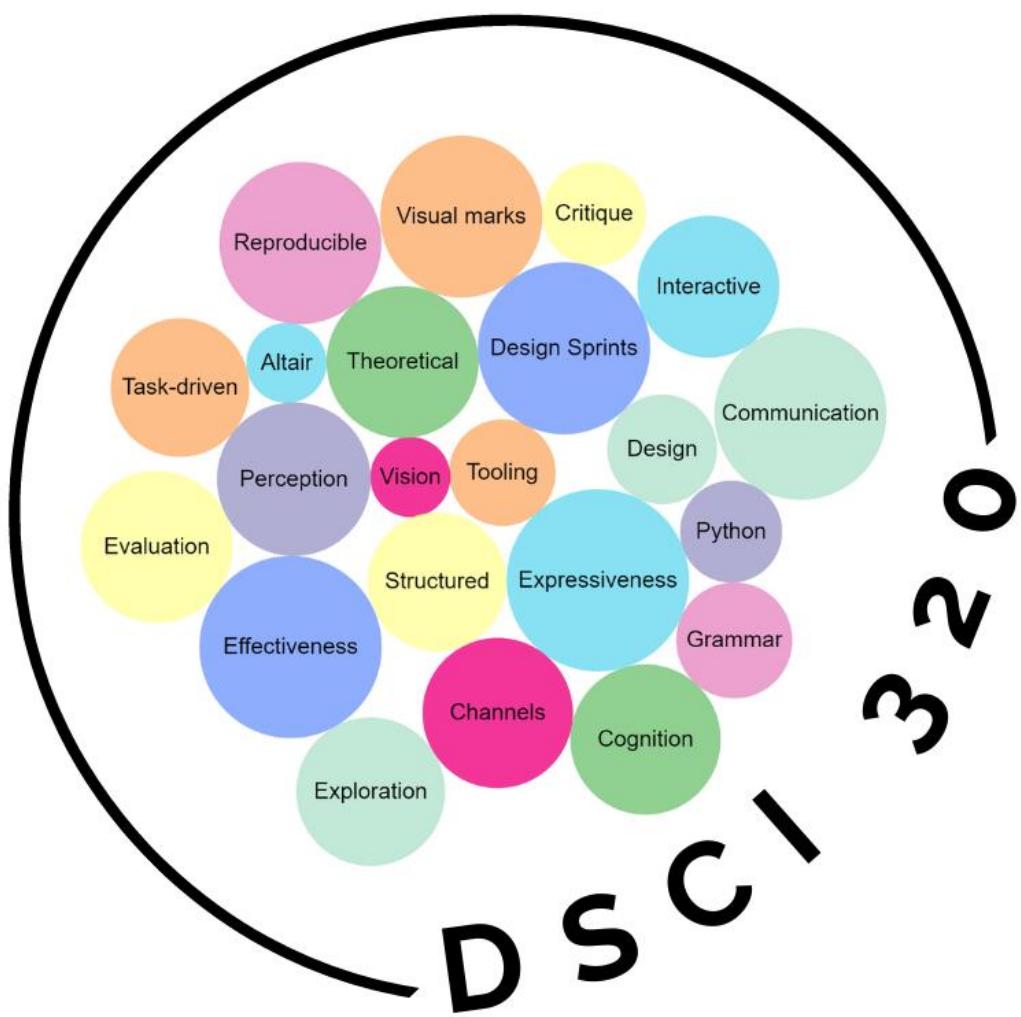
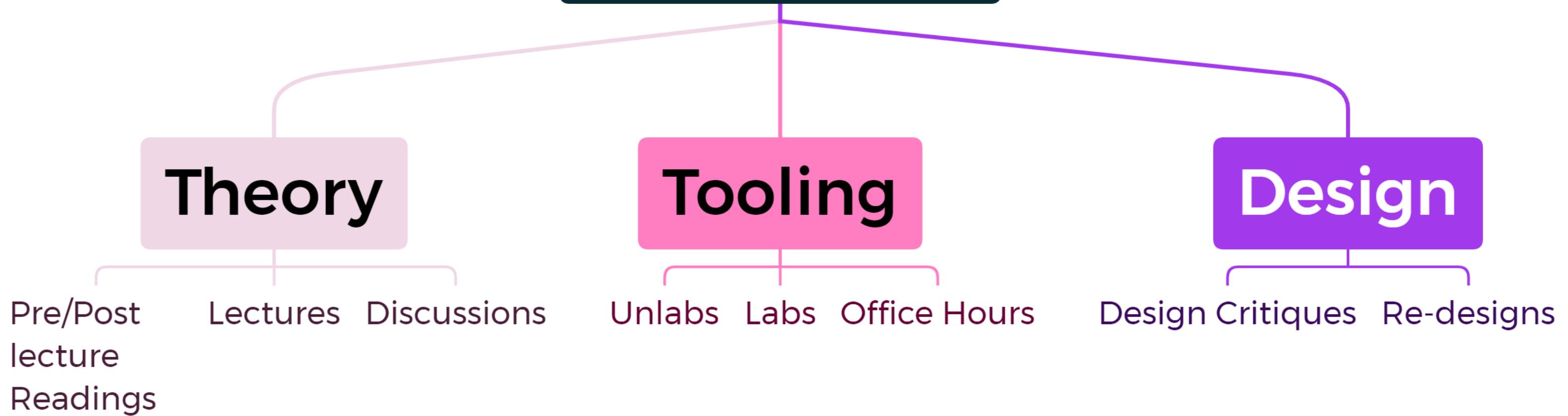


# Visualization for Data Science

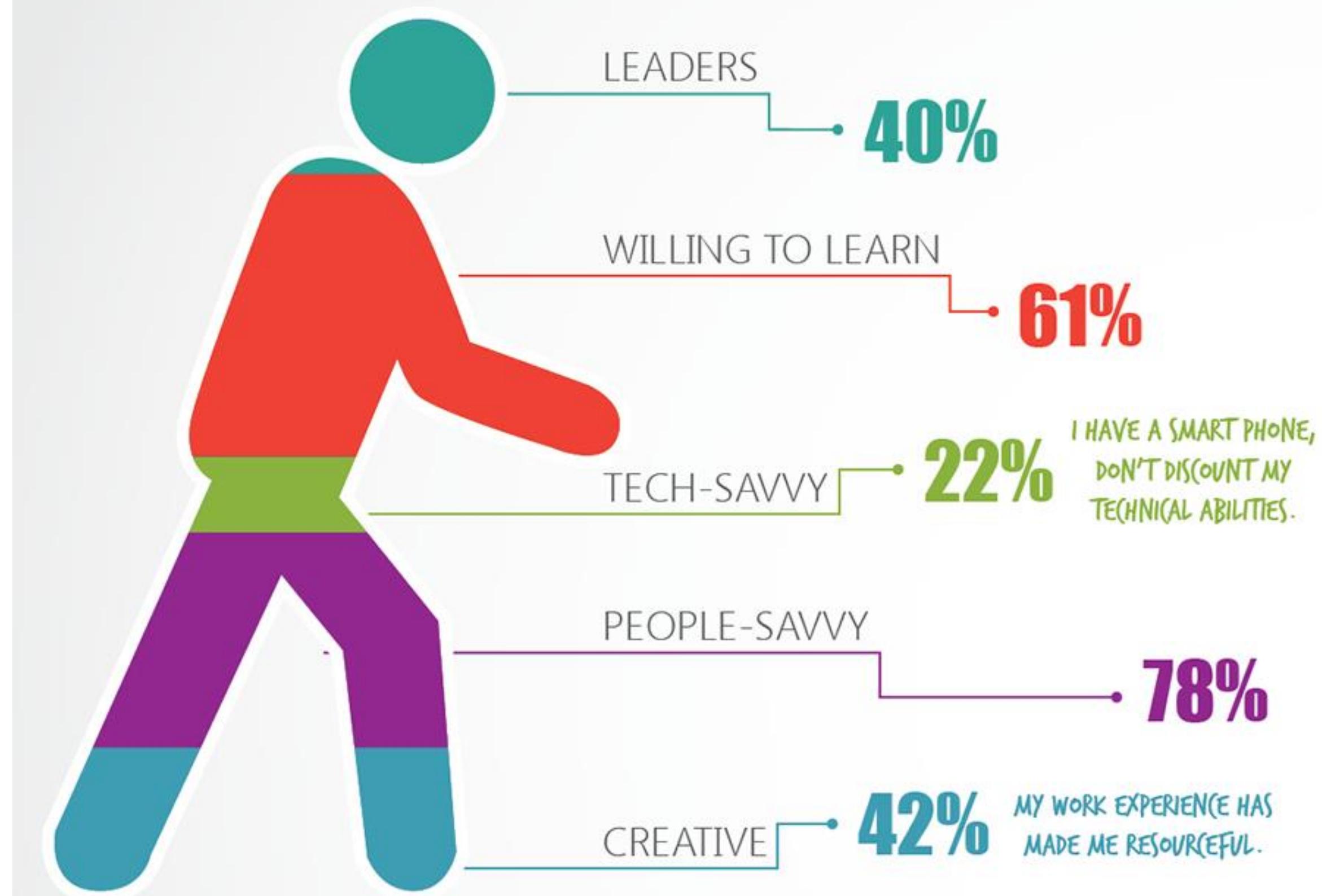
## Data Abstraction



# Viz4DSCI

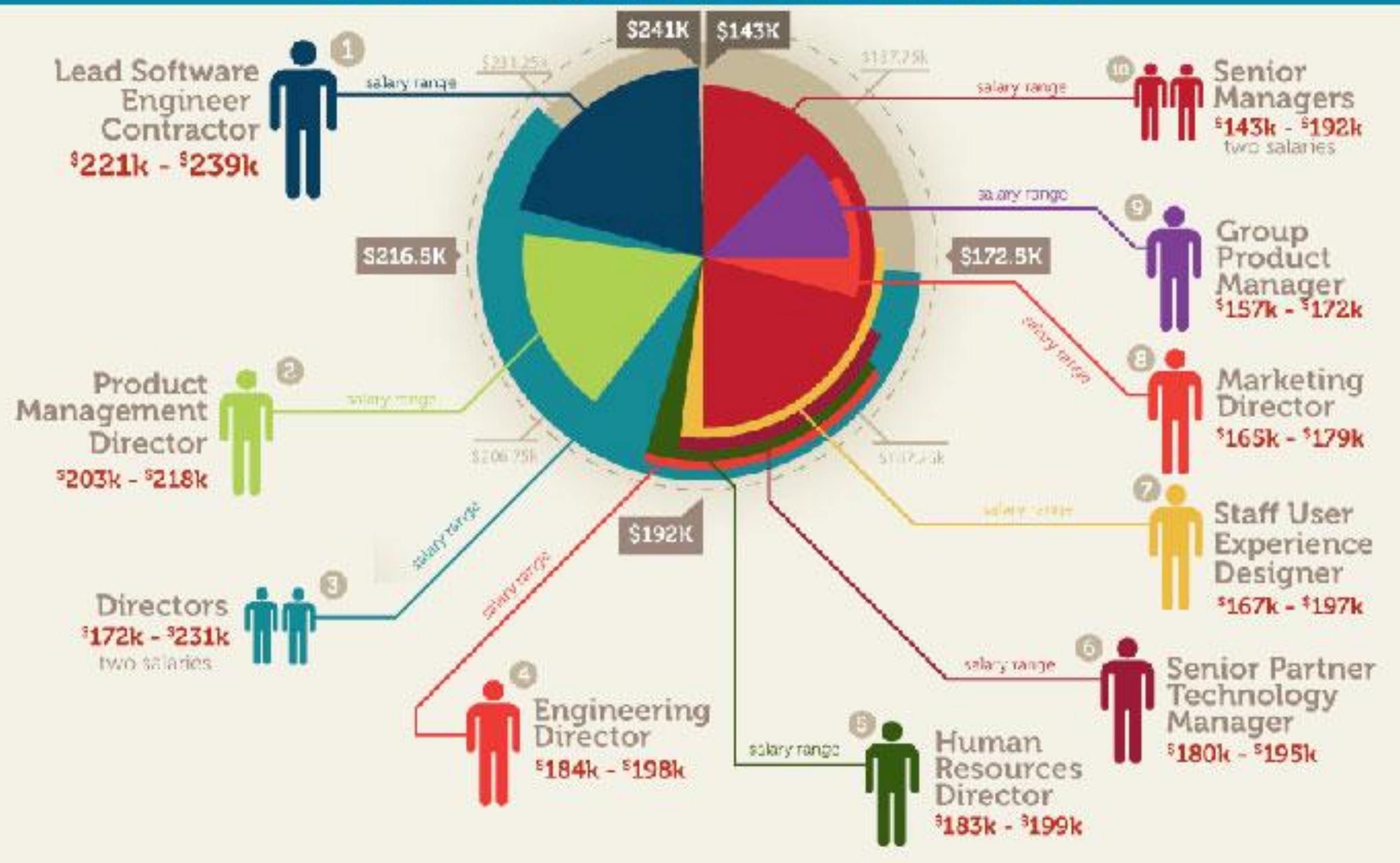


# HOW BABY BOOMERS DESCRIBE THEMSELVES



# top 10 salaries at Google™

RANGE FROM \$143,000 TO \$241,000 PER YEAR.



