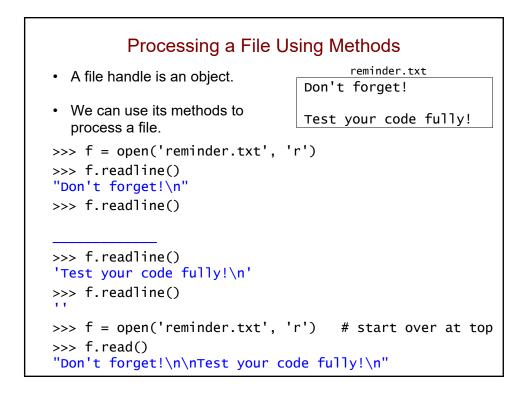


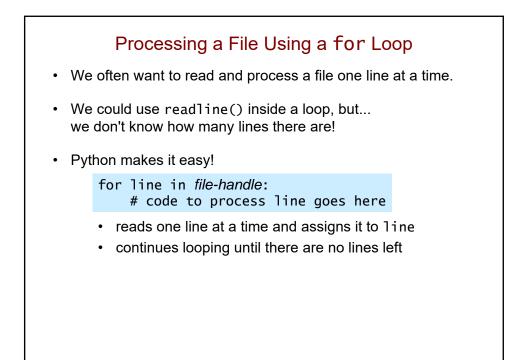
```
import sqlite3
filename = input('name of database file: ')
db = sqlite3.connect(filename)
cursor = db.cursor()
target_year = input('year to search for? ')
command = '''SELECT name, rating
             FROM Movie
             WHERE year = ?;'''
cursor.execute(command, [target_year])
count = 0
for row in cursor:
   print(row[0], '(' + row[1] + ')')
   count = count + 1
if count == 0:
    print('there are no movies from', target_year)
db.commit()
db.close()
```

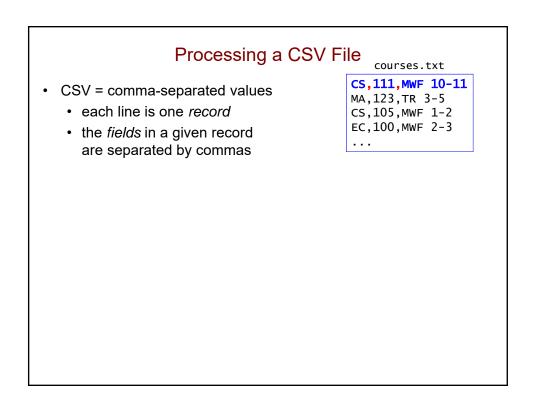


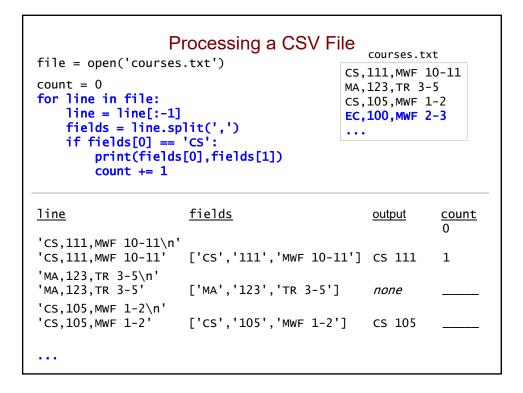


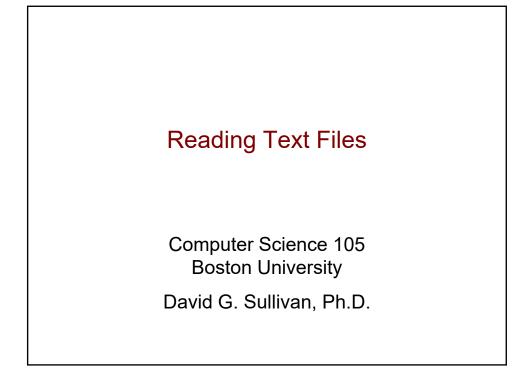


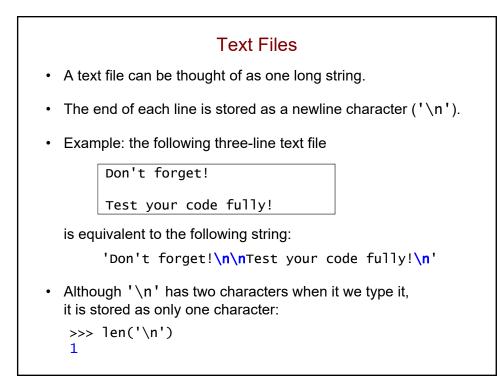


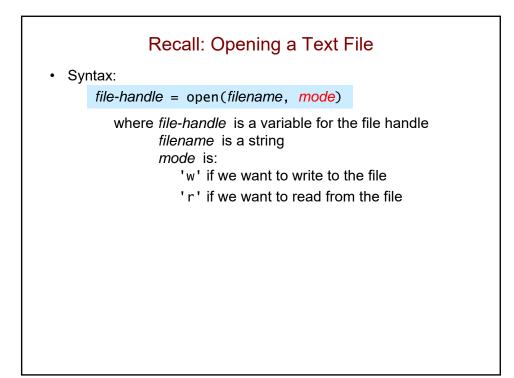


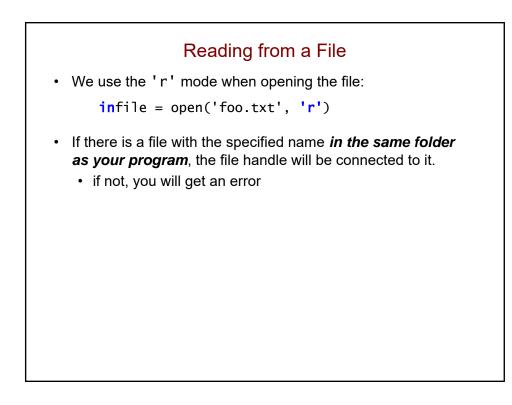


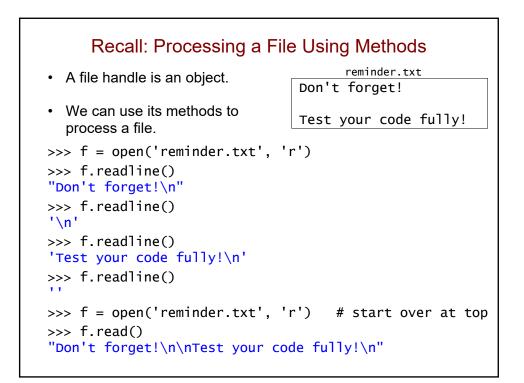


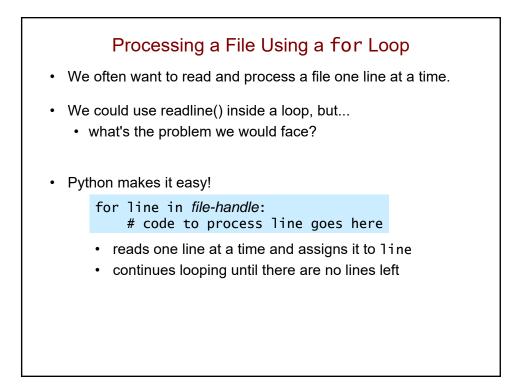


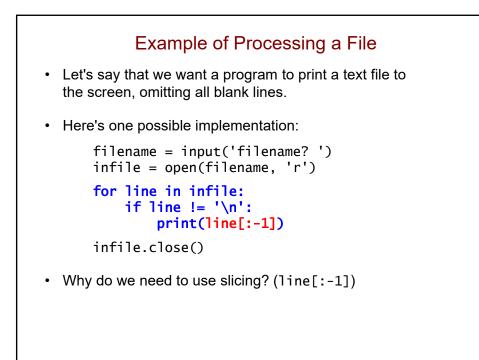




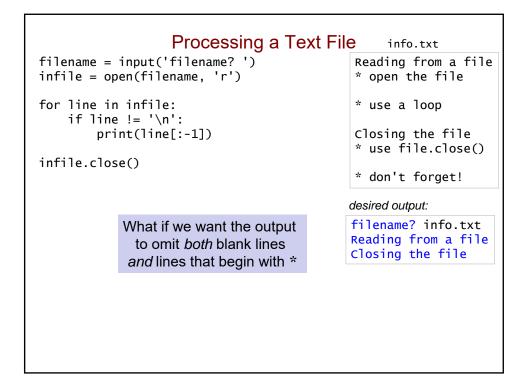


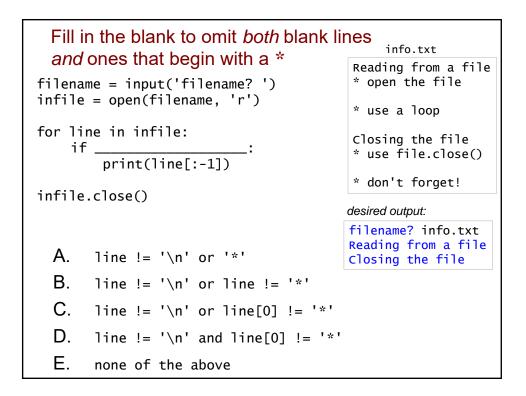


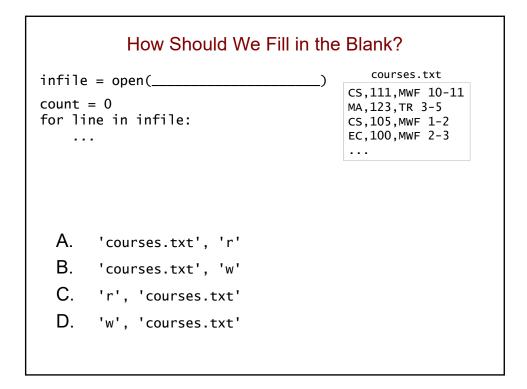


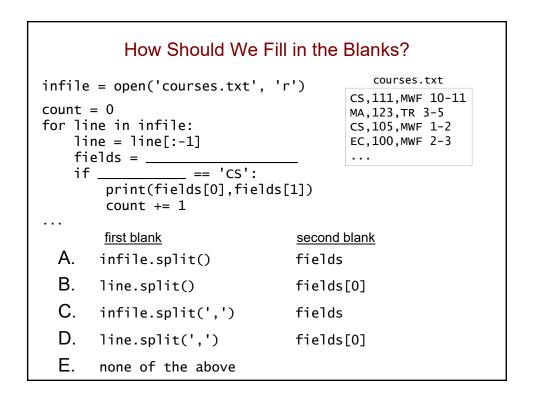


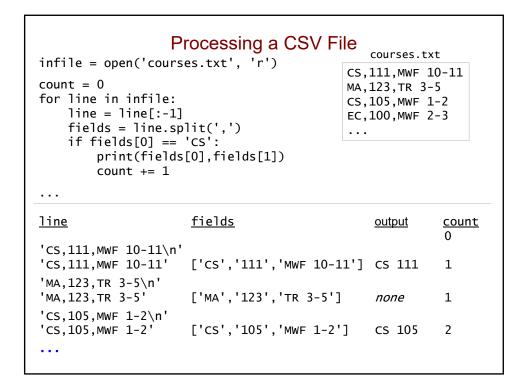
Proce filename = input('filenam infile = open(filename, '		ile info.txt Reading from a file * open the file
<pre>for line in infile: if line != '\n': print(line[:-1]) infile.close()</pre>		<pre>* use a loop Closing the file * use file.close() * don't forget!</pre>
<pre>line 'Reading from a file\n' '* open the file\n' '\n' '* use a loop\n' '\n' 'Closing the file\n' '* use file.close()\n' '\n' '\n' '* don't forget\n'</pre>	is the line printed? yes yes no yes no yes yes	output: filename? info.txt Reading from a file * open the file * use a loop Closing the file * use file.close()

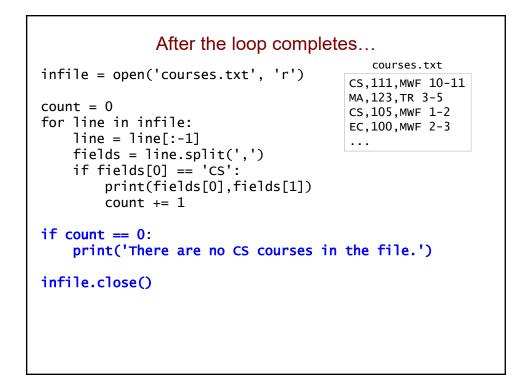


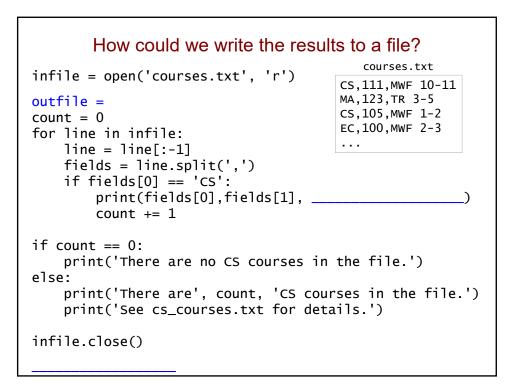


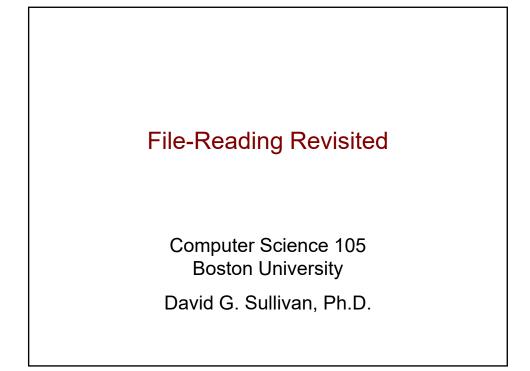


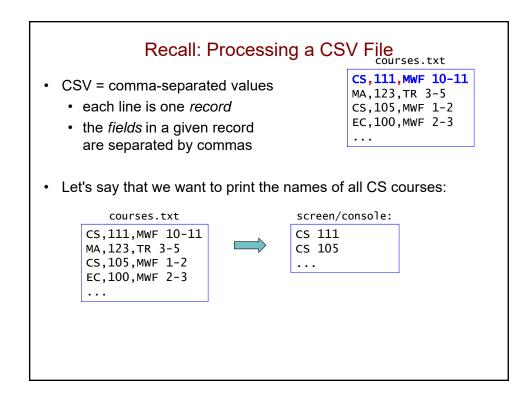


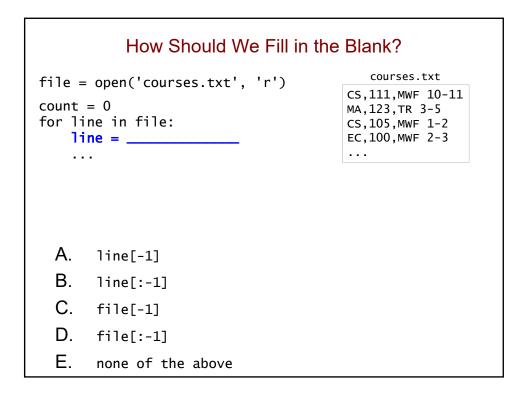


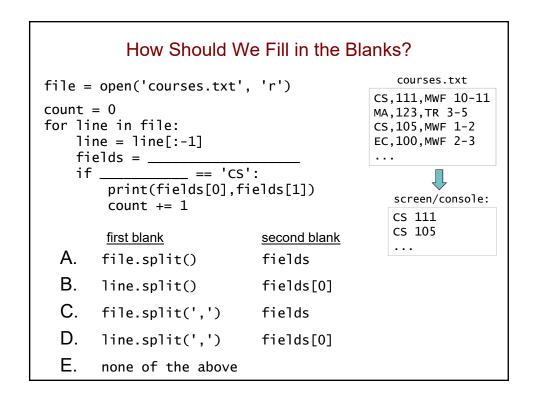


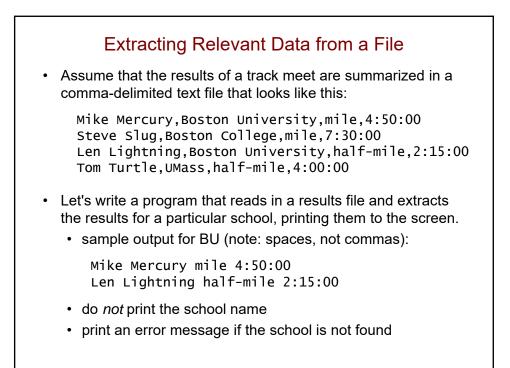




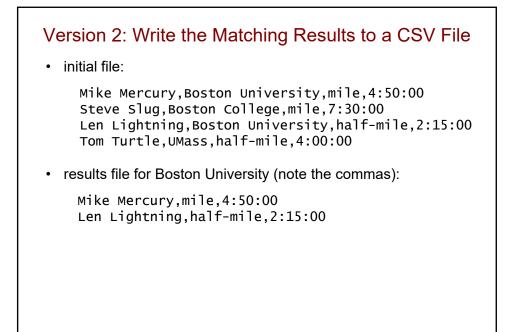




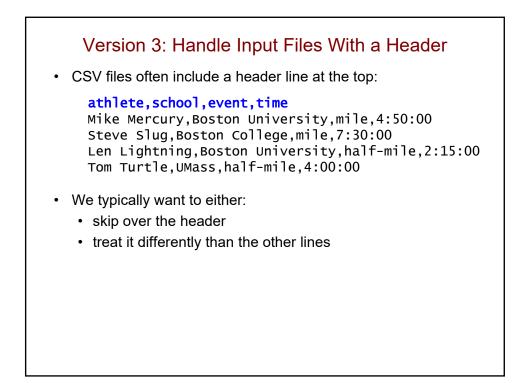


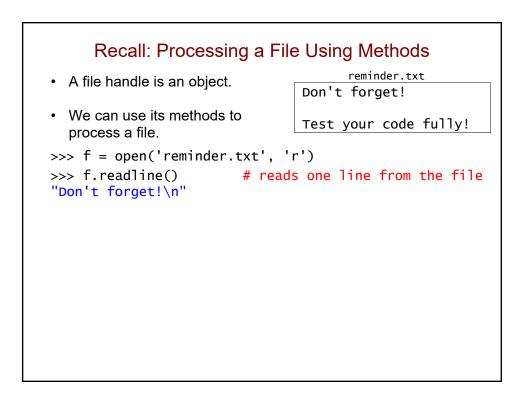


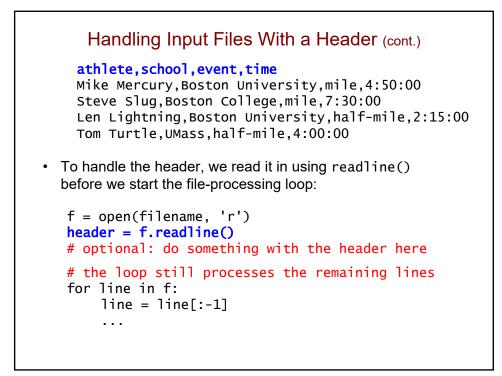
```
Extracting Relevant Data from a File
infilename = input('name of input file: ')
school = input('extract records of which school? ')
infile = open(infilename, 'r')
count = 0
for line in infile:
    line = line[:-1]
    fields = line.split(',')
    if fields[1] == school:
        count = count + 1
        print(fields[0], fields[2], fields[3])
if count == 0:
    print('No results for', school, 'in', infilename)
infile.close()
```

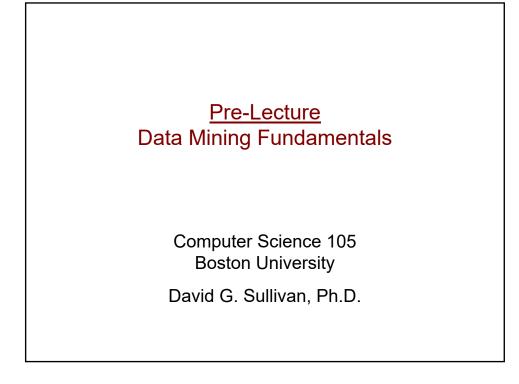


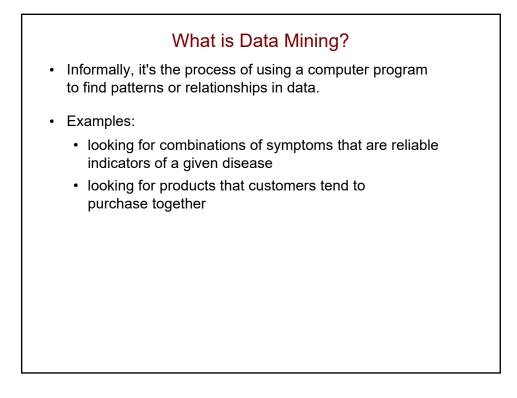
```
What changes do we need to make?
infilename = input('name of input file: ')
outfilename = input('name of output file: ')
school = input('extract records of which school? ')
infile = open(infilename, 'r')
outfile = _____
count = 0
for line in infile:
    line = line[:-1]
    fields = line.split(',')
    if fields[1] == school:
        count = count + 1
        print(fields[0] + ', ' + fields[2] + ', ' + fields[3],
                                    _)
if count == 0:
    print('No results for', school, 'in', infilename)
else:
    print(count, 'results written to', outfilename)
infile.close()
<u>outfile.close()</u>
```







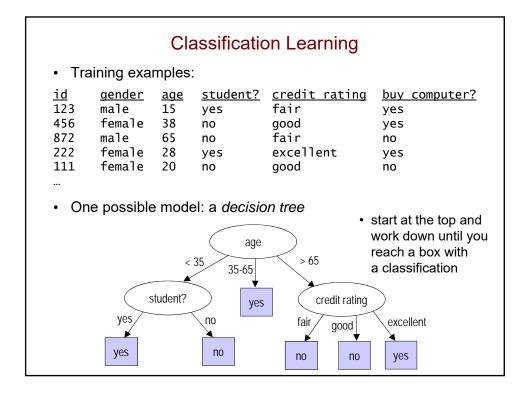


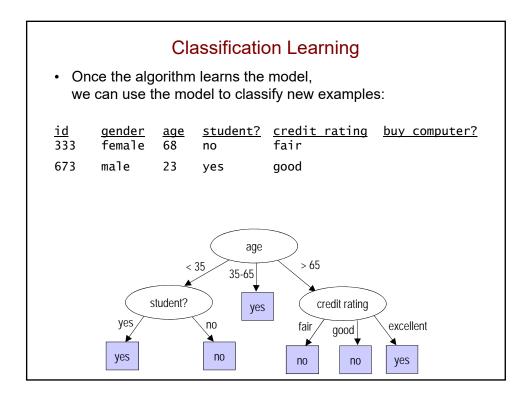


Machine Learning

- In data mining, we apply an algorithm that "learns" something about the data.
- These algorithms are machine-learning algorithms.
- We're ultimately going to consider three different types of machine learning:
 - · classification learning
 - association learning
 - numeric estimation

	Classification Learning						
	 Classification learning involves learning how to classify objects/entities on the basis of their characteristics. 						
•	 example: learning to determine whether a customer is likely to buy a computer in the next year (Yes/No). 						
	 We give the algorithm a set of <i>training examples</i> that have already been classified. 						
<u>id</u> 123	<u>gender</u> male	<u>age</u> 15	student?	<u>credit rating</u> fair	buy computer?		
456		-	yes no	good	yes ves		
872		65	no	fair	no		
	222 female 28 yes excellent yes						
	111 female 20 no good no						
	;						
	 The algorithm produces a <i>model</i> that can be used to classify other examples. 						





