# Data and Image Models

### Maneesh Agrawala
### Jessica Hullman

**CS 294-10: Visualization**

**Fall 2014**

# Last Time: The Purpose of Visualization

# Three functions of visualizations

**Record information**

- Photographs, blueprints, …

**Support reasoning about information (analyze)**

- Process and calculate
- Reason about data
- Feedback and interaction

**Convey information to others (present)**

- Share and persuade
- Collaborate and revise
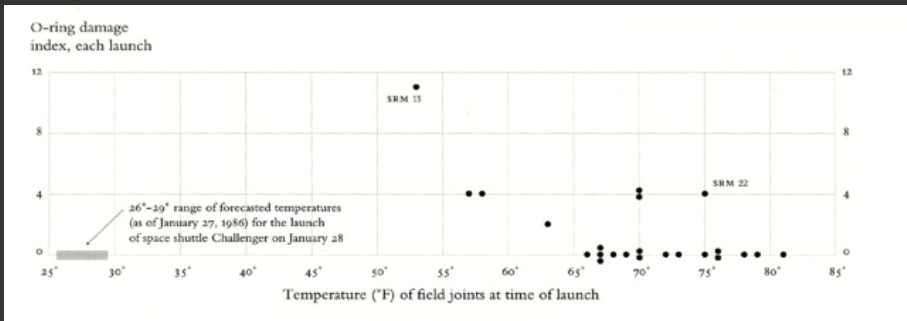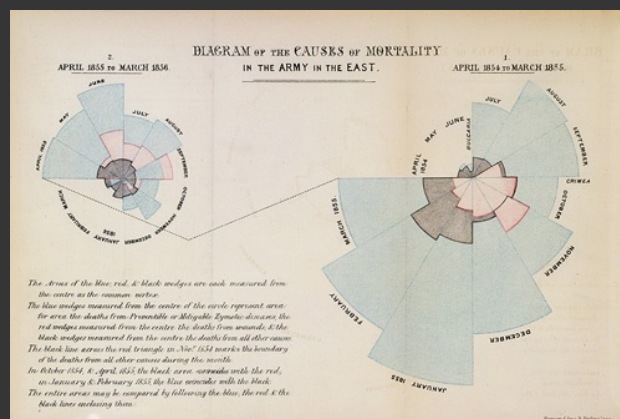- Emphasize important aspects of data

# Record information



Gallop, Bay Horse "Daisy" [Muybridge 1884-86]

# Analysis: Challenger



O-ring damage
index, each launch

Temperature (°F) of field joints at time of launch

26°–29° range of forecasted temperatures (as of January 27, 1986) for the launch of space shuttle Challenger on January 28
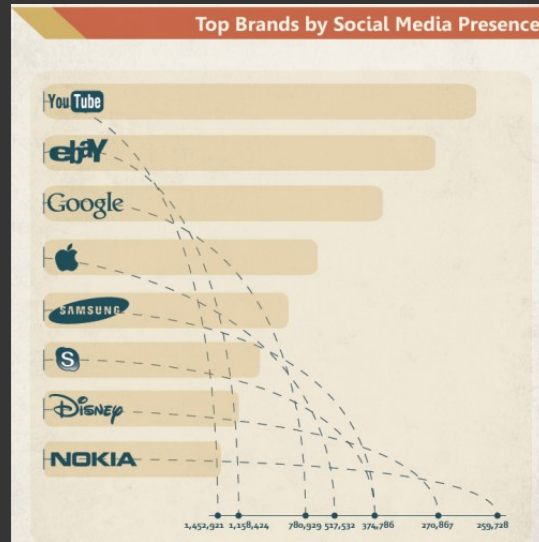
Visualizations drawn by Tufte show how low temperatures damage O-rings [Tufte 97]

# Communicate: War Deaths



Crimean War Deaths [Nightingale 1858]

## Confuse: Top Brands

**Top Brands by Social Media Presence**

You Tube

ebaY

Google



SAMSUNG

S

Disney

NOKIA

1,452,921  1,158,424  780,929 527,532 374,786  270,867  259,728

from wtfviz.net

---

# Announcements

**Auditors, *please* enroll in the class (1 unit, P/NP)**
- Requirements: Come to class and participate (online as well)
- Requirements: Assignment 1