# Learning Goals

- Describe the difference between how the phrases “dataset types”, “data types” are used in vis. Literature as opposed to programming
- Describe the characteristics of data
- Differentiate between the different types of data and dataset types



**Data Types** are the fundamental units in which observed phenomena are represented. The structural or mathematical interpretation of data

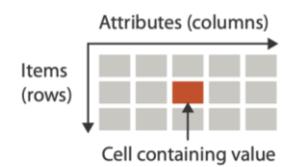
## Data Types

→ Items    → Attributes    → Links    → Positions    → Grids

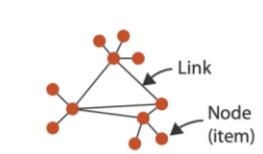
**Dataset** is any collection of information that is the target of analysis

### Dataset Types

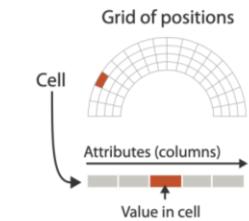
→ Tables



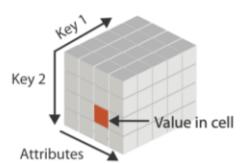
→ Networks



→ Fields (Continuous)



→ Multidimensional Table



→ Trees



→ Geometry (Spatial)



### Data and Dataset Types

#### Tables

Items  
Attributes

Items (nodes)  
Links  
Attributes

#### Networks & Trees

Grids  
Positions  
Attributes

#### Fields

Items  
Positions  
Attributes

#### Geometry

Items  
Positions

#### Clusters, Sets, Lists

Items

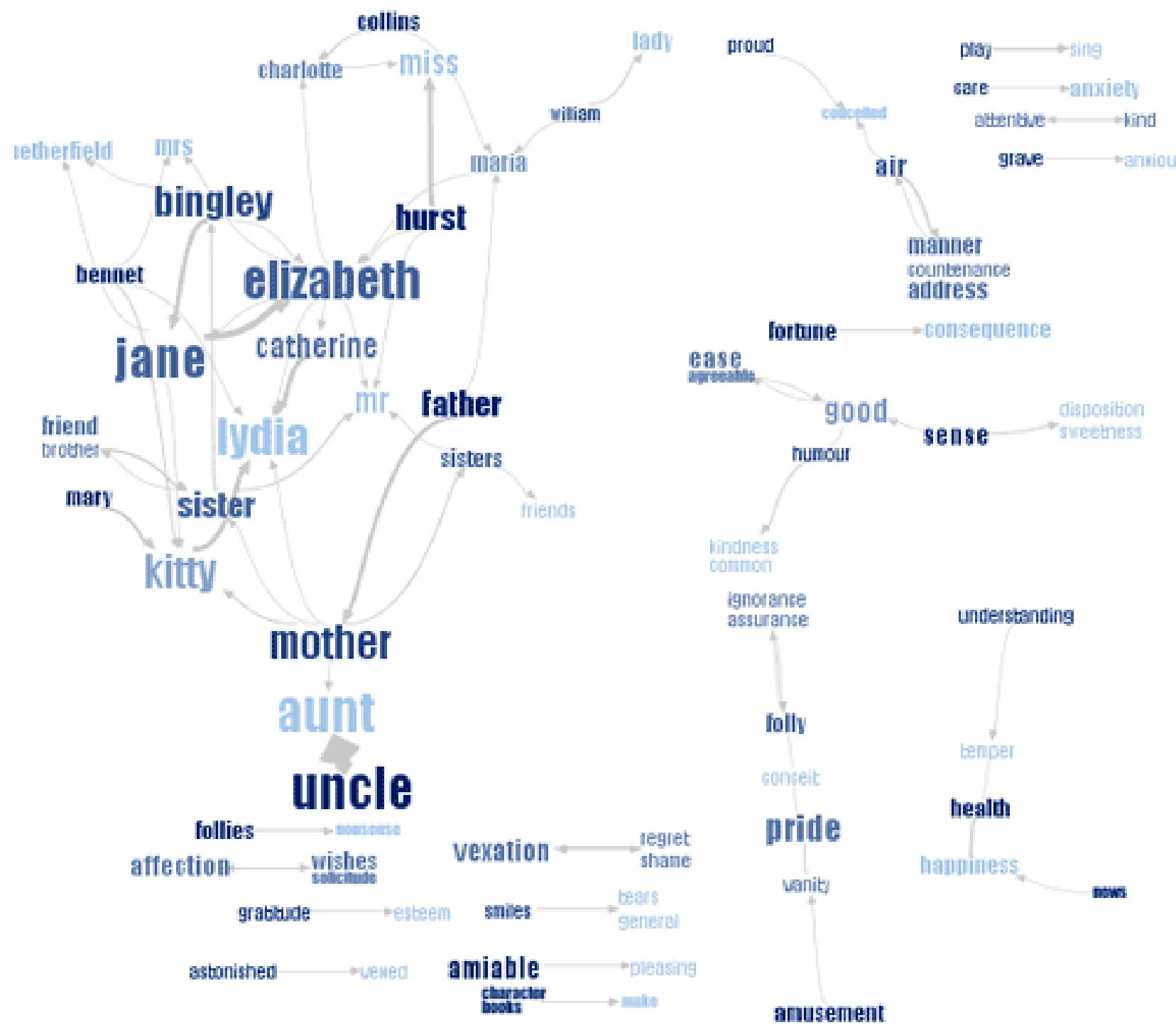
# Characterization

Data is characterized by its

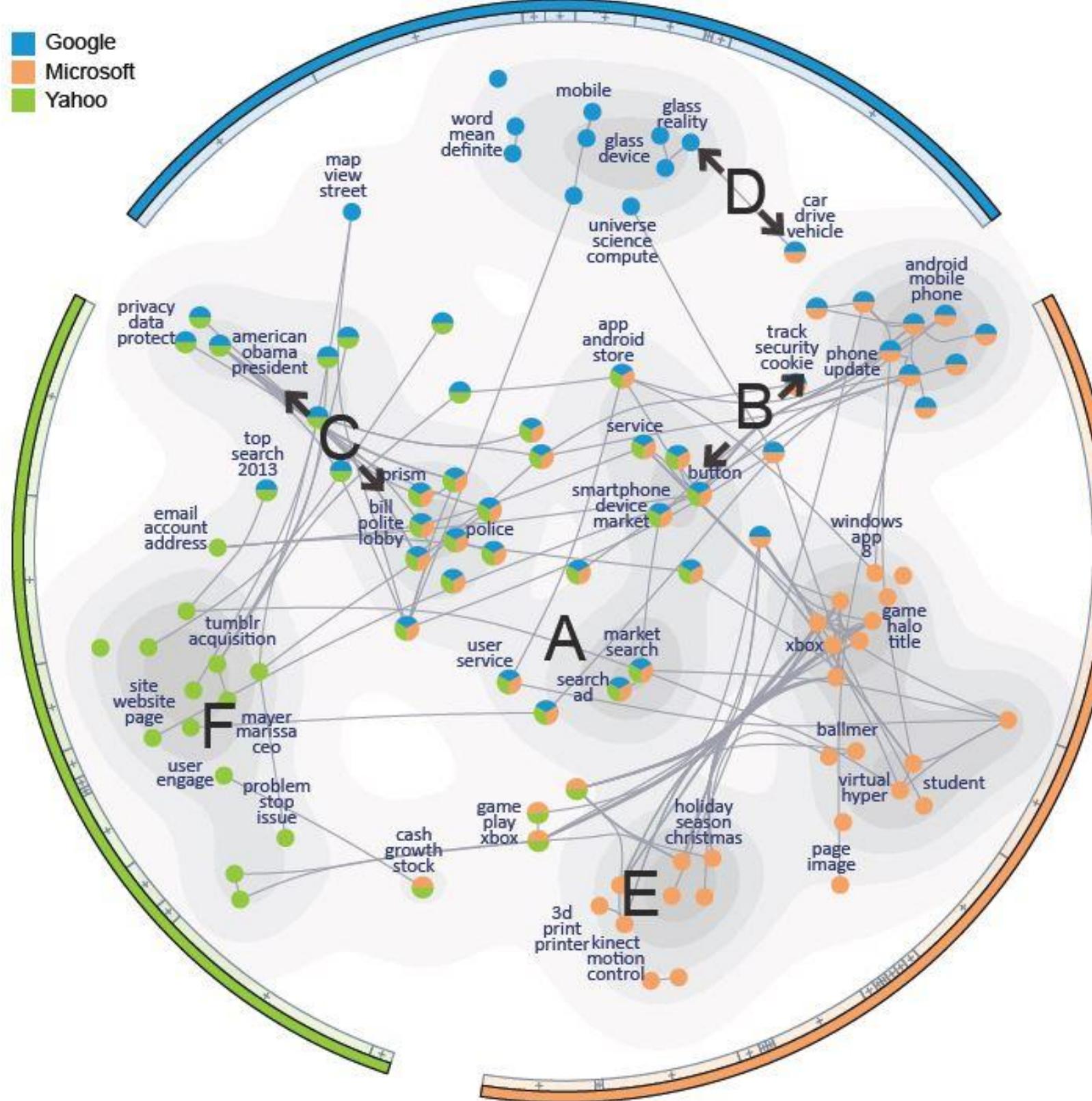
- Size (volume)
- Speed at which it generated (velocity)
- Quality (veracity)
- Structure



