

Information Visualization
Data Abstraction

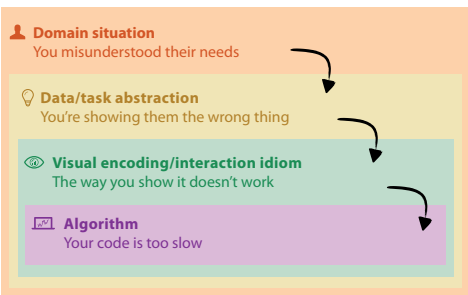
Tamara Munzner
Department of Computer Science
University of British Columbia

Lect 2, 9 Jan 2020

http://www.cs.ubc.ca/~tmm/courses/436V-20

Different threats to validity at each level

- cascading effects downstream



What does data mean?

- 14, 2.6, 30, 30, 15, 100001
What does this sequence of six numbers mean?
two points far from each other in 3D space?
two points close to each other in 2D space, with 15 links between them, and a weight of 100001 for the link?
something else??
Basil, 7, S, Pear
What about this data?
food shipment of produce (basil & pear) arrived in satisfactory condition on 7th day of month
Basil Point neighborhood of city had 7 inches of snow cleared by the Pear Creek Limited snow removal service
lab rat Basil made 7 attempts to find way through south section of maze, these trials used pear as reward food

Items & Attributes

- item: individual entity, discrete
eg patient, car, stock, city
"independent variable"
attribute: property that is measured, observed, logged...
eg height, blood pressure for patient
eg horsepower, make for car
"dependent variable"

attributes: name, age, shirt size, fave fruit

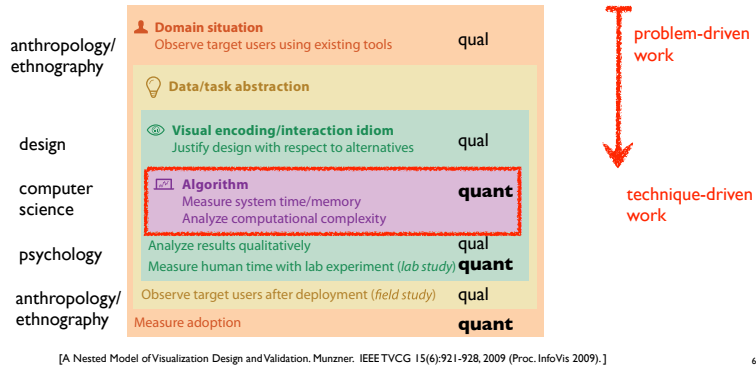
Name	Age	Shirt Size	Favorite Fruit
Amy	8	S	Apple
Basil	7	S	Pear
Clara	9	M	Durian
Desmond	13	L	Elderberry
Ernest	12	L	Peach
Fanny	10	S	Lychee
George	9	M	Orange
Hector	8	L	Loquat
Ida	10	M	Pear
Amy	12	M	Orange

item: person

Nested Model

Interdisciplinary: need methods from different fields at each level

- mix of qual and quant approaches (typically)



Now what?

- semantics: real-world meaning

Amy	8	S	Apple
Basil	7	S	Pear
Clara	9	M	Durian
Desmond	13	L	Elderberry
Ernest	12	L	Peach
Fanny	10	S	Lychee
George	9	M	Orange
Hector	8	L	Loquat
Ida	10	M	Pear
Amy	12	M	Orange

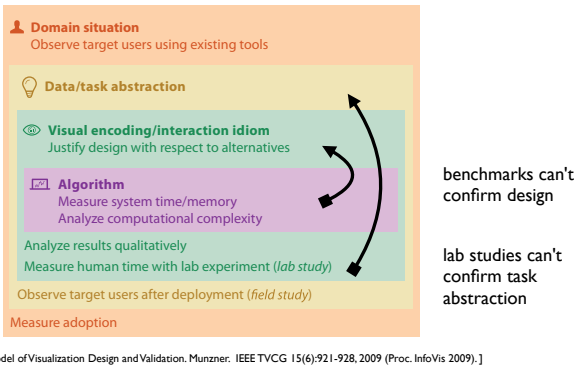
Other data types

- links
express relationship between two items
eg friendship on facebook, interaction between proteins
- positions
spatial data: location in 2D or 3D
pixels in photo, voxels in MRI scan, latitude/longitude
- (grids)
sampling strategy for continuous data

How to evaluate a visualization: So many methods, how to pick?

