# DATA & IMAGE MODELS

CS 448B | Fall 2023

MANEESH AGRAWALA

1

# The big picture

**task**
questions, goals,
assumptions

**data**
physical data type
conceptual data type

**domain**
metadata
semantics
conventions

**mapping**
visual encoding

**processing**
**algorithms**

**image**
graphical marks
visual attrs/channels

7

## TODAY

### Learning Objectives

1. Identify *properties* of data and images

2. Decide how to *encode data using visual attributes/channels*

3. Define concepts of *expressiveness* and *effectiveness*

4. Develop *automated chart design* algorithm

8

## DATA

9

http://walthickey.com/2017/10/18/whats-the-best-halloween-candy/

10

# DATA TABLE

**Halloween Candy Power Ranking Dataset**

| | competitorname | chocolate | fruity | caramel | peanutyalmondy | nougat | crispedricewafer | hard | bar | pluribus | sugarpercent | pricepercent | winpercent |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | competitorname | chocolate | fruity | caramel | peanutyalmondy | nougat | crispedricewafer | hard | bar | pluribus | sugarpercent | pricepercent | winpercent |
| 2 | 100 Grand | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | .73199999 | .86000001 | 66.971725 |
| 3 | 3 Musketeers | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | .60399997 | .51099998 | 67.602936 |
| 4 | One dime | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | .011 | .116 | 32.261086 |
| 5 | One quarter | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | .011 | .51099998 | 46.116505 |
| 6 | Air Heads | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | .90600002 | .51099998 | 52.341465 |
| 7 | Almond Joy | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | .465 | .76700002 | 50.347546 |
| 8 | Baby Ruth | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | .60399997 | .76700002 | 56.914547 |
| 9 | Boston Baked Beans | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | .31299999 | .51099998 | 23.417824 |
| 10 | Candy Corn | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | .90600002 | .32499999 | 38.010963 |
| 11 | Caramel Apple Pops | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | .60399997 | .32499999 | 34.517681 |
| 12 | Charleston Chew | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | .60399997 | .51099998 | 38.975037 |
| 13 | Chewey Lemonhead Fruit Mix | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | .73199999 | .51099998 | 36.017628 |

11

**Dataset**

**Data Field**

**Halloween Candy Power Ranking Dataset**

| | competitorname | chocolate | fruity | caramel | peanutyalmondy | nougat | crispedricewafer | hard | bar | pluribus | sugarpercent | pricepercent | winpercent |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | competitorname | chocolate | fruity | caramel | peanutyalmondy | nougat | crispedricewafer | hard | bar | pluribus | sugarpercent | pricepercent | winpercent |
| 2 | 100 Grand | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | .73199999 | .86000001 | 66.971725 |
| 3 | 3 Musketeers | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | .60399997 | .51099998 | 67.602936 |
| 4 | One dime | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | .011 | .116 | 32.261086 |
| 5 | One quarter | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | .011 | .51099998 | 46.116505 |
| 6 | Air Heads | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | .90600002 | .51099998 | 52.341465 |
| 7 | Almond Joy | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | .465 | .76700002 | 50.347546 |
| 8 | Baby Ruth | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | .60399997 | .76700002 | 56.914547 |
| 9 | Boston Baked Beans | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | .31299999 | .51099998 | 23.417824 |
| 10 | Candy Corn | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | .90600002 | .32499999 | 38.010963 |
| 11 | Caramel Apple Pops | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | .60399997 | .32499999 | 34.517681 |
| 12 | Charleston Chew | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | .60399997 | .51099998 | 38.975037 |
| 13 | Chewey Lemonhead Fruit Mix | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | .73199999 | .51099998 | 36.017628 |

**Data Item/Observation**

**Cell Value**

https://fivethirtyeight.com/videos/the-ultimate-halloween-candy-power-ranking/

12

# TIDY DATA [Wickham 2014]

How do rows and columns, match up with data fields, and observations?

