

## 02-1 Data models

Monday, March 13, 2023 17:33



2-1 Data  
models - ...

# Data Models

Sungahn Ko



1

## Outline

- Foundational topics for course
- Data types
- Data representations
- Data models
- Tables



2

# Data

- So far we have looked at many examples of visualization
- Ignored the fundamentals
- It all starts with **data**
- Good definition?
  - Collected from the world
  - Represents the world (somehow)
  - Models something that is interesting (?)



3

# Data

- Data is just an abstraction of a real phenomenon
- Corollaries:
  - Visualization is only as good as the data
  - Visualizations can be misleading
  - Good data is important!



4

## Data and Datasets

- Data is everywhere!
- Almost all of it is unstructured (95%)
  - Images
  - Video
  - Sound
  - Log files
  - Text
  - Web pages
- Need regular and structured datasets to analyze and visualize this data
- Often we must do this ourselves!



5

## Existing Structured Data

- Resources exist that collect data on the Web
- [Data.gov](#)
  - US Federal government dataset collection
- UCB data:
  - [http://vis.berkeley.edu/courses/cs294-10-fa14/wiki/index.php/Online\\_Datasets](http://vis.berkeley.edu/courses/cs294-10-fa14/wiki/index.php/Online_Datasets)
- Public Data Portal Korea:
   
<https://www.data.go.kr/>



6

## Deriving Structured Data: Wrangler (CHI 2011)

- <http://vis.stanford.edu/wrangler/>

*Check it !!!!!*

- Demo!
  - <https://vimeo.com/19185801>



7

## Data Models

- How to capture and structure our data?
- Often use three types of entities:
  - Objects
    - Actual items of interest
    - Different types possible
    - Example: people on Facebook
  - Relations
    - Connections between objects
    - Examples: friend relationships between people
  - Attributes
    - Characteristics of objects and relations
    - Property of an entity
    - Example: age, gender, color of object



8

## Relational Data Model

- Records in a **data table**
- Structured from amenable to analysis and visualization
- Fixed-length tuples (attributes)
- Each column (attribute) has a domain (type)
- Relational databases also allow relations between cases (often through related tables) – not our focus today



9

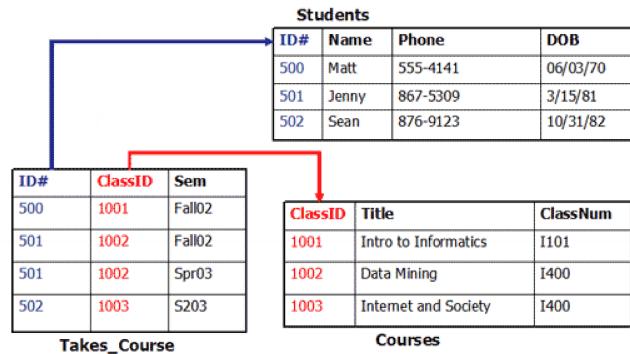
## Relational Algebra

- Manipulating relation data models
- Formalized in the standardized SQL language
  - Standard Query Language
- Selection (SELECT)
- Projection (WHERE)
- Sorting (ORDER BY)
- Aggregation (GROUP BY, SUM, MIN, ...)
- Set operations (UNION, INTERSECT...)
- Join (INNER JOIN, ...)



10

## Example: Relational Data



UNIST  
UNIVERSITY OF  
NATIONAL INSTITUTE OF  
SCIENCE AND TECHNOLOGY

11

# Variable Types

- **Nominal** • ?? data (labels) [*Category data*]
    - Supports only equality (same or different)
    - Examples: gender, car brand, fruit, bus number
  - **Ordinal** • ?? data (ordered) [*Integer data*]
    - Obeys the < relation, ordered set
    - Examples: days of week, Fresh/Soph./Junior/Senior
  - **Quantitative** • ?? data [*Real-number data*]
    - Supports arithmetic operations
    - Interval (zero arbitrary)
      - Example: Dates, location
    - Ratio (zero fixed)
      - Example: age, temperature, stock value

