

# MSiA-413 Introduction to Databases and Information Retrieval

## Lecture 1 Overview; Integer Representation

Instructor: Nikos Hardavellas

Slides adapted from Steve Tarzia

## Course Summary

- Learn how to handle real-world, **complex**, **messy** data with SQL relational databases
  - A powerful foundational technology
  - Like a filesystem, but better
    - (easy queries, indexing, concurrency, crash tolerance)
- Roughly speaking “Data Science” is:
  - Data management (*this course!*)
  - Statistics (e.g., IECS-304 Statistical Methods for Data Mining, EECS-349 Machine Learning)
  - Visualization (e.g., PSYCH-245 Presenting Ideas and Data)

You’ll learn to answer questions (about the past) using complex data sets.

# Things you cannot do with Excel and Matlab

- Model complex data relationships
  - Spreadsheets and matrices are very limiting formats
  - Just have records with attributes
  - Can't model one-to-many and many-to-many relationships
  - Multiple spreadsheets / multiple matrices for different types of data are possible
    - ...but, linking them is difficult
- Enforce data integrity constraints
  - Spreadsheet cells can have all kinds of weird data
  - Matlab matrices cannot easily handle anything other than numbers
- Spreadsheets are terrible for large datasets!
  - [ Excel forum on Reddit ]
  - Crashing on 10K-100K rows of data
  - Freezing on changes
  - Calculations take minutes
- Keep data and analysis separate

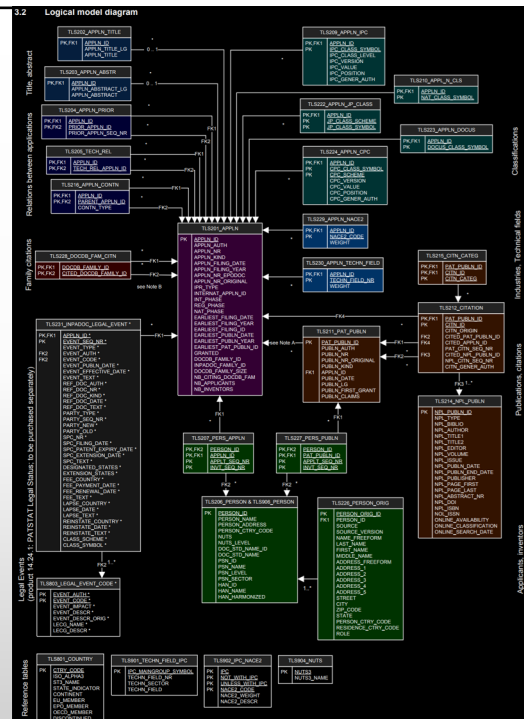
3

	A	B	C	D	E	F	G
1							
2							
3	Company Name	Invoice Date	Delivery Date	Amounts			
4	Jenny	01.09.2007	1900/01/00	2057			
5		01.11.2007	1900/01/00	2669			
6	Jenny Total			2669			
7	Sam	1900/01/01	1900/01/00	1426			
8		1998/01/01	01.01.1998	1185			
9		1998/01/01	1900/01/00	2359			
10		1998/01/01	01.06.1998	1886			
11		1998/01/01	1900/01/00	2359			
12		2000/07/01	01.07.2000	2486			
13	Sam Total			9342			
14							
15				Page 1			
16							
17	Peter	2000/01/01	1900/01/00	2385			
18		1975/04/01	1900/01/00	0			
19		2000/04/01	1900/01/00	0.000			
20		2005/06/01	1900/01/00	7 293.07			
21		1993/07/01	1900/01/00	42 717.42			
22		1993/07/01	01.07.1993	55 872.63			
23		01.08.2000	1900/01/00	40 176.80			
24		01.09.2000	1900/01/00	1585			
25		01.10.2001	1900/01/00	1384			
26		01.10.2004	01.10.2004	01518			
27		01.10.2007	01.10.2007	2057			

4

## PATSTAT: European Patent Office's International Patent Database

- 29 cross-referenced tables
- 6 DVDs of data
- 119GB of CSV files after unzipping



5

## Difficulties in plain Python, R, C++, Java, etc.

- Working with data that is larger than the computer's RAM (*scalability*)
- Keeping your data around after your program finishes (*persistence*)
- Efficiently searching through lots of data (*indexing*)
- Easily filtering and summarizing data (*querying*)
- Sharing data between multiple applications (*concurrency*)

6

## The Goal: Easy & Clean Descriptive Analytics

Answer a wide variety of complex questions using the same database:

- Where did our 10 biggest customers in 2007 live?

```
SELECT customer.name, customer.city, customer.province FROM
  customer JOIN order ON order.customer == customer.id
           JOIN order_item ON order_item.order == order.id
           JOIN product ON order_item.product = product.id
WHERE order.placed >= "2007-01-01" AND order.placed < "2008-01-01"
GROUP BY customer.id
ORDER BY SUM(order_item.qty * product.price) DESC
LIMIT 10;
```

- How many widgets are left in stock?
- What is the average price of the chairs we sell?