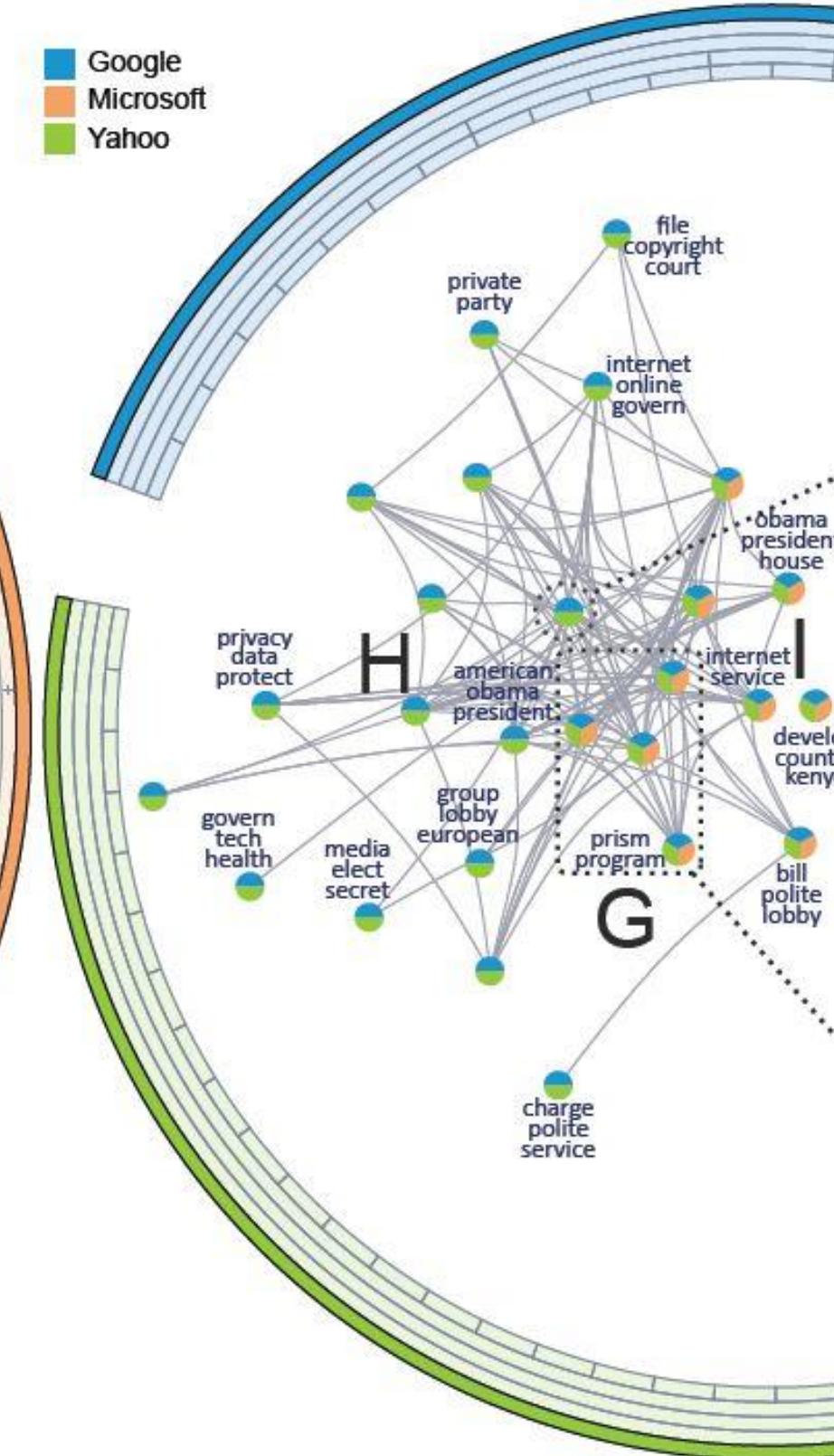
# Text Example



<http://hint.fm/projects/phrasenet/>



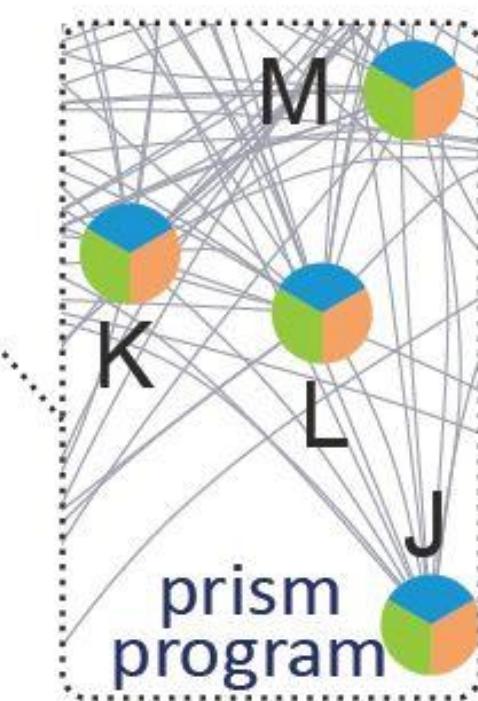
(a)



(1)



(c)



(d)

<http://shixiali.com/publications/TopicPanorama/index.htm>  
[http://shixiali.com/publications/TopicPanorama/video\\_eng.mp4](http://shixiali.com/publications/TopicPanorama/video_eng.mp4)

# What does data mean?

14, 2.6, 30, 30, 15, 100001

What does this sequence of six numbers mean?

# 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??

# What does data mean?

Basil, 7, S, Pear

What about this data?

# What does data mean?

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

# 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

# Semantics

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

# 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

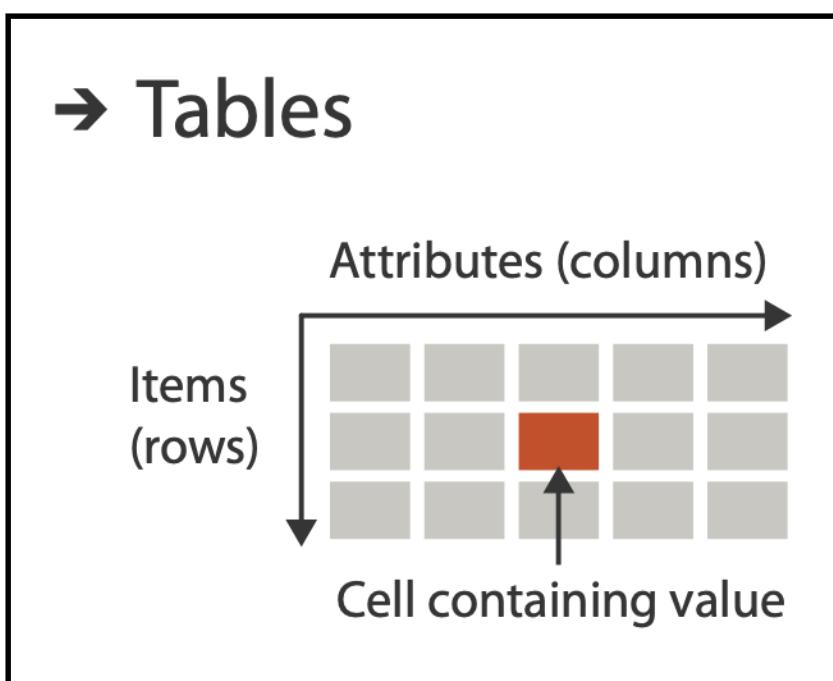
item: person

# Dataset types

## Tables

flat table

- one item per row
- each column is attribute
- cell holds value for item-attribute pair



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

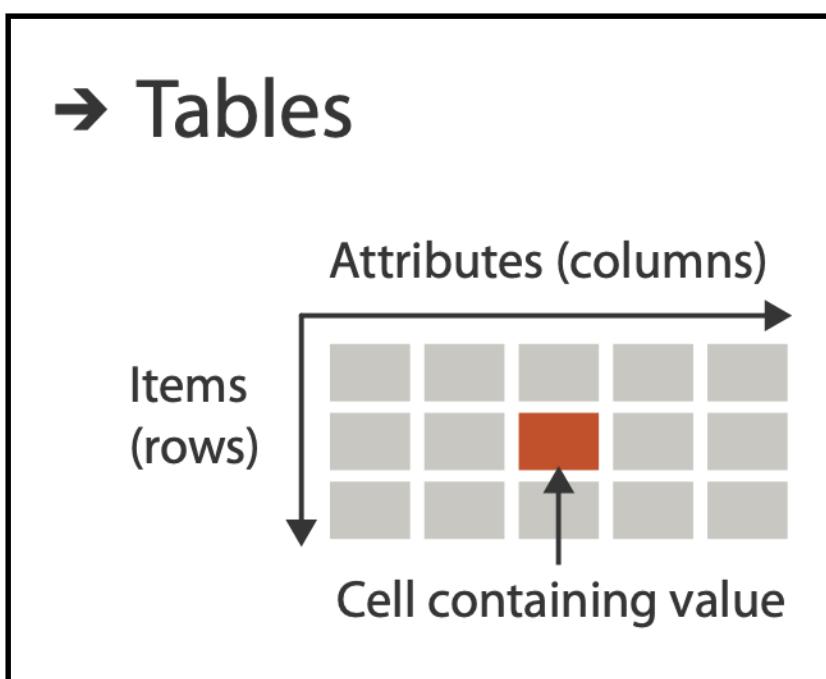
item: person

# Dataset types

## Tables

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

# Pulse Check

What is A?

What is B?

What is C?