In *tidy data*
1. Each field forms a column
2. Each observation forms a row
3. Each type of observational unit forms a table

Flexible starting point for analysis, transformation, and visualization

13

**Data models** are formal descriptions

**Math:** Sets with operations on them

**Example:** integers with + and × operators

**Conceptual models** are mental constructions

Include semantics and support reasoning

**Examples** (data vs. conceptual)

1D floats vs. temperature

3D vector of floats vs. spatial location

14

# DATA MODEL

| | competitorname | chocolate | fruity | caramel | peanutyalmondy | nougat | crispedricewafer | hard | bar | pluribus | sugarpercent | pricepercent | winpercent |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | | | | | | | | | | | | |
| 2 | 100 Grand | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | .73199999 | .86000001 | 66.971725 |
| 3 | 3 Musketeers | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | .60399997 | .51099998 | 67.602936 |
| 4 | One dime | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | .011 | .116 | 32.261086 |
| 5 | One quarter | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | .011 | .51099998 | 46.116505 |
| 6 | Air Heads | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | .90600002 | .51099998 | 52.341465 |
| 7 | Almond Joy | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | .465 | .76700002 | 50.347546 |
| 8 | Baby Ruth | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | .60399997 | .76700002 | 56.914547 |
| 9 | Boston Baked Beans | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | .31299999 | .51099998 | 23.417824 |
| 10 | Candy Corn | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | .90600002 | .32499999 | 38.010963 |
| 11 | Caramel Apple Pops | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | .60399997 | .32499999 | 34.517681 |
| 12 | Charleston Chew | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | .60399997 | .51099998 | 38.975037 |
| 13 | Chewey Lemonhead Fruit Mix | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | .73199999 | .51099998 | 36.017628 |

| **string** | **bool** | **bool** | **bool** | **bool** | **bool** | **bool** | **bool** | **bool** | **bool** | **float** | **float** | **float** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

**How is data stored in the database?**

https://github.com/fivethirtyeight/data/tree/master/candy-power-ranking

15

# CONCEPTUAL MODEL

| Header | Description |
|---|---|
| chocolate | Does it contain chocolate? |
| fruity | Is it fruit flavored? |
| caramel | Is there caramel in the candy? |
| peanutalmondy | Does it contain peanuts, peanut butter or almonds? |
| nougat | Does it contain nougat? |
| crispedricewafer | Does it contain crisped rice, wafers, or a cookie component? |
| hard | Is it a hard candy? |
| bar | Is it a candy bar? |
| pluribus | Is it one of many candies in a bag or box? |
| sugarpercent | The percentile of sugar it falls under within the data set. |
| pricepercent | The unit price percentile compared to the rest of the set. |
| winpercent | The overall win percentage according to 269,000 matchups. |

https://github.com/fivethirtyeight/data/tree/master/candy-power-ranking

16

# CONCEPTUAL MODEL

| Header | Description |
|---|---|
| chocolate | Does it contain chocolate? |
| fruity | Is it fruit flavored? |
| caramel | Is there caramel in the candy? |
| peanutalmondy | Does it contain peanuts or almonds? |
| nougat | Does it contain nougat? |
| crispedricewafer | Does it contain crisped rice or cookies? |
| hard | Is it a hard candy? |
| bar | Is it a candy bar? |
| pluribus | Is it one of many candies in a bad? |
| sugarpercent | The percentile of sugar (across dataset) |
| pricepercent | The unit price percentile (across dataset) |
| winpercent | The overall win percentage in 269K contests |

**Domain specific understanding about the data**

**Supports analysis and reasoning**

https://github.com/fivethirtyeight/data/tree/master/candy-power-ranking

17

# DATA TYPES

**N - Nominal (labels)**
Fruits: Apples, oranges, …
Operations: **=, ≠**

**O - Ordered**
Quality of meat: Grade A, AA, AAA
Operations: **=, ≠, <, >**

**Q - Interval (location of zero arbitrary)**
Dates: Jan, 19, 2016; Loc.: (LAT 33.98, LON -118.45)
Like a geometric point. Cannot compare directly
Only differences (i.e. intervals) may be compared
Operations **=, ≠, <, >, -**

**Q - Ratio (location of zero fixed)**
Physical measurement: Length, Mass, …
Counts and amounts
Like a geometric vector, origin is meaningful
Operations: **=, ≠, <, >, -, ÷**

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

18

# NOMINAL, ORDINAL, QUANTITATIVE

| Header | Description | |
|--------|-------------|---|
| competitorname | Name of candy | N |
| chocolate | Does it contain chocolate? | N (maybe O) |
| fruity | Is it fruit flavored? | N (maybe O) |
| caramel | Is there caramel in the candy? | N (maybe O) |
| peanutalmondy | Does it contain peanuts or almonds? | N (maybe O) |
| nougat | Does it contain nougat? | N (maybe O) |
| crispedricewafer | Does it contain crisped rice or cookies? | N (maybe O) |
| hard | Is it a hard candy? | N (maybe O) |
| bar | Is it a candy bar? | N (maybe O) |
| pluribus | Is it one of many candies in a bad? | N (maybe O) |
| sugarpercent | The percentile of sugar (across dataset) | Q-Ratio |
| pricepercent | The unit price percentile (across dataset) | Q-Ratio |
| winpercent | The overall win percentage in 269K contests | Q-Ratio |

https://github.com/fivethirtyeight/data/tree/master/candy-power-ranking

19

# DATA TYPES

## DIMENSIONS

Dimensions are often the **independent** variables

*Dimensions* contain **qualitative values that describe the data item** (such as names, dates, or geographical data)

