# CSE 564
# Visualization & Visual Analytics

# Applications and Basic Tasks

## Klaus Mueller

Computer Science Department
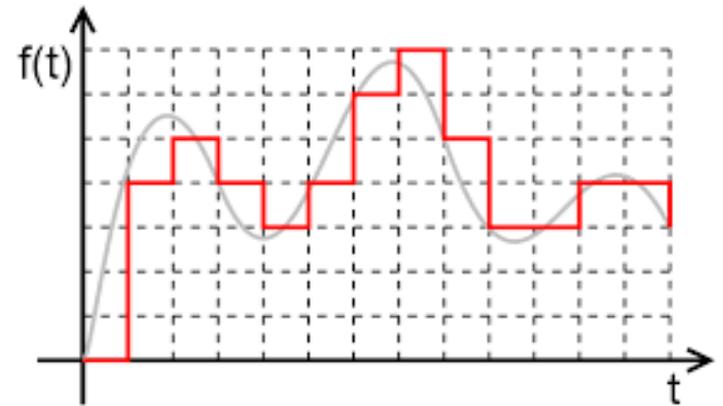Stony Brook University

| Lecture | Topic | Projects |
|---|---|---|
| 1 | Intro and logistics | |
| 2 | Basic visualizations and tasks, data types, examples, ethical considerations | |
| 3 | Data preparation (cleaning, imputation, data set integration) | |
| 4 | AI-assisted coding for VIS applications (design, debugging, refactoring) | Project #1 out |
| 5 | Big data and data reduction (distance/sim metrics, intro to clustering) | |
| 6 | High-D data and dimension reduction (PCA, subspaces, correlation maps) | |
| 7 | Cluster analysis: numerical data, categorical data | |
| 8 | Perception and cognition (human visual system, color, contrast, bias) | Project #2(a) out |
| 9 | Visual design and aesthetics | |
| 10 | Visualization of multivariate and high-dimensional data: direct methods | |
| 11 | Visualization of multivariate and high-D data: projections & embeddings | |
| 12 | Visualization and AI: mutual support and capabilities (VIS4AI, AI4VIS) | Project #2(b) out |
| 13 | Principles of interaction: drive what is visualized, analyzed & how  (HCI4VIS) | |
| 14 | Visual analytics (VA), human-centered AI, mixed-initiative system | |
| 15 | Midterm #1 (tentative date) | |
| 16 | VA system design and evaluation, collaborative VA, uncertainty, provenance | |
| 17 | Midterm #1 discussion (tentative date) | Final proj. proposal call out |
| 18 | Visualization of hierarchical data | |
| 19 | Visualization of maps and data with geo-reference | |
| 20 | Visualization of graphs, networks (incl. derivation of causal networks) | Final project proposal due |
| 21 | Vis. of time-varying, time-series, streaming data, progressive visualization | |
| 22 | Visualization of text, LLMs, and semantic data | |
| 23 | Ed Tufte revisited: principles, critiques and limits, responsible visualization | |
| 24 | Design of effective infographics | Final proj. prelim report due |
| 25 | Foundations scientific and medical visualization, intro to volume rendering | |
| 26 | Scientific visualization | Bonus project out (Vol Ren) |
| 27 | Story telling with data, data journalism | |
| 28 | Midterm #2 (tentative date) | |
| Final | Final project demo on zoom (public) | All final proj. materials due |