A – Attribute

B – Item

C – Cell

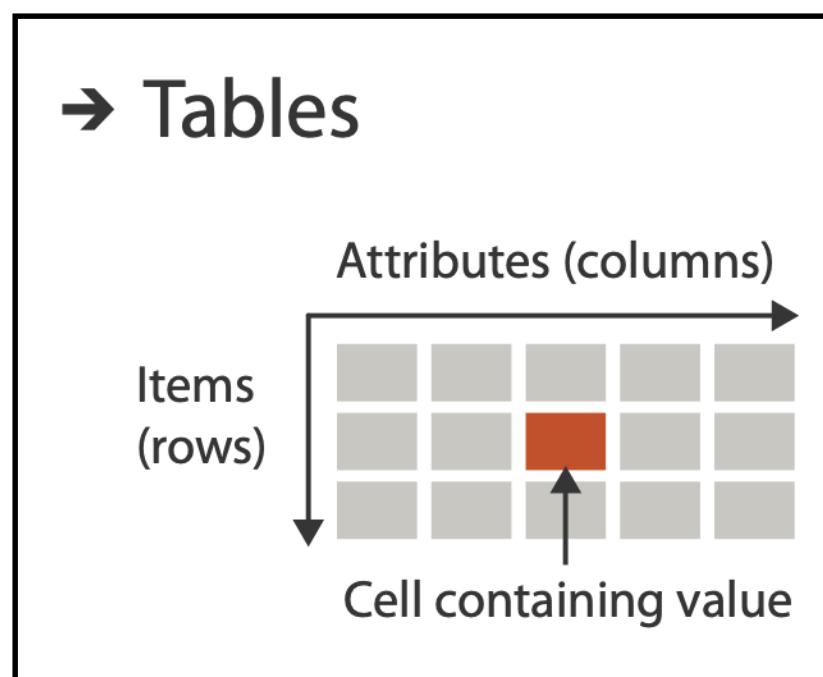
D – Data type

E – Dataset type

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

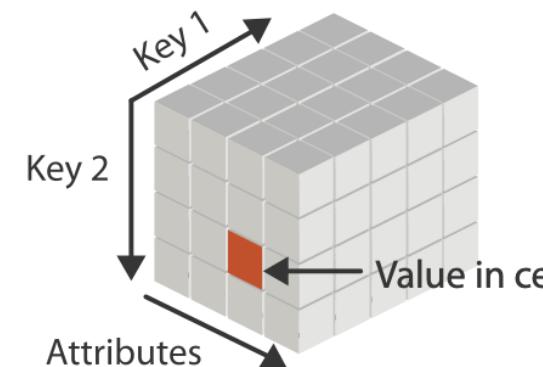
# Dataset types

## Tables



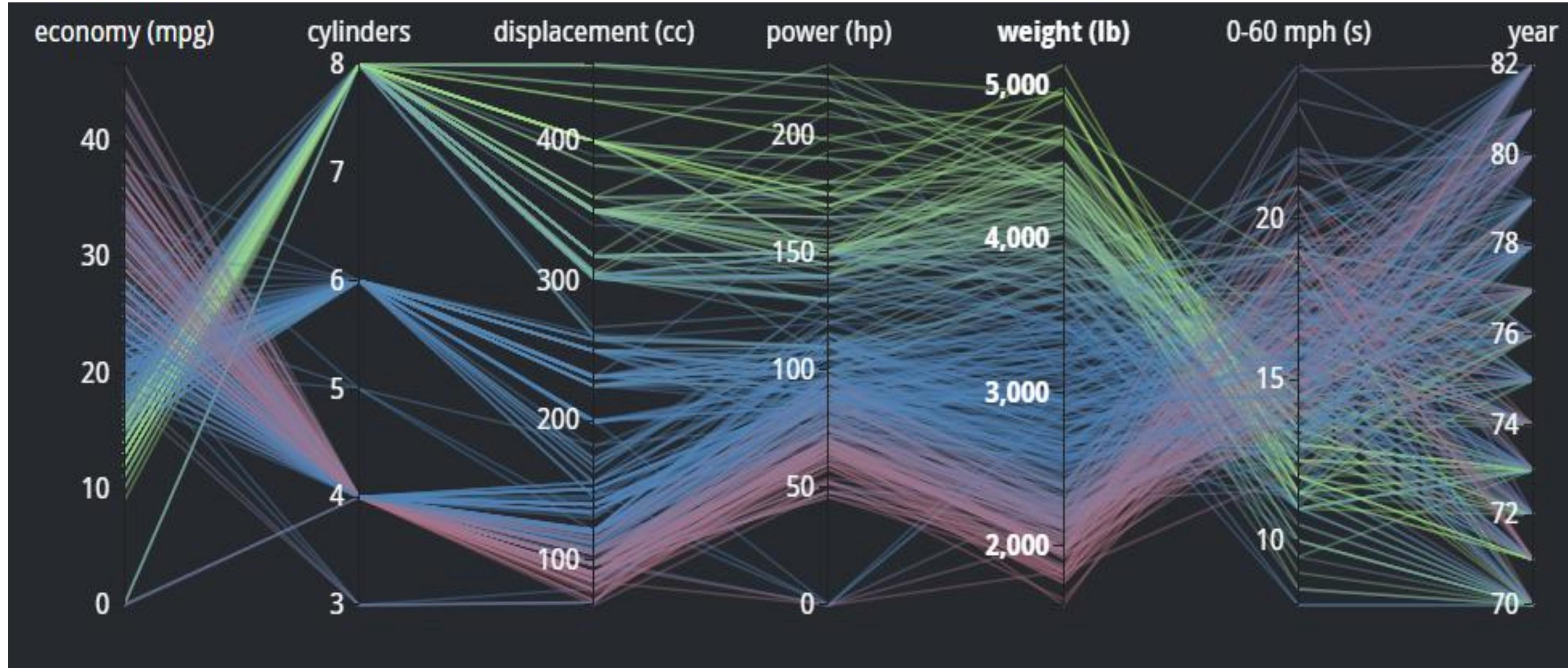
- multidimensional tables
  - indexing based on multiple keys
    - eg genes, patients

→ *Multidimensional Table*



	A	B	C	D	E	
1	A	B	C	D	E	
2	1	#	A	B	C	
3	2	1	#	A	B	
4	3	2	1	#1.2		
5	4	3	G	2	1500	529
6	5	4	L	T	GeneName	DESCRIPTION
7	6	5	P	4	LTF	TCGA-02-0001-01C-01R-0177-01
8	7	6	T	5	POSTN	TCGA-02-0003-01A-01R-0177-01
9	8	7	H	6	TMSL8	TCGA-02-0004-01A-01R-0298-01
10	9	8	R	7	HLA-DQA1	TCGA-02-0005-01A-01R-0298-01
11	10	9	S	8	RP11-35N6.1	TCGA-02-0006-01A-01R-0298-01
12	11	10	D	9	RP11-35N6.1	TCGA-02-0007-01A-01R-0298-01
13	12	11	I	10	STMN2	TCGA-02-0008-01A-01R-0298-01
14	13	12	I	11	DCX	TCGA-02-0009-01A-01R-0298-01
15	14	13	I	11	AGXT2L1	TCGA-02-0010-01A-01R-0298-01
16	15	14	I	12	IL13RA2	TCGA-02-0011-01A-01R-0298-01
17	16	15	M	13	SLN	TCGA-02-0012-01A-01R-0298-01
18	17	16	N	14	MEOX2	TCGA-02-0013-01A-01R-0298-01
19	18	17	F	15	COL11A1	TCGA-02-0014-01A-01R-0298-01
20	19	18	C	16	NNMT	TCGA-02-0015-01A-01R-0298-01
21	20	19	F	17	F13A1	TCGA-02-0016-01A-01R-0298-01
22	21	20	T	18	CXCL14	TCGA-02-0017-01A-01R-0298-01
					MBP	TCGA-02-0018-01A-01R-0298-01
					TF	TCGA-02-0019-01A-01R-0298-01
					KCND2	TCGA-02-0020-01A-01R-0298-01
						-1.777692395

# Visualizing Tables: Car Data Parallel Coordinates

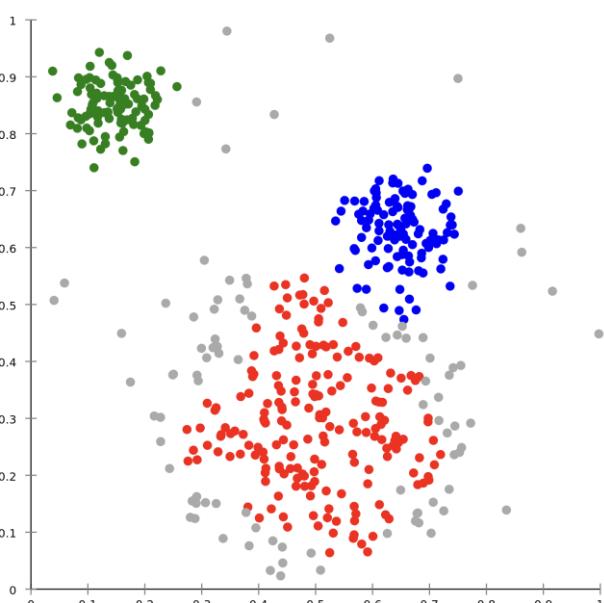
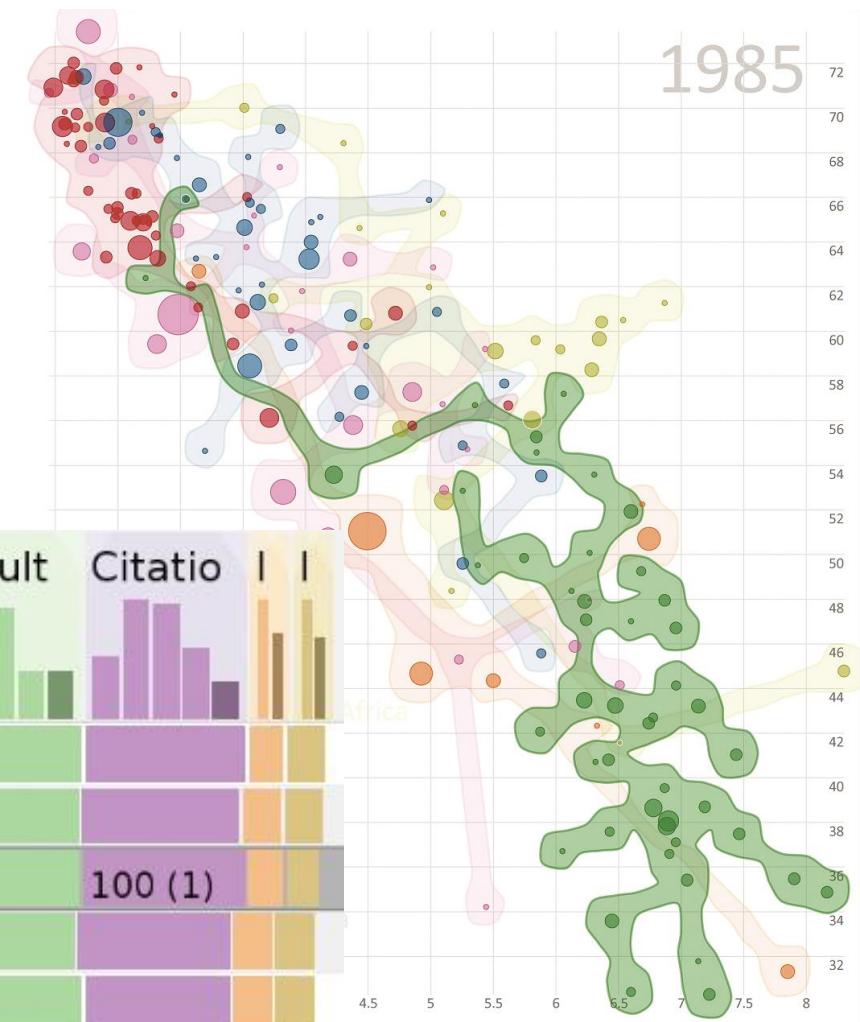
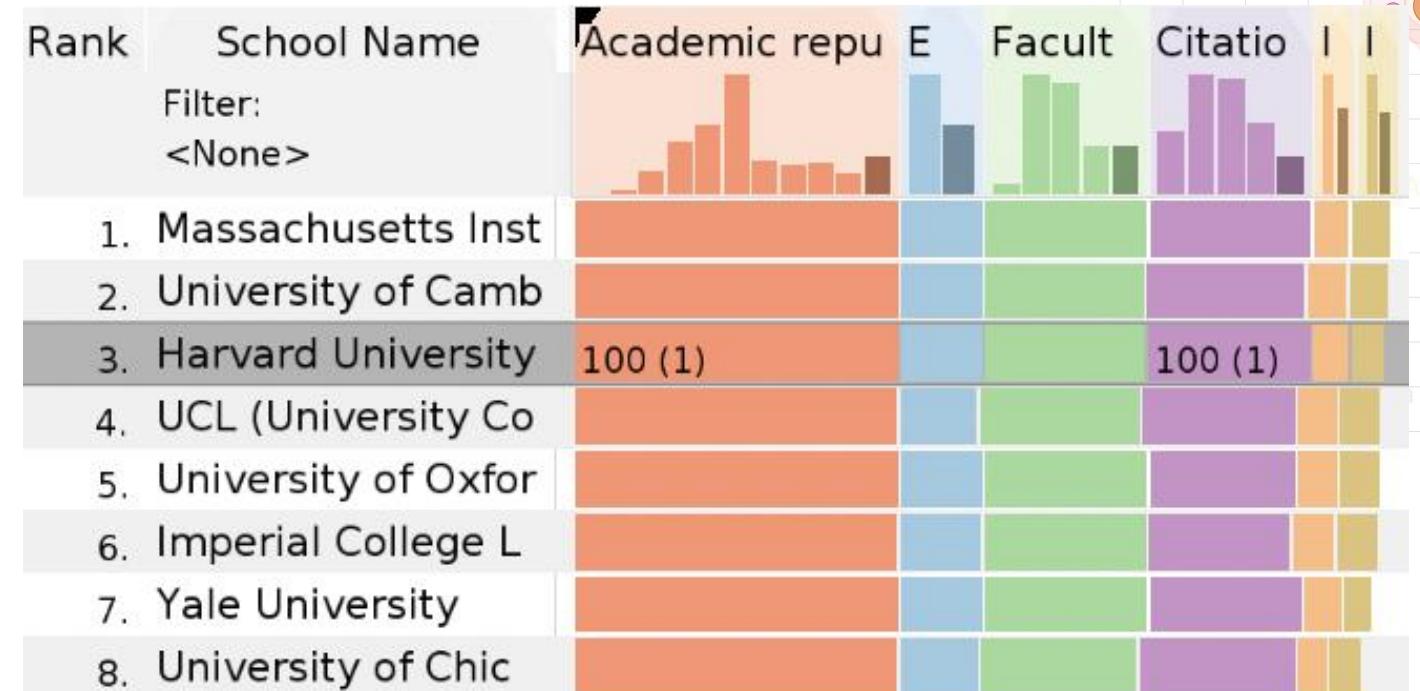