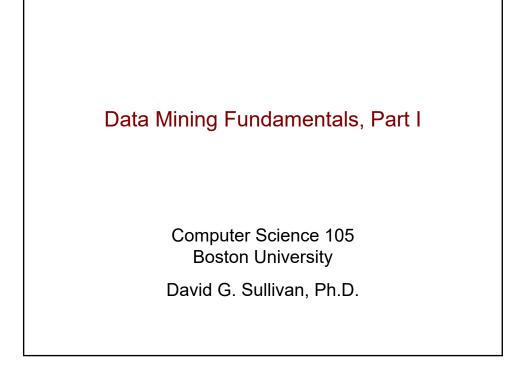
Some Terminology						
• Ea • Th	female female ch row in ch columr e attribute the <i>outpu</i>	65 28 20 the tra n is ref s can <i>t</i> attrib	yes no yes no ining data i erred to as be divided	an <i>attribute</i> . into two types: ne we want to de	buy computer? yes yes no yes no xample or instance.	
	input attril	butes 1	⇒ m	odel 🔿 of	utput attribute	

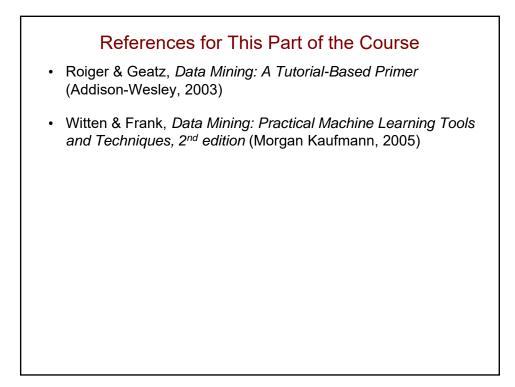
Some Terminology						
• Ea • Th	female ch row in t ch column e attribute	is ref s can	erred to as be divided	fair good fair excellent good	buy computer? yes yes no yes no xample or instance.	
 the <i>input</i> attributes – everything else 						
	stude credit ra		yes no	se est ang even the b	uy computer?	

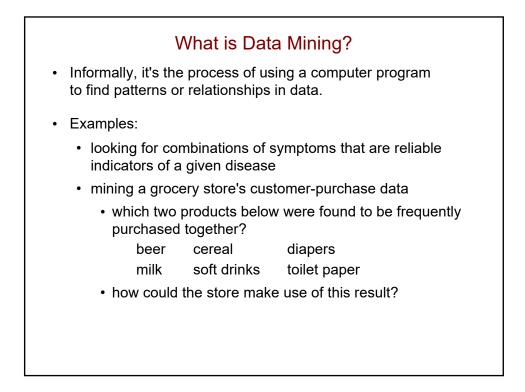
Nominal vs. Numeric <u>gender</u> student? credit rating buy computer? id <u>age</u> 123 male 15 yes fair yes 456 female 38 qood no yes 872 male 65 fair no no 222 female 28 yes excellent yes female 20 111 no good no ... Nominal attributes: • · have values that are "names" of categories • Numeric attributes: · have values that are numbers it makes sense to compare their values using < and > example: we could base predictions on whether age < 35 What about id? ٠

Classification Learning						
	female female e have a s	•	<u>student?</u> yes no yes no	<u>credit rating</u> fair good fair excellent good ute whose value	buy computer? yes yes no yes no	
determine/predict.That output attribute is <i>nominal</i>.						
• The input attributes can be either nominal or numeric.						

Numeric Estimation student? credit rating <u>gender</u> buy computer? id <u>age</u> 123 male 15 yes fair 0.75 0.90 456 female 38 no qood 872 male 65 no fair 0.23 222 female 28 yes excellent 0.68 female 20 0.37 111 no good • We have a single output attribute whose value we want to determine/predict. • That output attribute is *numeric*. • The input attributes can be either nominal or numeric.









Finding Patterns (cont.)

- In data mining:
 - · the data is stored in electronic form
 - the process is automated (at least in part) using a computer program
 - the program "mines" the data
 - "sifting through" it to try find something useful/valuable

Data Mining vs. Data Query

- Database queries in SQL are *not* the same thing as data mining.
- Queries allow us to extract factual information.
 - "shallow knowledge"
- In data mining, we attempt to extract patterns and relationships.
 - "hidden knowledge"

Machine Learning

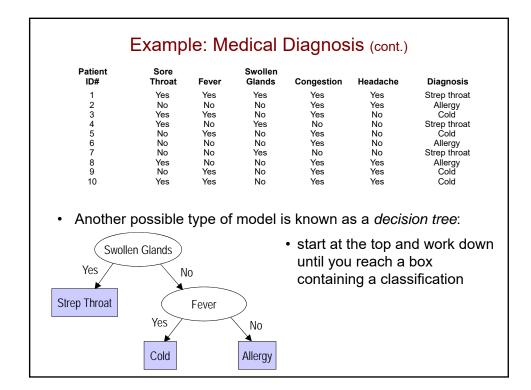
- In data mining, we apply an algorithm that "learns" something about the data.
 - a machine-learning algorithm
- We're ultimately going to consider three different types of machine learning:
 - classification learning
 - · association learning
 - numeric estimation

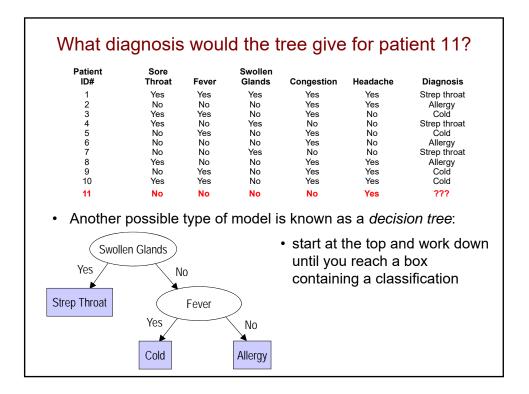
<list-item><list-item><list-item><list-item><list-item><list-item><list-item>

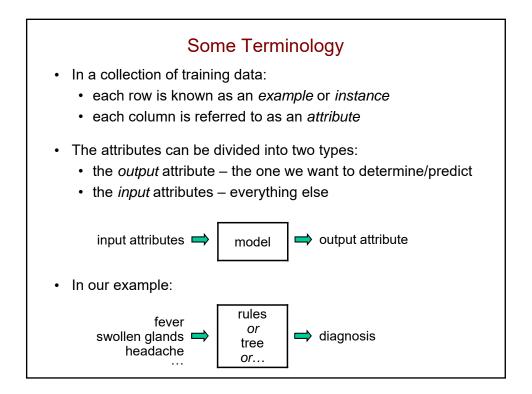
Sore					
Throat	Fever	Swollen Glands	Congestion	Headache	Diagnosis
Yes	Yes	Yes	Yes	Yes	Strep throa
No	No	No	Yes	Yes	Allergy
Yes	Yes	No	Yes	No	Cold
Yes	No	Yes	No	No	Strep throa
					Cold
					Allergy
					Strep throa
					Allergy Cold
Yes	Yes	No	Yes	Yes	Cold
	No Yes Yes No No Yes No	No No Yes Yes Yes No No Yes No No Yes No No Yes	NoNoNoYesYesNoYesNoYesNoYesNoNoNoNoNoNoYesYesNoNoYesNoNoNoYesNo	NoNoNoYesYesYesYesNoYesYesNoYesNoNoYesNoYesNoNoNoYesNoNoYesNoYesNoYesNoYesNoNoYesNoYesNoYesNoYesNoYes	NoNoNoYesYesYesYesNoYesNoYesNoYesNoNoNoYesNoYesNoNoYesNoYesNoNoNoNoYesNoNoNoYesNoNoYesNoNoYesYesNoNoYesNoYesNoYesNoYesYesNoYesNoYesYes



1 2 3 4 5	Yes No	Yes	V			Diagnosis
3 4			Yes	Yes	Yes	Strep throa
4		No	No	Yes	Yes	Allergy
	Yes	Yes	No	Yes	No	Cold
5	Yes	No	Yes	No	No	Strep throa
e	No No	Yes No	No No	Yes Yes	No No	Cold
6 7	NO	No	Yes	No	No	Allergy Strep throa
8	Yes	No	No	Yes	Yes	Allergy
9	No	Yes	No	Yes	Yes	Cold
10	Yes	Yes	No	Yes	Yes	Cold
patients i if sw then if sw	s a set c ollen Gl Diagnosi	of rules ands = s = St ands =	like the = Yes rep Thr = No an	e used for following: oat d Fever =	Ţ	

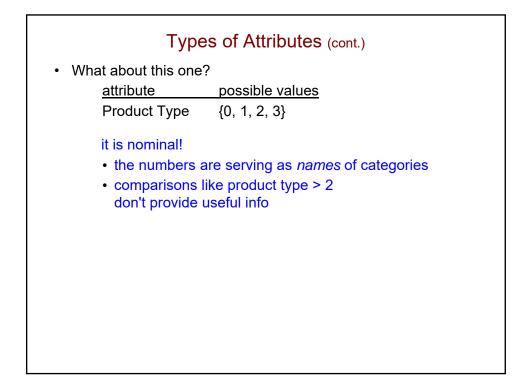


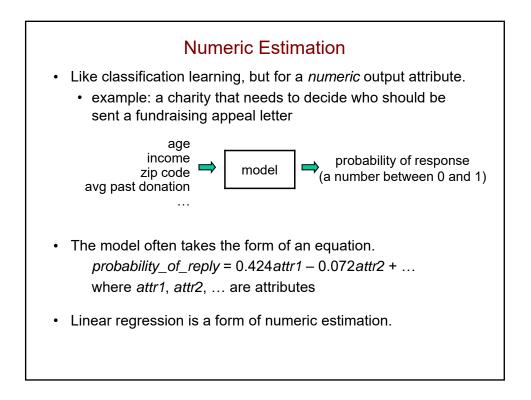


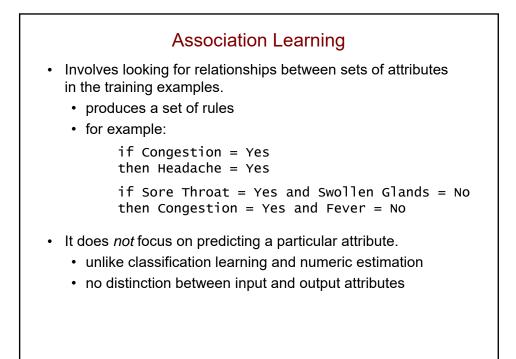


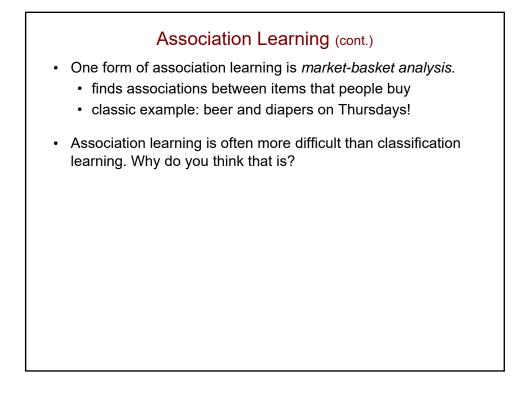
Types of Attributes
 Nominal attributes have values that are "names" of categories. there is a small set of possible values <u>attribute</u> possible values Fever {Yes, No} Diagnosis {Allergy, Cold, Strep Throat}
• In classification learning, the output attribute is always <i>nominal</i> .
Numeric attributes: A have values that are single numbers
 have values that are single numbers
 it makes sense to compare their values using < and >
 example: Body Temp
 each value is a single number like 98.0 or 101.5
 it could make sense to base our predictions on comparisons like Body Temp > 98.6

Туре	es of Attributes (cont.)
What about this one	?
attribute	possible values
Product Type	{0, 1, 2, 3}





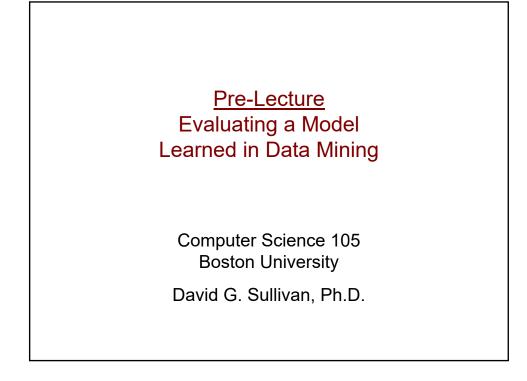


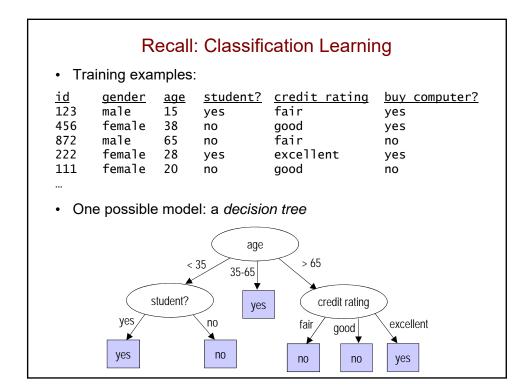


Customer ID	Account Type	Margin Account	Transaction Method	Trades/ Month	Sex	Age	Favorite Recreation	Annua Income
1005 1013	Joint Custodial	No No	Online Broker	12.5 0.5	F F	30–39 50–59	Tennis	40–59K
1013	Joint	No	Online	0.5 3.6	м	50–59 20–29	Skiing Golf	80–99k 20–39k
2110	Individual	Yes	Broker	22.3	Μ	30-39	Fishing	40-59k
1001	Individual	Yes	Online	5.0	М	40–49	Golf	60–79K
for eac	ch transa	action r	method.		<i>-</i> ave	eraye	trades/m	ontri
for eac	ch transa est appro	action r oach is	method.			erage	liades/iii	onun
for eac	ch transa	action r oach is	method.			erage	liades/iii	ontri
for eac	ch transa est appro databa	action r oach is se que	method.			Ū		Unin
for eac The be A. B.	ch transa est appro databa data m	action r oach is se que ining u	method. :: :ries sing class	sificatio	on le	arnin		ontri
for eac The be A.	ch transa est appro databa data m	action r oach is se que ining u	method. :: eries	sificatio	on le	arnin		ontri

Customer ID	Account Type	Margin Account	Transaction Method	Trades/ Month	Sex	Age	Favorite Recreation	Annua Income
1005 1013 1245 2110 1001	Joint Custodial Joint Individual Individual	No No No Yes Yes	Online Broker Online Broker Online	12.5 0.5 3.6 22.3 5.0	F M M M	30–39 50–59 20–29 30–39 40–49	Tennis Skiing Golf Fishing Golf	40–59K 80–99K 20–39K 40–59K 60–79K
			't know a n we use				rite recrea	ation,
	other fac	tors ca oach is	n we use				rite recrea	ition,
what o	other fac est appro databa	tors ca oach is se que	n we use	d to pr	edic	t it?		ition,
what o • The b A.	other fac est appro databa data m	tors ca oach is se que ining u	n we use :: ries	d to pr	edic on le	t it? arnin		ation,

Customer ID	Account Type	Margin Account	Transaction Method	Trades/ Month	Sex	Age	Favorite Recreation	
1005	Joint	No	Online	12.5	F	30–39	Tennis	40–59K
1013	Custodial	No	Broker	0.5	F	50-59	Skiing	80–99K
1245 2110	Joint Individual	No Yes	Online Broker	3.6 22.3	M M	20–29 30–39	Golf Fishing	20–39K 40–59K
1001	Individual	Yes	Online	5.0	M	40–49	Golf	40–39K 60–79K
what c	other fact	tors ca	't know a n we use	d to pr ∱	edic	t it?	K	 output is nomi
what c		tors ca	n we use	d to pr ↑ predictin	r edic g a sin	:t it? gle attril	► bute,	 output is nomi
what c	other fact	tors ca bach is	n we use	d to pr ∱	r edic g a sin mining	:t it? gle attril using ei	bute, ither	 output is nomi
what o	other fact est appro databas	tors ca bach is se que	n we use	d to pr f predicting so data r classifica	g a sin mining ation or	:t it? gle attril using ei estimat	bute, ither tion	 output is nomi
what c • The be A.	other fact est appro databas data m	tors ca bach is se que ining	n we use :: ries	d to pr predicting so data r classifica	g a sin nining ation or ation	st it? gle attril using ei estimat	bute, ither tion	creation, output is nomi so classificatio





Evaluating the Model For most non-trivial, real-world data sets, no learned model

- is likely to work perfectly on all possible examples.
- Our goal is *not* to create a model that perfectly matches the training data.
- Instead, we want a model that performs well on previously unseen examples.
 - we say that we want the model to generalize

٠

Test Examples To see how well a model generalizes, we typically withhold some of the available data as *test examples*. these examples are *not* used to train the model Let's assume we have data for 100 customers. all of the data is already classified use 90 examples to learn the model (the training data) use 10 examples to test the model (the test data)

Using th Test Exa		< 3!	age 35-65	> 65	
		student?	yes	Credit ra	iting
	yes yes]	no	fair good	excellent v no yes
				buy	computer?
<u>id</u> gen			<u>credit rat.</u>	<u>actual</u>	predicted
954 mal		no	good	yes	yes
888 fem		yes	good	no	yes
357 mal		yes	fair	yes	yes
245 fem 177 fem		no	excellent	no	no
523 mal		no no	good good	no	no
999 mal		no	good	yes no	no yes
126 fem		yes	fair	no	no
443 mal		yes	fair	yes	
747 fem	ale 47	no	excellent	no	

		Sumi	marizing	the Result	S	
					buy	computer?
<u>id</u>	<u>gender</u>	<u>age</u>	<u>student?</u>	<u>credit rat.</u>	<u>actual</u>	<u>predicted</u>
954	male	45	no	good	yes	yes
888	female	22	yes	good	no	yes
357	male	25	yes	fair	yes	yes
245	female	28	no	excellent	no	no
177	female	80	no	good	no	no
523	male	68	no	good	yes	no
999	male	37	no	good	no	yes
126	female		yes	fair	no	no
443	male	19	yes	fair	yes	yes
747	female	47	no	excellent	no	yes
	uracy of t or rate = 4		odel = 6 /10 40%	= 60%		
	blem: the being equ			all misclassific	cations	

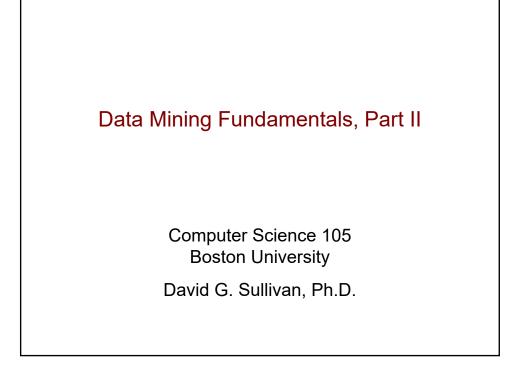
					buy	computer?
id	<u>gender</u>	<u>age</u>	<u>student?</u>	<u>credit rat.</u>	<u>actual</u>	<u>predicted</u>
954	male	45	no	good	yes	yes 🗲
888	female	22	yes	5		yes
357	male	25	yes	es fair		yes 🗲
245	female	28	no	no excellent		no
177	female	80	no	no good		no
523	male	68	no			no
999	male	37	no	good	no	yes
126	female	70	yes	fair	no	no
443	male	19	yes	fair	yes	yes 🗲
747	female	47	no	excellent	no	yes
			pred	dicted class		
			ye	s no		
a	ctual class	s: yes	3			
		no				

					buy	computer?
id	<u>gender</u>	<u>age</u>	<u>student?</u>	<u>credit rat.</u>	<u>actual</u>	predicted
954	male	45	no	good	yes	yes
888	female	22	yes			yes
357	male	25	yes	ves fair		yes
245	female	28	no			no
177	female	80	no	no good		no
523	male	68	no			no 🗲
999	male	37	no	good	no	yes
126	female	70	yes	fair	no	no
443	male	19	yes	fair	yes	yes
747	female	47	no	excellent	no	yes
			pred	dicted class		
			ye	s no		
a	ctual class	s: yes	3	1		
		no				

						buy	computer?	
id	<u>gender</u>	<u>age</u>	<u>student?</u>	<u>credi</u>	t rat.	<u>actual</u>	predicted	
954	male	45	no	good		yes	yes	
888	female	22	yes			no	yes 🗲	
357	male	25	yes	ves fair		yes	yes	
245	female	28	no	no excellent		no	no	
177	female	80	no			no	no	
523	male	68	no			yes	no	
999	male	37	no	- J		no	yes 🗲	
126	female	70	yes	fair		no	no	
443	male	male	19	yes	fair		yes	yes
747	female	47	no	excel	lent	no	yes ←	
			prec	dicted o	class			
			ye	S	no			
a	ctual class	s: yes	3		1			
		no	3					

	ι	Jsing	a Conf	usion Matri	x	
<u>id</u> 954 888 357 245 177 523 999 126 443	<u>gender</u> male female female female male male female female male	45 22 25 28 80 68 37 70	student? no yes yes no no no yes yes	<u>credit rat.</u> good good fair excellent good good fair fair	buy o actual yes no yes no yes no no yes	computer? predicted yes yes no no no yes no yes
747	female	47	no	excellent	no	yes
			prec yes	licted class		
a	ctual class	: yes	3	1		
		no	3		-	

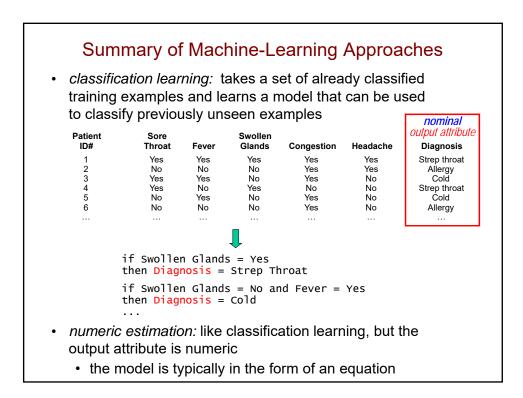
		Using	a Conf	usion Matri	x	
<u>id</u> 954 888 357 245 177 523 999 126 443 747	<u>gender</u> male female female female male female female female female	age 45 22 25 28 80 68 37 70 19 47	student? no yes yes no no no no yes yes no	<u>credit rat.</u> good good fair excellent good good fair fair excellent	buy c actual yes no yes no yes no yes no yes no	computer? predicted yes yes yes no no yes no yes yes
	ctual class the diago	no	ye: 3 3	dicted class s no 1 3 orrectly classif	the diagonal of the matrix ied examp	(

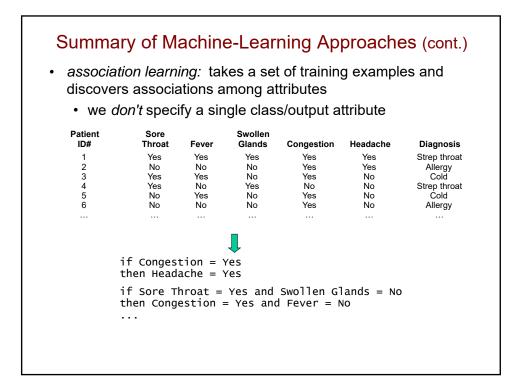


Rec	call: Types of Attributes
 there is typical 	s have values that are "names" of categories. ly a small set of possible values
<u>attribute</u> Fever	<u>possible values</u> {Yes, No}
	{Allergy, Cold, Strep Throat}
Body Temp	{below 96, 96-99, 99-102, above 102}
Numeric attributes	s have values that are <u>single</u> numbers.
	ly a wide range of possible values possible values
Body Temp	any single real number in 96.0-106.0 any single integer in \$15,000-250,000
 it makes sense \$210,000 > \$ 98.6 < 101.3 	e to <i>order/compare</i> their values 6125,000

Which of the attributes are numeric?