# Variable Types
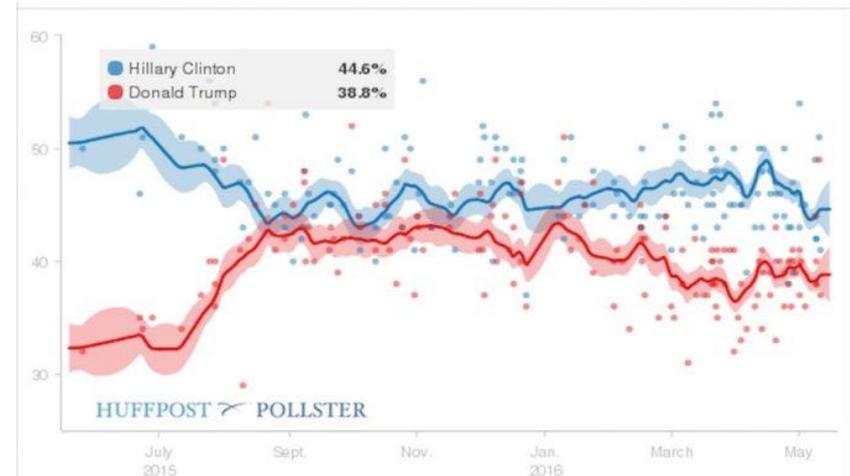
# NUMERICAL VARIABLES

## Numerical variables

- measure a **quantity** as a number
- like: 'how many' or 'how much'
- can be continuous (grey curve)
- or discrete (red steps)
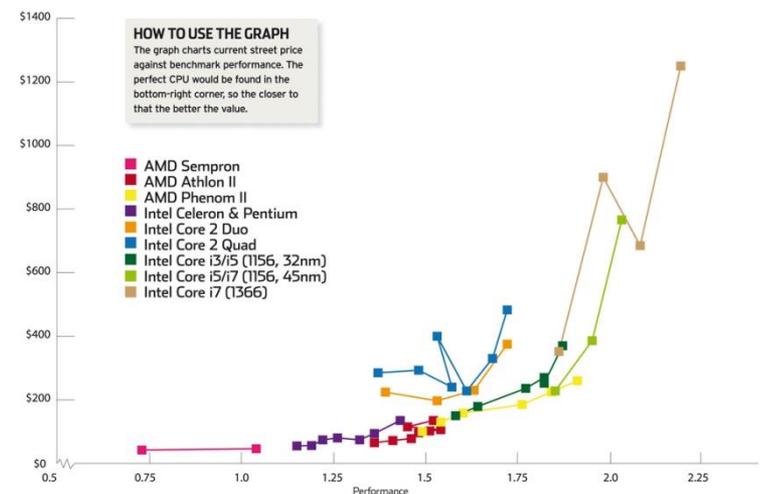
# NUMERICAL VARIABLES

Most often the x-axis is 'time'

- provides an intuitive & innate ordering of the data values
- the majority of people expect the x-axis to be 'time'



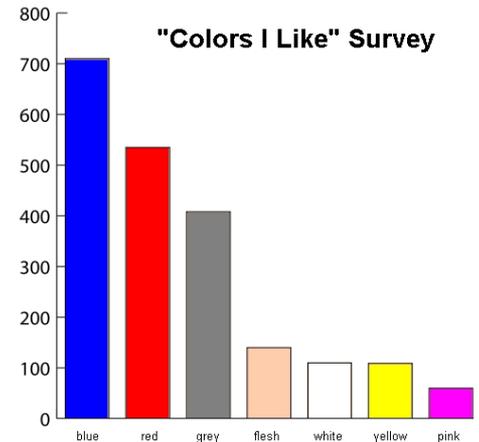But 'time' is not the only option

- engineers, statisticians, etc. will be receptive to this idea
- can you think of an example?

# Categorical Variables

Categorical variables

- describe a **quality** or characteristic
- like: 'what type' or 'which category'


"Colors I Like" Survey

- can be **ordinal** = ordered, ranked (distances need not be equal)
  – clothing size, academic grades, levels of agreement

- or **nominal** = not organized into a logical sequence
  – gender, business type, eye color, brand

# Categorical Variables

Usually plotted as bar charts or pie charts



Number of Colors in Bag of M&M Candies



Customer Satisfaction

??                                    ??

nominal

ordinal

but of course, you can plot either of them in either of these two representations

# Numbers are Good

But not everything is expressed in numbers

- images
- video
- text
- web logs
- ...

Do feature analysis to turn these abstract things into numbers

- a vector of numbers, to be concrete
- then apply your analysis as usual
- but keep the reference to the original data so you can return to the native domain where the analysis problem originated

# SENSOR DATA

## Characteristics

- often large scale
- time series



## Feature Analysis

- example: Motif discovery
- encode into 5D data vector



# features discovered in stream
feature [F1, F 2, F3, F4, F 5]
         [12, 3, 41, 12, 5]
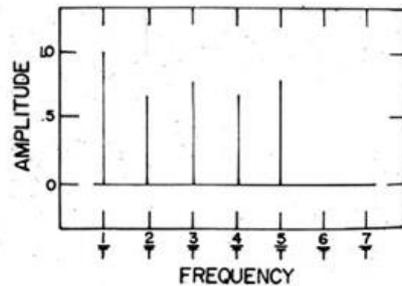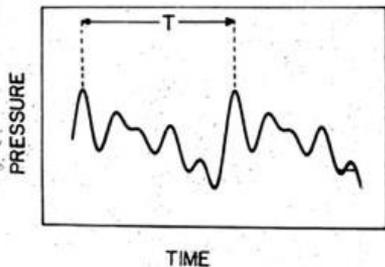