# Collections

how we group items

sets

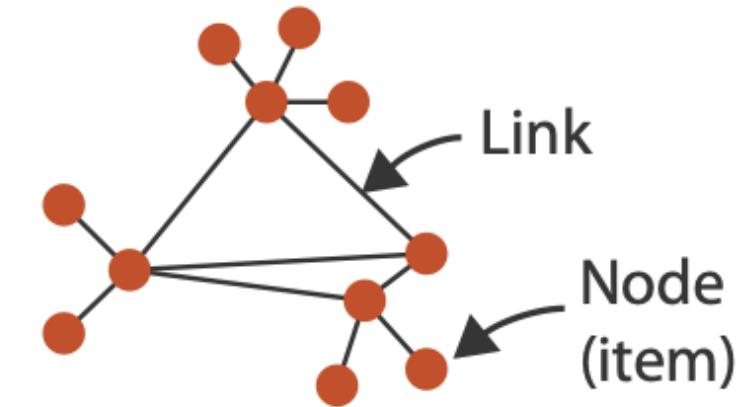
- unique items, unordered lists
- ordered, duplicates possible clusters
- groups of similar items



# Other data types

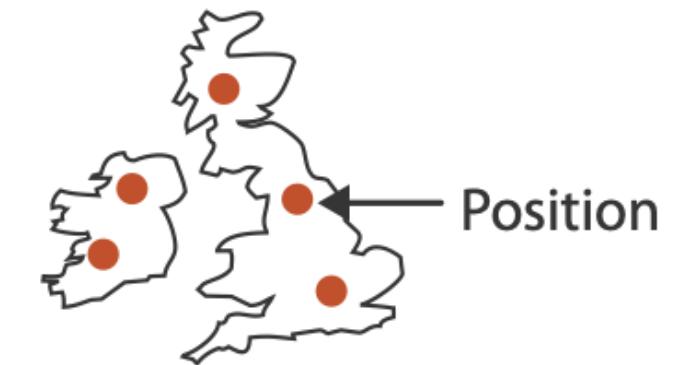
## Links

- express relationship between two items
- e.g. friendship on facebook, interaction between proteins



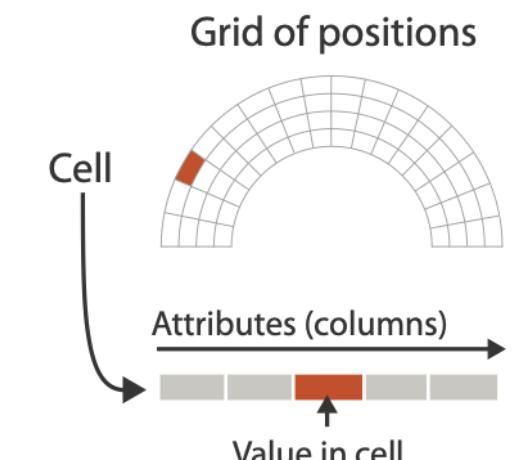
## Positions

- spatial data: location in 2D or 3D
- e.g. pixels in photo, voxels in MRI scan, latitude/longitude



## Grids

- sampling strategy for continuous data



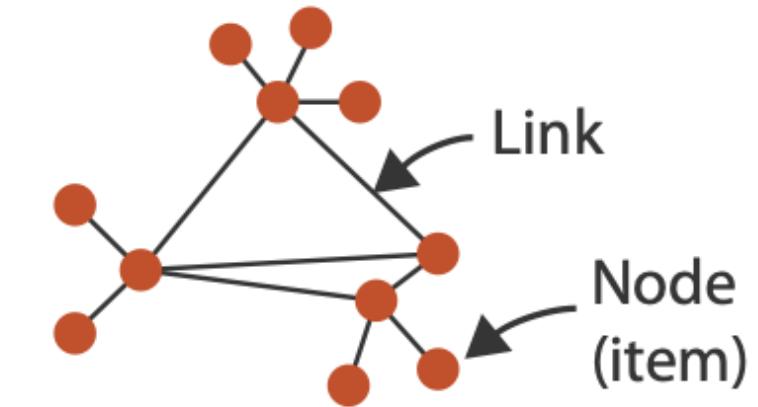
# Graphs/Networks

Used to express relationship between two items

Items (nodes) are connected with links.

Examples

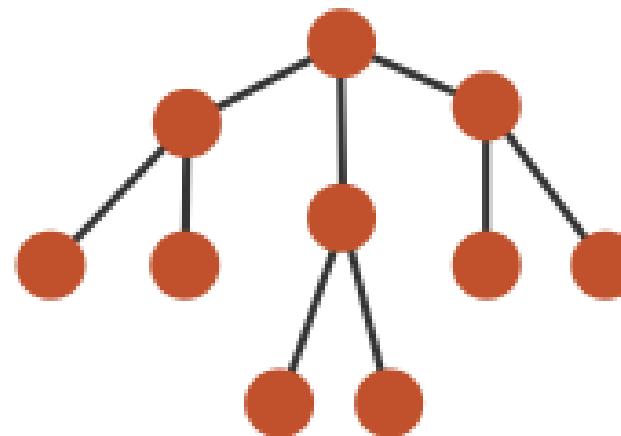
- social networks
- Power grids
- Road networks



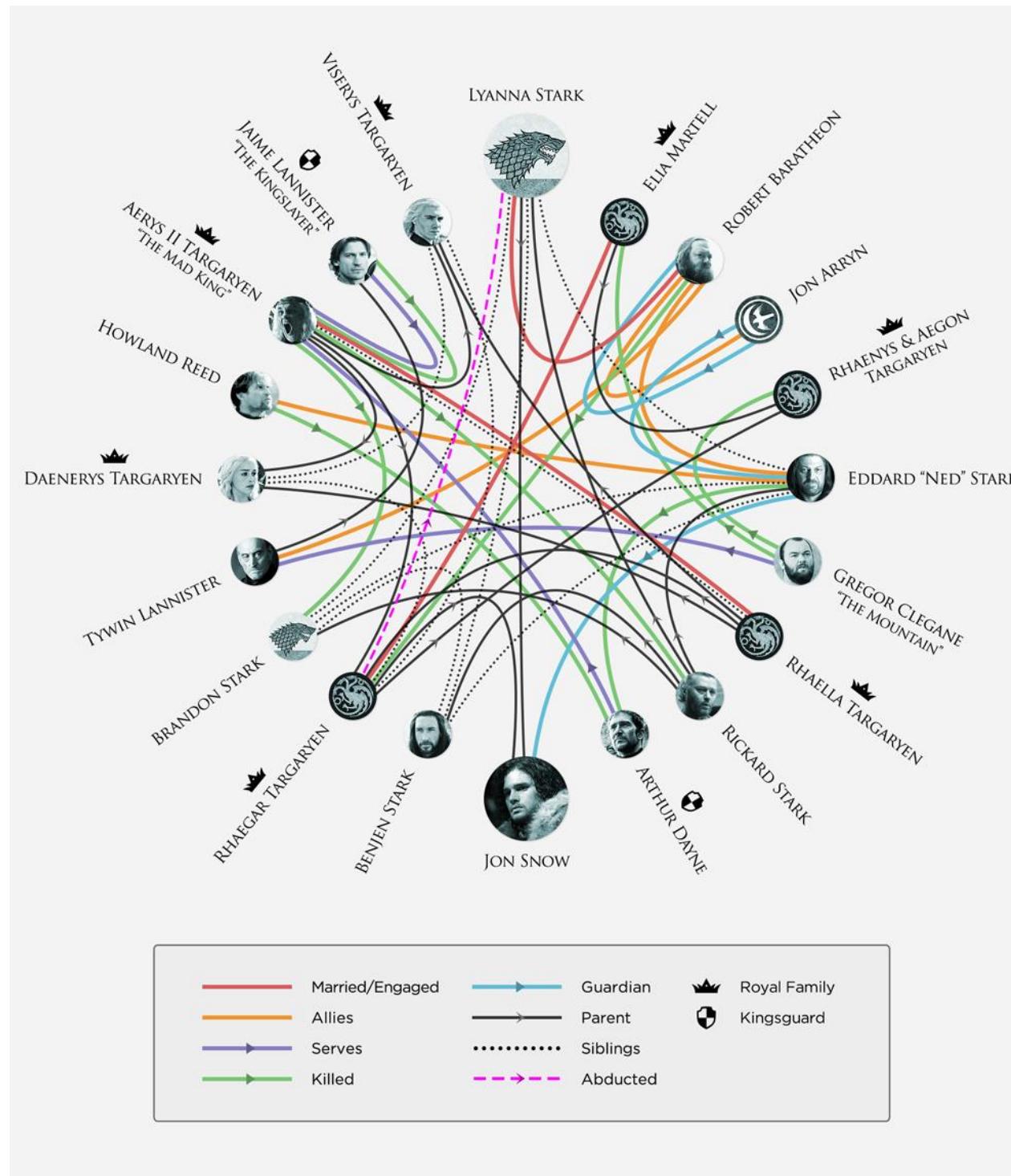
A tree is a subset of graphs

Basically a graph with no cycles.

Trees typically have roots and are directed.