Custor ID 1005 1013 1245	Type Joint Custodial	Margin Account No No No	Transaction Method Online Broker Online	Trades/ Month 12.5 0.5 3.6	Sex F F M	Age 30–39 50–59 20–29	Favorite Recreation Tennis Skiing Golf	Annual Income 40–59K 80–99K 20–39K	
2110 1001		Yes Yes	Broker Online	22.3 5.0	M	30–39 40–49	Fishing Golf	40–59K 60–79K	
А.	customer	ID							
В.	trades/mo	onth							
C.	age								
D.	two of the	above	Э						
E.	all of the a	above							

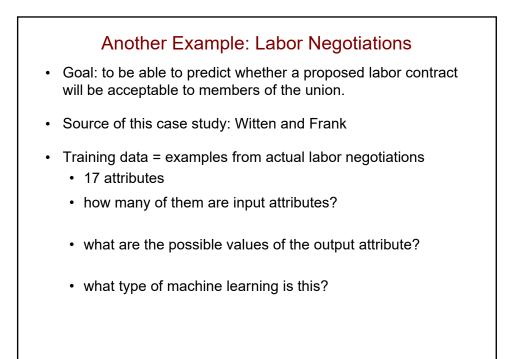


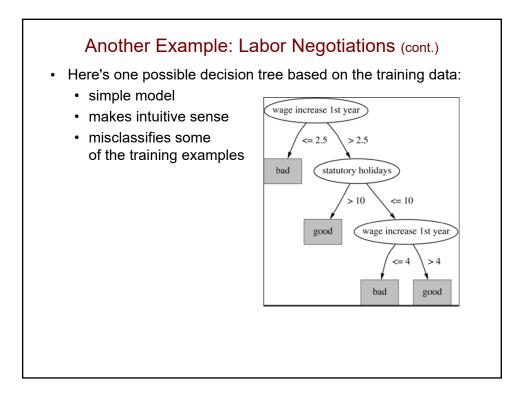


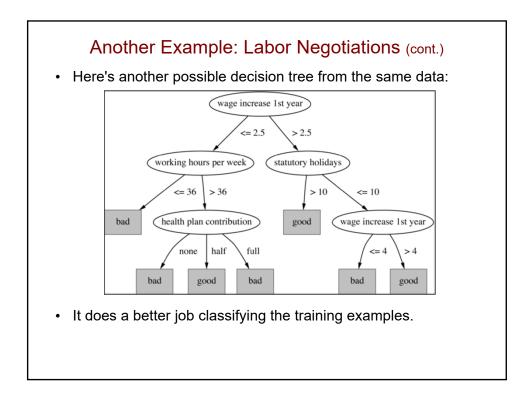
Customer ID		Aargin ccount	Transaction Method	Trades/ Month	Sex	Age	Favorite Recreation	Annua Income
1005 1013	Joint Custodial	No No	Online Broker	12.5 0.5	F F	30–39 50–59	Tennis Skiing	40–59K 80–99K
1245	Joint	No	Online	3.6	М	20-29	Golf	20–39K
2110 1001	Individual Individual	Yes Yes	Broker Online	22.3 5.0	M M	30–39 40–49	Fishing Golf	40–59K 60–79K
Α.	database	e quei	ries					
В.	data min	ing us	sing class	sificatio	on le	arnin	q	
		U	Ū				0	
							9	

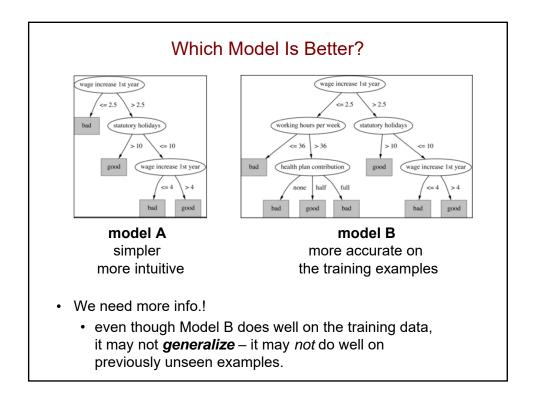
Customer ID	Account Type	Margin Account	Transaction Method	Trades/ Month	Sex	Age	Favorite Recreation	Annua Incom
1005 1013	Joint Custodial	No No	Online Broker	12.5 0.5	F F	30–39 50–59	Tennis Skiing	40–59k 80–99k
1245	Joint	No	Online	3.6	M	20-29	Golf	20–39k
2110	Individual	Yes	Broker	22.3	М	30–39	Fishing	40–59k
1001	Individual	Yes	Online	5.0	М	40–49	Golf	60–79k
accou	nt type, i	transad	to discove ction meth			•	between	
accou		transao oach is	ction meth			•	between	
accou The b	nt type, t est appro databa	transao oach is ise que	ction meth	nod, ar	nd aç	ge.		
accou The b A.	nt type, f est appro databa data m	transac oach is ise que iining u	ction meth :: eries	nod, ar sificatio	nd ao on le	ge. arnin		

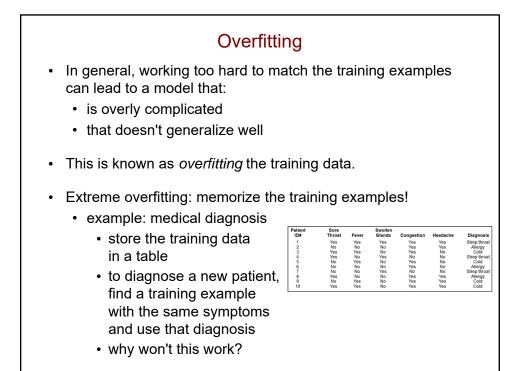
Customer ID	Account Type	Margin Account	Transaction Method	Trades/ Month	Sex	Age	Favorite Recreation	Annua Income
1005 1013 1245 2110 1001	Joint Custodial Joint Individual Individual	No No Yes Yes	Online Broker Online Broker Online	12.5 0.5 3.6 22.3 5.0	F F M M	30–39 50–59 20–29 30–39 40–49	Tennis Skiing Golf Fishing Golf	40–59K 80–99K 20–39K 40–59K 60–79K
	the annu		to know v <mark>ome</mark> .	vnich a	aund	ules	end to	
affect	•••••	ual inc	ome.	vnich a	und	utes i	end to	
affect	the annu	u <mark>al inc</mark> bach is	ome.	vnich a	uund	ules	iena to	
affect The be	the annu est appro databas	u <mark>al inc</mark> bach is se que	ome.					
affect • The bo	the annu est appro databas data mi	bach is se que ining u	ome. : ries	sificatio	on le	arnin		

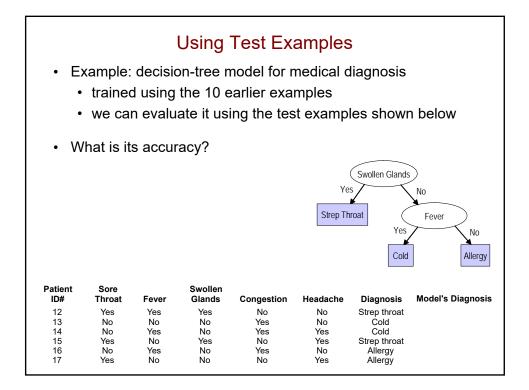












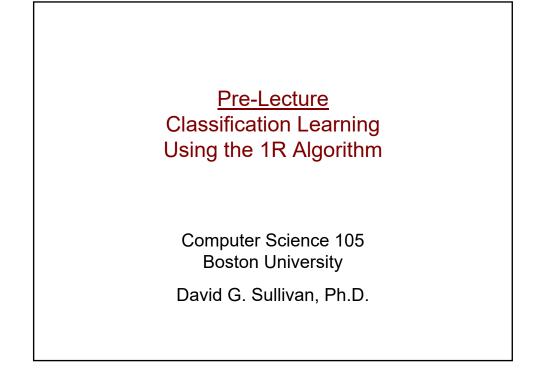
Evaluating Classification Learning Models Swollen Patient Sore Throat ID# Fever Glands Congestion Headache Diagnosis Model's Diagnosis 12 Yes Yes Yes No No Strep throat Strep throat 13 No No No Yes No Ċold Cold Allergy Cold 14 15 No Yes No Yes Yes Strep throat Allergy Yes Strep throat Yes No No Yes 16 17 No Yes No Yes No Ċold Allergy Yes No No No Yes Allergy The error rate of a model is the percentage of test examples ٠ that it misclassifies. in our example, the error rate = ____ error rate = 100 – accuracy Problem: these metrics treat all misclassifications as being equal. · this isn't always the case · example: more problematic to misclassify strep throat than to misclassify a cold or allergy

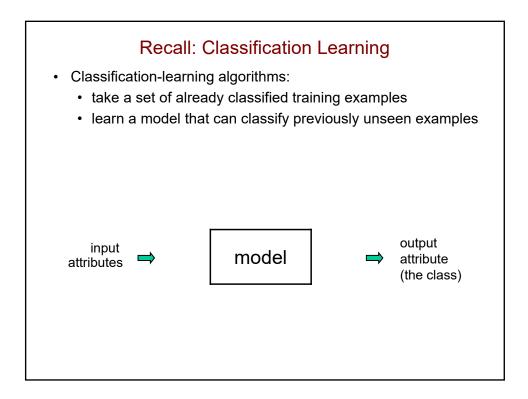
Patient ID#	Sore Throat	Fever	Swollen Glands	Congestion	Headache	Diagnosis	Model's Diagnosis
12 13	Yes No	Yes No	Yes No	No Yes	No No	Strep throat Cold	Strep throat
13	NO	Yes	No	Yes	Yes	Cold	Allergy Cold
15	Yes	No	Yes	No	Yes	Strep throat	Strep throat
16 17	No Yes	Yes No	No No	Yes No	No Yes	Allergy Allergy	Cold Allergy
	-		ore detaile onfusion r	•		ted class	
	-			•	predic	ted class	
w.	-	se a co		natrix:	predic	ted class	2
w.	e can us	se a co ass: c	onfusion r	natrix:	predic	ted class	2
w	e can us	se a co ass: c	onfusion r	<i>natrix</i> : cold	predic	ted class	2

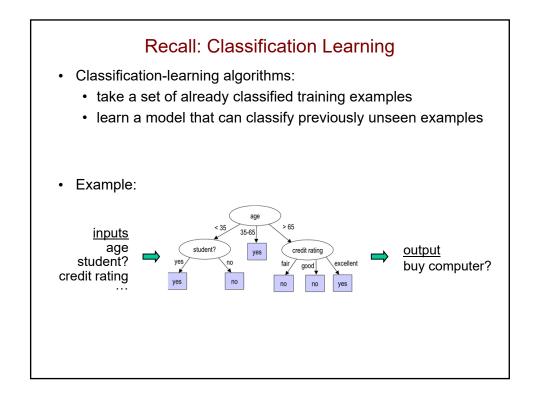
 Interpreting a Let's say that we had a lar that we obtained the follow 	rger numbe	r of test exar	
	•	predicted cla	iss
	cold	allergy	strep throat
 actual class: cold allergy strep throa what is the accuracy of total # of test cases = 1 	f the model	8 15 4 ?	7 3 33
 what is its error rate? 			

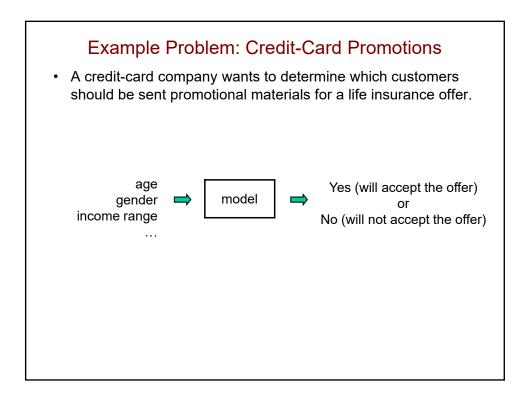
Interp	reting a Co	nfusion	Matrix (co	ont.)
		<u> </u>	predicted cla	ISS
		cold	allergy	strep throat
actual class:	cold allergy strep throat	25 6 5	8 15 4	7 3 33
 how many t 	test cases of s	strep throa	at are there?	
 how many a 	actual colds w	ere misdi	agnosed?	
what perce	ntage of actua	I colds w	ere correctly	diagnosed?

	w many act es the mod	el misc		
		cold	allergy	strep throat
actual class:	cold	5	3	2
	allergy	2	6	5
	strep throat	1	2	8

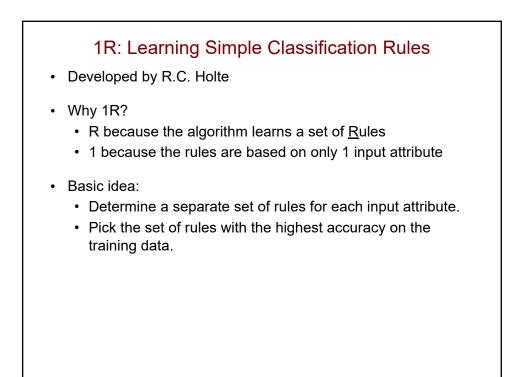


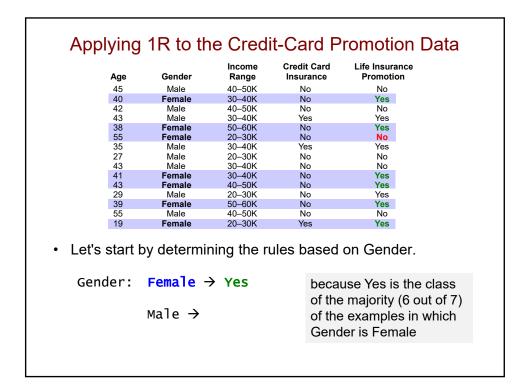


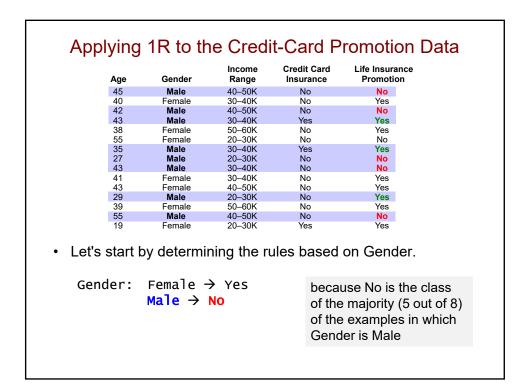




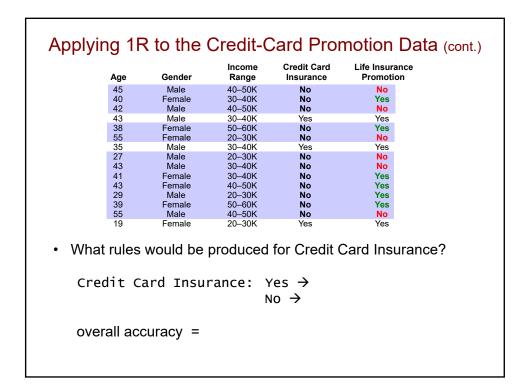
15 training	examples	:		
Age	Gender	Income Range	Credit Card Insurance	<i>class/</i> output attribute Life Insurance Promotion
45	Male	40–50K	No	Νο
40	Female	30–40K	No	Yes
42	Male	40–50K	No	No
43	Male	30–40K	Yes	Yes
38	Female	50–60K	No	Yes
55	Female	20–30K	No	No
35	Male	30–40K	Yes	Yes
27	Male	20–30K	No	No
43	Male	30–40K	No	No
41	Female	30–40K	No	Yes
43 29	Female Male	40–50K 20–30K	No No	Yes Yes
29 39	Female	20–30K 50–60K	No	Yes
39 55	Male	40–50K	No	No
19	Female	20–30K	Yes	Yes







Age	Gender	Income Range	Credit Card Insurance	Life Insurance Promotion
45	Male	40–50K	No	No
40	Female	30–40K	No	Yes
42	Male	40–50K	No	No
43	Male	30–40K	Yes	Yes
38	Female	50–60K	No	Yes
55	Female	20–30K	No	No
35	Male	30–40K	Yes	Yes
27	Male	20–30K	No	No
43	Male	30–40K	No	No
41	Female	30–40K	No	Yes
43	Female	40–50K	No	Yes
29	Male	20–30K	No	Yes
39	Female	50-60K	No	Yes
55	Male	40–50K	No	No
19	Female	20–30K	Yes	Yes
⁻ hus, we e	nd up with	the follow	wing rules b	ased on Gende
Gender:	Female -	→ Yes	(6 out of 7)	
	Male →	NO	(5 out of 8)	
overall acc	curacy =		<u>11</u> = 73% 15	



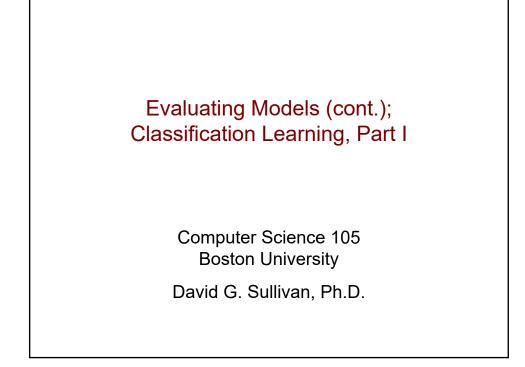
Age	Gender	Income Range	Credit Car Insurance	
45	Male	40–50K	No	No
40	Female	30–40K	No	Yes
42	Male	40–50K	No	No
43	Male	30–40K	Yes	Yes
38	Female	50–60K	No	Yes
55	Female	20–30K	No	No
35	Male	30–40K	Yes	Yes
27	Male	20–30K	No	No
43 41	Male Female	30–40K 30–40K	No No	No Yes
41	Female	30–40K 40–50K	No	Yes
29	Male	20–30K	No	Yes
39	Female	50-60K	No	Yes
55	Male	40–50K	No	No
19	Female	20–30K	Yes	Yes
nat rule		•		me Range?
come	Range: 2	0-30K →	No/Yes	(2 out of 4)
	-	0-40к →		(4 out of 5)
	-	0-50K →		(3 out of 4)
				· /
	5	0-60к →	Yes	(2 out of 2)

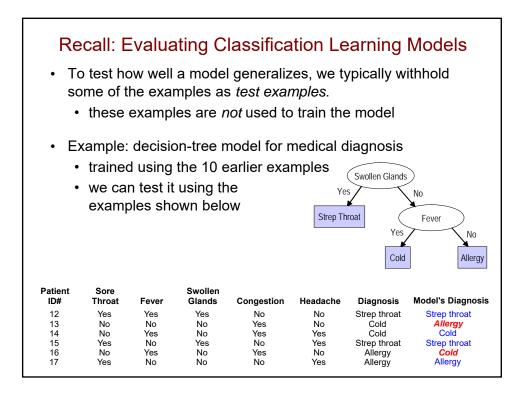
Age	Gender	Income Range	Credit Card Insurance	Life Insurance Promotion
45	Male	40–50K	No	No
40	Female	30–40K	No	Yes
42	Male	40–50K	No	No
43	Male	30–40K	Yes	Yes
38	Female	50–60K	No	Yes
55	Female	20–30K	No	No
35	Male	30–40K	Yes	Yes
27	Male	20–30K	No	No
43	Male	30–40K	No	No
41	Female	30–40K	No	Yes
43	Female	40–50K	No	Yes
29	Male	20–30K	No	Yes
39	Female	50-60K	No	Yes
55 19	Male Female	40–50K 20–30K	No Yes	No Yes
What rules Income R	ange: 20 30 40	produced)-30к →)-40к →)-50к →)-60к →	NO (2 Yes (4 NO (3	e Range? out of 4) out of 5) out of 4) out of 2)

 To hat of post 	ndle num		ributes,	we ne	ed to a			ange	÷	
• One a	pproach	. ,		•		es by ag b <i>inary</i> (2		split		
			Income	Cred	it Card	Life Insurance				
	Age (Gender	Range	n						
	45	Male	40–50K	1	١o	No				
		emale	1	١o	Yes					
	42	Male	40–50K		۱o	No				
	43	Male	30–40K	-	es	Yes				
		emale	50–60K	-	lo	Yes				
		Female 20–30K			lo	No				
	35 27	Male Male	30–40K 20–30K	-	'es lo	Yes No				
	43	Male	20–30K 30–40K		10 10	NO				
	43 Male 41 Female					Yes				
		Female	30–40K 40–50K		No	Yes				
	29	Male	20–30K		No	Yes				
	39 1	50-60K	1	lo	Yes					
	40–50K	1	۱o	No						
	19 Female			Y	es	Yes				
sort by A	ge: 19	27 29	35 38	39 40	41 42	43 43	43 45	55	55	
ADC 1	fe Ins: Y	NY	Y Y	ŶŶ	Y N	Y Y	N N	N	N	
<u> </u>	IC 1115. I	IN I	1 1	1 1	I IN	I Í		IN	IN	

	Handling Numeric Attributes (cont.)															
•	 Here's one possible binary split for age: 															
	Age:	19	27	29	35	38	39	40	41	42	43	43	43	45	55	55
	Life Ins:	Y	Ν	Y	Y	Y	Y	Y	Y	Ν	Y	Y	Ν	Ν	Ν	Ν
	 the corresponding rules are: 															
		Age: <= 39 → Yes > 39 → No						overall accuracy: 10/15 = 67%								
	 The following is one of the splits with the best overall accuracy: 															
	Age:	19	27	29	35	38	39	40	41	42	43	43	43	45	55	55
	Life Ins:	Y	Ν	Y	Y	Y	Y	Y	Y	Ν	Y	Y	Ν	Ν	Ν	Ν
	 the corresponding rules are: Age: <= 43 → Yes (9 out of 12) overall accuracy: 															
	> 43 \rightarrow No (3 out of 3) 12/15 = 80%															