Five different known motifs

# SENSOR DATA

## Characteristics

- often large scale
- time series

## Feature Analysis

- Fourier transform (FT, FFT)
- Wavelet transform (WT, FWT)

Fourier transform
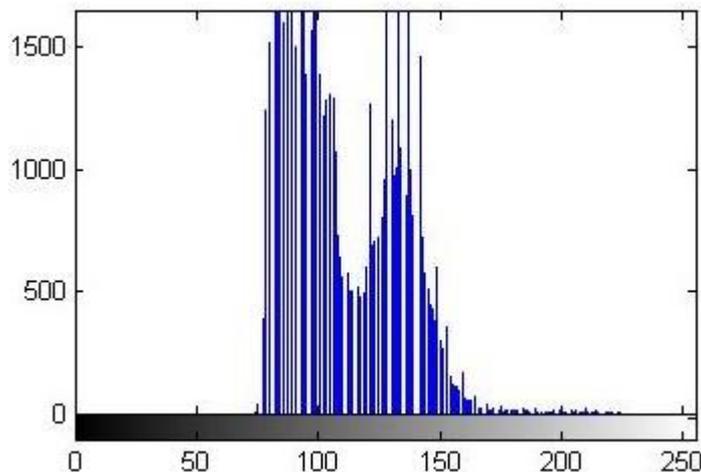Store spectrum into a vector

# IMAGE DATA

## Characteristics

- array of pixels
- representable as a vector of length [width x height]

## Feature Analysis

- value histograms
- encode into a 256-D vector
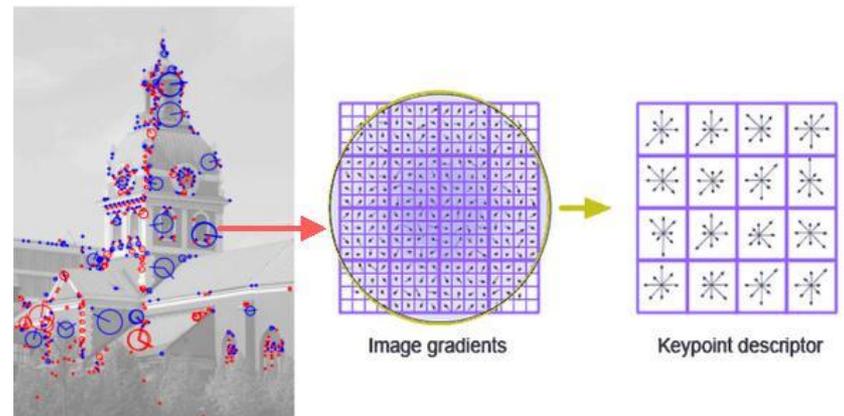
histograms



[0, 0, 0, ...., 10, ..., 1200, .....]

# Image Data

## Characteristics

- array of pixels
- representable as a vector of length [width x height]

## Feature Analysis

- value histograms
- gradient histograms
- FFT, FWT
- Scale Invariant Feature Transform (SIFT)
- Bag of Features (BoF)
- visual words

histograms





Image gradients → Keypoint descriptor

SIFT

# Video data

## Characteristics
- essentially a time series of images

## Feature Analysis
- many of the above techniques apply albeit extension is non-trivial

# Text Data

Characteristics
- often raw and unstructured

Feature analysis
- first step is to remove stop words and stem the data
- perform **named-entity recognition** to gain atomic elements
  - identify names, locations, actions, numeric quantities, relations
  - understand the structure of the sentence and complex events
- example:
  - Jim bought 300 shares of Acme Corp. in 2006.
  - [Jim]$_{Person}$ bought [300 shares] $_{Quantity}$ of [Acme Corp.]$_{Organiz.}$ in [2006]$_{Time}$
- distinguish between
  - application of grammar rules (old style, need experienced linguists)
  - statistical models (Google etc., need big data to build)