- Computational benchmarks?
quant: system performance, memory
- User study in lab setting?
quant: (human) time and error rates, preferences
qual: behavior/strategy observations
- Field study of deployed system?
quant: usage logs
qual: interviews with users, case studies, observations
- Analysis of results?
quant: metrics computed on result images
qual: consider what structure is visible in result images
- Justification of choices?
qual: perceptual principles, best practices

Mismatches: Common problem



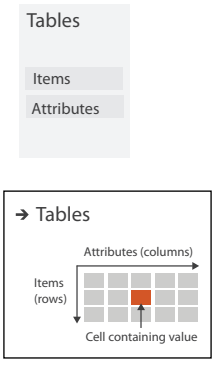
Now what?

- semantics: real-world meaning

Name	Age	Shirt Size	Favorite Fruit
Amy	8	S	Apple
Basil	7	S	Pear
Clara	9	M	Durian
Desmond	13	L	Elderberry
Ernest	12	L	Peach
Fanny	10	S	Lychee
George	9	M	Orange
Hector	8	L	Loquat
Ida	10	M	Pear
Amy	12	M	Orange

Dataset types

- flat table
one item per row
each column is attribute
cell holds value



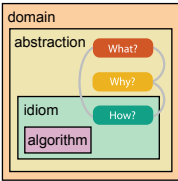
attributes: name, age, shirt size, fave fruit

Name	Age	Shirt Size	Favorite Fruit
Amy	8	S	Apple
Basil	7	S	Pear
Clara	9	M	Durian
Desmond	13	L	Elderberry
Ernest	12	L	Peach
Fanny	10	S	Lychee
George	9	M	Orange
Hector	8	L	Loquat
Ida	10	M	Pear
Amy	12	M	Orange

item: person

Nested model: Four levels of visualization design

- domain situation
who are the target users?
- abstraction
translate from specifics of domain to vocabulary of visualization
what is shown? data abstraction
why is the user looking at it? task abstraction
often must transform data, guided by task
- idiom
how is it shown?
visual encoding idiom: how to draw
interaction idiom: how to manipulate
- algorithm
efficient computation



[A Nested Model of Visualization Design and Validation. Munzner. IEEE TVCG 15(6):921-928, 2009 (Proc. InfoVis 2009).]
[A Multi-Level Typology of Abstract Visualization Tasks Brehmer and Munzner. IEEE TVCG 19(12):2376-2385, 2013 (Proc. InfoVis 2013).]

What: Data Abstraction

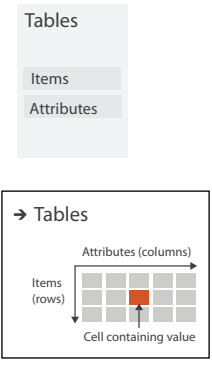
Now what?

- semantics: real-world meaning
- data types: structural or mathematical interpretation of data
item, link, attribute, position, (grid)
different from data types in programming!

Name	Age	Shirt Size	Favorite Fruit
Amy	8	S	Apple
Basil	7	S	Pear
Clara	9	M	Durian
Desmond	13	L	Elderberry
Ernest	12	L	Peach
Fanny	10	S	Lychee
George	9	M	Orange
Hector	8	L	Loquat
Ida	10	M	Pear
Amy	12	M	Orange

Dataset types

- flat table
one item per row
each column is attribute
cell holds value for item-attribute pair
unique key (could be implicit)



attributes: name, age, shirt size, fave fruit

ID	Name	Age	Shirt Size	Favorite Fruit
1	Amy	8	S	Apple
2	Basil	7	S	Pear
3	Clara	9	M	Durian
4	Desmond	13	L	Elderberry
5	Ernest	12	L	Peach
6	Fanny	10	S	Lychee
7	George	9	M	Orange
8	Hector	8	L	Loquat
9	Ida	10	M	Pear
10	Amy	12	M	Orange