**Class participation requirements**
- Complete readings before class
- In-class discussion
- Post at least 1 discussion substantive comment/question by 11am on day of lecture

**All, add yourself to participants page on the wiki**

**Class wiki**
http://vis.berkeley.edu/courses/cs294-10-fa14/wiki/

# Assignment 1: Visualization Design

Worldwide Disasters 1900-2008

Due by 11:59pm on Sep 9

# Data and Image Models

# The big picture

task

data
  physical type
    int, float, etc.
  abstract type
    nominal, ordinal, etc.

domain
  metadata
  semantics
  conceptual model

processing
algorithms

mapping
  visual encoding
  visual metaphor

image
  visual channel
  retinal variables

# Topics

**Properties of data or information**

**Properties of the image**

**Mapping data to images**

# Data

# Data models vs. Conceptual models

**Data models: low level descriptions of the data**
- Math: Sets with operations on them
- Example: integers with + and × operators

**Conceptual models: mental constructions**
- Include semantics and support reasoning

**Examples (data vs. conceptual)**
- (1D floats) vs. Temperature
- (3D vector of floats) vs. Space

# Taxonomy

- **1D (sets and sequences)**
- **Temporal**
- **2D (maps)**
- **3D (shapes)**
- **nD (relational)**
- **Trees (hierarchies)**
- **Networks (graphs)**

**Are there others?**

The eyes have it: A task by data type taxonomy for information
visualization [Schneiderman 96]

# Types of variables

## Physical types
- Characterized by storage format
- Characterized by machine operations

**Example:**
bool, short, int32, float, double, string, …

## Abstract types
- Provide descriptions of the data
- May be characterized by methods/attributes
- May be organized into a hierarchy

**Example:**
plants, animals, metazoans, …

# Nominal, ordinal and quantitative

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

**O - Ordered**
- Quality of meat: Grade A, AA, AAA

**Q - Interval (Location of zero arbitrary)**
- Dates: Jan, 19, 2006; Location: (LAT 33.98, LONG -118.45)
- Like a geometric point. Cannot compare directly
- Only differences (i.e. intervals) may be compared

**Q - Ratio (zero fixed)**
- Physical measurement: Length, Mass, Temp, …
- Counts and amounts
- Like a geometric vector, origin is meaningful

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

---

# Nominal, ordinal and quantitative

**N - Nominal (labels)**
- Operations: =, ≠

**O - Ordered**
- Operations: =, ≠, <, >, ≤, ≥

**Q - Interval (Location of zero arbitrary)**
- Operations: =, ≠, <, >, ≤, ≥, -
- Can measure distances or spans

**Q - Ratio (zero fixed)**
- Operations: =, ≠, <, >, ≤, ≥, -, ÷
- Can measure ratios or proportions

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

# From data model to N,O,Q data type

**Data model**

- 32.5, 54.0, -17.3, …
  - floats

**Conceptual model**

- Temperature

**Data type**

- Burned vs. Not burned (N)
- Hot, warm, cold (O)
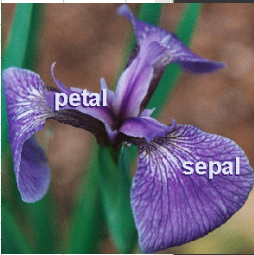- Continuous range of values (Q)



Iris Setosa        Iris Versicolor        Iris Virginica

Sepal and petal lengths and widths for three species of iris [Fisher 1936].

# Relational data model

Represent data as a **table** (*relation*)

Each **row** (*tuple*) represents a single record

  Each record is a fixed-length tuple

Each **column** (*attribute*) represents a single *variable*

  Each attribute has a *name* and a *data type*

A table's **schema** is the set of names and data types


A **database** is a collection of tables (relations)


# Relational algebra [Codd]

**Data transformations (SQL)**
- **Selection (WHERE) – restrict values**
- **Projection (SELECT) – choose subset of attributes**
- **Sorting (ORDER BY)**
- **Aggregation (GROUP BY, SUM, MIN, …)**
- **Set operations (UNION, …)**
- **Combine (INNNER JOIN, OUTER JOIN, …)**