## MEASURES

Measures are often the **dependent** variables

*Measures* contain numeric, **quantitative values that you can measure**. Measures can be aggregated (sum, count, average, std. deviation).

| | competitorname | chocolate | fruity | caramel | peanutyalmondy | nougat | crispedricewafer | hard | bar | pluribus | sugarpercent | pricepercent | winpercent |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | | | | | | | | | | | | |
| 2 | 100 Grand | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | .73199999 | .86000001 | 66.971725 |
| 3 | 3 Musketeers | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | .60399997 | .51099998 | 67.602936 |
| 4 | One dime | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | .011 | .116 | 32.261086 |

**NOTE: Distinction is not strict. The same variable may be treated either way depending on the task**

20

# DIMENSION OR MEASURE

| Header | Description |
|---|---|
| competitorname | Name of candy |
| chocolate | Does it contain chocolate? |
| fruity | Is it fruit flavored? |
| caramel | Is there caramel in the candy? |
| peanutalmondy | Does it contain peanuts or almonds? |
| nougat | Does it contain nougat? |
| crispedricewafer | Does it contain crisped rice or cookies? |
| hard | Is it a hard candy? |
| bar | Is it a candy bar? |
| pluribus | Is it one of many candies in a bad? |
| sugarpercent | The percentile of sugar (across dataset) |
| pricepercent | The unit price percentile (across dataset) |
| winpercent | The overall win percentage in 269K contests |

https://github.com/fivethirtyeight/data/tree/master/candy-power-ranking

21

# DIMENSION OR MEASURE

| Header | Description |
| --- | --- |
| competitorname | Name of candy |
| chocolate | Does it contain chocolate? |
| fruity | Is it fruit flavored? |
| caramel | Is there caramel in the candy? |
| peanutalmondy | Does it contain peanuts or almonds? |
| nougat | Does it contain nougat? |
| crispedricewafer | Does it contain crisped rice or cookies? |
| hard | Is it a hard candy? |
| bar | Is it a candy bar? |
| pluribus | Is it one of many candies in a bad? |
| sugarpercent | The percentile of sugar (across dataset) |
| pricepercent | The unit price percentile (across dataset) |
| winpercent | The overall win percentage in 269K contests |

https://github.com/fivethirtyeight/data/tree/master/candy-power-ranking

22

# DIMENSION OR MEASURE

| Header | Description |
| --- | --- |
| competitorname | Name of candy |
| chocolate | Does it contain chocolate? |
| fruity | Is it fruit flavored? |
| caramel | Is there caramel in the candy? |
| peanutalmondy | Does it contain peanuts or almonds? |
| nougat | Does it contain nougat? |
| crispedricewafer | Does it contain crisped rice or cookies? |
| hard | Is it a hard candy? |
| bar | Is it a candy bar? |
| pluribus | Is it one of many candies in a bad? |
| sugarpercent | The percentile of sugar (across dataset) |
| pricepercent | The unit price percentile (across dataset) |
| winpercent | The overall win percentage in 269K contests |