item: person

Table	A		B		C		S	T		U
	Order ID	Order Date	Order Priority	Product	Container	Product Base	Margin	Ship Date		
	3	10/14/06	5-Low	Large Box			0.8	10/21/06		
	6	2/21/08	4-Not Specified	Small Pack			0.55	2/22/08		
	32	7/16/07	2-High	Small Pack			0.79	7/17/07		
	32	7/16/07	2-High	Jumbo Box			0.72	7/17/07		
	32	7/16/07	2-High	Medium Box			0.6	7/18/07		
	32	7/16/07	2-High	Medium Box			0.65	7/18/07		
	35	10/23/07	4-Not Specified	Wrap Bag			0.52	10/24/07		
	35	10/23/07	4-Not Specified	Small Box			0.58	10/25/07		
	36	11/3/07	1-Urgent	Small Box			0.55	11/3/07		
	65	3/18/07	1-Urgent	Small Pack			0.49	3/19/07		
	66	1/20/05	5-Low	Wrap Bag			0.56	1/20/05		
	69	6/4/05	4-Not Specified	Small Pack			0.44	6/6/05		
	69	6/4/05	4-Not Specified	Wrap Bag			0.6	6/6/05		
	70	12/18/06	5-Low	Small Box			0.59	12/23/06		
	70	12/18/06	5-Low	Wrap Bag			0.82	12/23/06		
	96	4/17/05	2-High	Small Box			0.55	4/19/05		
	97	1/29/06	3-Medium	Small Box			0.38	1/30/06		
	129	11/19/08	5-Low	Small Box			0.37	11/28/08		
	130	5/8/08	2-High	Small Box			0.37	5/9/08		
	130	5/8/08	2-High	Medium Box			0.38	5/10/08		
	130	5/8/08	2-High	Small Box			0.6	5/11/08		
	132	6/11/06	3-Medium	Medium Box			0.6	6/12/06		
	132	6/11/06	3-Medium	Jumbo Box			0.69	6/14/06		
	134	5/1/08	4-Not Specified	Large Box			0.82	5/3/08		
	135	10/21/07	4-Not Specified	Small Pack			0.64	10/23/07		
	166	9/12/07	2-High	Small Box			0.55	9/14/07		
	193	8/8/06	1-Urgent	Medium Box			0.57	8/10/06		
	194	4/5/08	3-Medium	Wrap Bag			0.42	4/7/08		

	A		B		C		S		T		U
	Order Id	Order Date	Order	Priority	Product	Container	Product Base	Margin	Ship Date		
	3	10/14/06	5	Low		Large Box		0.8		10/21/06	
	6	2/21/08	4	Not Specified		Small Pack		0.55		2/22/08	
	32	7/16/07	2	High		Small Pack		0.79		7/17/07	
	32	7/16/07	2	High		Jumbo Box		0.72		7/17/07	
	32	7/16/07	2	High		Medium Box		0.6		7/18/07	
	32	7/16/07	2	High		Medium Box		0.65		7/18/07	
	35	10/23/07	4	Not Specified		Wrap Bag		0.52		10/24/07	
	35	10/23/07	4	Not Specified		Small Box		0.58		10/25/07	
	36	11/3/07	1	Urgent		Small Box		0.55		11/3/07	
	65	3/18/07	1	Urgent		Small Pack		0.49		3/19/07	
	66	1/20/05	5	Low		Wrap Bag		0.56		1/20/05	
	69	6/4/05	4	Not Specified		Small Pack		0.44		6/6/05	
	69	6/4/05	4	Not Specified		Wrap Bag		0.6		6/6/05	
	70	12/18/06	5	Low		Small Box		0.59		12/23/06	
	72	12/18/06	5	Low		Wrap Bag		0.82		12/23/06	
	96	4/17/05	2	High		Small Box		0.55		4/19/05	
	97	1/29/06	3	Medium		Small Box		0.38		1/30/06	
	129	11/19/08	5	Low		Small Box		0.37		11/28/08	
	130	5/8/08	2	High		Small Box		0.37		5/9/08	
	130	5/8/08	2	High		Medium Box		0.38		5/10/08	
	130	5/8/08	2	High		Small Box		0.6		5/11/08	
	132	6/11/06	3	Medium		Medium Box		0.6		6/12/06	
	132	6/11/06	3	Medium		Jumbo Box		0.69		6/14/06	
	134	5/1/08	4	Not Specified		Large Box		0.82		5/3/08	
	135	10/21/07	4	Not Specified		Small Pack		0.64		10/23/07	
	166	9/12/07	2	High		Small Box		0.55		9/14/07	
	193	8/8/06	1	Urgent		Medium Box		0.57		8/10/06	
	194	4/5/08	3	Medium		Wrap Bag		0.42		4/7/08	

# Dataset types

- multidimensional tables
  - indexing based on multiple keys
  - eg genes, patients

→ Tables

Attributes (columns)

Items (rows)

Cell containing value

→ Multidimensional Table

Key 1

Key 2

Attributes

Value in cell

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
2	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
3	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
4	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
5	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
6	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
7	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
8	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
9	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
10	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
11	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
12	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
13	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
14	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
15	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
16	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
17	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
18	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
19	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
20	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
21	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
22	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
23	1	1	1																							