# Text to Numeric Data

Create a term-document matrix

- turns text into a high-dimensional vector which can be compared
- use Latent Semantic Analysis (LSA) to derive a visualization
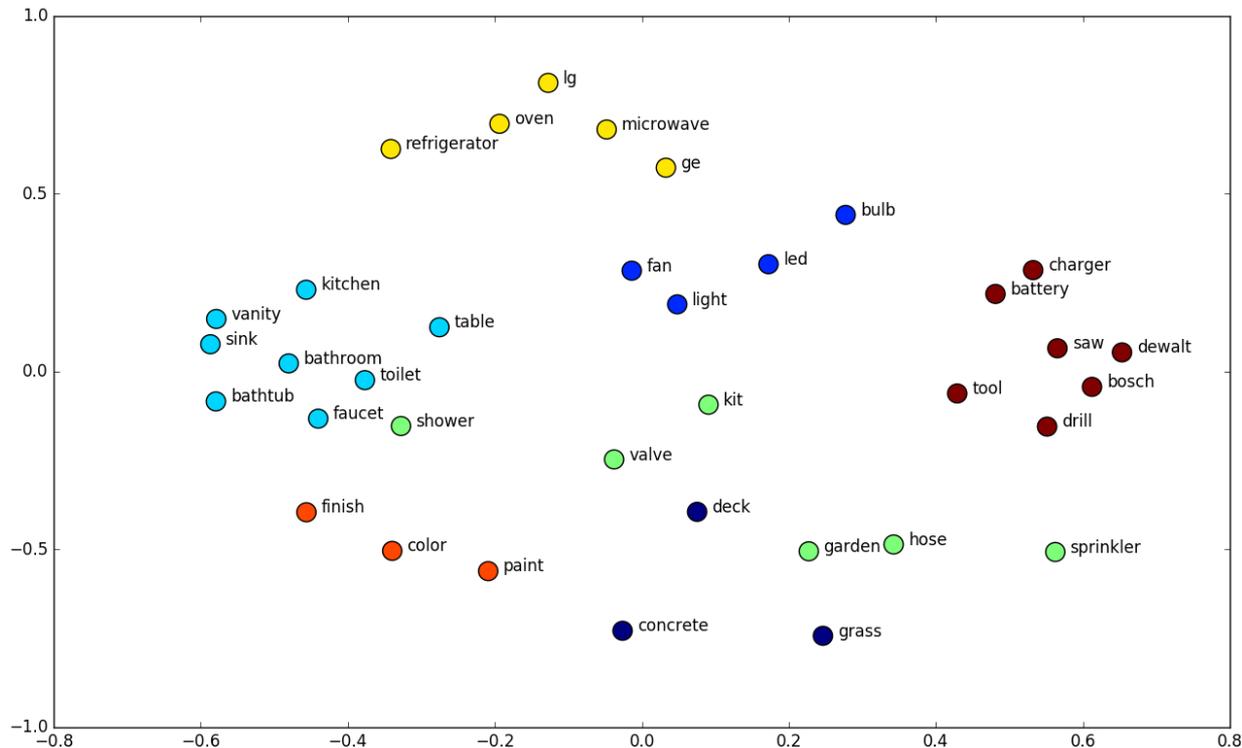


Term-Document Matrix



Word/document cluster

# WORD EMBEDDING

Train a shallow neural network (NN) on a corpus of text
- the NN weight vectors encode word similarity as a high-D vector
- use a 2D embedding technique to display

# Word Cloud

Maps the frequency of words in a corpus to size



https://www.jasondavies.com/wordcloud/

# OTHER DATA

## Weblogs

- typically represented as text strings in a pre-specified format
- this makes it easy to convert them into multidimensional representation of categorical and numeric attributes

## Network traffic

- characteristics of the network packets are used to analyze intrusions or other interesting activity
- a variety of features may be extracted from these packets
  - the number of bytes transferred
  - the network protocol used
  - IP ports used