# Statistical data model

**Variables or measurements**
**Categories or factors or dimensions**
**Observations or cases**

| Month | Control | Placebo | 300 mg | 450 mg |
|---|---|---|---|---|
| March | 165 | 163 | 166 | 168 |
| April | 162 | 159 | 161 | 163 |
| May | 164 | 158 | 161 | 153 |
| June | 162 | 161 | 158 | 160 |
| July | 166 | 158 | 160 | 148 |
| August | 163 | 158 | 157 | 150 |

**Blood Pressure Study (4 treatments, 6 months)**

# Dimensions and measures

**Dimensions:** Discrete variables describing data
- Dates, categories of values (independent vars)

**Measures:** Data values that can be aggregated
- Numbers to be analyzed (dependent vars)
- Aggregate as sum, count, average, std. deviation


# Dimensions and measures

**Independent vs. dependent variables**
- Example: $y = f(x,a)$
- Dimensions: Domain(x) $\times$ Domain(a)
- Measures: Range(y)

Image

# Visual language is a sign system


**Jacques Bertin**

**Images perceived as a set of signs**

**Sender encodes information in signs**

**Receiver decodes information from signs**

**Semiology of Graphics, 1967**

---

# Information in position



**1. A, B, C are distinguishable**
**2. B is 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

[Bertin, Semiology of Graphics, 1983]

# Visual variables

- Position (x 2)
- Size
- Value
- Texture
- Color
- Orientation
- Shape



**Note: Bertin does not consider 3D or time**
**Note: Card and Mackinlay extend the number of vars.**

# Information in color and value

**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**

---

# Bertins' "Levels of Organization"

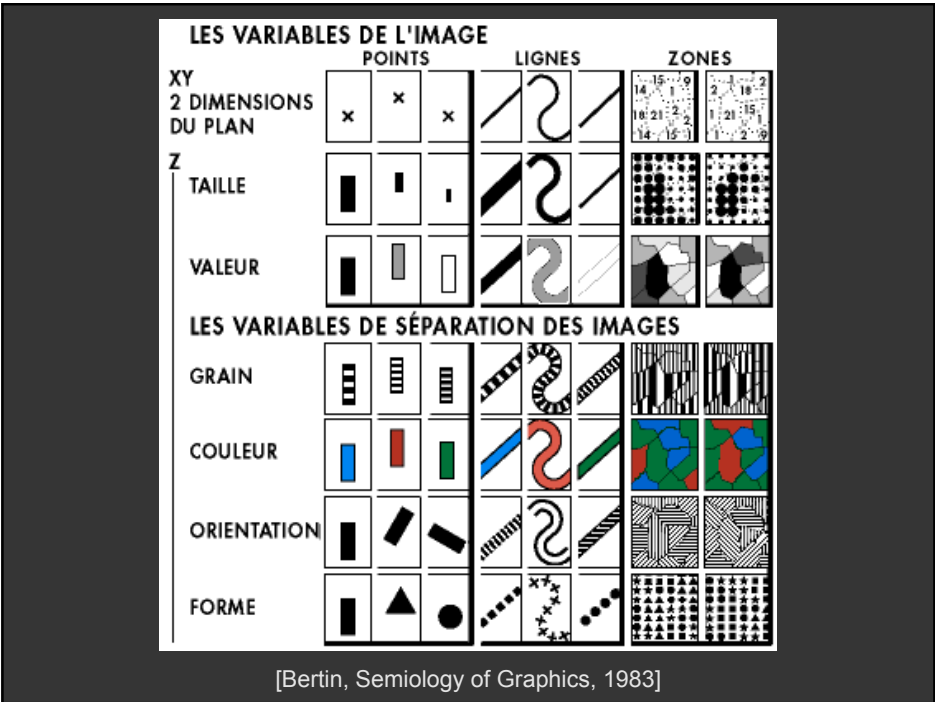| | | | |
|---|---|---|---|
| Position | N | O | Q |
| Size | N | O | Q |
| Value | N | O | Q |
| Texture | N | O | |
| Color | N | | |
| Orientation | N | | |
| Shape | N | | |

N  Nominal
O  Ordered
Q  Quantitative

Note: Q < O < N

Note: Bertin actually breaks visual variables down into differentiating (≠) and associating (≡)