https://github.com/fivethirtyeight/data/tree/master/candy-power-ranking

23

# U.S. CENSUS DATA

**People Count**: # 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 | marst | sex | people |
| 2 | 1850 | 0 | 0 | 1 | 1483789 |
| 3 | 1850 | 0 | 0 | 2 | 1450376 |
| 4 | 1850 | 5 | 0 | 1 | 1411067 |
| 5 | 1850 | 5 | 0 | 2 | 1359668 |
| 6 | 1850 | 10 | 0 | 1 | 1260099 |
| 7 | 1850 | 10 | 0 | 2 | 1216114 |
| 8 | 1850 | 15 | 0 | 1 | 1077133 |
| 9 | 1850 | 15 | 0 | 2 | 1110619 |
| 10 | 1850 | 20 | 0 | 1 | 1017281 |
| 11 | 1850 | 20 | 0 | 2 | 1003841 |
| 12 | 1850 | 25 | 0 | 1 | 862547 |
| 13 | 1850 | 25 | 0 | 2 | 799482 |
| 14 | 1850 | 30 | 0 | 1 | 730638 |
| 15 | 1850 | 30 | 0 | 2 | 639636 |
| 16 | 1850 | 35 | 0 | 1 | 588487 |
| 17 | 1850 | 35 | 0 | 2 | 505012 |
| 18 | 1850 | 40 | 0 | 1 | 475911 |
| 19 | 1850 | 40 | 0 | 2 | 428185 |
| 20 | 1850 | 45 | 0 | 1 | 384211 |
| 21 | 1850 | 45 | 0 | 2 | 341254 |
| 22 | 1850 | 50 | 0 | 1 | 321343 |
| 23 | 1850 | 50 | 0 | 2 | 286580 |
| 24 | 1850 | 55 | 0 | 1 | 194080 |
| 25 | 1850 | 55 | 0 | 2 | 187208 |
| 26 | 1850 | 60 | 0 | 1 | 174076 |

24

# CENSUS N, O, Q

**People Count**: Q-Ratio
**Year:** Q-Interval (maybe O)
**Age:** Q-Ratio (maybe O)
**Sex:** N
**Marital Status:** N

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | year | age | marst | sex | people |
| 2 | 1850 | 0 | 0 | 1 | 1483789 |
| 3 | 1850 | 0 | 0 | 2 | 1450376 |
| 4 | 1850 | 5 | 0 | 1 | 1411067 |
| 5 | 1850 | 5 | 0 | 2 | 1359668 |
| 6 | 1850 | 10 | 0 | 1 | 1260099 |
| 7 | 1850 | 10 | 0 | 2 | 1216114 |
| 8 | 1850 | 15 | 0 | 1 | 1077133 |
| 9 | 1850 | 15 | 0 | 2 | 1110619 |
| 10 | 1850 | 20 | 0 | 1 | 1017281 |
| 11 | 1850 | 20 | 0 | 2 | 1003841 |
| 12 | 1850 | 25 | 0 | 1 | 862547 |
| 13 | 1850 | 25 | 0 | 2 | 799482 |
| 14 | 1850 | 30 | 0 | 1 | 730638 |
| 15 | 1850 | 30 | 0 | 2 | 639636 |
| 16 | 1850 | 35 | 0 | 1 | 588487 |
| 17 | 1850 | 35 | 0 | 2 | 505012 |
| 18 | 1850 | 40 | 0 | 1 | 475911 |
| 19 | 1850 | 40 | 0 | 2 | 428185 |
| 20 | 1850 | 45 | 0 | 1 | 384211 |
| 21 | 1850 | 45 | 0 | 2 | 341254 |
| 22 | 1850 | 50 | 0 | 1 | 321343 |
| 23 | 1850 | 50 | 0 | 2 | 286580 |
| 24 | 1850 | 55 | 0 | 1 | 194080 |
| 25 | 1850 | 55 | 0 | 2 | 187208 |
| 26 | 1850 | 60 | 0 | 1 | 174076 |

25

# CENSUS DIM., MEAS.

| People Count: | Measure |
| Year: | Dimension |
| Age: | Depends! |
| Sex: | Dimension |
| Marital Status: | Dimension |

|  | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | year | age | marst | sex | people |
| 2 | 1850 | 0 | 0 | 1 | 1483789 |
| 3 | 1850 | 0 | 0 | 2 | 1450376 |
| 4 | 1850 | 5 | 0 | 1 | 1411067 |
| 5 | 1850 | 5 | 0 | 2 | 1359668 |
| 6 | 1850 | 10 | 0 | 1 | 1260099 |
| 7 | 1850 | 10 | 0 | 2 | 1216114 |
| 8 | 1850 | 15 | 0 | 1 | 1077133 |
| 9 | 1850 | 15 | 0 | 2 | 1110619 |
| 10 | 1850 | 20 | 0 | 1 | 1017281 |
| 11 | 1850 | 20 | 0 | 2 | 1003841 |
| 12 | 1850 | 25 | 0 | 1 | 862547 |
| 13 | 1850 | 25 | 0 | 2 | 799482 |
| 14 | 1850 | 30 | 0 | 1 | 730638 |
| 15 | 1850 | 30 | 0 | 2 | 639636 |
| 16 | 1850 | 35 | 0 | 1 | 588487 |
| 17 | 1850 | 35 | 0 | 2 | 505012 |
| 18 | 1850 | 40 | 0 | 1 | 475911 |
| 19 | 1850 | 40 | 0 | 2 | 428185 |
| 20 | 1850 | 45 | 0 | 1 | 384211 |
| 21 | 1850 | 45 | 0 | 2 | 341254 |
| 22 | 1850 | 50 | 0 | 1 | 321343 |
| 23 | 1850 | 50 | 0 | 2 | 286580 |
| 24 | 1850 | 55 | 0 | 1 | 194080 |
| 25 | 1850 | 55 | 0 | 2 | 187208 |
| 26 | 1850 | 60 | 0 | 1 | 174976 |