# Let's Look at Some Essential Graphical Representations
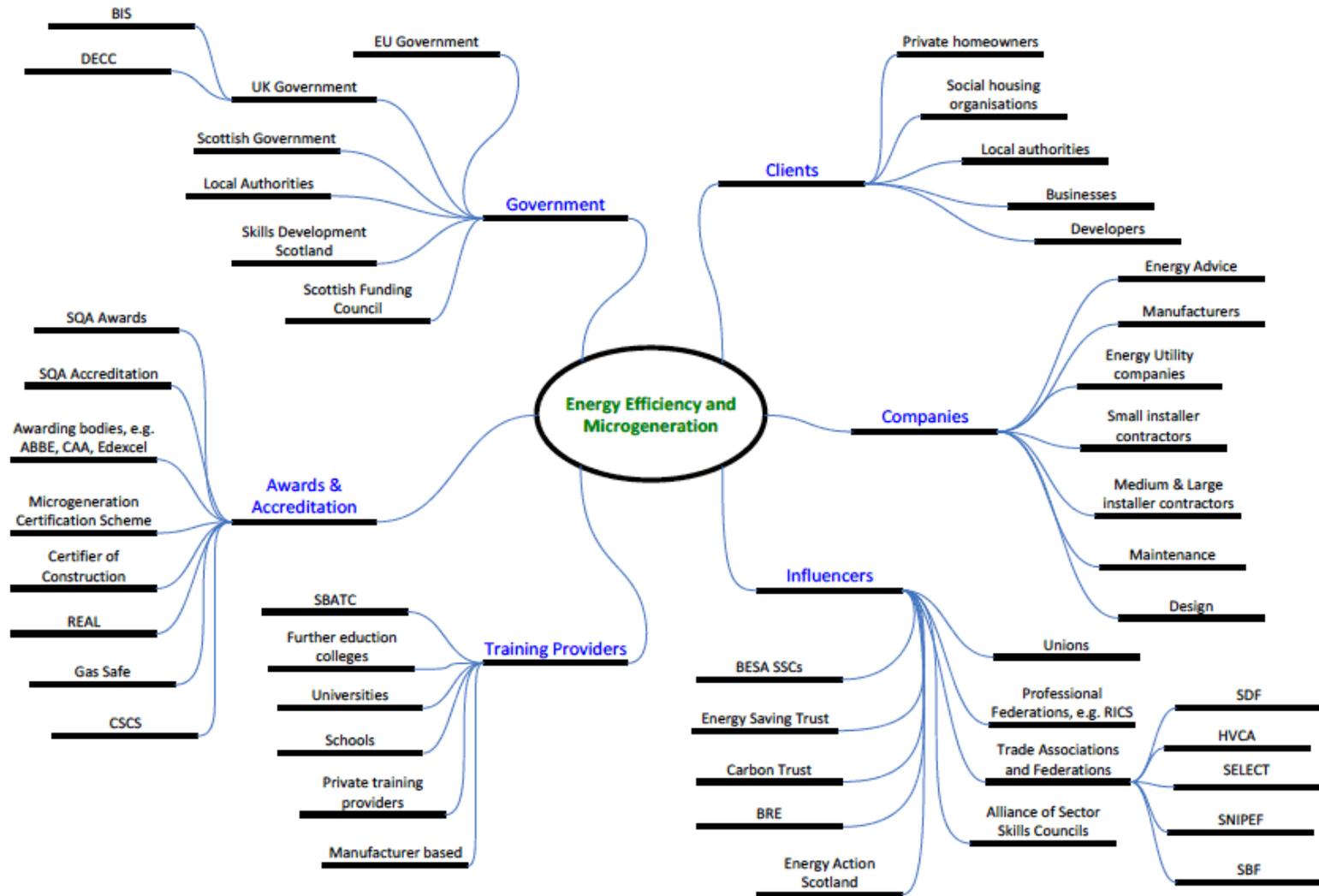
# And Do Some Advertising for D3

# Stakeholder Hierarchy

# FUNCTION CALL TREE

# More Complex Stakeholder Hierarchy



**Energy Efficiency and Microgeneration**

**Government**
- BIS
- DECC
- EU Government
- UK Government
- Scottish Government
- Local Authorities
- Skills Development Scotland
- Scottish Funding Council

**Clients**
- Private homeowners
- Social housing organisations
- Local authorities
- Businesses
- Developers

**Companies**
- Energy Advice
- Manufacturers
- Energy Utility companies
- Small installer contractors
- Medium & Large installer contractors
- Maintenance
- Design

**Awards & Accreditation**
- SQA Awards
- SQA Accreditation
- Awarding bodies, e.g. ABBE, CAA, Edexcel
- Microgeneration Certification Scheme
- Certifier of Construction
- REAL
- Gas Safe
- CSCS

**Training Providers**
- SBATC
- Further eduction colleges
- Universities
- Schools
- Private training providers
- Manufacturer based

**Influencers**
- BESA SSCs
- Energy Saving Trust
- Carbon Trust
- BRE
- Energy Action Scotland
- Unions
- Professional Federations, e.g. RICS
- Trade Associations and Federations
  - SDF
  - HVCA
  - SELECT
  - SNIPEF
  - SBF
- Alliance of Sector Skills Councils