Ge	Summary of 1 nder:Female → Yes Male → No	(6 out of 7)	overall accuracy: 11/15 = 73%
Cred.C	ard Ins: Yes \rightarrow Yes No \rightarrow No*	(3 out of 3) (6 out of 12)	overall accuracy: 9/15 = 60%
Income R	ange: 20-30K → No* 30-40K → Yes 40-50K → No 50-60K → Yes	(3 out of 4)	overall accuracy: 11/15 = 73%
	Age: <= 43 → Yes > 43 → No	(9 out of 12) (3 out of 3)	overall accuracy: 12/15 = 80%
	e the rules based on Age / on the training data, 1F	•	

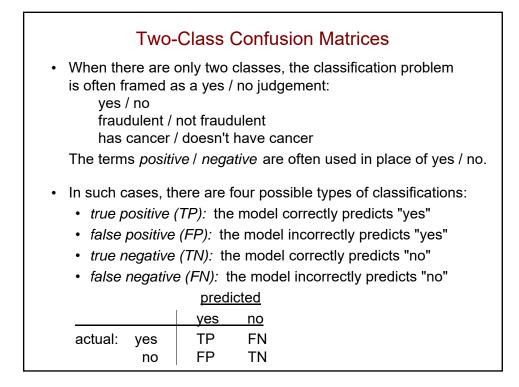




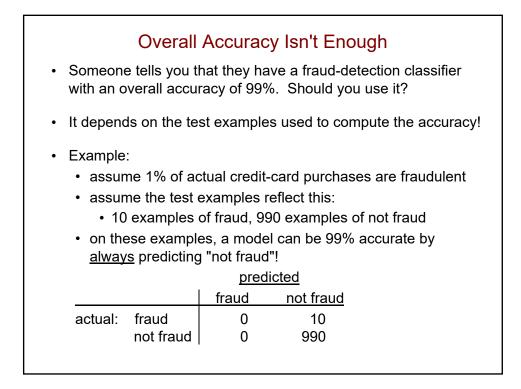
ID#	Sore Throat	Fever	Swollen Glands	Congestion	Headache	Diagnosis	Model's Diagnosi
12	Yes	Yes	Yes	No	No	Strep throat	Strep throat
13 14	No No	No Yes	No No	Yes Yes	No Yes	Cold Cold	Allergy Cold
15	Yes	No	Yes	No	Yes	Strep throat	Strep throat
16	No	Yes	No	Yes	No	Allergy	Cold
17	Yes	No	No	No	Yes	Allergy	Allergy
					predic	ted class	;
							_
_				cold			strep throat
-	actual cl	ass: c	old	cold			-
-	actual cl		old	<u>cold</u>			-
-	actual cl	a		<u>(1</u> 1			-
-	actual cl	a	llergy	<u>(1</u> 1			-

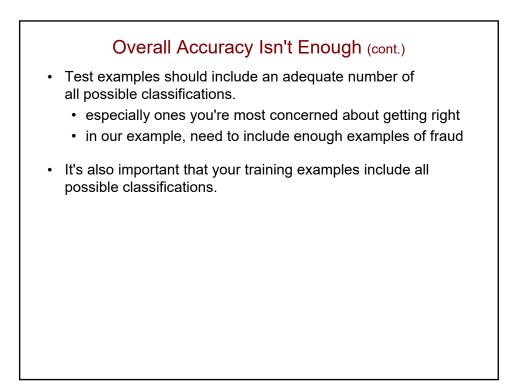
How many of	the predicti	ons of	allergy are	e incorrect?
			predicted cla	ISS
		cold	allergy	strep throat
actual class:	cold	5	3	2
	allergy	2	6	5
	strep throat	1	2	8

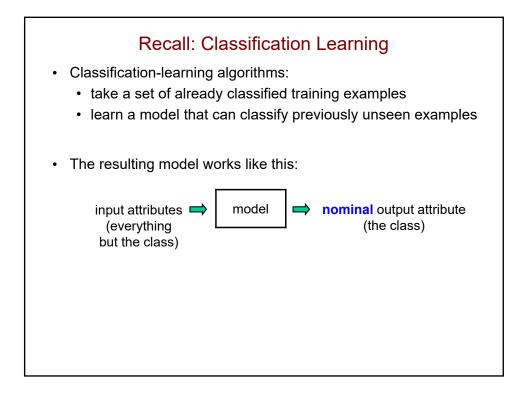
What is t	he overall a	accurac	y of the m	odel?
			predicted cla	<u>ss</u>
		cold	allergy	strep throat
actual class:	cold	5	3	2
	allergy	2	6	5
	strep throat	1	2	8
total # of test	examples = 3	4		

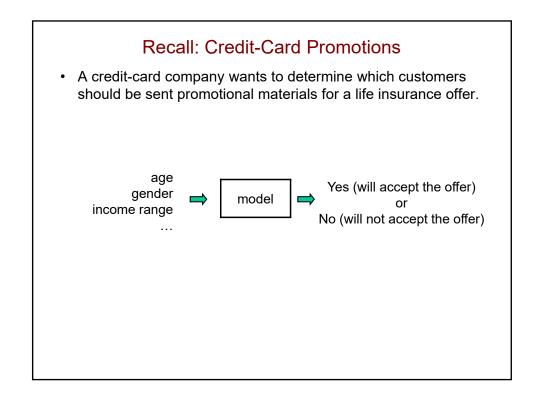


	Ŭ		ng Confusi credit-card fr	on Matrices
	wo differen lifferent mo		tion-learning	techniques and
Performa	ance on 400) test exam	nples:	
			by model A	overall accuracy =
		fraud	not fraud	(100+250)/400 = .875
actual:	fraud	100	10	
	not fraud	40	250	
		predicted	l by model B	overall accuracy =
		fraud	not fraud	(80+270)/400 = .875
actual:	fraud	80	30	
	not fraud	20	270	
which	model is b	etter?		









15 training	Gender	Income Range	Credit Card Insurance*	<i>class/</i> output attribute Life Insurance Promotion
45	Male	40–50K	No	No
40	Female	30–40K	No	Yes
42	Male	40–50K	No	No
43	Male	30–40K	Yes	Yes
38	Female	50-60K	No	Yes
55	Female	20-30K	No	No
35	Male	30–40K	Yes	Yes
27	Male	20-30K	No	No
43	Male	30–40K	No	No
41	Female	30-40K	No	Yes
43	Female	40-50K	No	Yes
29	Male	20-30K	No	Yes
39	Female	50–60K	No	Yes
55	Male	40-50K	No	No
19	Female	20–30K	Yes	Yes
specifying		e customer a	es/No attribute accepted a sir	

Recall: Summary	of 1R Results
Gender:Female → Yes Male → No	
Cred.Card Ins: Yes → Yes No → No*	(3 out of 3) overall accuracy: (6 out of 12) 9/15 = 60%
Income Range: 20-30K → No* 30-40K → Yes 40-50K → No 50-60K → Yes	(3 out of 4) $11/15 = 73%$
Age: <= 43 → Yes > 43 → No	(9 out of 12) overall accuracy: (3 out of 3) 12/15 = 80%
• 1R learned the above set of car	ndidate models.
 Because the rules based on Aga accuracy on the training data, 1 	e have the highest overall R selects them as the <i>final model</i> .

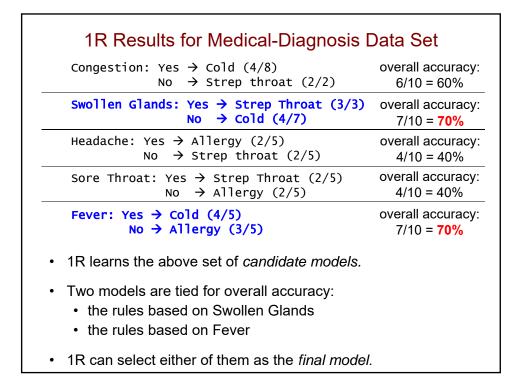
Patient ID#	Sore Throat	Fever	Swollen Glands	Congestion	Headache	Diagnosis
1	Yes	Yes	Yes	Yes	Yes	Strep throa
2	No	No	No	Yes	Yes	Allergy
3 4	Yes Yes	Yes No	No Yes	Yes No	No No	Cold Strep throa
5	No	Yes	No	Yes	No	Cold
5 6	No	No	No	Yes	No	Allergy
7	No	No	Yes	No	No	Strep throa
8	Yes	No	No	Yes	Yes	Allergy
9 10	No Yes	Yes Yes	No No	Yes Yes	Yes Yes	Cold Cold
We wa	nt to be ab	ole to di	agnose	new patie	nts	
	nt to be ab s the outpu		0	new patie	nts	
		ıt attribi	0	new patie	nts	
What is	s the outpu	ıt attribı D#	ute?	new patier	nts	
What is A. B.	s the outpu Patient I Swollen	ıt attribı D#	ute?	new patier	nts	
What is A. B. C.	s the outpu Patient I	ıt attribı D#	ute?	new patie	nts	
What is A. B.	s the outpu Patient I Swollen	ıt attribi D# Glands	ute?	new patier	nts	

ID#	Throat	Fever	Swollen Glands	Congestion	Headache	Diagnosis
1	Yes	Yes	Yes	Yes	Yes	Strep throat
2	No	No	No	Yes	Yes	Allergy
3	Yes	Yes	No	Yes	No	Cold
4	Yes	No	Yes	No	No	Strep throat
5	No	Yes	No	Yes	No	Cold
6	No	No	No	Yes	No	Allergy
7	No	No	Yes	No	No	Strep throat
8	Yes	No	No	Yes	Yes	Allergy
9 10	No Yes	Yes Yes	No No	Yes Yes	Yes Yes	Cold Cold
What is t	he outpu	ıt attribu	ute?			
Because	it is nom	ninal:				
• W0 D0		ificatio		a		
• we ne	ed class	ancatio	rieamin	iy		
		منامات ما	acrithm		100	
 1R is 	one pos	sible al	gonunn	we could ι	lse	

Patient ID#	Sore Throat	Fever	Swollen Glands	Congestion	Headache	Diagnosis
1	Yes	Yes	Yes	Yes	Yes	Strep throat
2	No	No	No	Yes	Yes	Allergy
3	Yes	Yes	No	Yes	No	Cold
4	Yes	No	Yes	No	No	Strep throat
5	No	Yes	No	Yes	No	Cold
6	No	No	No	Yes	No	Allergy
7	No	No	Yes	No	No	Strep throat
8	Yes	No	No	Yes	Yes	Allergy
9	No	Yes	No	Yes	Yes	Cold
10	Yes	Yes	No	Yes	Yes	Cold
•	tion			C		
A. Conges Yes No	→ Cold	p throat	:	• • •	estion: s → Colo → Allo	-

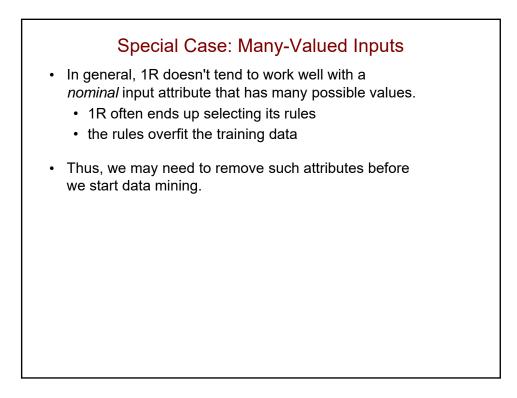
Patient ID#	Sore Throat	Fever	Swollen Glands	Congestion	Headache	Diagnosis
1	Yes	Yes	Yes	Yes	Yes	Strep throat
2	No	No	No	Yes	Yes	Allergy
3	Yes	Yes	No	Yes	No	Cold
4	Yes	No	Yes	No	No	Strep throa
5	No	Yes	No	Yes	No	Ċold
6	No	No	No	Yes	No	Allergy
7	No	No	Yes	No	No	Strep throat
8	Yes	No	No	Yes	Yes	Allergy
9	No	Yes	No	Yes	Yes	Cold
10	Yes	Yes	No	Yes	Yes	Cold
	en Gland → Stre → Alle	p Throa	at	U .	len Gland s → Colo → Stro	-
	en Gland → Stre → Cold	p Throa	at		e than one d be learn	

Patient ID#	Sore Throat	Fever	Swollen Glands	Congestion	Headache	Diagnosis
1	Yes No	Yes No	Yes No	Yes	Yes Yes	Strep throat Allergy
3 4 5	Yes Yes No	Yes No Yes	No Yes No	Yes No Yes	No No No	Cold Strep throat Cold
6 7	No No	No No	No Yes	Yes No	No No	Allergy Strep throat
8 9 10	Yes No Yes	No Yes Yes	No No No	Yes Yes Yes	Yes Yes Yes	Allergy Cold Cold
A. Headad Yes No	→ Cold	p throat		U .	ache: s → Allo → Colo	
B. Headad Yes	che: → Alle	rgy			than one be learn	

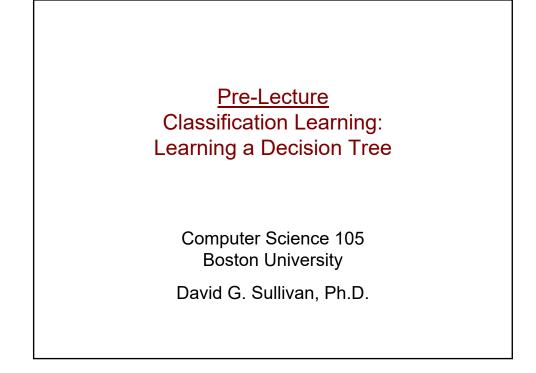


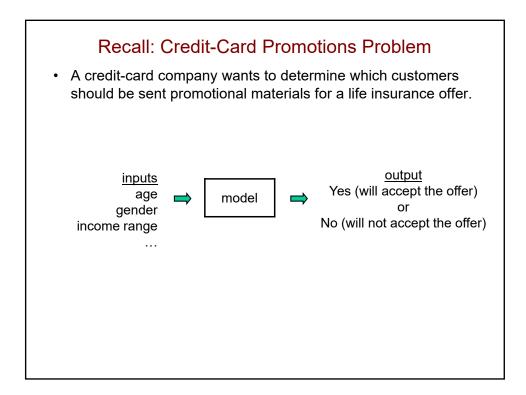
1YesYesYesYesYesStrep thr2NoNoNoYesYesAllergy3YesYesNoYesNoCold4YesNoYesNoNoStrep thr5NoYesNoYesNoCold6NoNoNoYesNoAllergy7NoNoYesNoNoStrep thr	ID#	Sore Throat	Fever	Swollen Glands	Congestion	Headache	Diagnosis
2 No No No Yes Yes Allergy 3 Yes Yes No Yes No Cold 4 Yes No Yes No No Strep thr 5 No Yes No Yes No Cold 6 No Yes No Yes No Allergy 7 No No Yes No No Strep thr 8 Yes No No Yes No Strep thr 8 Yes No No Yes Yes Allergy 9 No Yes No No Strep thr 10 Yes Yes No Yes Yes Cold f we learned rules based on Patient ID, what accuracy Yes Cold Yes Yes	1	Yes	Yes	Yes	-		Strep throa
4 Yes No Yes No No Strep thr. 5 No Yes No Yes No Cold 6 No No No Yes No Allergy 7 No No Yes No No Allergy 9 No Yes No Yes Yes Allergy 9 No Yes No Yes Yes Cold 10 Yes Yes No Yes Yes Cold f we learned rules based on Patient ID, what accuracy Strep thr. Strep thr. Strep thr.	2	No	No	No	Yes	Yes	
4 Yes No Yes No No Strep thr. 5 No Yes No Yes No Cold 6 No No No Yes No Allergy 7 No No Yes No No Allergy 9 No Yes No Yes Yes Allergy 9 No Yes No Yes Yes Cold 10 Yes Yes No Yes Yes Cold f we learned rules based on Patient ID, what accuracy Strep thr. Strep thr. Strep thr.	3	Yes	Yes	No	Yes		
5 No Yes No Yes No Cold 6 No No No Yes No Allergy 7 No No Yes No No Strep thrit 8 Yes No No Yes Yes Allergy 9 No Yes No Yes Yes Allergy 10 Yes Yes No Yes Yes Cold f we learned rules based on Patient ID, what accuracy	4	Yes	No	Yes	No	No	Strep throa
6 No No No Yes No Allergy 7 No No Yes No No Strep thr 8 Yes No No Yes Yes Allergy 9 No Yes No Yes Yes Allergy 10 Yes Yes No Yes Yes Cold f we learned rules based on Patient ID, what accuracy Interval Interval Interval Interval	5	No	Yes	No	Yes	No	
7 No No Yes No No Strep throws 8 Yes No No Yes Yes Allergy 9 No Yes No Yes Yes Cold 10 Yes Yes No Yes Yes Cold f we learned rules based on Patient ID, what accuracy	6	No	No	No	Yes	No	Allergy
8 Yes No No Yes Yes Aliergy 9 No Yes No Yes Yes Cold 10 Yes Yes No Yes Yes Cold f we learned rules based on Patient ID, what accuracy		No	No	Yes	No	No	Strep throa
9 No Yes No Yes Yes Cold 10 Yes Yes No Yes Yes Cold f we learned rules based on Patient ID, what accuracy		Yes	No	No	Yes	Yes	Allergy
f we learned rules based on Patient ID, what accuracy		No	Yes	No	Yes	Yes	
	10	Yes	Yes	No	Yes	Yes	Cold
				d on Pat	ient ID, wl	nat accura	асу
				d on Pat	ient ID, wl	nat accura	асу
				d on Pat	ient ID, wl	nat accura	асу
				d on Pat	ient ID, wl	nat accura	асу
				d on Pat	ient ID, wl	nat accura	асу
				d on Pat	ient ID, wl	nat accura	асу
				d on Pat	ient ID, wl	nat accura	асу

Yes No Yes Yes No	Yes No Yes	Yes No	Yes Yes	Yes	Strep throa
Yes Yes			Vec		
Yes	res			Yes	Allergy
	No	No Yes	Yes No	No No	Cold Strep throa
	Yes	No	Yes	No	Cold
No	No	No	Yes	No	Allergy
No	No	Yes	No	No	Strep throa
Yes	No	No	Yes	Yes	Allergy
	Yes	No		Yes	Cold
Yes	Yes	No	Yes	Yes	Cold
y have?	2 100% D: 1 2	%! → Str → All	ep thro ergy (1	at (1/1 /1)	•
		• •			
	Yes No Yes ned rule y have?	Yes No No Yes Yes Yes have? 100% ent ID: 1 2	Yes No No No Yes No Yes Yes No hed rules based on Pati y have? 100%! ent ID: $1 \rightarrow \text{Str}$ $2 \rightarrow \text{All}$	Yes No No Yes No Yes No Yes Yes Yes No Yes hed rules based on Patient ID, wh y have? 100%! ent ID: 1 → Strep thro	Yes No No Yes Yes No Yes No Yes Yes No Yes Yes Yes hed rules based on Patient ID, what accurate y have? 100%! ent ID: $1 \rightarrow \text{Strep throat (1/1)}$ $2 \rightarrow \text{Allergy (1/1)}$



ID#	Sore Throat	Fever	Swollen Glands	Congestion	Headache	Diagnosis
1	Yes	Yes	Yes	Yes	Yes	Strep throat
2	No	No	No	Yes	Yes	Allergy
3	Yes	Yes	No	Yes	No	Cold
4	Yes	No	Yes	No	No	Strep throat
5	No	Yes	No	Yes	No	Cold
6	No	No	No	Yes	No	Allergy
7	No	No	Yes	No	No	Strep throat
8	Yes	No	No	Yes	Yes	Allergy
9	No	Yes	No	Yes	Yes	Cold
10	_{Yes} are thre	^{Yes} ee poss	∾ ible clas	Yes	Yes Throat, (Cold
There Binary	are thre	ee poss es like f these	ible clas	ves ses: Strep roduce rule o1d) Throat, (^{Cold} Cold, Aller