**UNIST**  
UNIVERSITY OF  
SOUTHERN  
KOREA

S. S. Stevens, On the theory of scales of measurements, 1946

12

# Mathematical Operations

- N - Nominal (labels)
  - Operations:  $=, \neq$
- O – Ordered
  - Operations:  $=, \neq, <, >$
- Q - Interval (Location of zero arbitrary)
  - Operations:  $=, \neq, <, >, -, +$
  - Can measure distances or spans
- Q - Ratio (zero fixed)
  - Operations:  $=, \neq, <, >, -, +$
  - Can measure ratios or proportions



13

# Metadata

- Data about data (derived data)
- Describes:
  - Definition
  - Structure
  - Administration
- Examples:
  - Types of variables in data table
  - Language of a particular text
  - Dimensions, bit depth, timestamp for a photograph
- Metadata is often useful when treating data, and sometimes also for visualization!



14

# Data Dimensions

- Common dimensions: 1, 2, 3
  - 1 dimension – univariate
    - Temperature readings
  - 2 dimensions – bivariate
    - Positions on map (lat/long)
  - 3 dimensions – trivariate
    - Positions in space (3D)
- For more than 3 dimensions
  - Multivariate
  - Hypervariate



15

## Example: US Census Data

- People: # of people in group
- Year: 1850 – 2000 (every decade)
- Age: 0 – 90+
- Sex: Male, Female
- Marital Status: Single, Married, Divorced, ...
- 2348 data points

	A	B	C	D	E
1	year	age	marital	sex	people
2	1850	0	0	1	1559799
3	1850	0	0	2	1450376
4	1850	5	0	1	1411087
5	1850	5	0	2	1359668
6	1850	10	0	1	1310999
7	1850	10	0	2	1216114
8	1850	15	0	1	1077135
9	1850	15	0	2	1117109
10	1850	20	0	1	1317181
11	1850	20	0	2	1005841
12	1850	25	0	1	162147
13	1850	25	0	2	11144
14	1850	30	0	1	720588
15	1850	30	0	2	839898
16	1850	35	0	1	588497
17	1850	35	0	2	124612
18	1850	40	0	1	475911
19	1850	40	0	2	428185
20	1850	45	0	1	382111
21	1850	45	0	2	341194
22	1850	50	0	1	211143
23	1850	50	0	2	295980
24	1850	55	0	1	141290
25	1850	55	0	2	197209
26	1850	60	0	1	174976
27	1850	60	0	2	14144
28	1850	65	0	1	166527
29	1850	65	0	2	105584
30	1850	70	0	1	78677
31	1850	70	0	2	71743
32	1850	75	0	1	40834
33	1850	75	0	2	40229
34	1850	80	0	1	23159
35	1850	80	0	2	23849
36	1850	85	0	1	1886
37	1850	85	0	2	10511
38	1850	90	0	1	3529
39	1850	90	0	2	8389
40	1860	0	0	1	2120846
41	1860	0	0	2	2092162



16

## How to Represent Tabular Data?

- Standard answer in this course: **graphs!**
  - Statistical data graphics
  - Bar charts, line charts, pie charts, etc.
- There is a simpler way: **tables**
  - Also a graphical representation!
  - Textual representation
  - Useful for direct lookups
- When to use which format?
  - Tables: looking up individual values, precise data
  - Graphs: relationships, comparisons



18

## Side Note

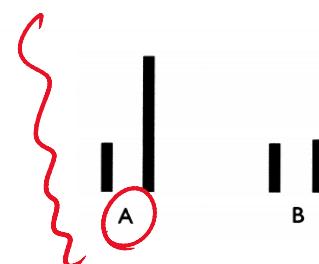
A: Which number is larger?

284      912

B: Which number is larger?

284      312

→ People answer A faster than B. Why?