7

## Database Management Systems

- E.g., Oracle, MS SQL Server, MySQL, PostgreSQL, (SQLite, Access)
- Often run on a remote, multi-user server
  - Typically you need to know the hostname and have a username and password.
- May be connected to one or more software applications or may be standalone.
- Client libraries exist for every common programming language
  - But you usually query the database using the SQL language

8

## Course Outline (subject to change)

- Data in detail
  - Numeric formats
    - Binary, integers, floats, precision
    - Dates and times
  - Text encodings
    - ASCII, UTF-8, special characters
  - Organizing data in files
    - CSV, XML, JSON
  - Messy data
    - Missing entries, fuzzy matching
- SQL relational databases
  - Data modelling
    - One-to-many, many-to-many relationships
    - Integrity & foreign key constraints
  - Structured Query Language (SQL)
    - Select, create table, update, delete
    - Joining tables
    - Subqueries & temporary tables
    - Indexes and execution plans
- Plus more if time permits

9

Questions about course  
content?

10

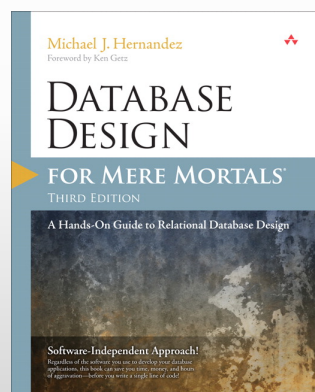
## Course Logistics

- READ THE SYLLABUS !!!

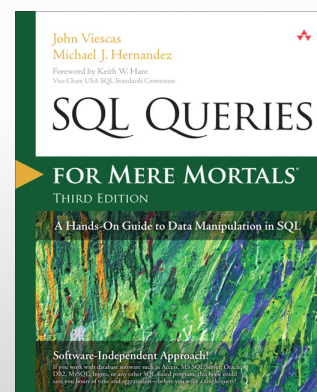
11

## Recommended Books

Hernandez “Database Design for Mere Mortals”



Viescas & Hernandez “SQL Queries for Mere Mortals”



12

Questions about logistics?

13

## Part 2: Integer Representation

14

## Computers store information in **binary**

- Ones and Zeros
- ...000100100001001001110011011010101010111100000...
- Called “bits,” meaning “**binary** digits”
- Why?
  - Simplicity
  - Noise robustness
  - By convention
- But how do we get meaning from a sequence of ones and zeros?

15

## Data is zeros and ones plus **context**

- An **encoding** defines what the zeros and ones represent
- “01000100011000010111010001100001” can represent:
  - The number 1,147,237,473 as an **integer**
  - The number 901.8184 as a **float**
  - The four letters “Data” in the **ASCII** character encoding
  - **This color** (at 37% transparency) in **RGBA**
  - 32 separate True or False values
- Any crazy encoding is possible, but there are some standards

16



# Integers

- Integers are the simplest of all data encodings
- Whole numbers only (no fractions)
- Numbers are represented directly in the “base two” positional notation
- The familiar “base ten” representation of numbers is just a convention due to the fact that humans have ten fingers
  - Still, Aztecs used base-20 (vigesimal)
  - Babylonians used base-60 (sexagesimal)
- What number base will octopuses evolve to use?