# Visualizing tables

# Dataset types

Tables

Items

Attributes

Networks & Trees

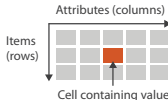
Items (nodes)

Links

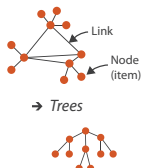
Attributes

- **network/graph**
  - nodes (vertices) connected by links (edges)
  - tree is special case: no cycles
    - often have roots and are directed

→ Tables



→ Networks

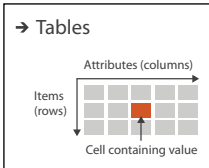


# Visualizing networks

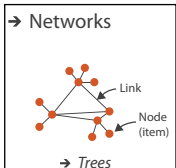
# Dataset types

Tables	Networks & Trees	Fields
Items	Items (nodes)	Grids
Attributes	Links	Positions
	Attributes	Attributes

→ Tables

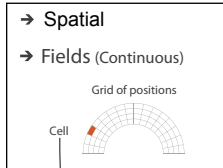


→ Networks



→ Spatial

→ Fields (Continuous)



23

## Spatial fields

- attribute values associated with cells
- cell contains value from continuous domain
  - eg temperature, pressure, wind velocity
- measured or simulated

→ Spatial

→ Fields (Continuous)

Grid of positions

Cell

Attributes (columns)

Value in cell

• weather station

• land

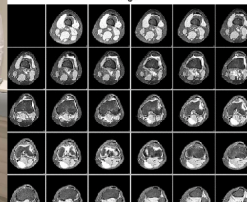
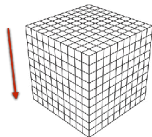
• population


• urban

500 km

## Spatial fields

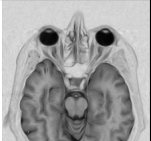
- attribute values associated with cells
- cell contains value from continuous domain
  - eg temperature, pressure, wind velocity
- measured or simulated
- beyond the scope of this class
  - sampling where attributes are measured
  - interpolation how to model attributes elsewhere
  - grid types





## Spatial fields

- attribute values associated with cells
- cell contains value from continuous domain
  - eg temperature, pressure, wind velocity
- measured or simulated
- beyond the scope of this class
  - sampling where attributes are measured
  - interpolation how to model attributes elsewhere
  - grid types, tensors



The image displays three distinct types of spatial fields. The top visualization is a scalar field, represented by a grayscale axial cross-section of a human brain, where each pixel's intensity corresponds to a scalar value. The middle visualization is a vector field, shown as a dense collection of small black arrows on a white background, representing the direction and magnitude of a vector quantity like wind velocity. The bottom visualization is a tensor field, depicted as a complex, multi-colored pattern of small dots or ellipses, where the color and shape of each element represent the properties of a tensor, such as stress or strain in a material.

# Dataset types

Tables

Items

Attributes

Networks & Trees

Items (nodes)

Links

Attributes

Fields

Grids

Positions

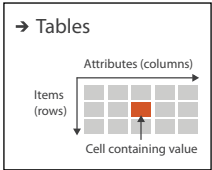
Attributes

Geometry

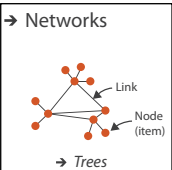
Items

Positions

→ Tables

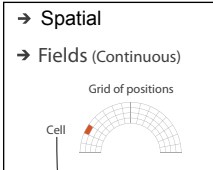


→ Networks

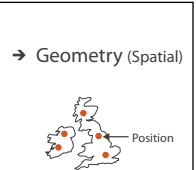


→ Spatial

→ Fields (Continuous)



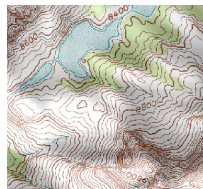

→ Geometry (Spatial)



27

## Geometry

- shape of items
- explicit spatial positions
- points, lines, curves, surfaces, regions
  - (volumes outside scope of class)
- boundary between computer graphics and visualization
  - graphics: geometry taken as given
  - vis: geometry is result of a design decision