*"The form of representation most appropriate for an artifact depends on the task to be performed" – D. A. Norman, 1993*



19

## Example: Tables and Graphs

Cancer site	Relative survival rate, % (SE)			
	5 years	10 years	15 years	20 years
Ducts, breast and pharynx	56.7 (1.3)	44.2 (1.4)	37.5 (1.6)	33.0 (1.8)
Oesophagus	14.2 (1.4)	7.9 (1.3)	7.7 (1.6)	5.4 (2.0)
Stomach	22.8 (0.8)	12.9 (1.0)	12.9 (1.0)	12.9 (1.0)
Colon	52.0 (0.8)	55.4 (1.0)	53.9 (1.2)	53.4 (1.6)
Rectum	52.6 (1.2)	55.2 (1.4)	53.8 (1.8)	49.2 (2.3)
Liver and intrahepatic bile duct	7.5 (1.1)	5.8 (1.2)	6.3 (1.9)	7.6 (2.0)
Pancreas	4.0 (0.8)	3.0 (0.9)	2.7 (0.8)	2.7 (0.8)
Lung	16.0 (0.8)	10.0 (1.0)	9.0 (1.1)	8.3 (1.3)
Lung and bronchus	15.0 (0.4)	10.6 (0.4)	8.1 (0.4)	6.5 (0.4)
Melanomas	89.0 (0.4)	86.7 (1.1)	85.5 (1.9)	82.8 (1.9)
Breast	84.7 (0.4)	82.0 (0.4)	80.5 (0.5)	78.0 (0.5)
Cervix uteri	70.5 (1.4)	64.1 (1.8)	58.8 (2.1)	60.0 (2.4)
Corpus uteri and uterus,	84.3 (1.0)	83.2 (1.3)	80.8 (1.7)	79.2 (2.0)
NSCLC	—	—	—	—
Ovary	55.0 (0.4)	49.3 (1.6)	49.9 (1.9)	49.6 (2.4)
Prostate	95.0 (0.4)	92.0 (0.4)	87.0 (1.1)	81.0 (2.0)
Testis	94.7 (1.1)	94.0 (1.3)	91.1 (1.8)	88.2 (2.3)
Thyroid	98.0 (0.4)	95.2 (0.9)	87.1 (1.4)	81.4 (2.0)
Tissue	94.7 (1.1)	94.0 (1.3)	91.1 (1.8)	88.2 (2.3)
Melanomas	89.0 (0.8)	86.7 (1.1)	83.5 (1.5)	82.8 (1.9)
Breast	86.4 (0.4)	78.3 (0.6)	71.3 (0.7)	65.0 (1.0)
Hodgkin's disease	85.1 (1.7)	79.8 (2.0)	73.8 (2.4)	67.1 (2.8)
Corpus uteri, uterus	84.3 (1.0)	83.2 (1.3)	80.8 (1.7)	79.2 (2.0)
Urinary bladder	80.0 (1.4)	74.0 (1.8)	70.0 (2.1)	64.0 (2.4)
Cervix,uteri	70.5 (1.4)	64.1 (1.8)	62.8 (2.1)	60.0 (2.4)
Larynx	68.8 (2.1)	56.7 (2.5)	45.8 (2.8)	37.8 (3.1)
Rectum	62.6 (1.2)	55.2 (1.4)	51.8 (1.6)	49.2 (2.3)
Kidney, renal pelvis	61.8 (1.3)	54.4 (1.6)	49.8 (1.8)	47.3 (2.6)
Colon	61.7 (0.8)	55.4 (1.0)	53.9 (1.2)	52.3 (1.6)
Non-Hodgkin's	59.0 (1.3)	52.0 (1.6)	47.0 (1.7)	47.0 (1.7)
Oral cavity,pharynx	56.7 (1.3)	44.2 (1.4)	37.5 (1.6)	33.0 (1.8)
Ovary	55.0 (1.3)	49.3 (1.6)	49.9 (1.9)	49.6 (2.4)
Leukemia	42.5 (1.2)	32.4 (1.3)	29.7 (1.8)	26.2 (1.7)
Brain, nervous system	32.0 (1.4)	29.2 (1.5)	27.6 (1.7)	26.1 (1.9)
Multiple myeloma	29.5 (1.6)	12.7 (1.5)	7.0 (1.5)	4.8 (1.5)
Stomach	29.0 (1.4)	16.0 (1.4)	8.0 (1.4)	5.9 (1.4)
Lung and bronchus	15.0 (0.4)	10.6 (0.4)	8.1 (0.4)	6.5 (0.4)
Esophagus	14.2 (1.4)	7.9 (1.3)	7.7 (1.6)	5.4 (2.0)
Liver,bile duct	7.5 (1.1)	5.8 (1.2)	6.3 (1.5)	7.6 (2.0)
Pancreas	4.0 (0.5)	3.0 (1.0)	2.7 (0.6)	2.7 (0.8)