(drawing from <http://drawingpencilarts.com/realistic-octopus-drawing/>)



## Integers in detail

Decimal  $137_{\text{ten}}$

$$\begin{array}{r}
 1 \quad 3 \quad 7 \\
 \underline{\times 10^2} \quad \underline{\times 10^1} \quad \underline{\times 10^0} \quad \leftarrow \text{powers of 10} \\
 100 + 30 + 7 = 137
 \end{array}$$

Binary  $10001001_{\text{two}} = 137_{\text{ten}}$

$$\begin{array}{r}
 1 \quad 0 \quad 0 \quad 0 \quad 1 \quad 0 \quad 0 \quad 1 \\
 \underline{\times 2^7} \quad \underline{\times 2^6} \quad \underline{\times 2^5} \quad \underline{\times 2^4} \quad \underline{\times 2^3} \quad \underline{\times 2^2} \quad \underline{\times 2^1} \quad \underline{\times 2^0} \quad \leftarrow \text{powers of 2} \\
 128 + 0 + 0 + 0 + 8 + 0 + 0 + 1 = 137
 \end{array}$$

## Simple binary integers

$$1_{\text{ten}} = 1_{\text{two}}$$

$$2_{\text{ten}} = 10_{\text{two}}$$

$$4_{\text{ten}} = 100_{\text{two}}$$

$$8_{\text{ten}} = 1000_{\text{two}}$$

$$16_{\text{ten}} = 10000_{\text{two}}$$

$$32_{\text{ten}} = 100000_{\text{two}}$$

$$64_{\text{ten}} = 1000000_{\text{two}}$$

$$128_{\text{ten}} = 10000000_{\text{two}}$$

$$3_{\text{ten}} = 11_{\text{two}}$$

$$7_{\text{ten}} = 111_{\text{two}}$$

$$15_{\text{ten}} = 1111_{\text{two}}$$

$$31_{\text{ten}} = 11111_{\text{two}}$$

$$63_{\text{ten}} = 111111_{\text{two}}$$

$$127_{\text{ten}} = 1111111_{\text{two}}$$

$$255_{\text{ten}} = 11111111_{\text{two}}$$

19

There are only **10** types of people in this world... those who understand binary and those who don't.

(Stop and practice)

20

## Binary tricks

- Remember the first eight powers of two:
  - 2, 4, 8, 16, 32, 64, 128, 256
- Remember that  $2^{10} = 1024 \approx 1000$ 
  - Lets you estimate the number of binary digits in a decimal integer:  
Every three decimal digits gives about ten binary digits
  - $2^{20} \approx 1$  million
  - $2^{30} \approx 1$  billion
- Remember the important large powers of two:
  - $2^8 = 256$
  - $2^{16} \approx 64$  thousand
  - $2^{32} \approx 4$  billion
  - $2^{64} \approx$  really big ( $\approx 18$  quintillion, or in CS parlance: 16 exa-...)

21

## Addition in binary

$$4 + 7 = 11$$

1 ← carry

$$\begin{array}{r} 4 \\ + 7 \\ \hline 11 \end{array}$$

$$100 + 111 = 1011$$

1 ← carry

$$\begin{array}{r} 100 \\ + 111 \\ \hline 1011 \end{array}$$

22

## More binary addition

$$63 + 98 = 161$$

$$\begin{array}{r}
 1\ 1 \leftarrow \text{carry} \\
 6\ 3 \\
 +\ 9\ 8 \\
 \hline
 1\ 6\ 1
 \end{array}$$

$$11111 + 110010 = 1010001$$

$$\begin{array}{r}
 1\ 1\ 1\ 1\ 1 \leftarrow \text{carry} \\
 1\ 1\ 1\ 1\ 1 \\
 +\ 1\ 1\ 0\ 0\ 1\ 0 \\
 \hline
 1\ 0\ 1\ 0\ 0\ 0\ 1
 \end{array}$$

23

## Subtraction: addition's tricky pal

$$161 - 98 = 63$$

$$\begin{array}{r}
 \rightsquigarrow \quad \rightsquigarrow \text{borrow} \\
 \cancel{0}1\ \cancel{1}5\ \cancel{6}\ 11 \\
 -\quad\quad 9\quad 8 \\
 \hline
 \quad\quad 6\quad 3
 \end{array}$$

$$1010001 - 110010 = 11111$$

$$\begin{array}{r}
 \rightsquigarrow \quad \rightsquigarrow \quad \rightsquigarrow \quad \rightsquigarrow \quad \rightsquigarrow \text{borrow} \\
 \cancel{0}1\ \cancel{1}0\ \cancel{1}0\ \cancel{1}0\ \cancel{1}0\ \cancel{1}0\ 1 \\
 -\quad\quad 1\quad 1\quad 0\quad 0\quad 1\quad 0 \\
 \hline
 \quad\quad\quad 1\quad 1\quad 1\quad 1\quad 1
 \end{array}$$