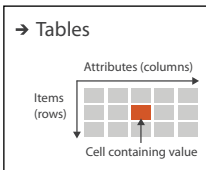


28

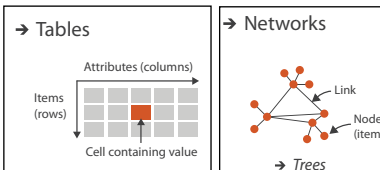
# Dataset types

Tables	Networks & Trees	Fields	Geometry	Clusters, Sets, Lists
Items	Items (nodes)	Grids	Items	Items
Attributes	Links	Positions	Positions	
	Attributes	Attributes		

→ Tables

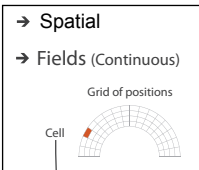


→ Networks

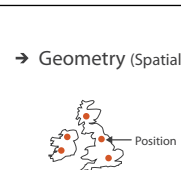


→ Spatial

→ Fields (Continuous)



→ Geometry (Spatial)



29

# Collections

- how we group items
- sets
  - unique items, unordered
- lists
  - ordered, duplicates possible
- clusters
  - groups of similar items

Rank	School Name	Academic repu	E	Facult	Citatio	I
Filters: <None>						
1.	Massachusetts inst					
2.	University of Camb					
3.	Harvard University	100 (1)			100 (1)	
4.	UCL (University Co					
5.	University of Oxfor					
6.	Imperial College L					
7.	Yale University					
8.	University of Chic					

# Dataset and data types

→ Data and Dataset Types

Tables	Networks & Trees	Fields	Geometry	Clusters, Sets, Lists
Items	Items (nodes)	Grids	Items	Items
Attributes	Links	Positions	Positions	
	Attributes	Attributes		

→ Data Types

→ Items    → Attributes    → Links    → Positions    → Grids

31

# Attribute types

- which classes of values & measurements?
- categorical (nominal)
  - compare equality
  - no implicit ordering
- ordered
  - ordinal
    - less/greater than defined
  - quantitative
    - meaningful magnitude
    - arithmetic possible

➡ Attribute Types

➔ Categorical

➔ Ordered

➔ Ordinal

➔ Quantitative

The diagram illustrates the hierarchy of attribute types. At the top is 'Attribute Types', which branches into 'Categorical' and 'Ordered'. 'Categorical' is further divided into 'Ordinal' and 'Quantitative'. 'Ordinal' is represented by three t-shirts of increasing size, and 'Quantitative' is represented by three horizontal lines of increasing length.

Table

A	B	C	S	T	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack	0.55	2/22/08
32	7/16/07	2-High	Small Pack	0.79	7/17/07
32	7/16/07	2-High	Jumbo Box	0.72	7/17/07
32	7/16/07	2-High	Medium Box	0.6	7/18/07
32	7/16/07	2-High	Medium Box	0.65	7/18/07
35	10/23/07	4-Not Specified	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1/20/05	5-Low	Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Specified	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Specified	Wrap Bag	0.6	6/6/05
70	12/18/06	5-Low	Small Box	0.59	12/23/06
70	12/18/06	5-Low	Wrap Bag	0.82	12/23/06
96	4/17/05	2-High	Small Box	0.55	4/19/05
97	1/29/06	3-Medium	Small Box	0.38	1/30/06
129	11/19/08	5-Low	Small Box	0.37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06
194	4/5/08	3-Medium	Wrap Bag	0.42	4/7/08

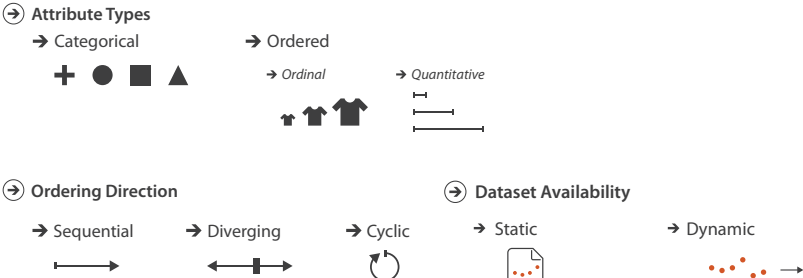
categorical  
ordinal  
quantitative

A	B	C	S	T	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack	0.55	2/22/08
32	7/16/07	2-High	Small Pack	0.79	7/17/07
32	7/16/07	2-High	Jumbo Box	0.72	7/17/07
32	7/16/07	2-High	Medium Box	0.6	7/18/07
32	7/16/07	2-High	Medium Box	0.65	7/18/07
35	10/23/07	4-Not Specified	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1/20/05	5-Low	Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Specified	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Specified	Wrap Bag	0.6	6/6/05
70	12/18/06	5-Low	Small Box	0.59	12/23/06
70	12/18/06	5-Low	Wrap Bag	0.82	12/23/06
96	4/17/05	2-High	Small Box	0.55	4/19/05
97	1/29/06	3-Medium	Small Box	0.38	1/30/06
129	11/19/08	5-Low	Small Box	0.37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06

Quiz: What kind of variable?