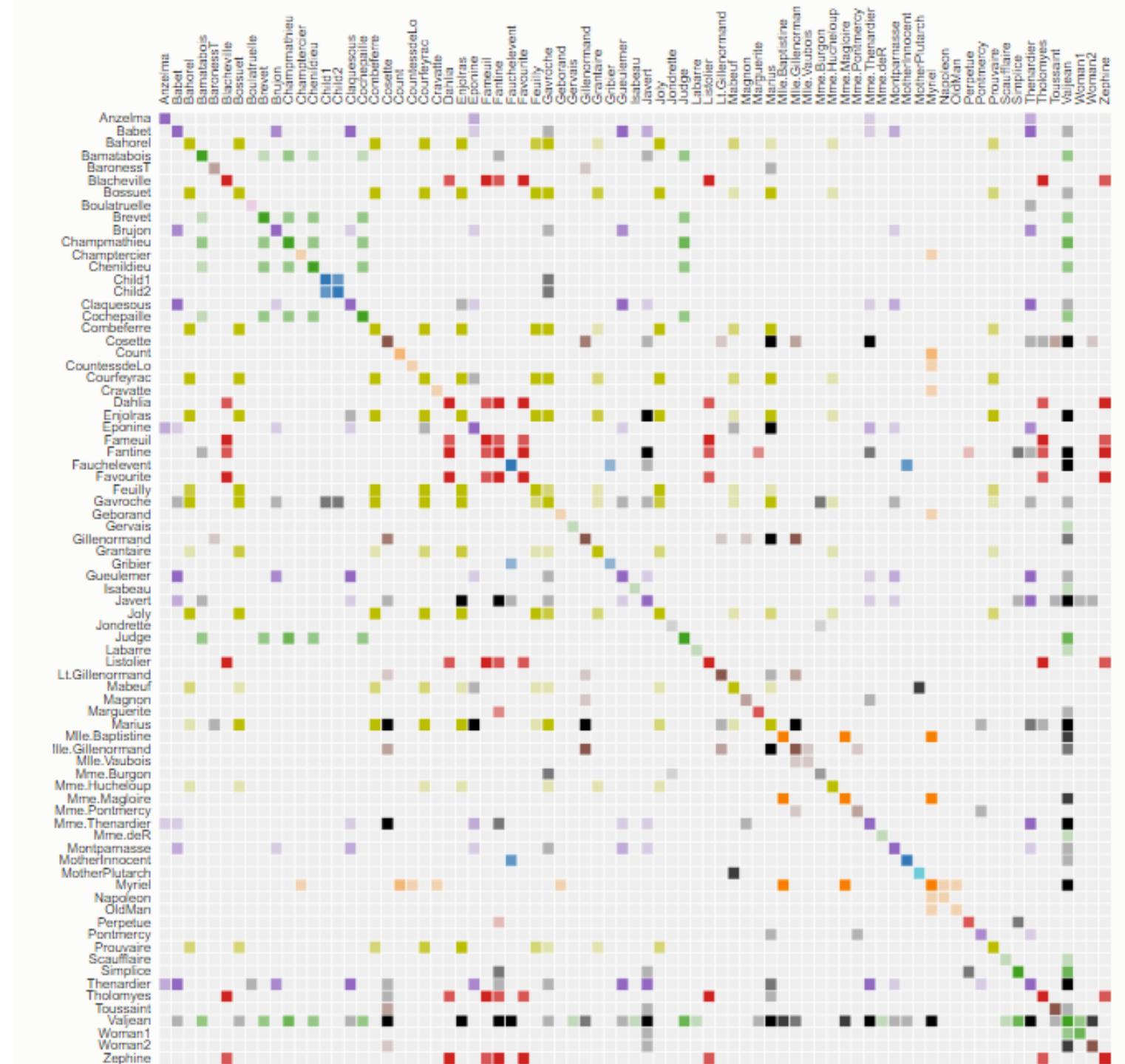


# Visualizing Networks/Graphs



<https://imgur.com/PG0963W>

# Les Misérables Co-occurrence



<https://bostocks.org/mike/miserables/>

# Visualizing Trees

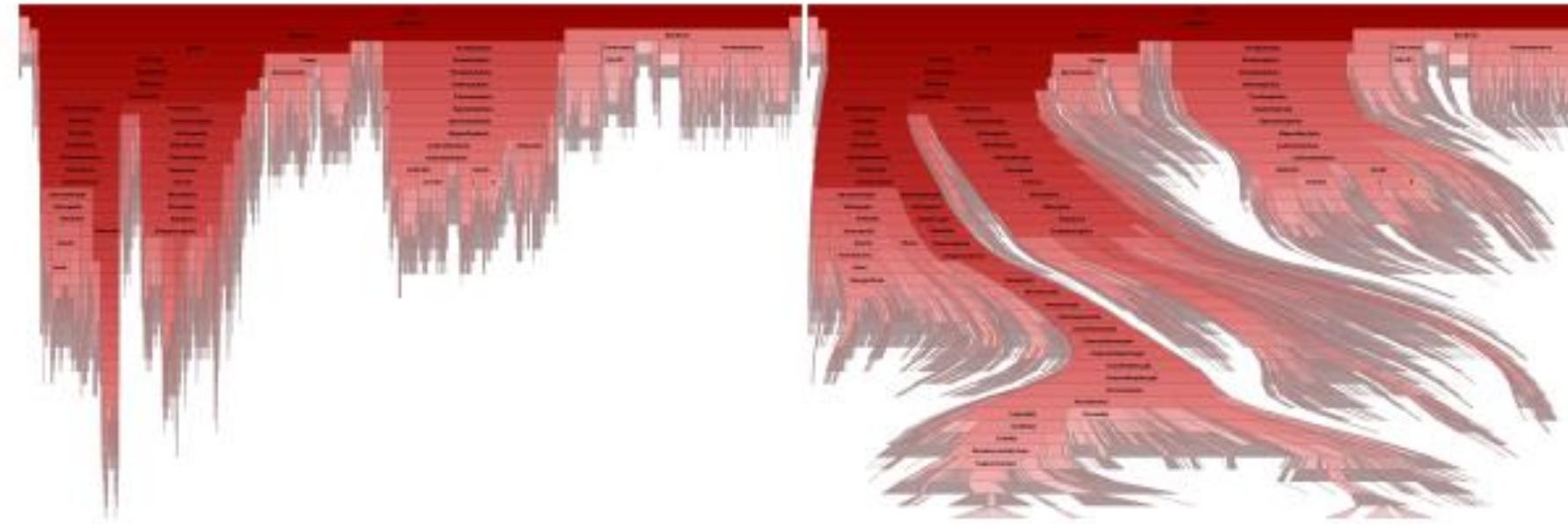


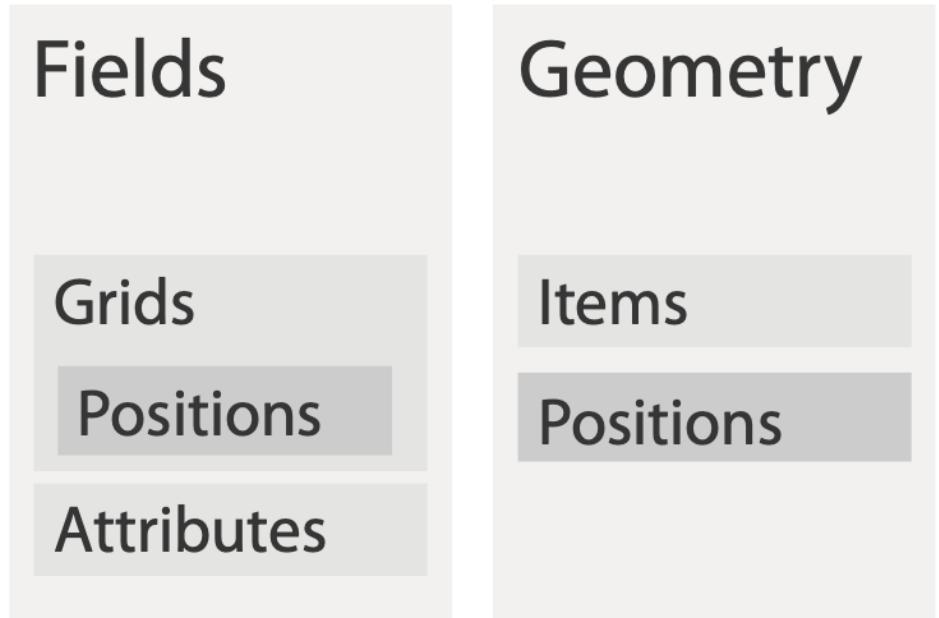
Figure 1: Two visualizations of the NCBI taxonomy dataset [13, 29]: Standard icicle plot (left), space-reclaiming icicle plot (right).



<https://ieeexplore.ieee.org/document/9086292>

<https://observablehq.com/@d3/zoomable-sunburst>

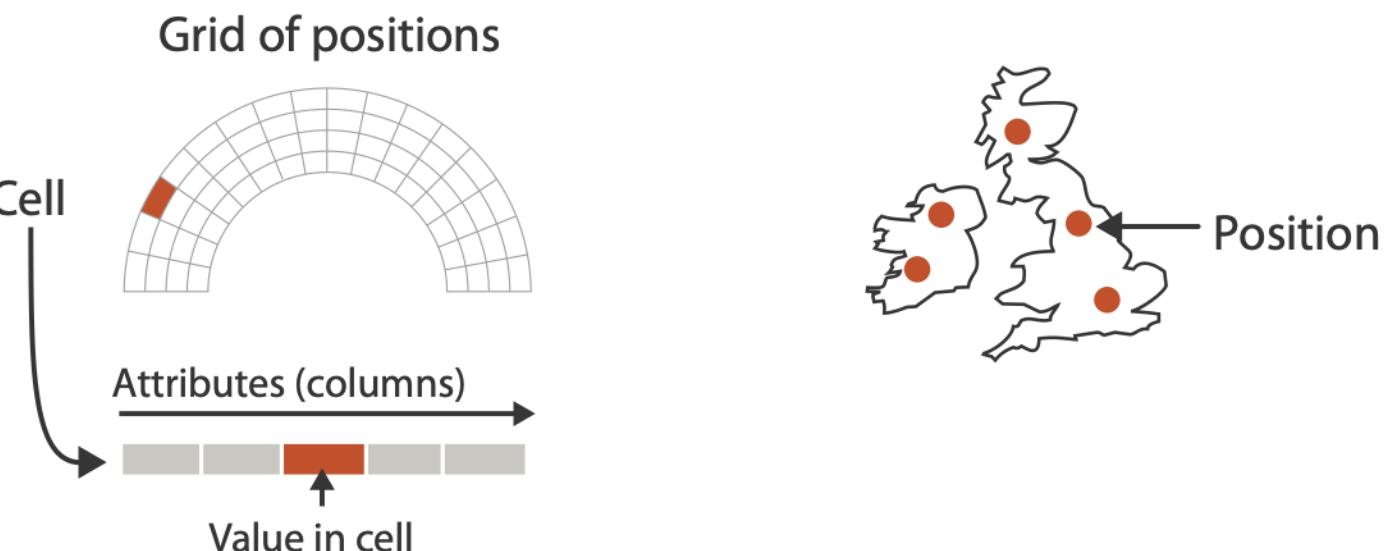
# Dataset types



→ Spatial

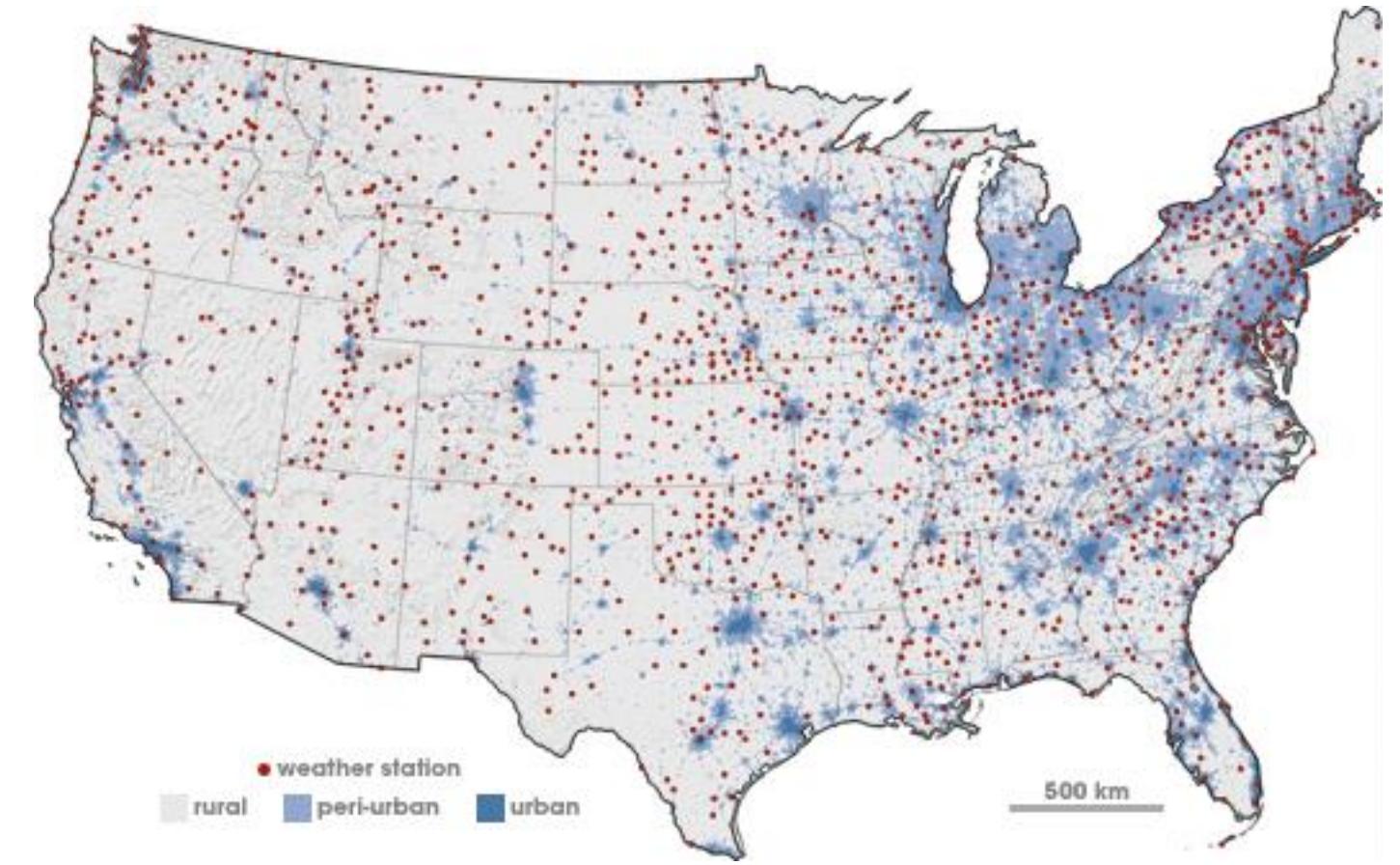
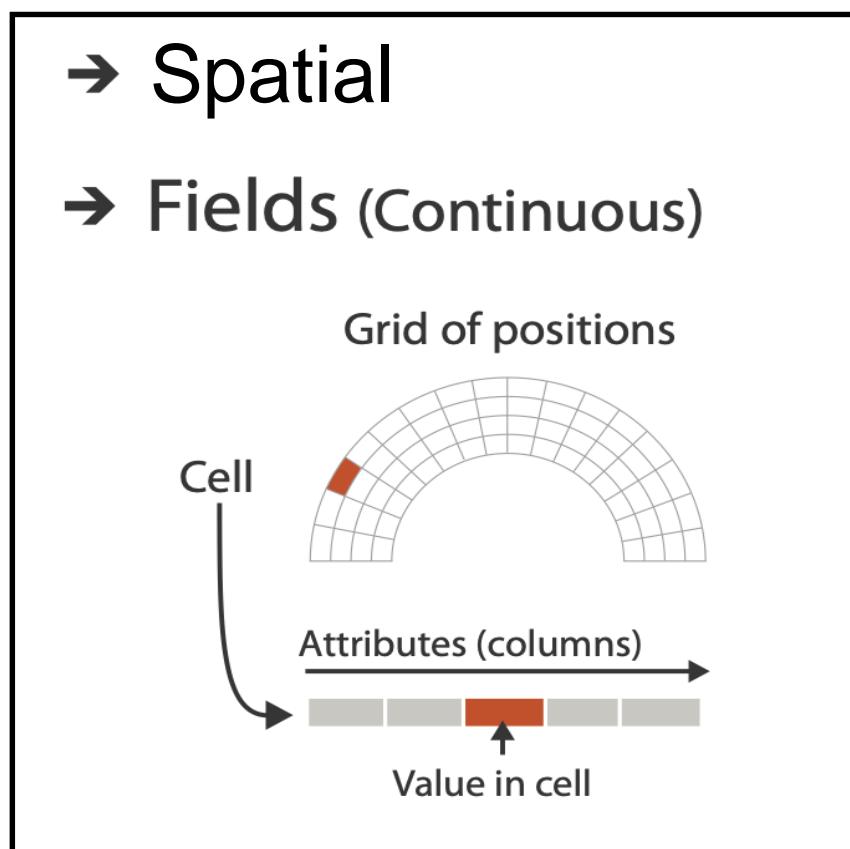
→ Fields (Continuous)