26

# DATA TABLES & TRANSFORMATIONS

27

# RELATIONAL ALGEBRA [Codd 1970] / SQL

**Operations on data tables: table(s) in, table out**

Projection (SELECT) – select a set of columns

Selection (WHERE) – filter rows

Sorting (ORDER BY) – order rows

Aggregation (GROUP BY, SUM, MIN, ...)

    partition rows into groups and summarize

Combination (JOIN, UNION, ...)

    integrate data from multiple tables

| ID | Name | Population | Med. Income |
|----|------|-----------|-------------|
| 100 | Valley East | 3,200 | 45,000 |
| 101 | Val Therese | 4,125 | 48,000 |
| 102 | Capreol | 2,109 | 39,000 |
| 103 | Eastwood | 4,500 | 43,500 |
| 104 | Lynnwood | 3,459 | 42,000 |
| 105 | Kingsway | 3,443 | 55,000 |
| 106 | Prince Anne | 2,986 | 52,500 |
| 107 | Whitefish | 1,998 | 39,000 |

Attributes — Primary key — Cardinality — Tuple — Attribute value

28

---

# RELATIONAL ALGEBRA [Codd 1970] / SQL

Projection (SELECT) – select a set of columns

`select day, stock`

| day | stock | price |
|-----|-------|-------|
| 10/3 | AMZN | 957.10 |
| 10/3 | MSFT | 74.26 |
| 10/4 | AMZN | 965.45 |
| 10/4 | MSFT | 74.69 |

→

| day | stock |
|-----|-------|
| 10/3 | AMZN |
| 10/3 | MSFT |
| 10/4 | AMZN |
| 10/4 | MSFT |

29

# RELATIONAL ALGEBRA [Codd 1970] / SQL

Selection (WHERE) – filter rows

```
select * where price > 100
```

| day | stock | price |
|-----|-------|-------|
| 10/3 | AMZN | 957.10 |
| 10/3 | MSFT | 74.26 |
| 10/4 | AMZN | 965.45 |
| 10/4 | MSFT | 74.69 |

➡

| day | stock | price |
|-----|-------|-------|
| 10/3 | AMZN | 957.10 |
| 10/4 | AMZN | 965.45 |

30

# RELATIONAL ALGEBRA [Codd 1970] / SQL

Sorting (ORDER BY) – order records

```
select * order by stock
```

| day | stock | price |
|-----|-------|-------|
| 10/3 | AMZN | 957.10 |
| 10/3 | MSFT | 74.26 |
| 10/4 | AMZN | 965.45 |
| 10/4 | MSFT | 74.69 |

➡

| day | stock | price |
|-----|-------|-------|
| 10/3 | AMZN | 957.10 |
| 10/4 | AMZN | 965.45 |
| 10/3 | MSFT | 74.26 |
| 10/4 | MSFT | 74.69 |

31

# RELATIONAL ALGEBRA [Codd 1970] / SQL

Aggregation (GROUP BY, SUM, MIN, …)

```
select stock min(price) group by stock
```

| day | stock | price |
|-----|-------|-------|
| 10/3 | AMZN | 957.10 |
| 10/3 | MSFT | 74.26 |
| 10/4 | AMZN | 965.45 |
| 10/4 | MSFT | 74.69 |

→

| stock | min(price) |
|-------|-----------|
| AMZN | 957.10 |
| MSFT | 74.26 |

32

# RELATIONAL ALGEBRA [Codd 1970] / SQL

Combination (JOIN) multiple tables together

| day | stock | price |
|-----|-------|-------|
| 10/3 | AMZN | 957.10 |
| 10/3 | MSFT | 74.26 |
| 10/4 | AMZN | 965.45 |
| 10/4 | MSFT | 74.69 |