# Hierarchies

Questions you might have

- how large is each group of stakeholders (or function)?
  - tree with quantities
- what fraction is each group with respect to the entire group?
  - partition of unity
- how is information disseminated among the stakeholders (or functions)?
  - information flow
- how close (or distant) are the individual stakeholders (functions) in terms of some metric?
  - force directed layout

# Invoke Nature

More scalable tree, and natural with some randomness

http://animateddata.co.uk/lab/d3-tree/

# Collapsible Tree

A standard tree, but one that is scalable to large hierarchies

http://mbostock.github.io/d3/talk/20111018/tree.html

# Zoomable Partition Layout

A tree that is scalable and has partial partition of unity

http://mbostock.github.io/d3/talk/20111018/partition.html

# Sunburst

More space efficient since it's radial, has partial partition of unity

https://observablehq.com/@kerryrodden/sequences-sunburst

# Bubble Charts

No hierarchy information, just quantities

https://observablehq.com/@d3/bubble-chart

# Circle Packing

Quantities and containment, but not partition of unity

http://mbostock.github.io/d3/talk/20111116/pack-hierarchy.html

# Treemap

Quantities, containment, and full partition of unity

http://mbostock.github.io/d3/talk/20111018/treemap.html

# Chord Diagram

Relationships among group fractions, not necessarily a tree

https://observablehq.com/@d3/chord-diagram

# Hierarchical Edge Bundling

Relationships of individual group members, also in terms of quantitative measures such as information flow

http://mbostock.github.io/d3/talk/20111116/bundle.html

# Collapsible Force Layout

Relationships within organization members expressed as distance and proximity

[http://mbostock.github.io/d3/talk/20111116/force-collapsible.html](http://mbostock.github.io/d3/talk/20111116/force-collapsible.html)

# Voronoi Tessellation

Shows the closest point on the plane for a given set of points... and a new point via interaction

https://observablehq.com/collection/@d3/d3-delaunay

# Data Type Conversions and Transformation

# Numeric to Categorical Data: Discretization (1)

Solution 1:

- divide the numeric attribute values into φ **equi-width** ranges
- each range/bucket has the same width
- example: customer age



- what is lost here?

# PROBLEM WITH EQUI-WIDTH HISTOGRAM

Age ranges of customers could be unevenly distributed within a bin
- this could be an interesting anomaly

Solution 2:

- divide the numeric attribute values into φ **equi-depth** ranges
- same number of samples in each bin
- (again) example: customer age:



- what is the disadvantage here?
- extra storage needed: must store the start/end value for each bin

# NUMERIC TO CATEGORICAL DATA: DISCRETIZATION (3)

Solution 3:

- what if all the bars have seemingly the same height
- or are dominated by one large peak



- switch to log scaling of the y-value

# OTHER TRANSFORMATIONS



- *none*: $x^* = x$ (leaves points unchanged)
- *half*: $x^* = x/2$ (squeezes all points together)
- *square*: $x^* = x^2$ (pulls points toward left of frame)
- *square root*: $x^* = \sqrt{x}$ (mildly pulls points toward right of frame)
- *log*: $x^* = log(x)$ (strongly pulls points toward right of frame)
- *inverse*: $x^* = 1/x$ (reverses scale and squeezes points into left of frame)
- *logit*: $x^* = (log(x/(1-x)) + 10)/20$ (squeezes points toward middle of frame)
- *sigmoid*: $x^* = 1/(1 + exp(-20x + 10))$ (expands points away from middle of frame)

Legend:
- None
- Half
- Square
- Sqrt
- Log
- Inverse
- Logit
- Sigmoid

Dang and Wilkinson,
"Transforming Scagnostics to
Reveal Hidden Features", TVCG 2014

# Infinite Zooms

# Prelude: How To Smooth a Discrete Signal?

Slide a window across the signal

- stop at each discrete sample point
- average the original data points that fall into the window
- store this average value at the sample point
- move the window to the next sample point
- repeat

# Zooming into Images



300%
ZOOM

Pixel replication creates jagged edges,
but looks sharp

300%
ZOOM

Smoothing looks more
natural, but blurs detail

PS3

# Anti-Aliasing Via Smoothing



Pixel replication creates jagged edges, but looks sharp

Smoothing looks more natural, but blurs detail

# The Solution

What's the underlying problem?