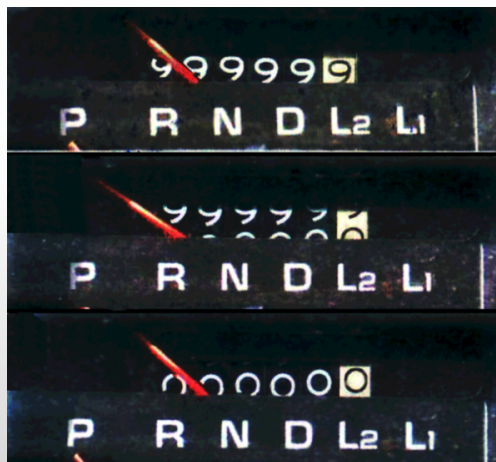
24

## What about negative integers?

- **Signed integers** can represent both positive and negative integers
- We need an extra bit to represent the sign of the number
- But we don't just use a simple sign bit
- We use **two's complement** to represent negative numbers, because it
  - Simplifies the computer's addition and subtraction circuitry, and
  - And it has just one representation of zero
- Negative numbers “roll over” from the top of the binary range.

25

## Works like an old-style car odometer



26

## Two's complement for three-bit numbers

3: 011  
 2: 010  
 1: 001  
 0: 000  
 -1: 111 ← rollover  
 -2: 110  
 -3: 101  
 -4: 100  
     ↑  
     sign bit

$$1 - 2 = 1 + (-2) = -1$$

$$001 + 110 = 111$$

Subtraction is done in the exact same way as addition!

No need to learn how to “borrow.”

27

## Subtraction works just like addition!

No need to learn how to “borrow.”

Just negate the second number and add.

$$3 - 2 = 3 + (-2) =$$

1 1 ← carry  
 0 1 1  
 + 1 1 0  
 ———  
 0 0 1 ← our answer!

We ignore the final carry because it falls outside of the 3-bits we are working with. That's how we roll-over between negative and positive.

3: 011  
 2: 010  
 1: 001  
 0: 000  
 -1: 111  
 -2: 110  
 -3: 101  
 -4: 100

28

## Two's complement negation

To negate a number:

- **Flip** all the bits. Ones become zeros and zeros become ones.
- **Add one**

For example -3

- Start with the bits for three: **011**
- Flip the bits: **100**
- Add one: **101**

3: 011  
2: 010  
1: 001  
0: 000  
-1: 111  
-2: 110  
-3: 101  
-4: 100

29

## Negating with Complement and Increment

- Claim: Following holds for 2's complement (when defined)
  - $\sim x + 1 = -x$
- Complement
  - Observation:  $\sim x + x = 1111\dots 11_2 = -1$

$$\begin{array}{r}
 x \quad 10011101 \\
 + \quad \sim x \quad 01100010 \\
 \hline
 -1 \quad 11111111
 \end{array}$$

- Increment

- $\sim x + 1 = \sim x + x - x + 1 = \cancel{1} - x + \cancel{1} = -x$
- $\sim x + 1 = -x$

30

## Overflow: when numbers don't fit

For example,  $2 + 2 = 4$

4 cannot be represented in a three-bit **signed** integer.

What happens when we try this addition?

**1** ← carry  
0 1 0  
+ 0 1 0  
1 0 0 ← answer looks like -4!

- The computer will throw an **exception** if the signs of the operands were the same, but the sign of the result is different.
  - positive + negative cannot overflow
  - positive + positive should give a positive
  - negative + negative should give a negative
- Remember that the left-most bit indicates the sign.

**3:** 011  
**2:** 010  
**1:** 001  
**0:** 000  
**-1:** 111  
**-2:** 110  
**-3:** 101  
**-4:** 100

31

## Examples with 4 and 8 bits

4-bit is between -8 and 7

8-bit is between -128 and 127

(Stop and practice)

32



## Reading Assignment and Practice

- Read “Representing Numbers in Computers” at <http://www.stat.berkeley.edu/~nolan/stat133/Spr04/chapters/representations.pdf>
- Practice converting numbers to and from binary
- Practice binary addition and subtraction (check with online tools)
- Browse data sets from Kaggle.com
  - Don't forget to click the “Data” tab
- Watch this video to see how addition is actually implemented in hardware  
<https://www.youtube.com/watch?v=1I5ZMmrOfnA>  
Search YouTube for “PBS ALU”