→

| day | stock | price | min |
|-----|-------|-------|-----|
| 10/3 | AMZN | 957.10 | 957.10 |
| 10/3 | MSFT | 74.26 | 74.26 |
| 10/4 | AMZN | 965.45 | 957.10 |
| 10/4 | MSFT | 74.69 | 74.26 |

| stock | min |
|-------|-----|
| AMZN | 957.10 |
| MSFT | 74.26 |

```
select t.day, t.stock, t.price, a.min
from table as t, aggregate as a
where t.stock = a.stock
```

33

14

**Original**

| YEAR | AGE | MARST | SEX | PEOPLE |
|------|-----|-------|-----|--------|
| 1850 | 0 | 0 | 1 | 1,483,789 |
| 1850 | 5 | 0 | 1 | 1,411,067 |
| 1860 | 0 | 0 | 1 | 2,120,846 |
| 1860 | 5 | 0 | 1 | 1,804,467 |
| . . . | | | | |

**Pivoted or Cross-Tabulation**

| AGE | MARST | SEX | 1850 | 1860 | . . . |
|-----|-------|-----|------|------|-------|
| 0 | 0 | 1 | 1,483,789 | 2,120,846 | . . . |
| 5 | 0 | 1 | 1,411,067 | 1,804,467 | . . . |
| . . . | | | | | |

Which format might we prefer? Why?

38

# ANNOUNCEMENTS

40

15

# CLASS PARTICIPATION REQUIREMENTS

**Complete** required **readings** and **notebooks** before class

**Attend class** and be a part of the in-class discussion

**Post** at least 1 discussion substantive comment/question per week

Due by 8pm the following Sunday
1 free pass for the quarter

**Class home page**
https://magrawala.github.io/cs448b-fa23/

41

# READING/NOTEBOOK/LECTURE RESPONSES

**Good responses typically exhibit one or more**
    **Critiques** of arguments made in the papers/lectures
    **Analysis** of implications or future directions for ideas in readings/lectures
    **Insightful questions** about the readings/lectures

**Responses should not be summaries**
    Should be substantive (1-2 paragraphs is typical)

42

# OBSERVABLE NOTEBOOKS / VEGA-LITE



**Vega-Lite** is a *declarative* API for programming visualizations

**Do the exercises** (fork notebook)

**Lec. on Wed 10/4 will assume you have done the 1st three notebooks**

43

# ASSIGNMENT 1: VISUALIZATION DESIGN
## Due TODAY

Design a static visualization for a data set
You must choose the message you want to convey. What question(s) do you want to answer? What insight do you want to communicate?

**Data: Stanford Undergraduate Majors**

Stanford University publishes a variety of datasets through the Stanford Institutional Rsearch & Decision Support website. They have published a data table containing information about the number of Stanford undergraduates obtaining a Bachelor's degree in 75 different fields of study from 2003 to 2022. We have filtered and wrangled this data to the top 10 fields of study by cummulative degrees conferred over the time period to produce a dataset with the following attributes:

- **Year:** Academic year between 2003 and 2022. (Academic years run July-June so Year=2003 covers July 2002 to June 2003.)
- **FieldOfStudy:** Field in which degree was obtained.
- **Count:** Number of students earning a Bachelor's degree.

The extracted dataset is available in csv format: TopFieldsStanfordBachelors.csv

44

# ASSIGNMENT 2:  EXP. DATA ANALYSIS
## Due 10/16   11:30am

Use Tableau or Vega-Lite to formulate & answer data questions

First steps
- Step 1: Pick domain & data
- Step 2: Pose questions
- Step 3: Profile data
- Iterate as needed

Create visualizations
- See different views of data
- Refine questions

Author a report
- Screenshots of most insightful views (8+)
- Include titles and captions for each view

Indexed Gas Prices by Region Over Time

When prices peak, the West Coast gets hit hardest while the Gulf Coast almost always gets hit softest.

Where
- Central Atlantic
- East Coast
- Gulf Coast
- Lower Atlantic
- Midwest
- New England
- Rocky Mountain
- West Coast

45

IMAGE

46

# MARKS & VISUAL ATTRs

**Marks:** **geometric primitives**

points          lines          areas

**Visual Attributes:** **control mark appearance**

Position (2x)

Size

Value

Texture

Color

Orientation

Shape

*Semiology of Graphics*
*J. Bertin, 1967*

47

# CODING INFORMATION IN POSITION

+
C
+
+
B
+
A

1. A, B, C are distinguishable
2. Three points are colinear: B between A and C
3. BC is twice as long as AB

∴ Encode quantitative variables

**"Resemblance, order and proportional are the three signfields in graphics." - Bertin**

49

# CODING INFORMATION IN COLOR

**Value is perceived as ordered**

∴ Encode ordinal variables (O)