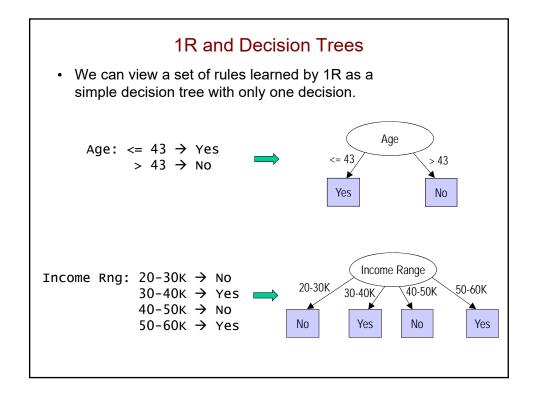


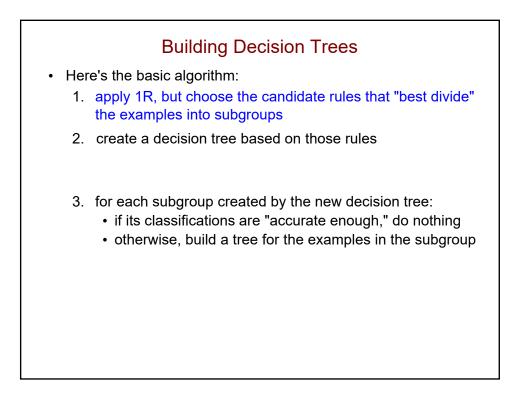


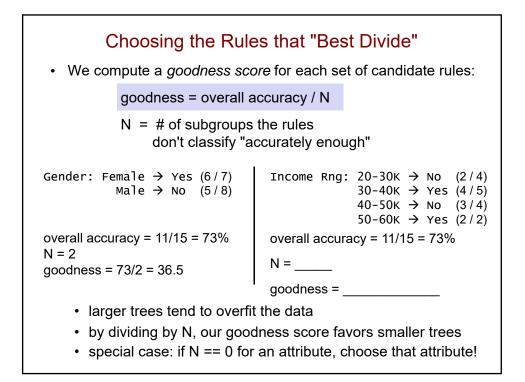
IncomeCredit CardLife InsuranceAgeGenderRangeInsurancePromotion45Male40–50KNoNo40Female30–40KNoYes42Male40–50KNoNo43Male30–40KYesYes38Female50–60KNoNo35Female20–30KNoNo35Male30–40KYesYes27Male20–30KNoNo43Male30–40KNoNo43Male30–40KNoNo43Female20–30KNoNo43Male30–40KNoYes43Female40–50KNoYes43Female40–50KNoYes39Female40–50KNoYes39Female40–50KNoYes55Male40–50KNoYes					class/ output attribute
40 Female 30–40K No Yes 42 Male 40–50K No No 43 Male 30–40K Yes Yes 38 Female 50–60K No Yes 55 Female 20–30K No No 35 Male 30–40K Yes Yes 27 Male 20–30K No No 43 Male 30–40K No No 43 Male 30–40K No No 43 Male 30–40K No No 41 Female 30–40K No Yes 43 Female 40–50K No Yes 29 Male 20–30K No Yes 39 Female 50–60K No Yes 55 Male 40–50K No No	Age	Gender			
42 Male 40–50K No No 43 Male 30–40K Yes Yes 38 Female 50–60K No Yes 55 Female 20–30K No No 35 Male 30–40K Yes Yes 27 Male 20–30K No No 43 Male 30–40K No No 43 Male 30–40K No No 43 Female 30–40K No Yes 43 Female 30–40K No Yes 43 Female 40–50K No Yes 29 Male 20–30K No Yes 39 Female 55 Male 40–50K No	45	Male	40–50K	No	No
43 Male 30-40K Yes Yes 38 Female 50-60K No Yes 55 Female 20-30K No No 35 Male 30-40K Yes Yes 27 Male 20-30K No No 43 Male 30-40K Yes Yes 41 Female 30-40K No Yes 43 Female 40-50K No Yes 29 Male 20-30K No Yes 39 Female 40-50K No Yes 55 Male 20-30K No Yes	40	Female	30–40K	No	Yes
38 Female 50–60K No Yes 55 Female 20–30K No No 35 Male 30–40K Yes Yes 27 Male 20–30K No No 43 Male 30–40K No No 41 Female 30–40K No Yes 43 Female 30–40K No Yes 29 Male 20–30K No Yes 39 Female 20–30K No Yes 39 Female 20–60K No Yes 55 Male 40–50K No No	42	Male	40–50K	No	No
55 Female 20–30K No No 35 Male 30–40K Yes Yes 27 Male 20–30K No No 43 Male 30–40K No No 41 Female 30–40K No Yes 43 Female 40–50K No Yes 43 Female 40–50K No Yes 39 Female 50–60K No Yes 55 Male 40–50K No No	43	Male	30–40K	Yes	Yes
35 Male 30-40K Yes Yes 27 Male 20-30K No No 43 Male 30-40K No No 41 Female 30-40K No Yes 43 Female 40-50K No Yes 43 Female 20-30K No Yes 29 Male 20-30K No Yes 39 Female 50-60K No Yes 55 Male 40-50K No No	38	Female	50–60K	No	Yes
27 Male 20–30K No No 43 Male 30–40K No No 41 Female 30–40K No Yes 43 Female 40–50K No Yes 43 Female 40–50K No Yes 29 Male 20–30K No Yes 39 Female 50–60K No Yes 55 Male 40–50K No No	55	Female	20–30K	No	No
43 Male 30-40K No No 41 Female 30-40K No Yes 43 Female 40-50K No Yes 43 Female 40-50K No Yes 29 Male 20-30K No Yes 39 Female 50-60K No Yes 55 Male 40-50K No No		Male		Yes	
41 Female 30–40K No Yes 43 Female 40–50K No Yes 29 Male 20–30K No Yes 39 Female 50–60K No Yes 55 Male 40–50K No No					
43 Female 40–50K No Yes 29 Male 20–30K No Yes 39 Female 50–60K No Yes 55 Male 40–50K No No					
29 Male 20–30K No Yes 39 Female 50–60K No Yes 55 Male 40–50K No No					
39 Female 50–60K No Yes 55 Male 40–50K No No					
55 Male 40–50K No No					
19 Female 20–30K Yes Yes	19	Female	20–30K	Yes	Yes

No → No (6 out of 12) 9/15 = 60% Income Range: 20-30K → No (2 out of 4) 30-40K → Yes (4 out of 5) 40-50K → No (3 out of 4) 50-60K → Yes (2 out of 2) Age: <= 43 → Yes (9 out of 12) overall accurace	Gender: Female → Yes Male → No	· · · ·	overall accuracy: 11/15 = 73%
$30-40 \times \rightarrow \text{Yes} (4 \text{ out of } 5) \qquad \text{overall accurac} \\ 40-50 \times \rightarrow \text{No} \qquad (3 \text{ out of } 4) \qquad 11/15 = 73\% \\ 50-60 \times \rightarrow \text{Yes} \qquad (2 \text{ out of } 2) \qquad \text{overall accurac} \\ \text{Age: } <= 43 \rightarrow \text{Yes} \qquad (9 \text{ out of } 12) \qquad \text{overall accurac} \\ \end{array}$		()	•
	30-40K → Yes 40-50K → No	(4 out of 5) (3 out of 4)	overall accuracy: 11/15 = 73%
	Age: <= 43 → Yes > 43 → No		

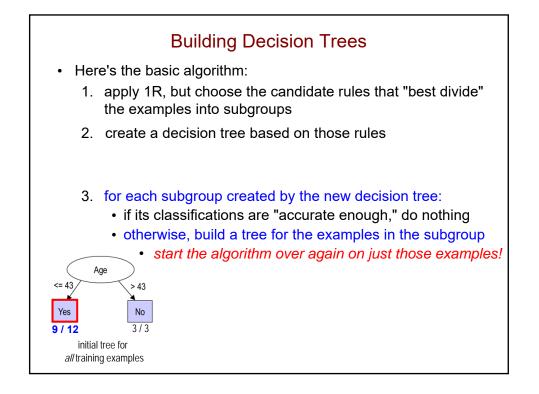
• When building a decision tree, we need to consider other factors.



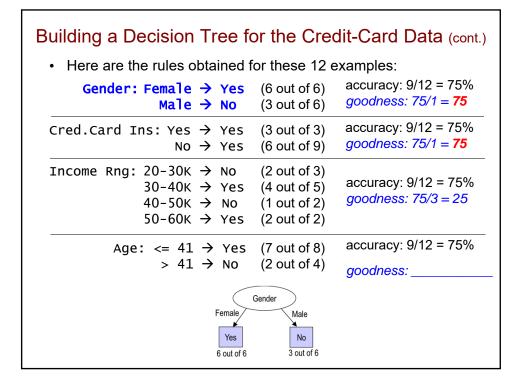


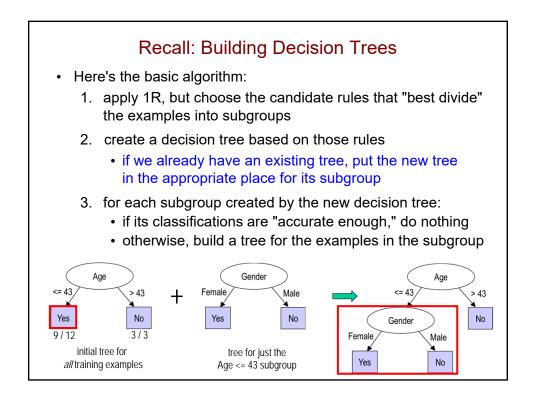


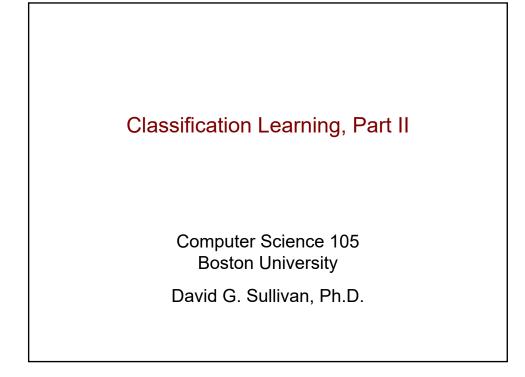
Here are tl	he rules v	veo	obtaine	ed for each at	ttribute using 1R:
Gender				(6 out of 7) (5 out of 8)	accuracy: 11/15 = 73% goodness: 73/2 = 36.5
Cred.Card I				(3 out of 3) (6 out of 12)	accuracy: 9/15 = 60% goodness: 60/1 = 60
Income Rng:	20-30K 30-40K 40-50K 50-60K	${\rightarrow}$	Yes No	(3 out of 4)	accuracy: 11/15 = 73% goodness: 73/3 = 24.3
Age	e: <= 43 > 43		Yes No	· · · ·	accuracy: 12/15 = 80% goodness: 80/1 = 80
			> (Age	

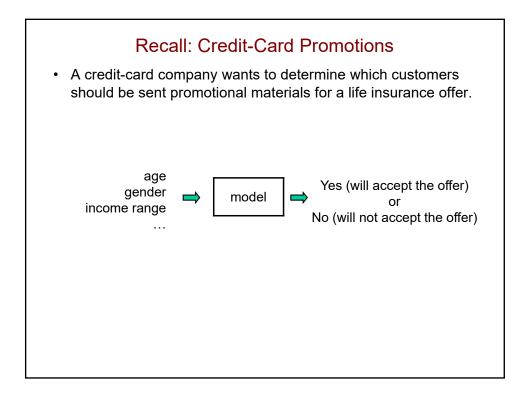


	Repeating the Algorithm on a Subgroup Here are the 12 examples in the Age <= 43 subgroup:								
Age 40 42 43 38 35 27 43 41 43 29 39 39 19	Gender Female Male Female Male Female Female Female Female	Income Range 30–40K 40–50K 30–40K 30–40K 30–40K 30–40K 40–50K 20–30K 50–60K 20–30K	Credit Card Insurance No Yes No Yes No No No No No Yes	Life Insurance Promotion Yes No Yes Yes No No Yes Yes Yes Yes Yes Yes					

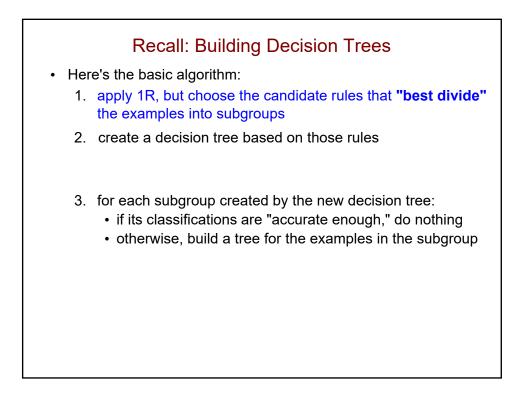








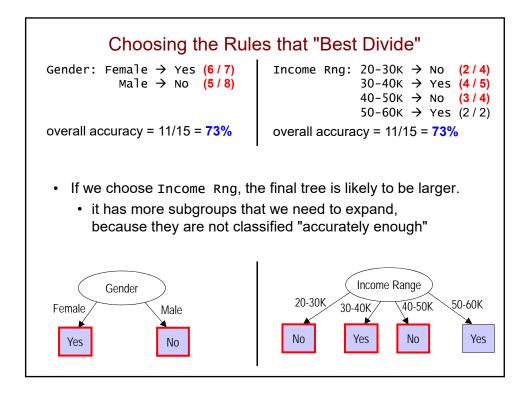
Age	Gender	Income Range	Credit Card Insurance*	class/ output attribute Life Insurance Promotion
45	Male	40–50K	No	No
40	Female	30–40K	No	Yes
42	Male	40-50K	No	No
43	Male	30–40K	Yes	Yes
38	Female	50–60K	No	Yes
55	Female	20–30K	No	No
35	Male	30–40K	Yes	Yes
27	Male	20–30K	No	No
43	Male	30–40K	No	No
41	Female	30–40K	No	Yes
43	Female	40–50K	No	Yes
29	Male	20–30K	No	Yes
39	Female	50–60K	No	Yes
55	Male	40–50K	No	No
19	Female	20–30K	Yes	Yes
specifying		e customer a	es/No attribut accepted a si	

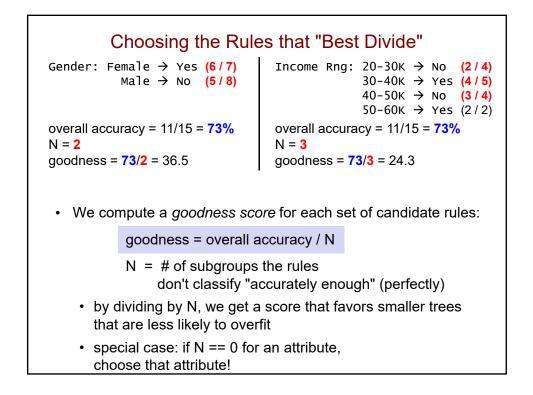


Choosing the Rules that "Best Divide"

Gender: Female → Yes (6/7) Male → No (5/8)	Income Rng: $20-30K \rightarrow No$ (2/4) $30-40K \rightarrow Yes$ (4/5)
	40-50K → No (3/4) 50-60K → Yes (2/2)
overall accuracy = 11/15 = 73%	overall accuracy = 11/15 = 73%

Precall: Overfitting In general, working too hard to match the training examples can lead to a model that: is overly complicated that doesn't generalize well This is known as *overfitting* the training data. The larger a decision tree gets, the more likely it is to overfit. its rules/decisions are based on smaller and smaller subgroups of training data

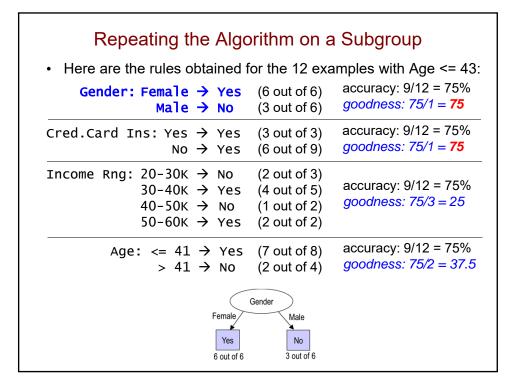


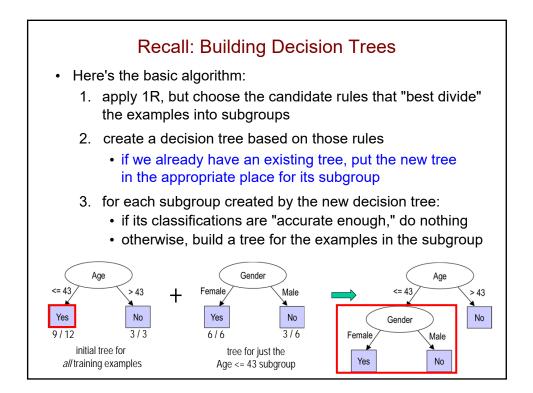


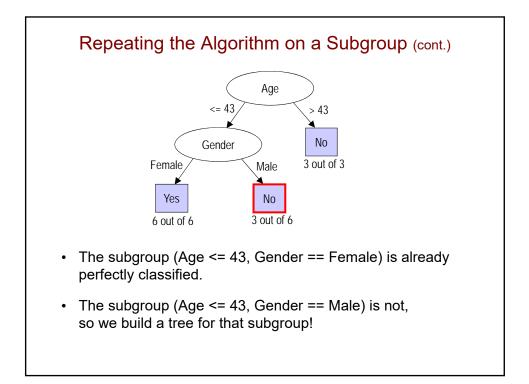
Consider these rules based on different training data. Which set of rules has the highest goodness score?

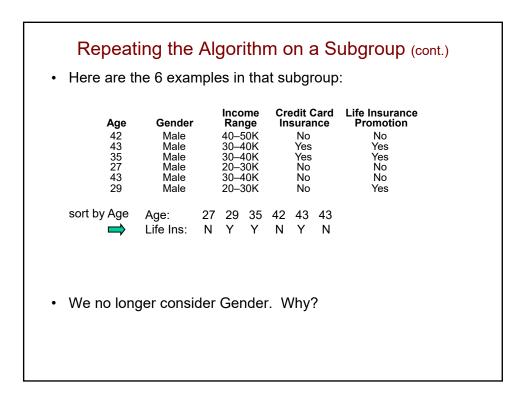
Gender: Female \rightarrow Yes (6 out of 7) Α. accuracy: 8/10 = 80% Male \rightarrow No (2 out of 3) goodness: 80/2 = 40 **B**. Cred.Card Ins: Yes \rightarrow Yes (3 out of 3) accuracy: 7/10 = 70% $NO \rightarrow NO$ (4 out of 7) goodness: ? C. Income Rng: 20-30K \rightarrow No (3 out of 3) accuracy: 8/10 = 80% 30-40K → Yes (2 out of 2) goodness: ? 40-50K → No* (1 out of 2) 50-60K → Yes (2 out of 3) D. more than one of the above

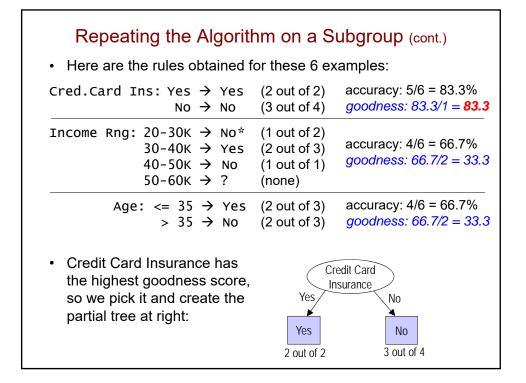
Here are the rules we obta	ned for each a	ttribute using 1R:
Gender:Female → Yes Male → No	6 out of 7) (5 out of 8)	accuracy: 11/15 = 73% goodness: 73/2 = 36.5
Cred.Card Ins: Yes → Yes No → No	· /	accuracy: 9/15 = 60% goodness: 60/1 = 60
Income Rng: 20-30K → No 30-40K → Yes 40-50K → No 50-60K → Yes	(4 out of 5) (3 out of 4)	accuracy: 11/15 = 73% goodness: 73/3 = 24.3
Age: <= 43 → Ye > 43 → No	,	accuracy: 12/15 = 80% goodness: 80/1 = 80
<= 43 Yes 9 out of	Age > 43	

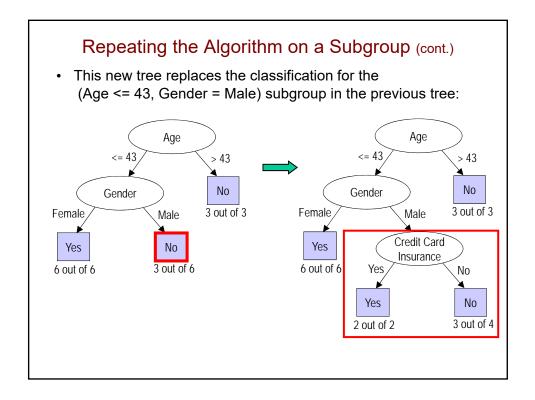


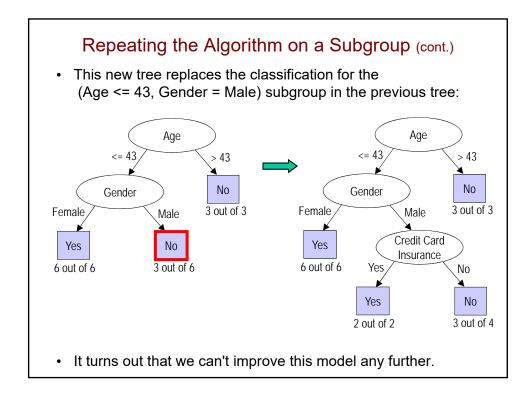








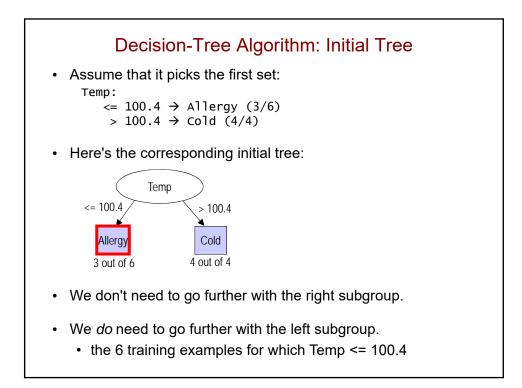


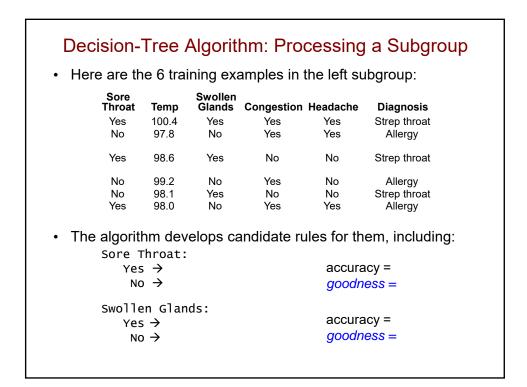


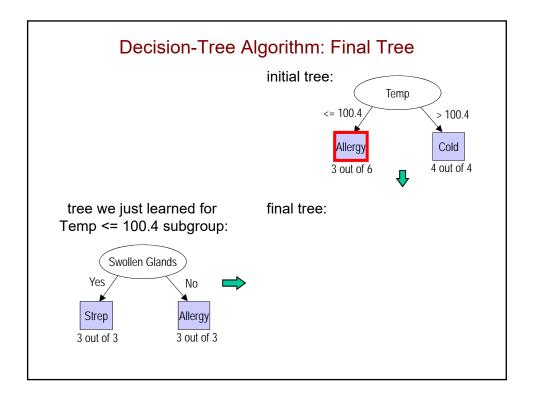
10 11				viously, bu ne person'	•	Yes/No) mperature
Patient ID#	' Sore Throat	Temp	' Swollen Glands	Congestion		Diagnosis
1	Yes	100.4	Yes	Yes	Yes	Strep throat
2	No	97.8	No	Yes	Yes	Allergy
3	Yes	101.2	No	Yes	No	Cold
4	Yes	98.6	Yes	No	No	Strep throat
5	No	102.0	No	Yes	No	Ċold
6	No	99.2	No	Yes	No	Allergy
7	No	98.1	Yes	No	No	Strep throat
8	Yes	98.0	No	Yes	Yes	Allergy
9	No	102.5	No	Yes	Yes	Cold
10	Yes	100.7	No	Yes	Yes	Cold
•	ain wan atients.	it to lea	rn a mo	del that all	ows us to	diagnose

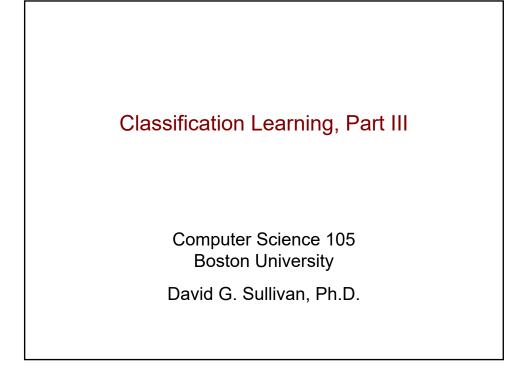
2No97.8NoYesYesAllergy3Yes101.2NoYesNoCold4Yes98.6YesNoNoStrep throat5No102.0NoYesNoCold6No99.2NoYesNoAllergy	Patient ID#	Sore Throat	Temp	Swollen Glands	Congestion	Headache	Diagnosis
4Yes98.6YesNoNoStrep throat5No102.0NoYesNoCold6No99.2NoYesNoAllergy7No98.1YesNoNoStrep throat8Yes98.0NoYesYesAllergy9No102.5NoYesYesCold10Yes100.7NoYesYesCold	1	Yes	100.4	Yes	Yes	Yes	Strep throat
4Yes98.6YesNoNoStrep throat5No102.0NoYesNoCold6No99.2NoYesNoAllergy7No98.1YesNoNoStrep throat8Yes98.0NoYesYesAllergy9No102.5NoYesYesCold10Yes100.7NoYesYesCold	2	No	97.8	No	Yes	Yes	Allergy
5No102.0NoYesNoČold6No99.2NoYesNoAllergy7No98.1YesNoNoStrep throat8Yes98.0NoYesYesAllergy9No102.5NoYesYesCold10Yes100.7NoYesYesCold	3	Yes	101.2	No	Yes	No	
6No99.2NoYesNoAllergy7No98.1YesNoNoStrep throat8Yes98.0NoYesYesAllergy9No102.5NoYesYesCold10Yes100.7NoYesYesCold		Yes	98.6	Yes	No	No	Strep throat
7No98.1YesNoNoStrep throat8Yes98.0NoYesYesAllergy9No102.5NoYesYesCold10Yes100.7NoYesYesCold	5	No	102.0	No	Yes	No	Ċold
8 Yes 98.0 No Yes Yes Allergy 9 No 102.5 No Yes Yes Cold 10 Yes 100.7 No Yes Yes Cold		No	99.2	No	Yes	No	Allergy
9 No 102.5 No Yes Yes Cold 10 Yes 100.7 No Yes Yes Cold							Strep throat
10 Yes 100.7 No Yes Yes Cold							
	-						•
ls Patient ID# numeric or nominal?	10	Yes	100.7	No	Yes	Yes	Cold
	ls Pati	ent ID#	numeri	ic or non	ninal?		
What is the accuracy of rules based on Patient ID#?	What i	s the ac	curacy	of rules	based on	Patient II	D#?
 we get one rule for each ID, which correctly classifies its example. This is overfitting! 		0			-		
In general, we should:		oral we	should	4.			