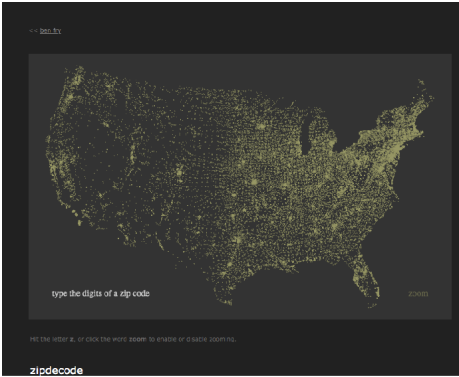
- 50 meter race times
- college major
- Amazon rating for product
- product name

Other data concerns



Hierarchical data

- multi-level structure
  - space
  - time
  - others
- example: zipcode



<https://benfry.com/zipcode/>

Data abstraction: Three operations

- translate from domain-specific language to generic visualization language
- identify dataset type(s), attribute types
- identify cardinality
  - how many items in the dataset?
  - what is cardinality of each attribute?
    - number of levels for categorical data
    - range for quantitative data
- consider whether to transform data
  - guided by understanding of task

Data vs conceptual models

- data model
  - mathematical abstraction
    - sets with operations, eg floats with \* / - +
    - variable data types in programming languages
- conceptual model
  - mental construction (semantics)
  - supports reasoning
  - typically based on understanding of tasks [stay tuned, next week]
- data abstraction process relies on conceptual model
  - for transforming data if needed

Data vs conceptual model, example

- data model: floats
  - 32.52, 54.06, -14.35, ...
- conceptual model
  - temperature
- multiple possible data abstractions
  - continuous to 2 significant figures: quantitative
    - task: forecasting the weather
  - hot, warm, cold: ordinal
    - task: deciding if bath water is ready
  - above freezing, below freezing: categorical
    - task: decide if I should leave the house today

Derived attributes

- derived attribute: compute from originals
  - simple change of type
  - acquire additional data
  - complex transformation
- more on this next time

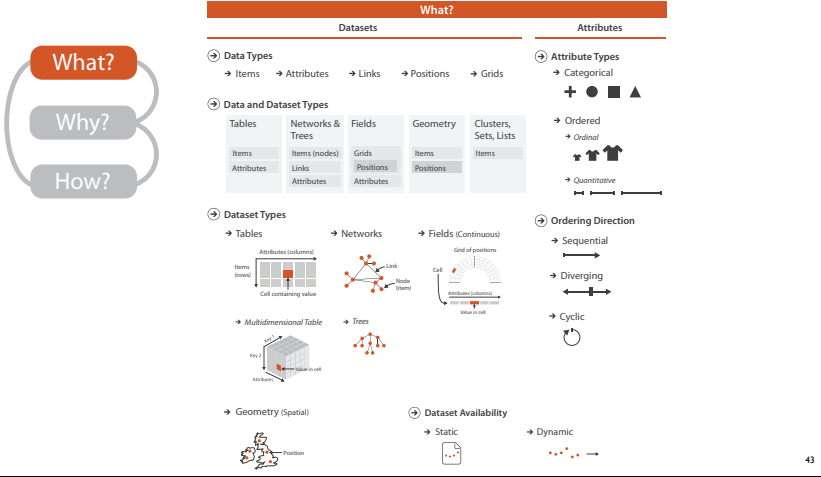


Data abstraction practice

- 2018 Central Park Squirrel Census

<https://www.thesquirrelcensus.com/>

<https://data.cityofnewyork.us/Environment/2018-Central-Park-Squirrel-Census-Squirrel-Data/vfnx-vebw>



Todo this week

- D3 videos to watch this week
  - refresher only if you need it: JS/HTML [90 min]
  - Intro to HTML/CSS/SVG [35 min]
  - Intro to D3.js [45 min]
- Quiz 1 to do this week, due by Fri Jan 10, 8am
- remember, no in-person labs this week!
- Foundations Exercise 1 out today (Thu Jan 9)
  - due Wed Jan 15

Credits

- Visualization Analysis and Design (Ch 2)
- Alex Lex & Miriah Meyer, <http://dataviscourse.net/>