→ Geometry (Spatial)



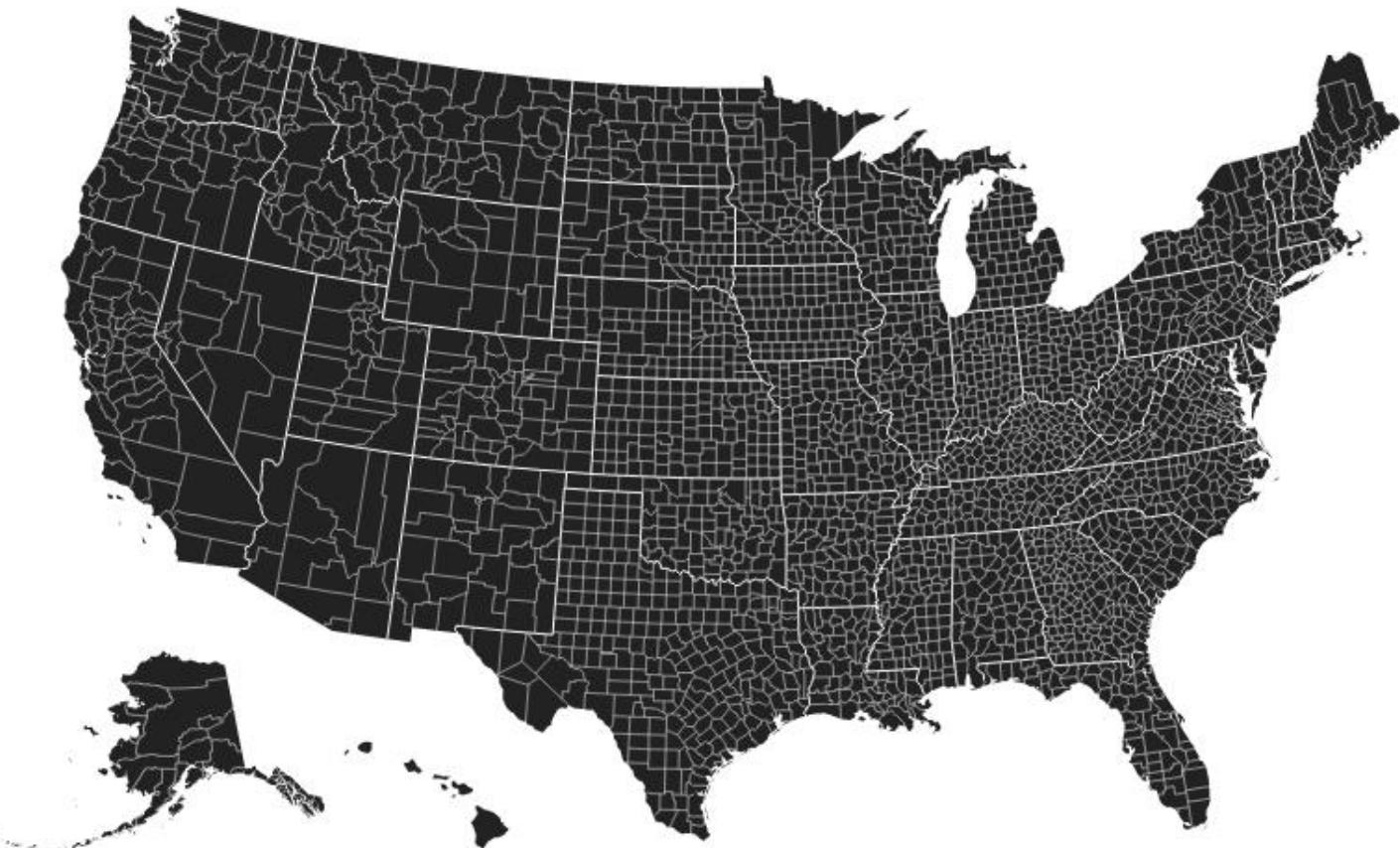
# Spatial fields

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



# Geometry

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



# Attribute Types

# 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	

## Data Types

→ Items

→ Attributes

→ Links

→ Positions

→ Grids

# 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

→ Categorical

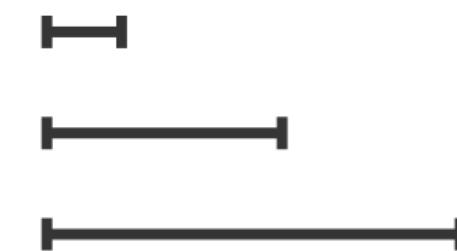


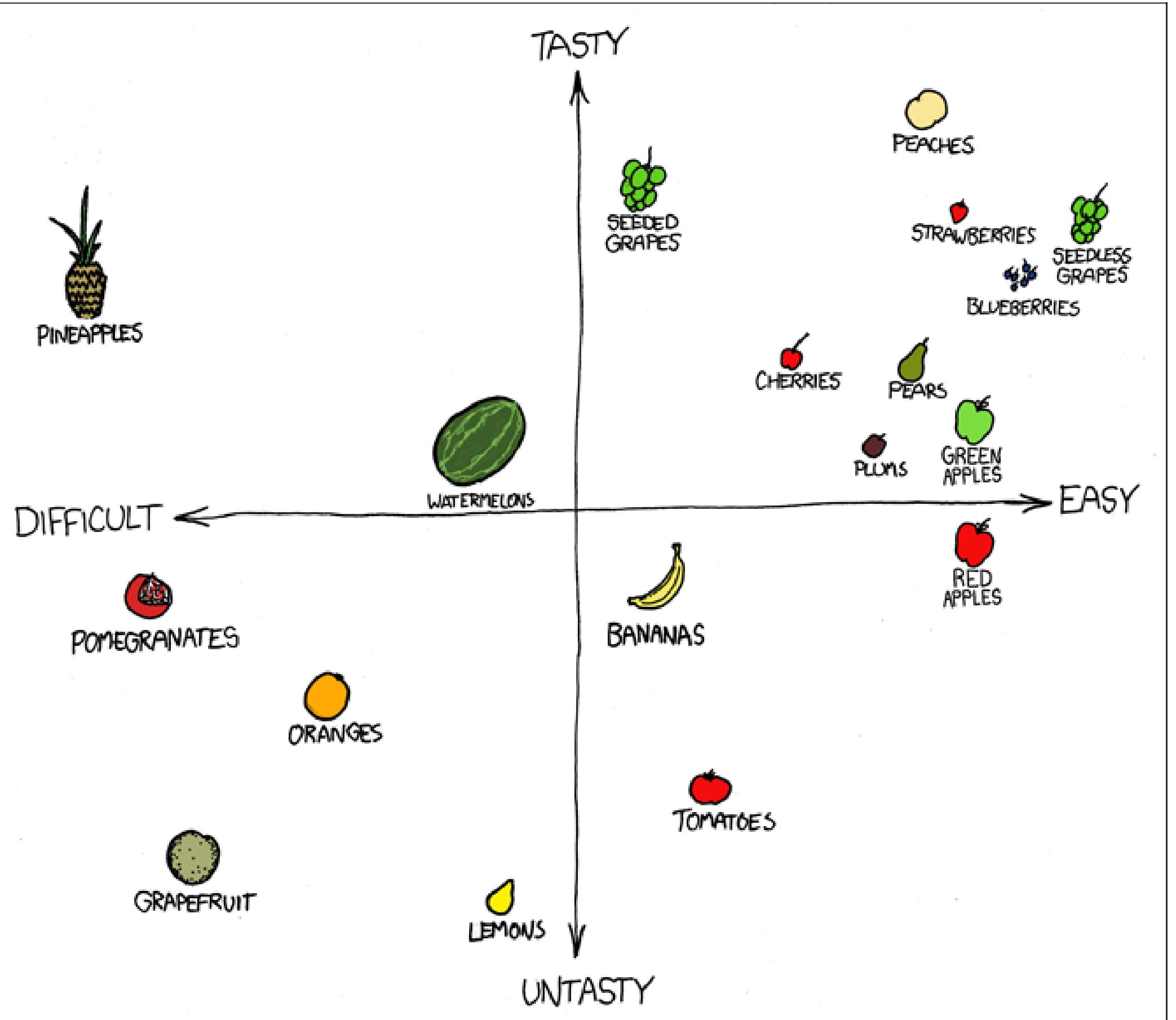
→ Ordered

→ *Ordinal*



→ *Quantitative*





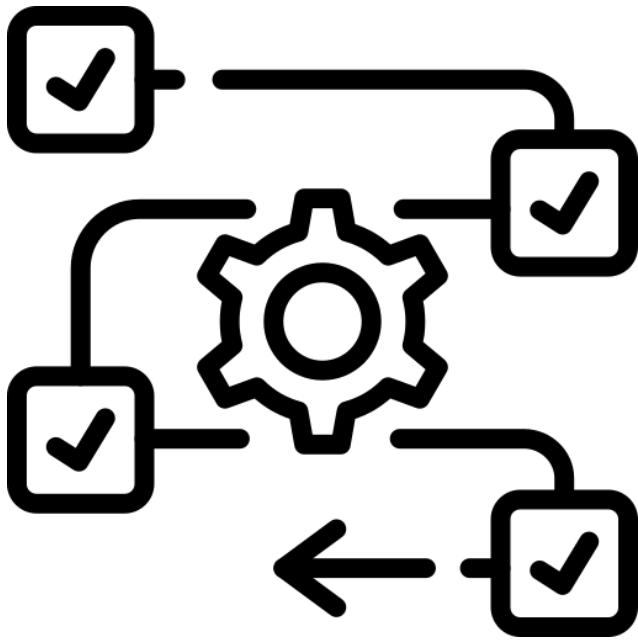
1. What is an item
2. What is an attribute
3. What is a feature
4. What is the semantics
5. What is the keys
6. What is type for Order ID
7. What is the type for Order Date
8. What is the type for Order Priority
9. What is the type for Product Container
10. What is the type for Product Base margin
11. What is the type for Ship Date