∴ Encode continuous variables (Q) [not as well]

**Hue is normally perceived as unordered**

∴ Encode nominal variables (N) using color

50

---

# BERTIN'S "LEVELS OF ORGANIZATION"

| | N | O | Q |
|---|---|---|---|
| **Position** | N | O | Q |
| **Size** | N | O | Q |
| **Value** | N | O | Q |
| **Texture** | N | o | |
| **Color** | N | | |
| **Orientation** | N | | |
| **Shape** | N | | |

N  Nominal
O  Ordered
Q  Quantitative

**Note: Q ⊂ O ⊂ N**

51

# VISUAL ENCODING

52

---

# ENCODINGS: MAP DATA to MARK ATTRIBUTES

mark: rect
data → size (height)

53

## ENCODINGS: MAP DATA to MARK ATTRIBUTES



mark: rect
data → size (height)

mark: points
data$_1$ → x-pos
data$_2$ → y-pos

54

## ENCODINGS: MAP DATA to MARK ATTRIBUTES



mark: rect
data → size (height)

mark: points
data$_1$ → x-pos
data$_2$ → y-pos

mark: points
data$_1$ → x-pos
data$_2$ → y-pos
data$_3$ → color

55

## ENCODINGS: MAP DATA to MARK ATTRIBUTES



mark: rect
data → size (height)

mark: points
data$_1$ → x-pos
data$_2$ → y-pos

mark: points
data$_1$ → x-pos
data$_2$ → y-pos
data$_3$ → color

mark: points
data$_1$ → x-pos
data$_2$ → y-pos
data$_3$ → color
data$_4$ → size

56

# DECONSTRUCTIONS

57

**Commercial and Political Atlas [Playfair 1786/1801]**

Time (Q) → x-position
Exports/Imports Values (Q) → y-position
Exports/Imports (N, O) → color
Balance for/against (Q) → area (maybe length??)
Balance for/against (N,O) → color

60



https://finviz.com/map.ashx

61

# ENCODINGS

market cap (Q) → rectangle size

mkt sector (N), mkt cap (Q) → rect. pos.

loss vs. gain (N, O) → color hue

magnitude of loss or gain (Q) → color value

https://finviz.com/map.ashx



62

# MINARD's MARCH on MOSCOW



**Figurative Map of the successive losses of men of the French army during the Russian Campaign 1812–1813**

63

# MARK COMPOSITION

 temperature (Q) → y-position

+

longitude (Q), time (O) → x-position

―――――――――――――――――

=



temp across space & time (Q x Q,O)

# MARK COMPOSITION

**+** | latitude (Q) → y-position

**+** longitude (Q) → x-position

army size (Q) → width

────────────────────────

**=**

army position (Q x Q) and army size (Q)

66

---

latitude (Q)

longitude (Q)

army size (Q)

temperature (Q)

longitude (Q), time(O)

67

68



# FORMALIZING DESIGN

69

# COMBINATORICS OF ENCODINGS

**Challenge:**

Assume k visual attributes/channels and n data fields

Pick the best encoding from the exponential number of possibilities $(n+1)^k$

70

# PRINCIPLES

**Challenge**

Assume k visual attributes/channels and n data fields

Pick the best encoding from the exponential number of possibilities $(n+1)^k$

**Principle of Consistency**

Properties of image (visual variables) should match properties of data

**Principle of Importance Ordering**

Encode most important information in the most effective way

71

# EXPRESSIVENESS CRITERIA [Mackinlay 1986]

## Expressiveness

A set of facts is expressible in a visual language if the sentences (i.e., the visualizations) in the language express *all* the facts in the set of data, and *only* the facts in the data.
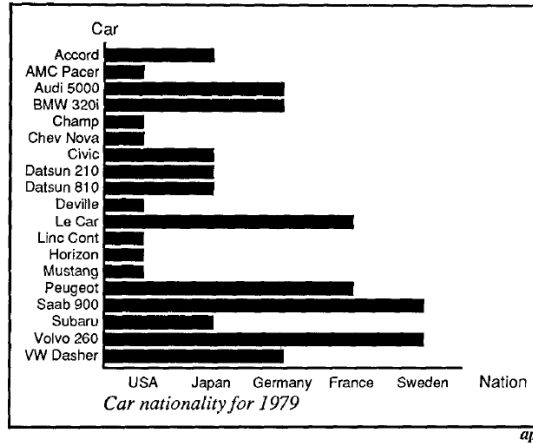
72

# CANNOT EXPRESS ALL THE FACTS

## Horizontal dot plot

A one-to-many (1 → N) relation cannot be expressed in a single horizontal dot plot because multiple tuples are mapped to the same position



73

# EXPRESSES FACTS NOT IN THE DATA



Car

Car nationality for 1979

Fig. 11. Incorrect use of a bar chart for the *Nation* relation. The lengths of the bars suggest an ordering on the vertical axis, as if the USA cars were longer or better than the other cars, which is not true for the *Nation* relation.