 Here are some of the candidate rul Temp: <= 100.4 → Allergy (3/6) 100.4 → cald (4(4)) 	les for the initial tree: accuracy = $7/10 = 70\%$ goodness = $70 / 1 = 7$
> 100.4 → Cold (4/4)	900011633 = 707 7 = 7
Swollen Glands: Yes → Strepthroat (3/3) No → Cold (4/7)	accuracy = 7/10 = 70% goodness = ?
Congestion:	
Yes \rightarrow Cold (4/8)	accuracy = 6/10 = 609
No \rightarrow Strep throat (2/2)	goodness = ?
• The other candidates are not as go	ood as these three.
• The decision-tree algorithm could p	pick which of these?
5	







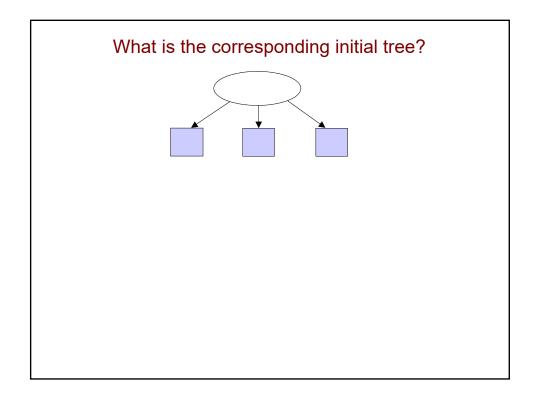


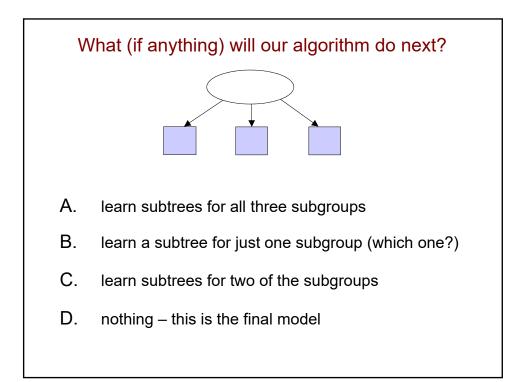
Color	Height	Stripes	Texture	Poisonous
Purple	Tall	Yes	Rough	Yes
Purple	Tall	Yes	Smooth	Yes
Red	Short	Yes	Hairy	No
Purple	Short	No	Smooth	No
Blue	Short	Yes	Hairy	Yes
Red	Tall	No	Rough	No
Blue	Tall	Yes	Smooth	Yes
Blue	Short	Yes	Rough	Yes
Red	Short	No	Smooth	No
Purple	Short	No	Hairy	Yes
Purple	Tall	No	Smooth	No

	What rul	es would	1R learn	based o	n Color?
	Color	Height	Stripes	Texture	Poisonous
	Purple	Tall	Yes	Rough	Yes
	Purple	Tall	Yes	Smooth	Yes
	Red	Short	Yes	Hairy	No
	Purple	Short	No	Smooth	No
	Blue	Short	Yes	Hairy	Yes
	Red	Tall	No	Rough	No
	Blue	Tall	Yes	Smooth	Yes
	Blue	Short	Yes	Rough	Yes
	Red	Short	No	Smooth	No
	Purple	Short	No	Hairy	Yes
	Purple	Tall	No	Smooth	No
A.		→ Yes → Yes → No	C.	Color: Purple Red Blue	 → Yes → No → Yes
В.	_	→ Yes → No → No	D.	more thar could be l	n one of these earned

	Wha	at would t	he final 1	R model	be?
	Color	Height	Stripes	Texture	Poisonous
	Purple	Tall	Yes	Rough	Yes
	Purple	Tall	Yes	Smooth	Yes
	Red	Short	Yes	Hairy	No
	Purple	Short	No	Smooth	No
	Blue	Short	Yes	Hairy	Yes
	Red	Tall	No	Rough	No
	Blue	Tall	Yes	Smooth	Yes
	Blue	Short	Yes	Rough	Yes
	Red	Short	No	Smooth	No
	Purple	Short	No	Hairy	Yes
	Purple	Tall	No	Smooth	No
A.	Red \rightarrow	Yes (3/5) No (3/3) Yes (3/3)	(9/11)	Texture: Rough Smooth Hairy	 → ? → ? → ?
В.		??	D.	Stripes: Yes → No →	
Ε.	two or more	of the model	ls are tied		

	Which set of rules would b our decision-tree alg	
A.	Color: Purple → Yes (3/5) Red → No (3/3) Blue → Yes (3/3)	accuracy: 9/11 = 82% goodness:
B.	Height: Tall → Yes (3/5) Short → No* (3/6)	accuracy: 6/11 = 55% goodness:
C.	Texture: Rough → Yes (2/3) Smooth → No (3/5) Hairy → Yes (2/3)	accuracy: 7/11 = 64% goodness:
D.	Stripes: Yes → Yes (5/6) No → No (4/5)	accuracy: 9/11 = 82% goodness:
E.	more than one of the above could	be chosen



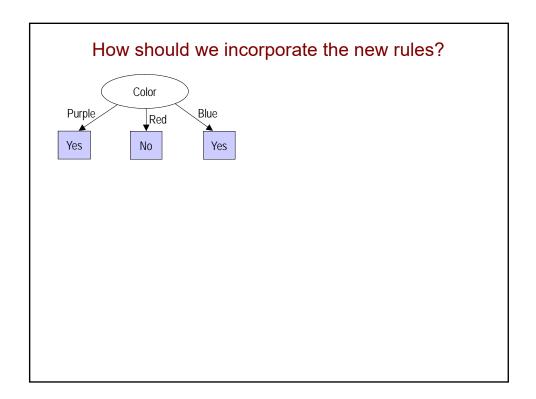


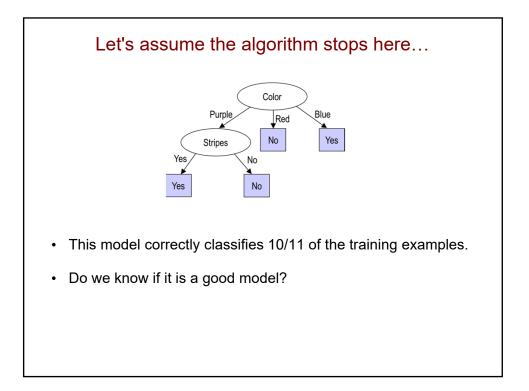
PurpleTallYesRoughYesPurpleTallYesSmoothYesPurpleTallYesSmoothYesRedShortYesHairyNoPurpleShortNoSmoothNoBlueShortYesHairyYesRedTallNoRoughNoBlueTallYesSmoothYesBlueShortYesRoughYesBlueShortYesRoughYesRedShortNoSmoothNoPurpleShortNoHairyYes		Height	Stripes	Texture	Poisonous
RedShortYesHairyNoPurpleShortNoSmoothNoBlueShortYesHairyYesRedTallNoRoughNoBlueTallYesSmoothYesBlueShortYesRoughYesBlueShortYesRoughYesBlueShortNoSmoothNo	Purple	•	•	Rough	
PurpleShortNoSmoothNoBlueShortYesHairyYesRedTallNoRoughNoBlueTallYesSmoothYesBlueShortYesRoughYesBlueShortYesRoughYesRedShortNoSmoothNo	Purple	Tall	Yes	Smooth	Yes
BlueShortYesHairyYesRedTallNoRoughNoBlueTallYesSmoothYesBlueShortYesRoughYesRedShortNoSmoothNo	Red	Short	Yes	Hairy	No
RedTallNoRoughNoBlueTallYesSmoothYesBlueShortYesRoughYesRedShortNoSmoothNo	Purple	Short	No	Smooth	No
BlueTallYesSmoothYesBlueShortYesRoughYesRedShortNoSmoothNo	Blue	Short	Yes	Hairy	Yes
BlueShortYesRoughYesRedShortNoSmoothNo	Red	Tall	No	Rough	No
RedShortNoSmoothNo	Blue	Tall	Yes	Smooth	Yes
	Blue	Short	Yes	Rough	Yes
Purple Short No Hairy Yes	Red	Short	No	Smooth	No
	Purple	Short	No	Hairy	Yes
Purple Tall No Smooth No	Purple	Tall	No	Smooth	No

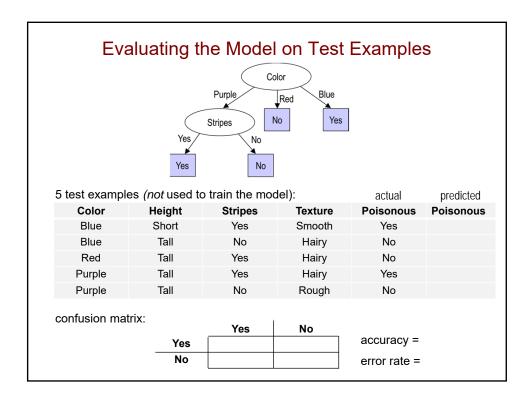
Color	Height	Stripes	Texture	Poisonous
Purple	Tall	Yes	Rough	Yes
Purple	Tall	Yes	Smooth	Yes
Purple	Short	No	Smooth	No
Purple	Short	No	Hairy	Yes
Purple	Tall	No	Smooth	No

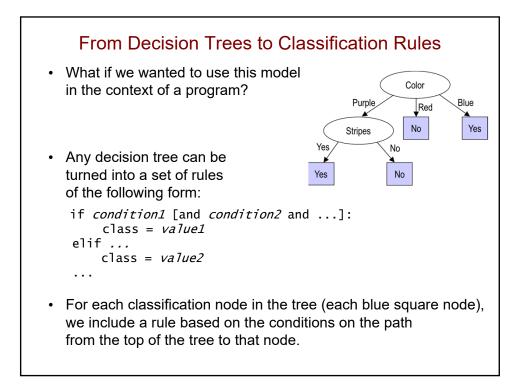
Purple Purple Purple	Height Tall Tall	Stripes Yes	Rough	
Purple	Tall		rtougn	Yes
		Yes	Smooth	Yes
	Short	No	Smooth	No
Purple	Short	No	Hairy	Yes
Purple	Tall	No	Smooth	No

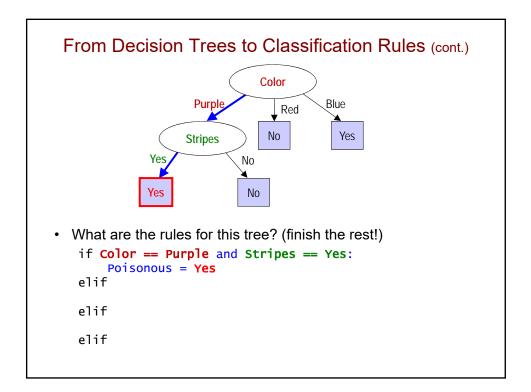
	What rul	es will be	selected	for this s	ubgroup?
•	just the Purp	ole subgroup):		
	Color	Height	Stripes	Texture	Poisonous
	Purple	Tall	Yes	Rough	Yes
	Purple	Tall	Yes	Smooth	Yes
	Purple	Short	No	Smooth	No
	Purple	Short	No	Hairy	Yes
	Purple	Tall	No	Smooth	No
Α.	Height: Tall → Short →		accurad goodne		Generate all three sets of candidate rules,
В.	Texture: Rough → Smooth → Hairy →		accurae goodne		and then determine which will be selected!
C.	Stripes: Yes → No →		accurae goodne		D. more than one could be selected

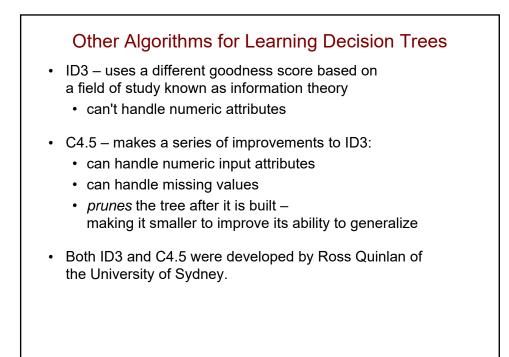


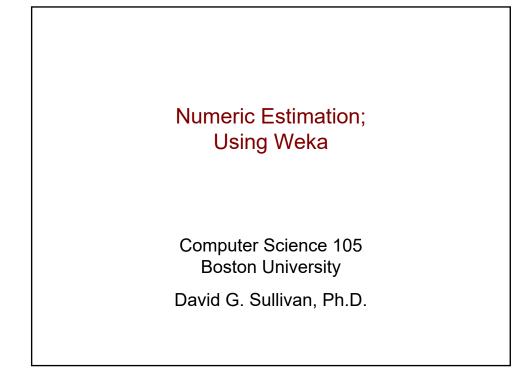


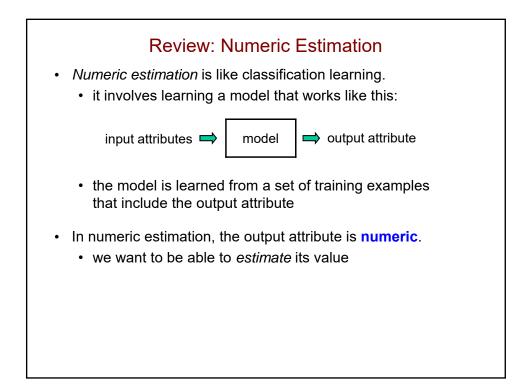


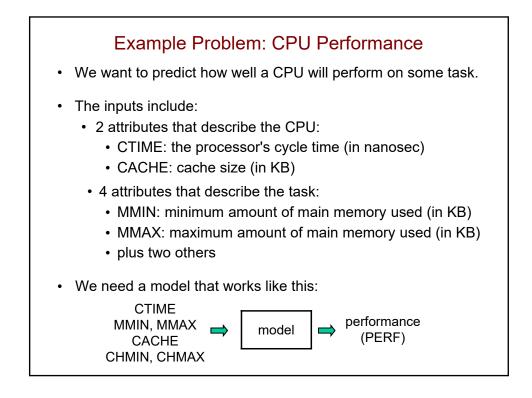




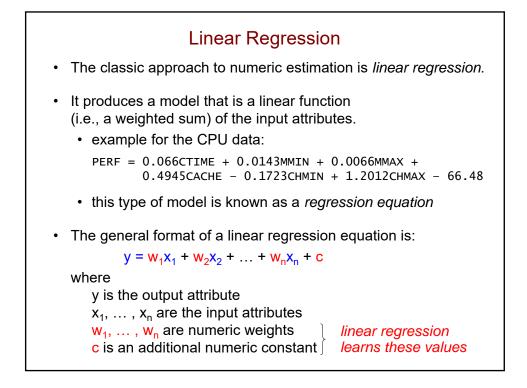


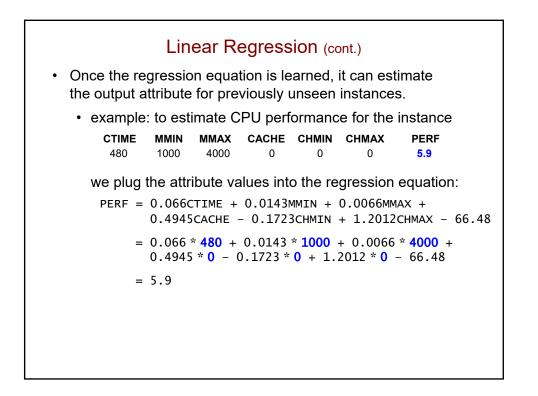


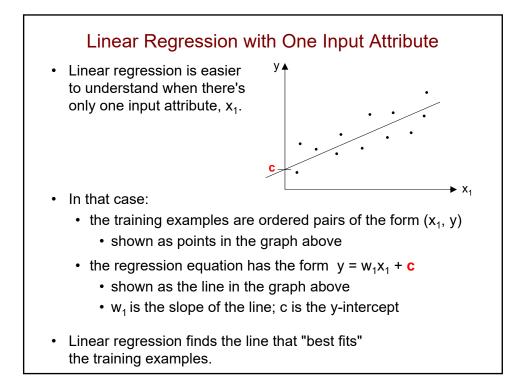


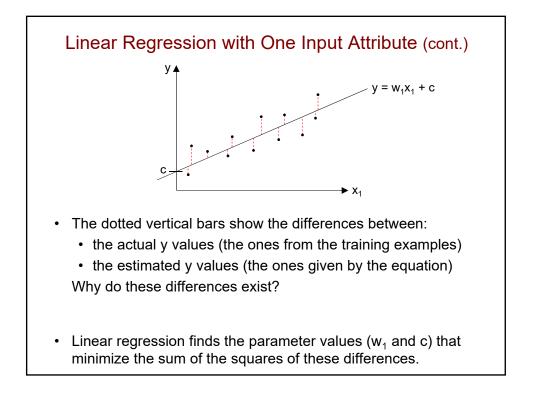


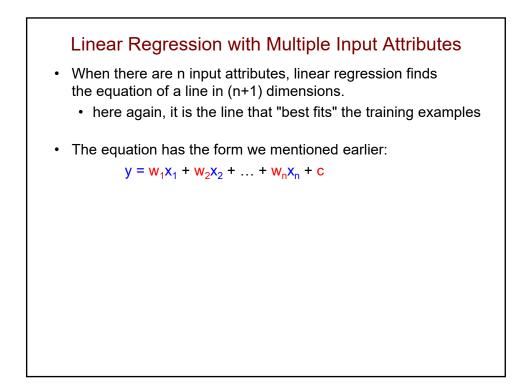
	Exan	nple F	Proble	m: CPI	U Per	formar	ICE (cont.)
•	There are	e 209 tr	aining e	example	s. Here	e are five	e of them:
							class/
	input attri	ibutes:					output attribute
	CTIME	MMIN	MMAX	CACHE	CHMIN	CHMAX	PERF
	125	256	6000	256	16	128	198
	29 29	8000 8000	32000 32000	32 32	8	32 32	269 172
	125	2000	8000	0	8 2	14	52
	480	512	8000	32	0	0	67

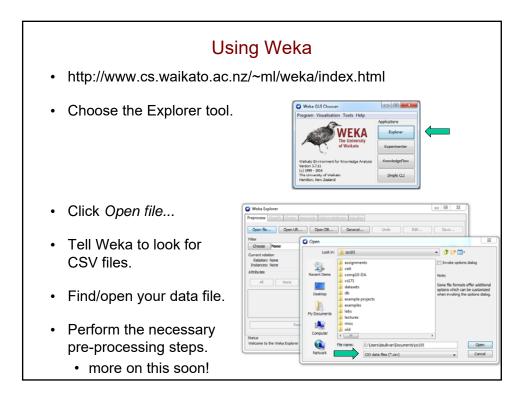


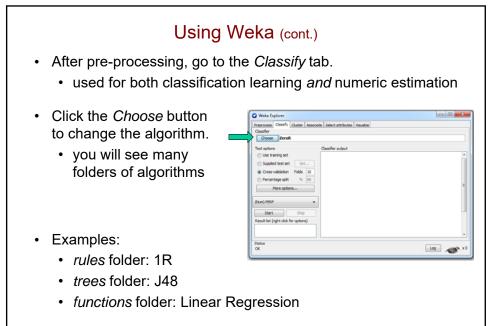




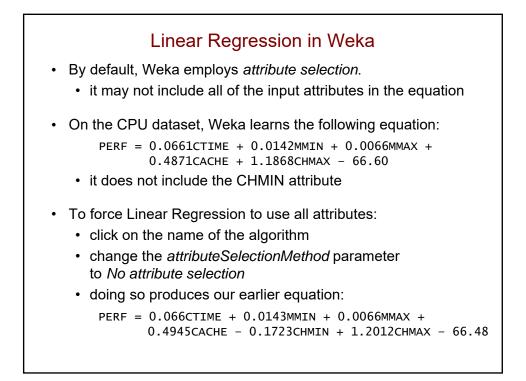








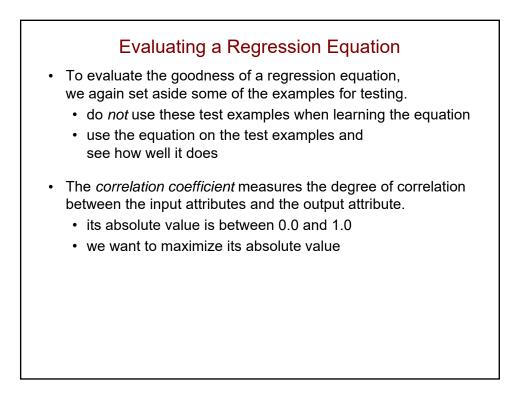
• Feel free to try algorithms that we haven't discussed in lecture!



The Coefficients in Linear Regression

PERF = 0.066CTIME + 0.0143MMIN + 0.0066MMAX + 0.4945CACHE - 0.1723CHMIN + 1.2012CHMAX - 66.48

- Notes about the coefficients:
 - · what do the signs of the coefficients mean?
 - what about their magnitudes?