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

City	Condition	Temperature
<a href="#">Calgary</a>	Mainly Clear	-4°C
<a href="#">Charlottetown</a>	Light Snowshower	-6°C
<a href="#">Edmonton</a>		-7°C
<a href="#">Fredericton</a>	Clear	-9°C
<a href="#">Halifax</a>	Light Snow	-6°C
<a href="#">Iqaluit</a>	Clear	-28°C
<a href="#">Montréal</a>	Mainly Clear	-9°C
<a href="#">Ottawa (Kanata - Orléans)</a>	Mainly Clear	-10°C
<a href="#">Prince George</a>	Fog	-6°C
<a href="#">Québec</a>	Mainly Clear	-15°C
<a href="#">Regina</a>	Mist	-13°C
<a href="#">Saskatoon</a>	Mist	-11°C
<a href="#">St. John's</a>	Mostly Cloudy	-5°C
<a href="#">Thunder Bay</a>	Light Snow	0°C
<a href="#">Toronto</a>	Cloudy	-1°C
<a href="#">Vancouver</a>	Mainly Clear	7°C
<a href="#">Victoria</a>	Mainly Clear	5°C
<a href="#">Whitehorse</a>	Mostly Cloudy	-12°C
<a href="#">Winnipeg</a>	Mist	-6°C
<a href="#">Yellowknife</a>	Light Snow	-13°C

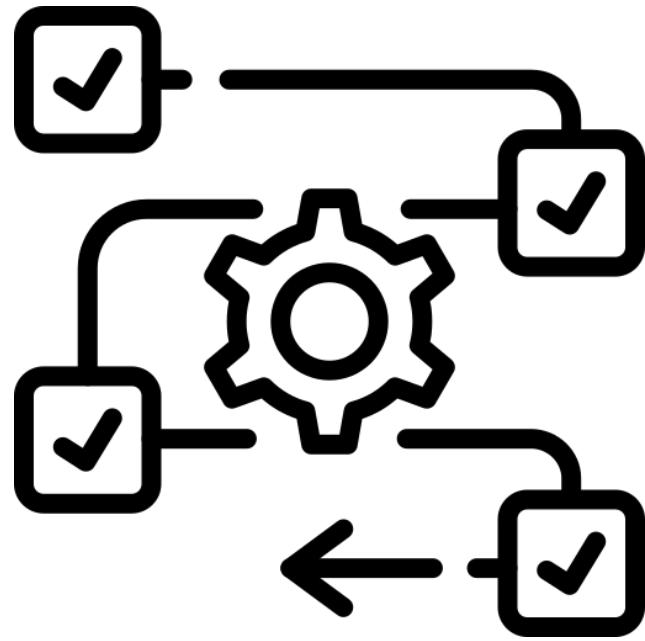
# Data abstraction: Three operations

- 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



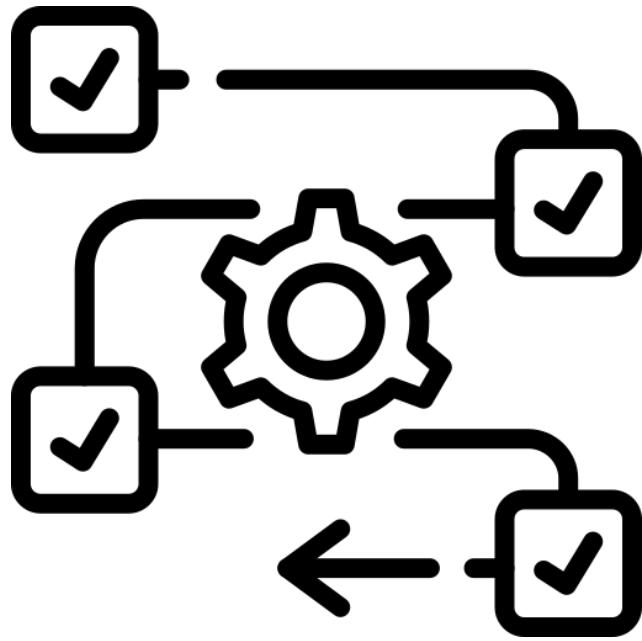
# Data vs conceptual models

- 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!]
- 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



# What?

## Datasets

## Attributes

### → Data Types

→ Items    → Attributes    → Links    → Positions    → Grids

### → Attribute Types

→ Categorical



### → Data and Dataset Types

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

→ Ordered

→ Ordinal

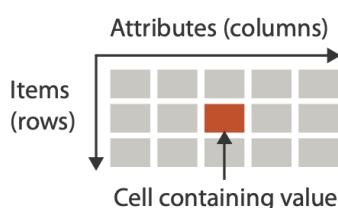


→ Quantitative

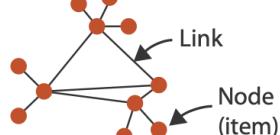


### → Dataset Types

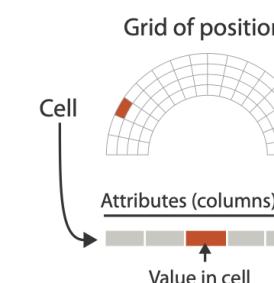
→ Tables



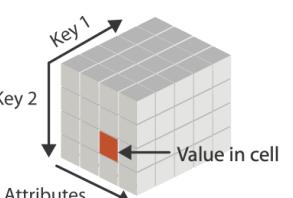
→ Networks



→ Fields (Continuous)



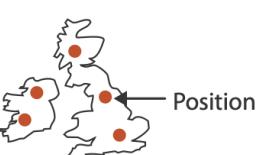
→ Multidimensional Table



→ Trees



→ Geometry (Spatial)



### → Ordering Direction

→ Sequential



→ Diverging



→ Cyclic



### → Dataset Availability

→ Static



→ Dynamic



What?

Why?

How?

# Design Sprint

What kind of questions can we answer?

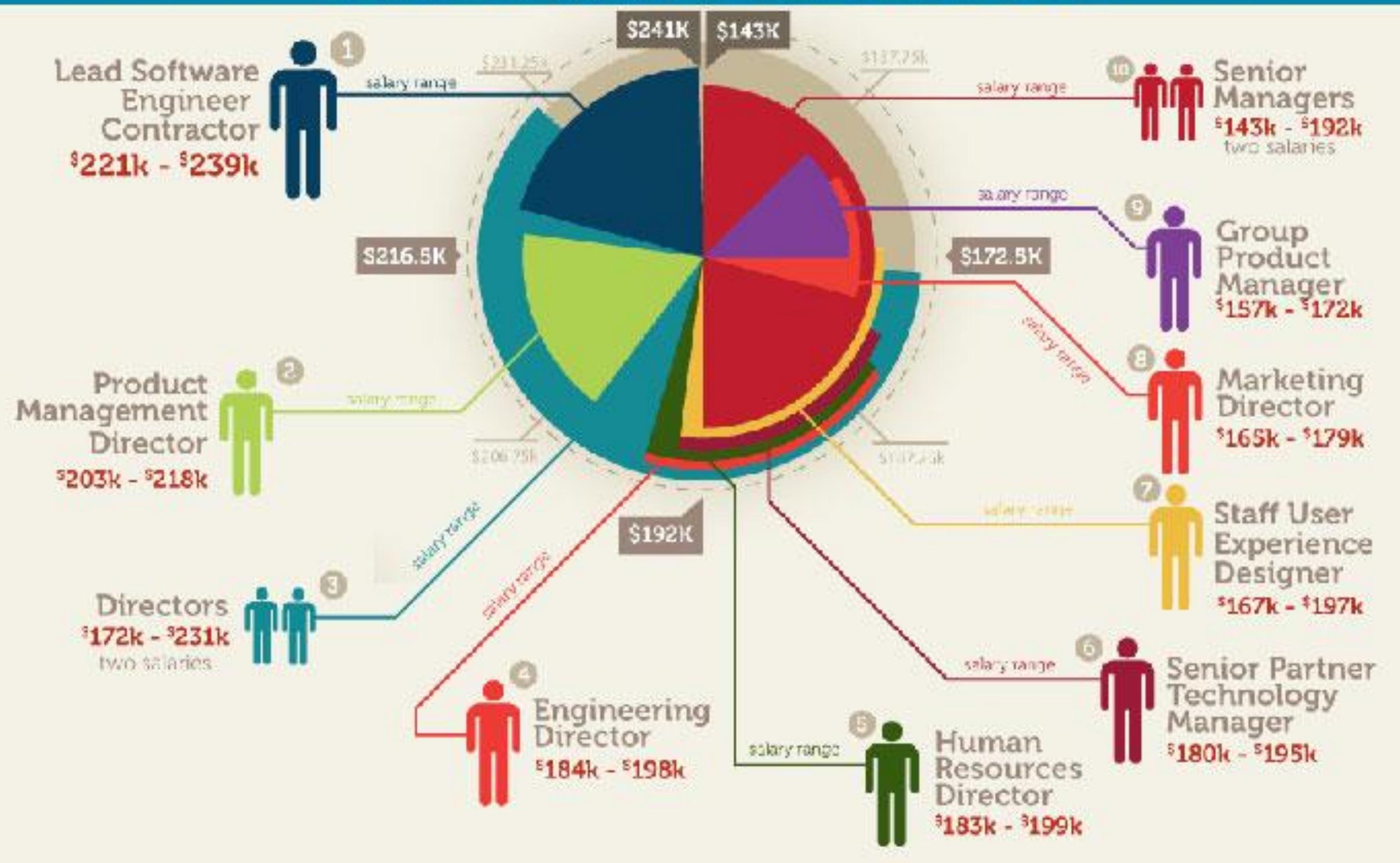
What is the data?

Who might be the audience?

City	Condition	Temperature
<a href="#">Calgary</a>	Mainly Clear	-4°C
<a href="#">Charlottetown</a>	Light Snowshower	-6°C
<a href="#">Edmonton</a>		-7°C
<a href="#">Fredericton</a>	Clear	-9°C
<a href="#">Halifax</a>	Light Snow	-6°C
<a href="#">Iqaluit</a>	Clear	-28°C
<a href="#">Montréal</a>	Mainly Clear	-9°C
<a href="#">Ottawa (Kanata - Orléans)</a>	Mainly Clear	-10°C
<a href="#">Prince George</a>	Fog	-6°C
<a href="#">Québec</a>	Mainly Clear	-15°C
<a href="#">Regina</a>	Mist	-13°C
<a href="#">Saskatoon</a>	Mist	-11°C
<a href="#">St. John's</a>	Mostly Cloudy	-5°C
<a href="#">Thunder Bay</a>	Light Snow	0°C
<a href="#">Toronto</a>	Cloudy	-1°C
<a href="#">Vancouver</a>	Mainly Clear	7°C
<a href="#">Victoria</a>	Mainly Clear	5°C
<a href="#">Whitehorse</a>	Mostly Cloudy	-12°C
<a href="#">Winnipeg</a>	Mist	-6°C
<a href="#">Yellowknife</a>	Light Snow	-13°C

# top 10 salaries at Google™

RANGE FROM \$143,000 TO \$241,000 PER YEAR.



# Learning Goals

- Describe the difference between how the phrases “dataset types”, “data types” are used in vis. Literature as opposed to programming
- Describe the characteristics of data
- Differentiate between the different types of data and dataset types