Length is interpreted as encoding a quantitative value

74

# EFFECTIVENESS CRITERIA [Mackinlay 1986]

## Effectiveness

A visualization is more effective than another visualization if the information conveyed by one visualization is more readily *perceived* than the information in the other visualization.

Subject of the Perception Lecture

75

# MACKINLAY'S RANKING



Conjectured *effectiveness* of encodings by data type

76

# AUTOMATIC CHART DESIGN [Mackinlay 1986]

**APT** – "A Presentation Tool"

**User formally specifies data model and type**
Input: list of data variables ordered by importance

**APT searches over the design space**
Tests expressiveness of each visual encoding (rule-based)
Generates encodings that pass test
Rank by perceptual effectiveness criteria
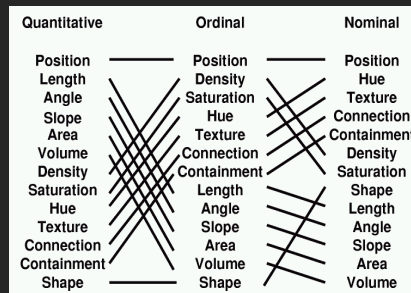
**Outputs *most effective* visualization**

77

# APT [Mackinlay 1986]

Encode most important data using highest ranking visual variable for the data type

| Price | Mileage | Weight | Repair |
|-------|---------|--------|--------|
| 13,500 | 22 | 3000 | great |
| 7,200 | 31 | 1500 | ok |
| 11,300 | 12 | 4200 | terrible |
| ... | ... | ... | ... |

→
1. **Price (Q)**
2. **Mileage (Q)**
3. **Weight (Q)**
4. **Repair (N)**



**mark: lines**

→
**Price (Q) → y-pos**
**Mileage (Q) → x-pos**
**Weight (Q) → size**
**Repair (N) → color**

Automating the design of graphical presentation of relational information
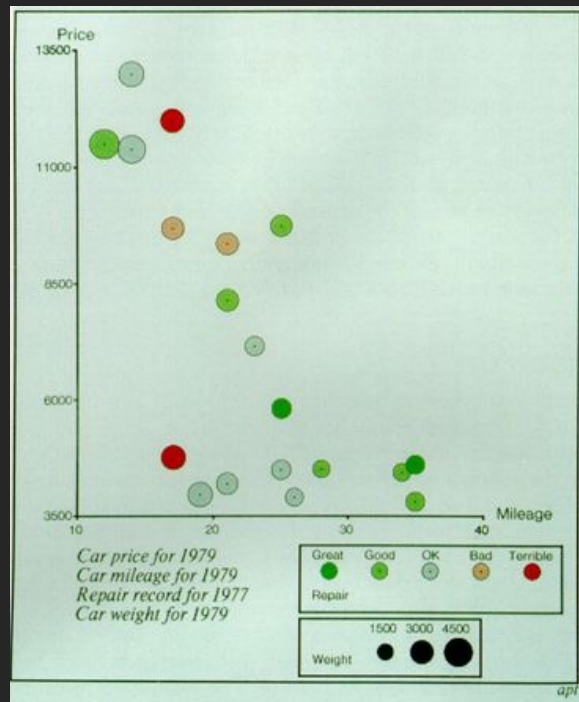J. Mackinlay, 1986

79

---

# APT [Mackinlay 1986]

Automatically generated chart for cars data

**Cars Data**
1. Price (Q)
2. Mileage (Q)
3. Weight (Q)
4. Repair (Q)



80

# LIMITATIONS

**Does not cover many visualization techniques**

Networks, maps, diagrams

Also, 3D, animation, illustration, …

**Does not consider interaction**

**Does not consider semantics or conventions**

**Assumes single visualization as output**

81

# SUMMARY

**Formal specification**

Data model: tidy data, N,O,Q types

Image model: marks, visual attributes/channels

Encodings map data to mark attributes/channels

**Choose *expressive* and *effective* encodings**

Rule-based test of expressiveness

Perceptual effectiveness rankings

82