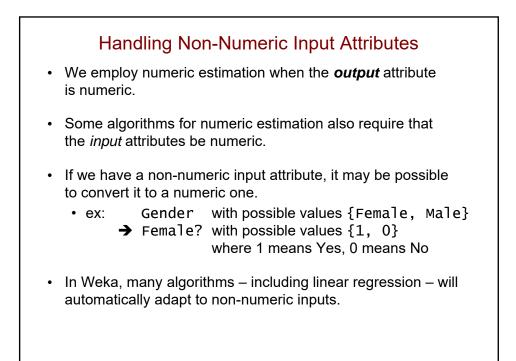


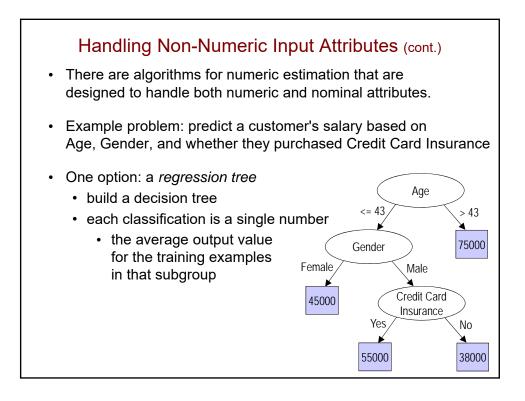
Simple Linear Regression

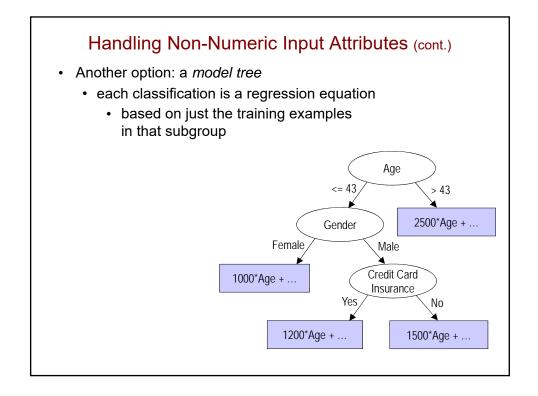
- This algorithm in Weka creates a regression equation that uses only one of the input attributes.
 - · even when there are multiple inputs
 - like 1R, but for numeric estimation
- We can use it as a baseline.
 - · determine the correlation coefficient of its model
 - if a more complex model has a lower correlation coefficient, don't use it!
 - (we can use 1R in a similar way when doing classification learning)
- It also gives insight into which of the input attributes has the largest impact on the output.

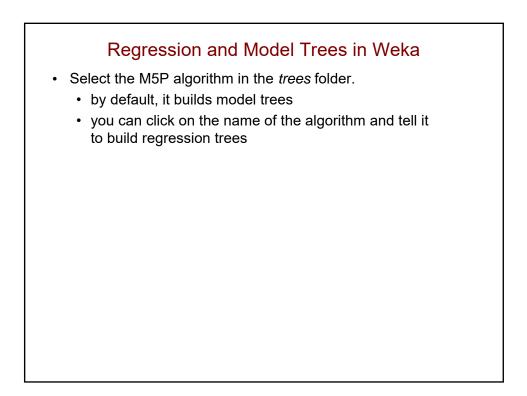


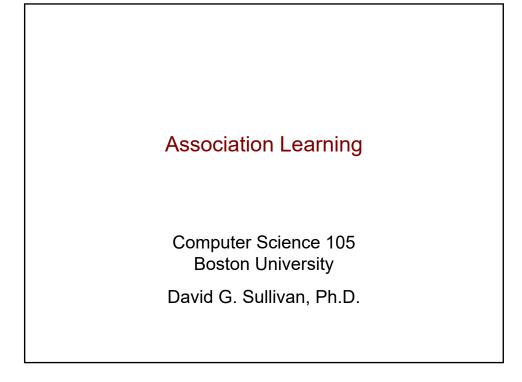
- A. Numeric estimation produces a model that predicts the value of a single output attribute.
- B. The model produced by numeric estimation does not need to use all of the input attributes.
- C. In order to perform numeric estimation, the input attributes must be numeric.
- D. A numeric-estimation model is learned from a set of training examples that include values for the output attribute.

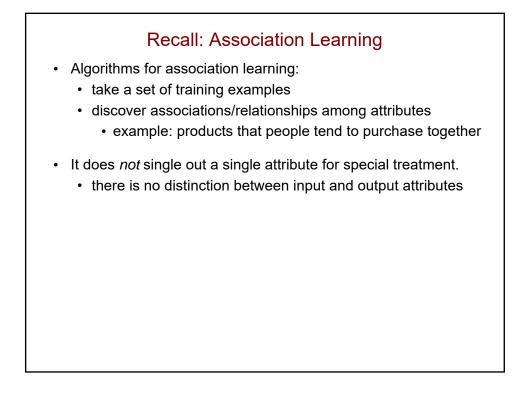




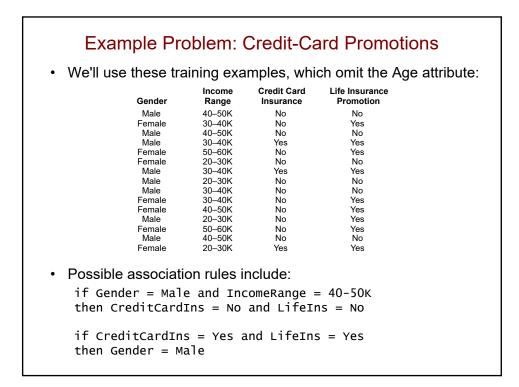


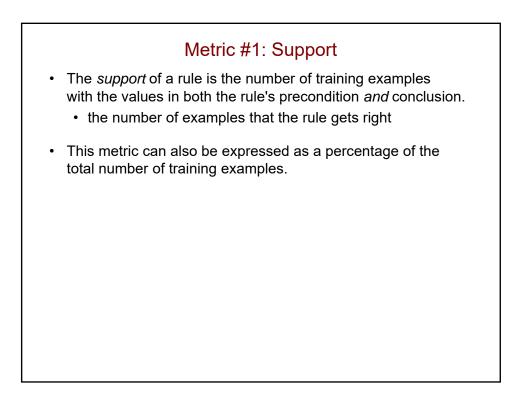


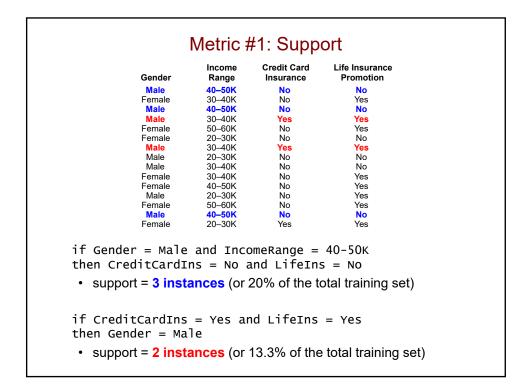




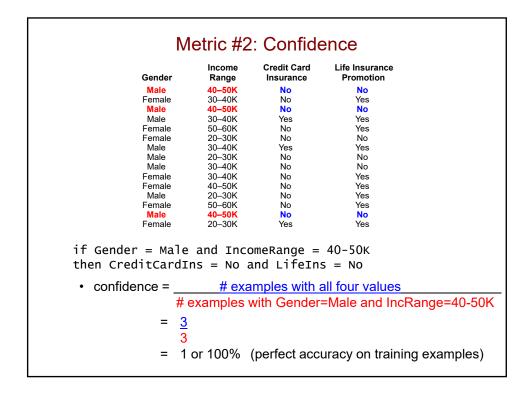
The Converse of a Rule
 The converse of a rule is obtained by swapping the precondition and conclusion.
 example: here's one rule: if PurchaseDiapers = Yes then PurchaseBeer = Yes
its converse is: if PurchaseBeer = Yes then PurchaseDiapers = Yes
 The converse of a rule is <u>not</u> necessarily true. example: this rule is true: if name = 'Perry Sullivan' then yearBorn = 2000
its converse is not! if yearBorn = 2000 then name = 'Perry Sullivan'







	Metric #2: Confidence
	e of a rule provides a measure of a rule's accuracy predicts the values in the conclusion.
	question: if the precondition of the rule holds, that the conclusion also holds?
Here's the forn	nula:
confidence =	# examples with the values in the support the support
	# examples with the values in just the precondition

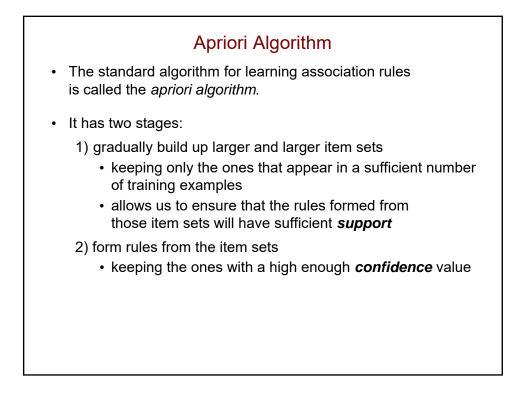


Gender	Income Range	Credit Card Insurance	Life Insurance Promotion
Male	40–50K	No	No
Female	30–40K	No	Yes
Male	40–50K	No	No
Male	30–40K	Yes	Yes
Female	50–60K	No	Yes
Female	20–30K	No	No
Male	30–40K	Yes	Yes
Male	20–30K	No	No
Male	30–40K	No	No
Female	30–40K	No	Yes
Female	40–50K	No	Yes
Male	20–30K	No	Yes
Female	50–60K	No	Yes
Male	40–50K	No	No
Female	20–30K	Yes	Yes
if CreditCardIns then Gender = Ma		nd LifeIns	= Yes
 confidence = 	# exa	moles with	all three values
# (examples	with CreditC	ardIns=Yes, LifeIns=Ye
= 2			
= <u>2</u> 3			
0			

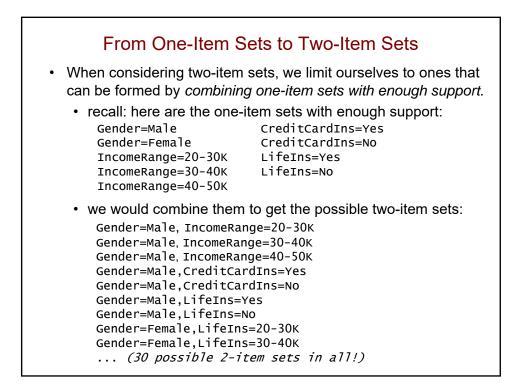
Gender	Income Range	Credit Card Insurance	Life Insurance Promotion
Male	40–50K	No	No
Female	30–40K	No	Yes
Male	40–50K	No	No
Male	30–40K	Yes	Yes
Female	50–60K	No	Yes
Female	20–30K	No	No
Male	30–40K	Yes	Yes
Male	20–30K	No	No
Male	30–40K	No	No
Female	30–40K	No	Yes
Female	40–50K	No	Yes
Male	20–30K	No	Yes
Female Male	50-60K	No No	Yes
Female	40–50K 20–30K	Yes	No Yes
remaie	20-30K	tes	tes
F LifeIns = Yes nen Gender = Fen support =	nale and	CreditCar	dIns = No
confidence =			

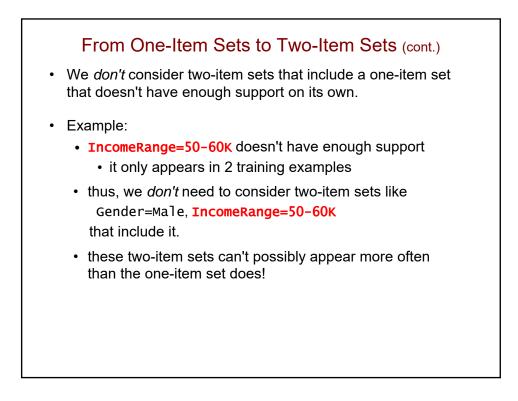
Learning Association Rules
 For a given dataset, there are a large number of association rules that could be learned.
example:
if CreditCardIns = Yes and LifeIns = Yes and IncomeRange = 20-30K then Gender = Female
has a confidence of 100%, but it is only based on
a single example (i.e., its support = 1)
 To cut down the number of rules that we consider, we limit ourselves to ones with sufficient support.
 Of these rules, we keep the most accurate ones – the ones with a confidence value that is above some minimum value.

	Income	Credit Card	Life Insurance
Gender	Range	Insurance	Promotion
Male	40–50K	No	No
Female	30–40K	No	Yes
Male	40–50K	No	No
Male	30–40K	Yes	Yes
Female	50–60K	No	Yes
Female	20–30K	No	No
Male	30–40K	Yes	Yes
Male	20–30K	No	No
Male	30–40K	No	No
Female	30–40K	No	Yes
Female	40–50K	No	Yes
Male	20–30K	No	Yes
Female Male	50–60K 40–50K	No No	Yes No
Female	20–30K	Yes	Yes
			alues that appears
one or more traini	ng examp	nes.	
 example: the it 	em set c	reditCardI	ns=Yes, LifeIns=Ye
appears in 3 tr			·
	•	•	
	d to forma	two difforos	nt rules with support

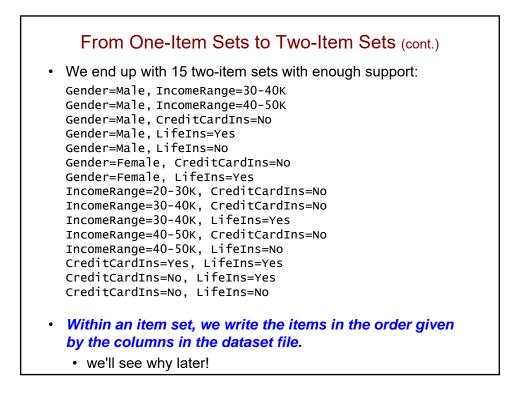


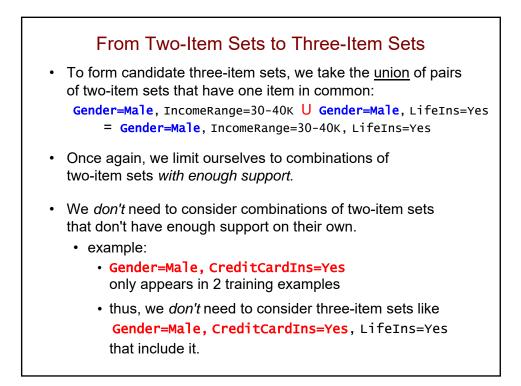
First S	tage: E	Building It	tem Sets
Gender	Income Range	Credit Card Insurance	Life Insurance Promotion
Male	40–50K	No	No
Female	30–40K	No	Yes
Male	40–50K	No	No
Male	30–40K	Yes	Yes
Female	50-60K	No	Yes
Female	20–30K	No	No
Male	30–40K	Yes	Yes
Male	20–30K	No	No
Male	30–40K	No	No
Female	30–40K	No	Yes
Female	40–50K	No	Yes
Male	20–30K	No	Yes
Female	50-60K	No	Yes
Male Female	40–50K 20–30K	No	No
Female	20-30K	Yes	Yes
Assume we want it			·
We get 9 one-item	sets that	t meet this of	criterion:
Gender=Male		CreditCa	ardīns=Yes
		0.00.000	
Gender=Female	2	Creattea	ardIns=No
IncomeRange=2	20-30K	LifeIns=	=Yes
IncomeRange=3			-No
		LITETU3-	-110
IncomeRange=4	10-50K		
 everything but In 	comeRan	ge=50-60к , \	which is in only 2 examples

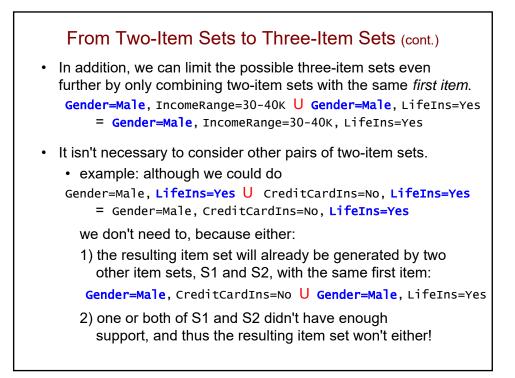


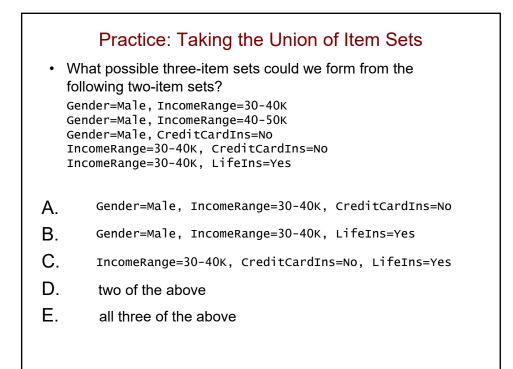


	r Range	Insurance	Promotion
Male	40–50K	No	No
Female	30–40K	No	Yes
Male	40–50K	No	No
Male	30–40K	Yes	Yes
Female	50–60K	No	Yes
Female	20–30K	No	No
Male	30–40K	Yes	Yes
Male	20–30K	No	No
Male	30–40K	No	No
Female	30–40K	No	Yes
Female		No	Yes
Male	20–30K	No	Yes
Female	00 0011	No	Yes
Male	40–50K	No	No
Female	20–30K	Yes	Yes
with enough suexample: we	pport.	Gender=Ma	d only keep the ones le, IncomeRange=20-30









 For our dataset, there are 56 three 	ee-item sets in all:
Gender=Female, Income=20-30K, CredCardIns=Yes	Gender=Male, Income=20-30K, CredCardIns=Yes
Gender=Female, Income=20-30K, CredCardIns=No	Gender=Male, Income=20-30K, CredCardIns=No
Gender=Female, Income=20-30K, LifeIns=Yes	Gender=Male, Income=20-30K, LifeIns=Yes
Gender=Female, Income=20-30K, LifeIns=No	Gender=Male, Income=20-30K, LifeIns=No
Gender=Female, Income=30-40K, CredCardIns=Yes	Gender=Male, Income=30-40K, CredCardIns=Yes
Gender=Female, Income=30-40K, CredCardIns=No	Gender=Male, Income=30-40K, CredCardIns=No
Gender=Female, Income=30-40K, LifeIns=Yes	Gender=Male, Income=30-40K, LifeIns=Yes
Gender=Female, Income=30-40K, LifeIns=No	Gender=Male, Income=30-40K, LifeIns=No
Gender=Female, Income=40-50K, CredCardIns=Yes	Gender=Male, Income=40-50K, CredCardIns=Yes
Gender=Female, Income=40-50K, CredCardIns=No	Gender=Male, Income=40-50K, CredCardIns=No
Gender=Female, Income=40-50K, LifeIns=Yes	Gender=Male, Income=40-50K, LifeIns=Yes
Gender=Female, Income=40-50K, LifeIns=No	Gender=Male, Income=40-50K, LifeIns=No
Gender=Female, Income=50-60K, CredCardIns=Yes	Gender=Male, Income=50-60K, CredCardIns=Yes
Gender=Female, Income=50-60K, CredCardIns=No	Gender=Male, Income=50-60K, CredCardIns=No
Gender=Female, Income=50-60K, LifeIns=Yes	Gender=Male, Income=50-60K, LifeIns=Yes
Gender=Female, Income=50-60K, LifeIns=No Gender=Female. CredCardIns=Yes. LifeIns=Yes	Gender=Male, Income=50-60K, LifeIns=No Gender=Male. CredCardIns=Yes. LifeIns=Yes
<pre>Gender=Female, CredCardIns=Yes, LifeIns=No Gender=Female, CredCardIns=No, LifeIns=Yes</pre>	Gender=Male, CredCardIns=Yes, LifeIns=No Gender=Male, CredCardIns=No, LifeIns=Yes
Gender=Female, CredCardIns=NO, LifeIns=Yes	Gender=Male, CredCardIns=No, LifeIns=Yes Gender=Male, CredCardIns=No, LifeIns=No
Income=20-30K, CredCardIns=Yes, LifeIns=Yes	Income=40-50K, CredCardIns=Yes, LifeIns=Yes
<pre>Income=20-30K, CredCardIns=Yes, LifeIns=No Income=20-30K, CredCardIns=No, LifeIns=Yes</pre>	<pre>Income=40-50K, CredCardIns=Yes, LifeIns=No Income=40-50K, CredCardIns=No, LifeIns=Yes</pre>
Income=20-30K, CredCardIns=No, LifeIns=Yes	Income=40-50K, CredCardIns=No, LifeIns=Yes
Income=30-40K. CredCardIns=No, LifeIns=No	Income=50-60K, CredCardIns=No, LifeIns=No
Income=30-40K, CredCardIns=Yes, LifeIns=No	Income=50-60K, CredCardIns=Yes, LifeIns=Yes
Income=30-40K, CredCardIns=Yes, LifeIns=No	Income=50-60K, CredCardIns=Yes, LifeIns=Yes
Income=30-40K, CredCardIns=No, LifeIns=No	Income=50-60K, CredCardIns=No, LifeIns=No

	Gender	Income Range	Credit Card Insurance	Life Insurance Promotion
	Male	40–50K	No	No
	Female	30–40K	No	Yes
	Male	40–50K	No	No
	Male	30–40K	Yes	Yes
	Female	50–60K	No	Yes
	Female	20–30K	No	No
	Male	30–40K	Yes	Yes
	Male	20–30K	No	No
	Male	30–40K	No	No
	Female	30–40K	No	Yes
	Female	40–50K	No	Yes
	Male	20–30K	No	Yes
	Female	50-60K	No	Yes
	Male	40–50K	No	No
	Female	20–30K	Yes	Yes
	•		e-items sets east 3 exam	s, only 5 have suffici ples:
Gender=M Gender=M Gender=M	Male, Inco Male, Creo Female, Cr	omeRange= litCardIn reditCard	40-50K, Lit s=No, Life Ins=No, Lit	Ins=No

From Three-Item Sets to Four-Item Sets

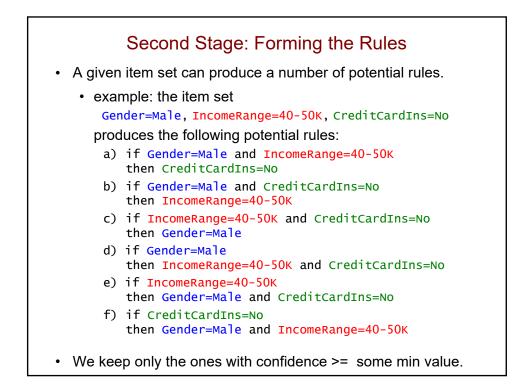
```
Gender=Male, IncomeRange=40-50K, CreditCardIns=No
Gender=Male, IncomeRange=40-50K, LifeIns=No
Gender=Male, CreditCardIns=No, LifeIns=No
Gender=Female, CreditCardIns=No, LifeIns=Yes
IncomeRange=40-50K, CreditCardIns=No, LifeIns=No
```