# Encoding rules

---

# Univariate data

| | A | B | C | |
|---|---|---|---|---|
| 1 | | | | measure |

# Univariate data

**A B C**

measure

**1**

7
5
3
1

A  B  C  D  E

Tukey box plot

low | Middle 50% | high

Mean

0                                        20

A      B      C      D

# Bivariate data

**A B C**

**1**
**2**

C

B         D              F

A                E

Scatter plot is common

20

# Trivariate data

|   | A | B | C |
|---|---|---|---|
| 1 |   |   |   |
| 2 |   |   |   |
| 3 |   |   |   |

3D scatter plot is possible



---

# Three variables

**Two variables [x,y] can map to points**
- Scatterplots, maps, …

**Third variable [z] must use …**
- Color, size, shape, …

## Large design space (visual metaphors)



[Bertin, Graphics and Graphic Info. Processing, 1981]

## Multidimensional data

**How many variables can be depicted in an image?**

|   | A | B | C |
|---|---|---|---|
| 1 |   |   |   |
| 2 |   |   |   |
| 3 |   |   |   |
| 4 |   |   |   |
| 5 |   |   |   |
| 6 |   |   |   |
| 7 |   |   |   |
| 8 |   |   |   |

## Multidimensional data

How many variables can be depicted in an image?

| | A | B | C |
|---|---|---|---|
| 1 | | | |
| 2 | | | |
| 3 | | | |
| 4 | | | |
| 5 | | | |
| 6 | | | |
| 7 | | | |
| 8 | | | |

*"With up to three rows, a data table can be constructed directly as a single image … However, an image has only three dimensions.  And this barrier is impassible."*                    **Bertin**

# Deconstructions

# Stock chart  from the late 90s



---

# Stock chart  from the late 90s



- Time → x-position (Q, linear)
- Price → y-position (Q, linear)

# Playfair 1786



Exports and Imports to and from DENMARK & NORWAY from 1700 to 1780.
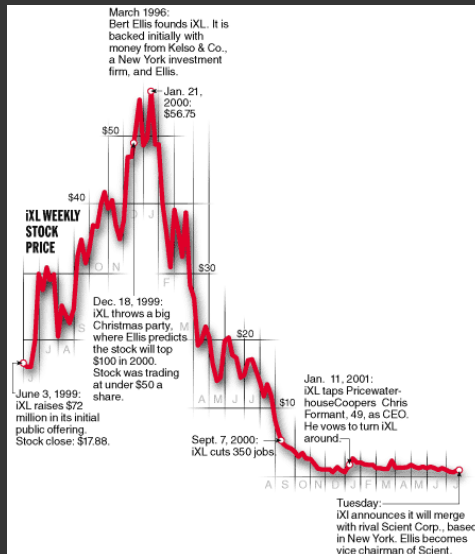
# Playfair 1786



Exports and Imports to and from DENMARK & NORWAY from 1700 to 1780.

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

# Minard 1869: Napoleon's march



# Single axis composition



[based on slide from Mackinlay]

# Mark composition

Temperature → y-position (Q, linear)

**+** longitude → x-position (Q, linear)

_____

**=** 

temp over longitude (Q x Q)

# Mark composition

Longitude → y-position (Q, linear)

**+** Latitude → x-position (Q, linear)

**+** army size → width (Q, linear)

_____

**=** 

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

longitude (Q, lin)

latitude (Q, lin)

army size (Q, lin)

temperature (Q, lin)

longitude (Q, lin)

# Minard 1869: Napoleon's march

**Depicts at least 5 quantitative variables**
**Any others?**

# Automated design
## Jock Mackinlay's APT 86



# Combinatorics of encodings

**Challenge:**

Assume 8 visual encodings and n data attributes

Pick the best encoding from the exponential number of possibilities $(n+1)^8$

**Principle of Consistency:**

The properties of the image (visual variables) should match the properties of the data

**Principle of Importance Ordering:**

Encode the most important information in the most effective way

# Mackinlay's expressiveness criteria

**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.

# Cannot express the facts

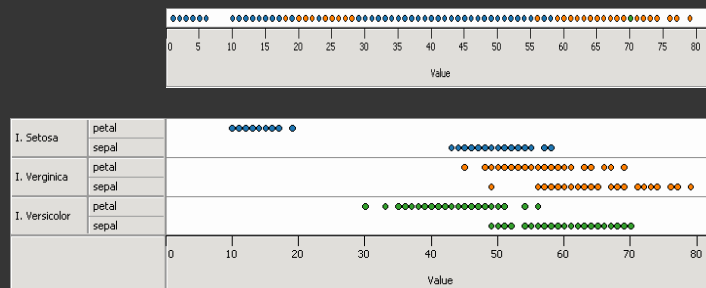A one-to-many (1 → N) relation cannot be expressed in a single horizontal dot plot because multiple tuples are mapped to the same position

# Expresses facts not in the data

**A length is interpreted as a quantitative value;**

**∴ Length of bar says something untrue about N data**



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.