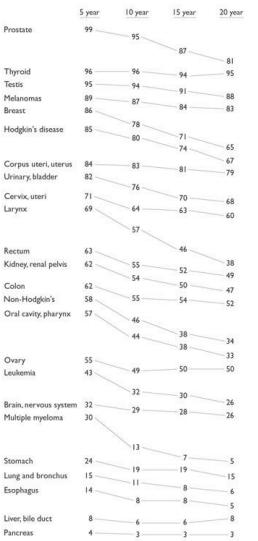
Table 4: Most recent period estimates of relative survival rates, by cancer site

	Estimates of relative survival rates, by cancer site			
	5 year	10 year	15 year	20 year
Prostate	99.0 0.4	95.2 0.9	87.1 1.4	81.0 2.0
Thyroid	96.0 0.4	94.0 0.4	91.4 0.4	89.1 0.4
Testis	97.1 1.1	94.0 1.3	91.1 1.8	88.2 2.3
Melanomas	89.0 0.8	86.7 1.1	83.5 1.5	82.8 1.9
Breast	86.0 0.4	78.3 0.6	71.3 0.7	65.0 1.0
Hodgkin's disease	85.1 1.7	79.8 2.0	73.8 2.4	67.1 2.8
Corpus uteri,uterus	84.3 1.0	83.2 1.3	80.8 1.7	79.2 2.0
Urinary bladder	80.0 1.4	74.0 1.8	70.0 2.1	64.0 2.4
Cervix,uteri	70.5 1.4	64.1 1.8	62.8 2.1	60.0 2.4
Larynx	68.8 2.1	56.7 2.5	45.8 2.8	37.8 3.1
Rectum	62.6 1.2	55.2 1.4	51.8 1.6	49.2 2.3
Kidney, renal pelvis	61.8 1.3	54.4 1.6	49.8 1.8	47.3 2.6
Colon	61.7 0.8	55.4 1.0	53.9 1.2	52.3 1.6
Non-Hodgkin's	59.0 1.3	52.0 1.6	47.0 1.7	47.0 1.7
Oral cavity,pharynx	56.7 1.3	44.2 1.4	37.5 1.6	33.0 1.8
Ovary	55.0 1.3	49.3 1.6	49.9 1.9	49.6 2.4
Leukemia	42.5 1.2	32.4 1.3	29.7 1.8	26.2 1.7
Brain, nervous system	32.0 1.4	29.2 1.5	27.6 1.7	26.1 1.9
Multiple myeloma	29.5 1.6	12.7 1.5	7.0 1.5	4.8 1.5
Stomach	29.0 1.4	16.0 1.4	8.0 1.4	5.9 1.4
Lung and bronchus	15.0 0.4	10.6 0.4	8.1 0.4	6.5 0.4
Esophagus	14.2 1.4	7.9 1.3	7.7 1.6	5.4 2.0
Liver,bile duct	7.5 1.1	5.8 1.2	6.3 1.5	7.6 2.0
Pancreas	4.0 0.5	3.0 1.0	2.7 0.6	2.7 0.8

"For [...] small data sets, usually a simple table shows the data more effectively than a graph, let alone a chartjunk graph." – E. R. Tufte, 2003



21



## Acknowledgments

- Dr. Niklas Elmquist, University of Maryland
- John Stasko, Georgia Tech
- Jeff Heer, University of Washington
- S. S. Stevens, On the theory of scales of measurements, *Science*, 103(2684):677-680, 1946.



22