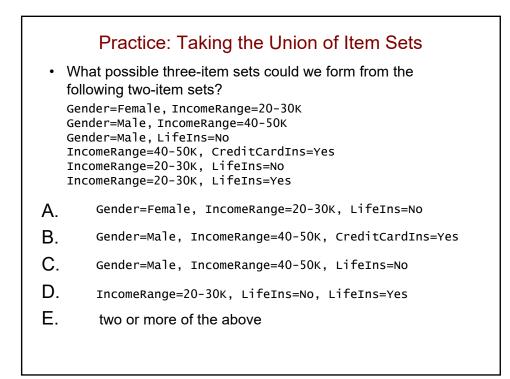
- To form potential four-item sets, we take the union of pairs of surviving three-item sets with the same first two items.
 - more generally, to form n-item sets, we take the union of pairs of (n – 1)-item sets with the same first n – 2 items
- We get only one potential four-item set:
 Gender=Male, IncomeRange=40-50K, CreditCardIns=No, LifeIns=No and it has enough support.
- There can't be any five-item sets (because there are only four attributes), so we're done building item sets!

```
Results of the First Stage

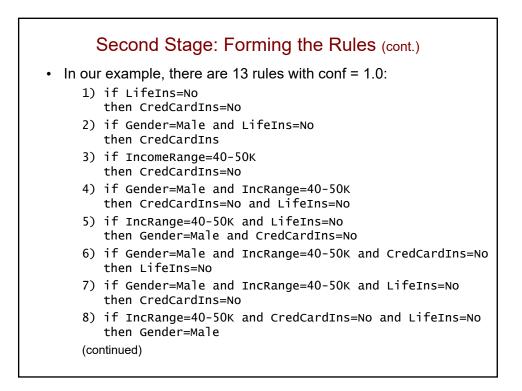
    Here are all item sets with two or more items and support >= 3:

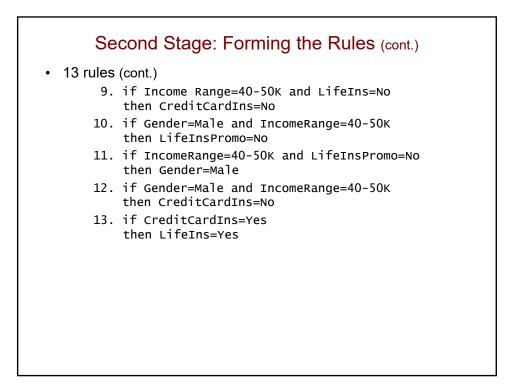
                                Gender=Female, CreditCardIns=No
Gender=Male, IncRange=30-40K
Gender=Male, IncRange=40-50K
                                Gender=Female, LifeIns=Yes
Gender=Male, CredCardIns=No
                                CredCardIns=Yes, LifeIns=Yes
Gender=Male, LifeIns=Yes
                                CredCardIns=No, LifeIns=Yes
Gender=Male, LifeIns=No
                                CredCardIns=No, LifeIns=No
IncRange=20-30K, CredCardIns=No
IncRange=30-40K, CredCardIns=No
IncRange=30-40K, LifeIns=Yes
IncRange=40-50K, CredCardIns=No
IncRange=40-50K, LifeIns=No
Gender=Male, IncRange=40-50K, CredCardIns=No
Gender=Male, IncRange=40-50K, LifeIns=No
Gender=Male, CredCardIns=No, LifeIns=No
Gender=Female, CredCardIns=No, LifeIns=Yes
IncRange=40-50K, CredCardIns=No, LifeIns=No
Gender=Male, IncRange=40-50K, CredCardIns=No, LifeIns=No
```

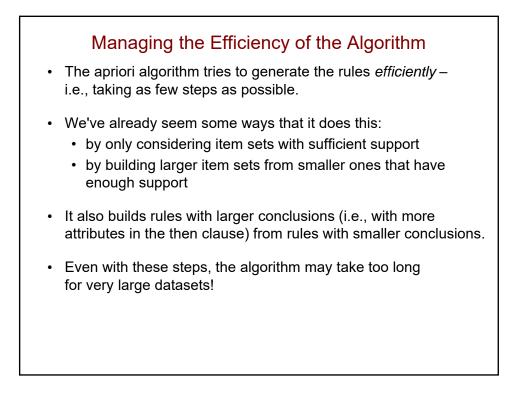




	Gender	Income Range	Credit Card Insurance	Life Insurance Promotion	
	Male Female Male Male	40–50K 30–40K 40–50K 30–40K	No No No Yes	No Yes No Yes	
	Female Female Male Male	50–60K 20–30K 30–40K 20–30K	No No Yes No	Yes No Yes No	Assume that we require a minimum confidence of 1.0
	Male Female Female Male Female Female	30–40K 30–40K 40–50K 20–30K 50–60K 40–50K 20–30K	No No No No No Yes	No Yes Yes Yes No Yes	
A.	if Gender=Male a then CreditCard	c support confide			
В.	if Gender=Male a then IncomeRange		support = ? confidence = ?		
C.	both rules would b	e kept			
D.	neither rule would	be kept			

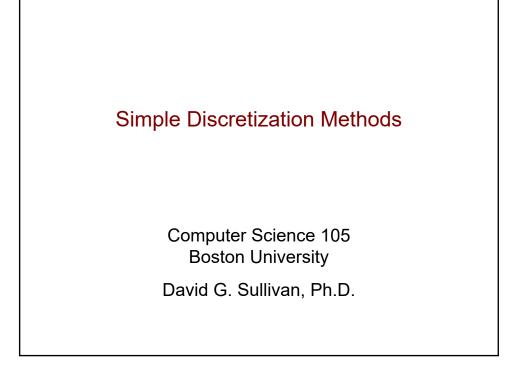


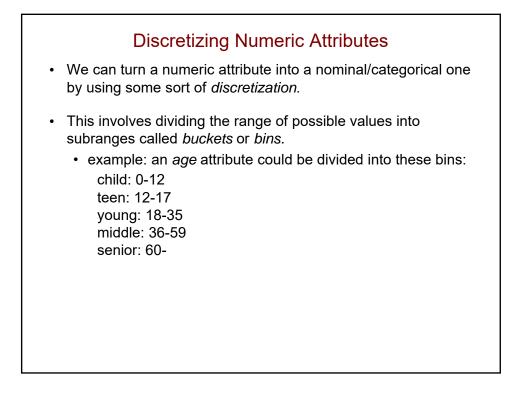


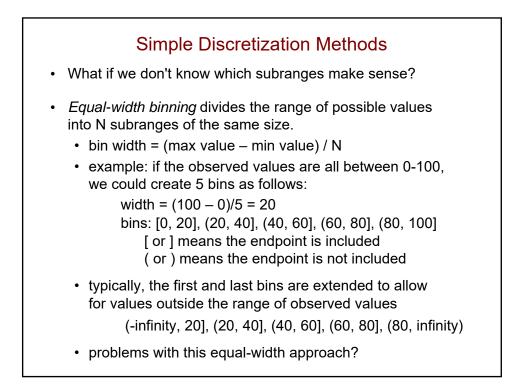


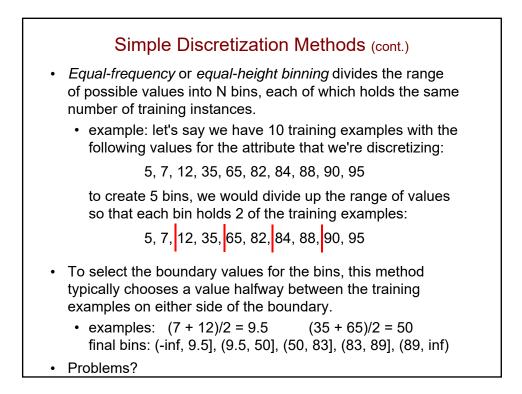
Managing the Efficiency of the Algorithm (cont.)

- To improve the efficiency even further, we can:
 - specify a large initial support value
 - the larger the support value, the sooner the first phase will finish
 - have the algorithm gradually decrease this support value and rerun the algorithm until it has generated enough rules
 - the *delta* parameter in Weka specifies how much the support should be decreased each time









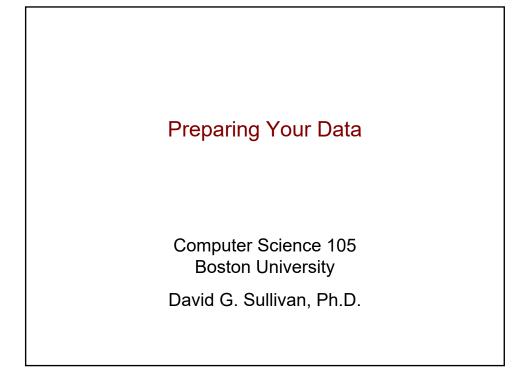
Discretization Example

• Let's say we have 8 training examples with the following values for Age:

17, 23, 35, 41, 51, 58, 70, 89

We want to discretize Age into 4 bins.

Which bins would be given by equal-height disc.? Let's say we have 8 training examples with the following values for Age: 17, 23, 35, 41, 51, 58, 70, 89 We want to discretize Age into 4 bins. A. (-infinity, 29], (29, 46], (46, 64], (64, infinity) B. [17, 29], (29, 46], (46, 64], (64, 89] C. (-infinity, 35], (35, 53], (53, 71], (71, infinity) D. [17, 35], (35, 53], (53, 71], (71, 89]



The Data Mining Process
 Key steps: assemble the data in the format needed for data mining typically a text file perform the data mining interpret/evaluate the results apply the results

Denormalization

- The data for a given entity (e.g., a customer) may be:
 - spread over multiple tables
 - spread over multiple rows within a given table
- To prepare for data mining, we need to *denormalize* the data.
 - multiple rows for a given entity → a single row

Denormalization										
 Example: finding associations between courses students take. 										
Student		Course				Enrolled				
id	name	name	start_time	end_time			course_name	credit_status		
12345678	Jill Jones	CS 105	13:00:00	14:00:00		12345678	CS 105	ugrad		
25252525	Alan Turing		09:30:00			25252525	CS 111	ugrad		
33566891	Audrey Chu		11:00:00			45678900	CS 460	grad		
45678900	Jose Delgado		16:00:00			33566891	CS 105	non-credit		
	Count Dracula		12:00:00 14:30:00			45678900	CS 510	grad		

Transforming the Data

- We may also need to reformat or transform the data.
 - we can use a Python program to do the reformatting!
- One reason for transforming the data: many machine-learning algorithms can only handle certain types of data.
 - some algorithms only work with *nominal* attributes attributes with a specified set of possible values
 - examples: {yes, no}
 - {strep throat, cold, allergy}
 - other algorithms only work with numeric attributes

Recall: Simple Discretization Methods

- · We've discussed two methods for discretization.
- *Equal-width binning* divides the range of possible values into N subranges of the same size.
- Equal-frequency or equal-height binning divides the range of possible values into N bins, each of which holds the same number of training instances.

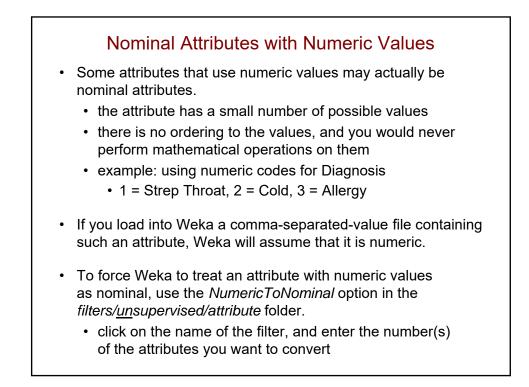
Discretization Example

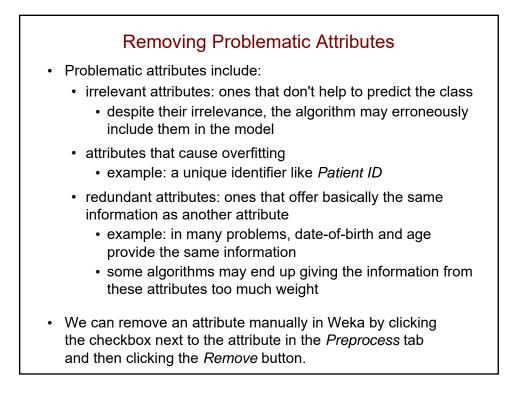
• Let's say we have 8 training examples with the following values for Age:

17, 23, 35, 41, 51, 58, 70, 89

We want to discretize Age into 4 bins.

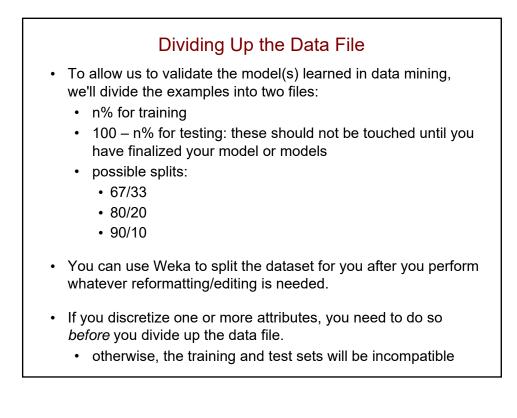
Discretization in Weka In Weka, you can discretize an attribute by applying the appropriate filter to it. After loading in the dataset in the *Preprocess* tab, click the *Choose* button in the *Filter* portion of the tab. For equal-width or equal-height, you choose the *Discretize* option in the *filters/<u>un</u>supervised/attribute* folder. by default, it uses equal-width binning to use equal-frequency binning instead, click on the name of the filter and set *useEqualFrequency* to *True*Another option: *Discretize* in *filters/supervised/attribute* folder attempts to learn *meaningful* cutoffs, based on your output

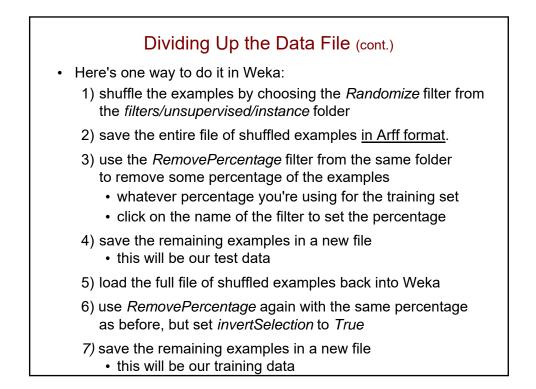


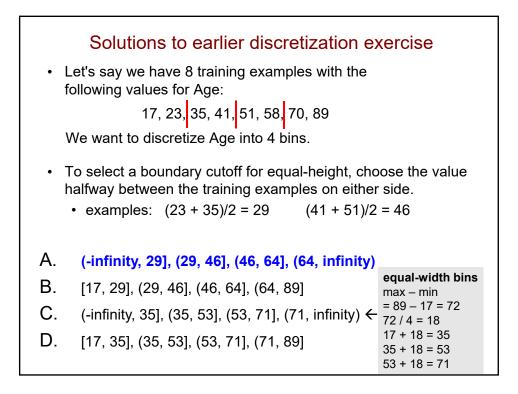


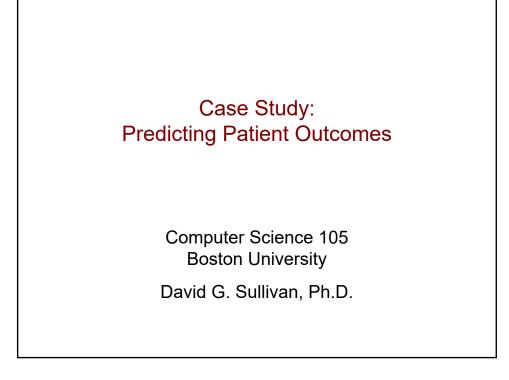
Undoing Preprocess Actions

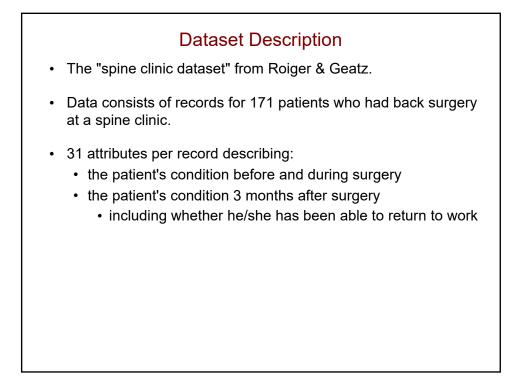
- In the Preprocess tab, the *Undo* button allows you to undo actions that you perform, including:
 - applying a filter to a dataset
 - · manually removing one or more attributes
- If you apply two filters without using *Undo* in between the two, the second filter will be applied to the results of the first filter.
- Undo can be pressed multiple times to undo a sequence of actions.

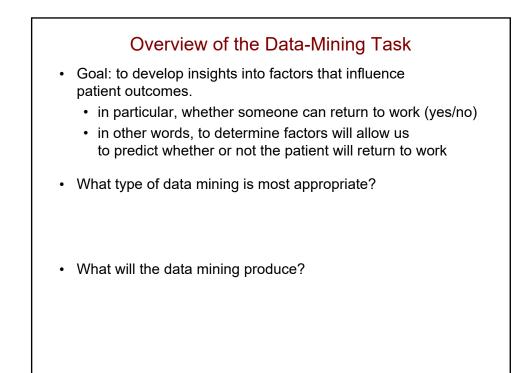


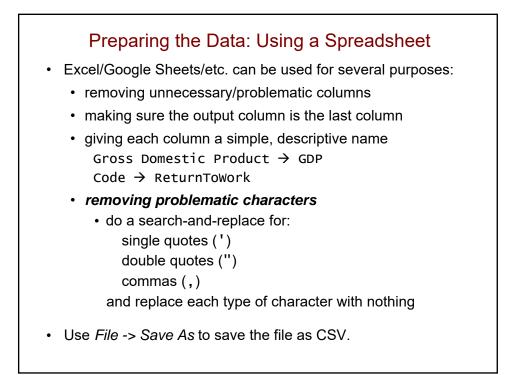


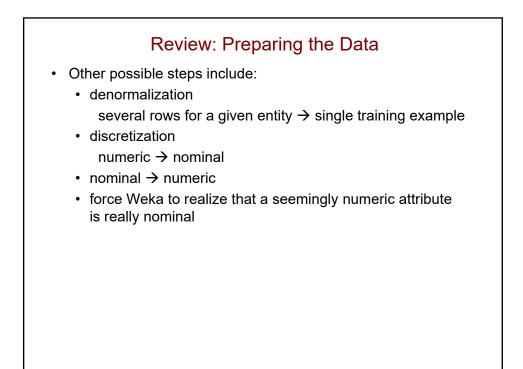


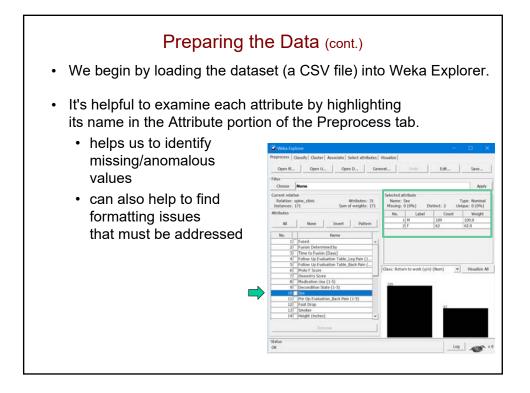






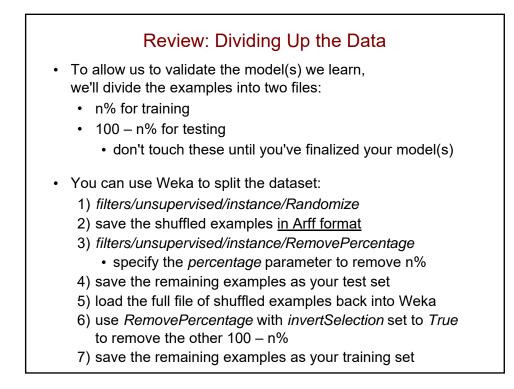






How many attributes should be removed/transformed?

Patient ID	Sex	# of Levels	Smoker (y/n)	Patient Type	Age	Return to Work (y/n)
1005	М	1	0	3100	30–39	1
1013	F	2	1	1400	50–59	0
1245	М	1	1	3100	20–29	••• 1
2110	F	3	0	2500	30–39	0
1001	F	2	1	1400	40–49	1



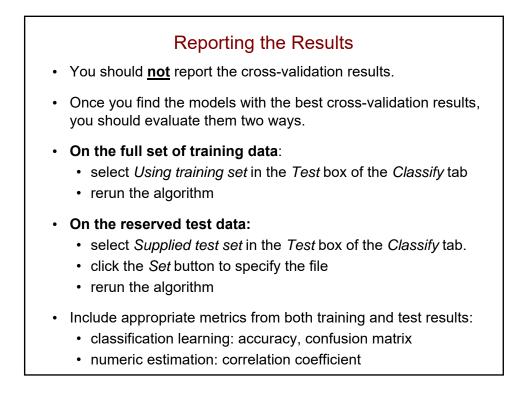
Experimenting with Different Techniques

- Use Weka to try different techniques on the training data.
- For each technique, examine:
 - the resulting model
 - the validation results
 - for classification models: overall accuracy, confusion matrix
 - · for numeric estimation models: correlation coefficient
- If the model is something you can interpret, make sure it seems reasonable.
- Try to improve the validation results by:
 - · changing the algorithm used
 - changing the algorithm's parameters

 Remember to Start with a Baseline For classification learning: 1R
 you can also use 0R to determine what % of your training data has the most common class value
For numeric estimation: <u>simple</u> linear regression
 Include the results of these baselines to put your other results in context.
 example: 80% accuracy isn't that impressive if 0R has 78% accuracy
 being honest about your results is better than making exaggerated claims!

Cross Validation

- When validating classification/estimation models, Weka performs *10-fold cross validation* by default:
 - 1) divides the training data into 10 subsets
 - 2) repeatedly does the following:
 - a) holds out one of the 10 subsets
 - b) builds a model using the other 9 subsets
 - c) tests the model using the held-out subset
 - 3) reports results that average the 10 models together
- We use cross validation when exploring possible models, because it gives a sense of how well the model will generalize.
- Note: the model reported in the output window is learned from <u>all</u> of the training examples.
 - · the cross-validation results do not actually evaluate it



Discussing the Results

- Your report should include more than just the numeric results.
- You should include an *intelligent discussion* of the results.
 - · compare training vs. test results
 - · how well do the models appear to generalize?
 - which attributes are included in the models?
 - for classification learning:
 - what do the confusion matrices tell you?
 - for numeric estimation:
 - which attributes have positive coefficients?
 - which have negative?
 - remember: the *magnitude* of the coefficients may <u>not</u> be significant
 - · are the models intuitive? why or why not?