[Mackinlay, APT, 1986]

---

# Mackinlay's effectiveness criteria

**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 perception lecture**

# Mackinlay's ranking

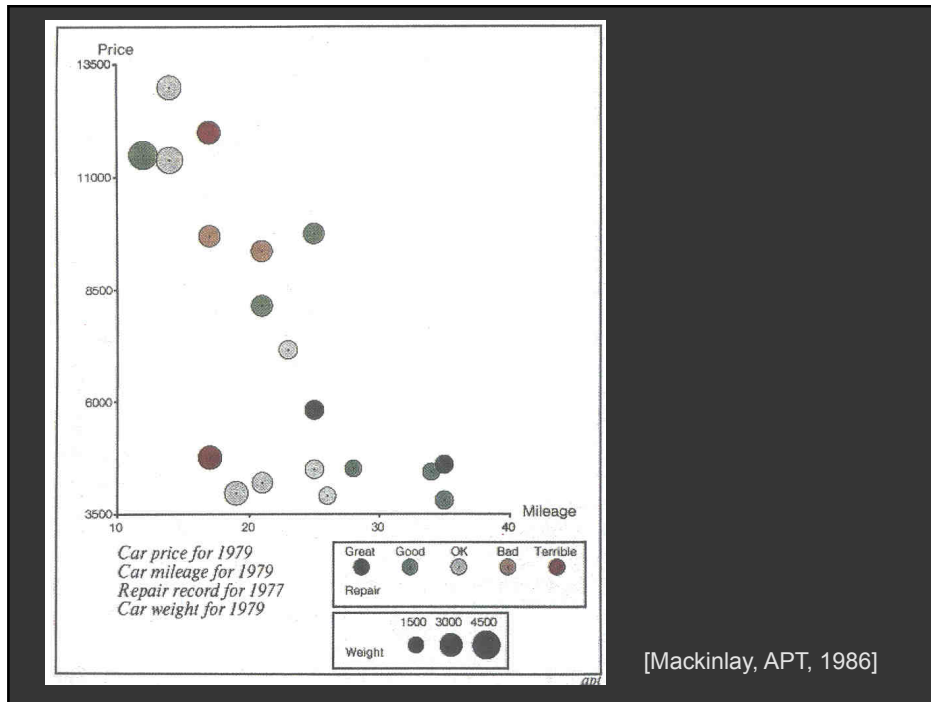| Quantitative | Ordinal | Nominal |
|---|---|---|
| Position | Position | Position |
| Length | Density | Hue |
| Angle | Saturation | Texture |
| Slope | Hue | Connection |
| Area | Texture | Containment |
| Volume | Connection | Density |
| Density | Containment | Saturation |
| Saturation | Length | Shape |
| Hue | Angle | Length |
| Texture | Slope | Angle |
| Connection | Area | Slope |
| Containment | Volume | Area |
| Shape | Shape | Volume |

Conjectured *effectiveness* of the encoding

# Mackinlay's design algorithm

- **User formally specifies data model and type**
- **APT searches over design space**
  - **Tests expressiveness of each visual encoding**
  - **Generates image for encodings that pass test**
  - **Tests perceptual effectiveness of resulting image**
- **Outputs most effective visualization**

[Mackinlay, APT, 1986]

# Limitations

## Does not cover many visualization techniques
- Bertin and others discuss networks, maps, diagrams
- They do not consider 3D, animation, illustration, photography, …

## Does not model interaction

# Summary

**Formal specification**
- Data model
- Image model
- Encodings mapping data to image

**Choose expressive and effective encodings**
- Formal test of expressiveness
- Experimental tests of perceptual effectiveness