- detail can't be refined upon zoom
- can just be replicated or blurred

raster graphics

vector graphics

The solution...

- represent detail as a function that can be mathematically refined
- replace raster graphics by vector graphics

# Scalable Vector Graphics (SVG)

# Photographs and Images in SVG

Vector graphics tends to have an "cartoonish" look



raster graphics                    vector graphics

# Photographs and Images in SVG

# D3 Uses SVG



The Wealth & Health of Nations

42.0%

# SMOOTHING FOR DE-NOISING

Filtering/smoothing also eliminates noise in the data

# LET'S TALK ABOUT BAR CHARTS

# BAR CHARTS CAN BE DATA SMOOTHERS

In some ways, bar charts reduce noise and uncertainties in the data

- the bins do the smoothing

Example:

- obesity over age (group)



SOURCE: Analysis of the 2007/08 Canadian Community Health Survey, Statistics Canada.



Gallup-Healthways Well-Being Index

GALLUP

# Categorical Bar Charts

Bar charts that hold categorical data

- smoothing by semantic grouping
- for example, Europe vs. {France, Spain, Italy, Germany, ...}

## Top Oil Reserves



MMbbl = one million barrels

# Bar Charts vs. Histograms

## Histograms

- bars show the frequency of numerical data
- quantitative data
- elements are grouped together, so that they are considered as ranges
- bars cannot be reordered
- width of bars need not be the same

## Bar charts

- uses bars to compare different categories of data
- comparison of discrete variables
- elements are taken as individual entities
- bars can be reordered
- width of bars need to be the same

# How many Bars in a Bar Chart

How many bars are too many (in a chart)

- if individual categories are the focus? 12 is a good rule
- if the overall trend is the important factor? 50 or even more
- eventually you can switch to a line chart



- sort bars by height and use 'other' to aggregate the bar chart tails into a single bar
- find a grouping that can semantically aggregate bars, for example aggregate countries into continents

more information

# Illusions & Ethical Considerations (Visualizations Can Deceive)

# Visualization Can be Deceptive

# Visualization Can be Deceptive

# Visualization Can be Deceptive



Count the number of black dots
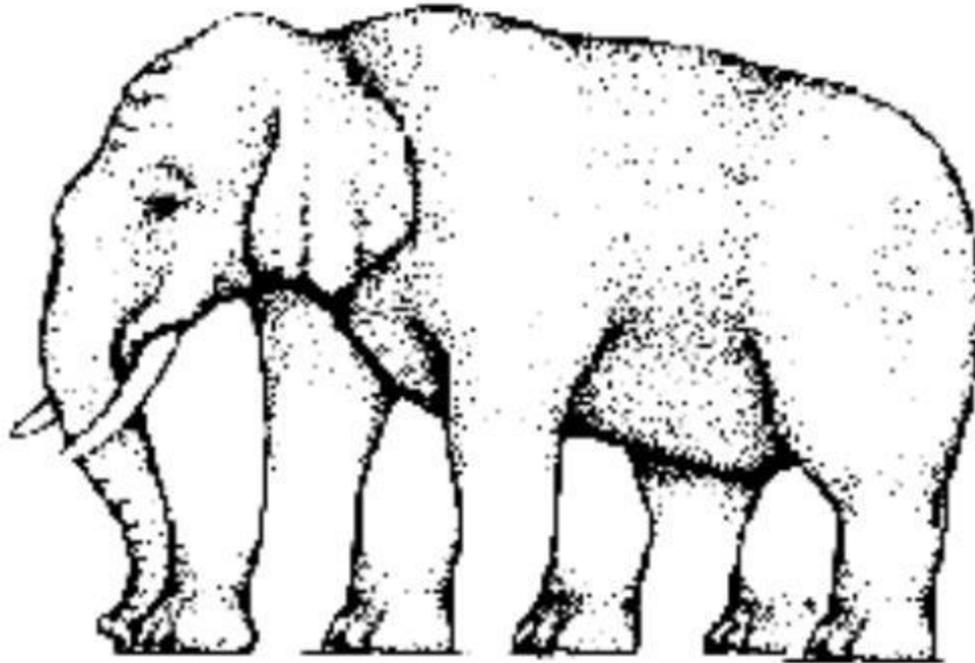
# Visualization Can be Deceptive

# Visualization Can be Deceptive



Are the horizontal lines parallel or do they slope?

# Visualization Can be Deceptive
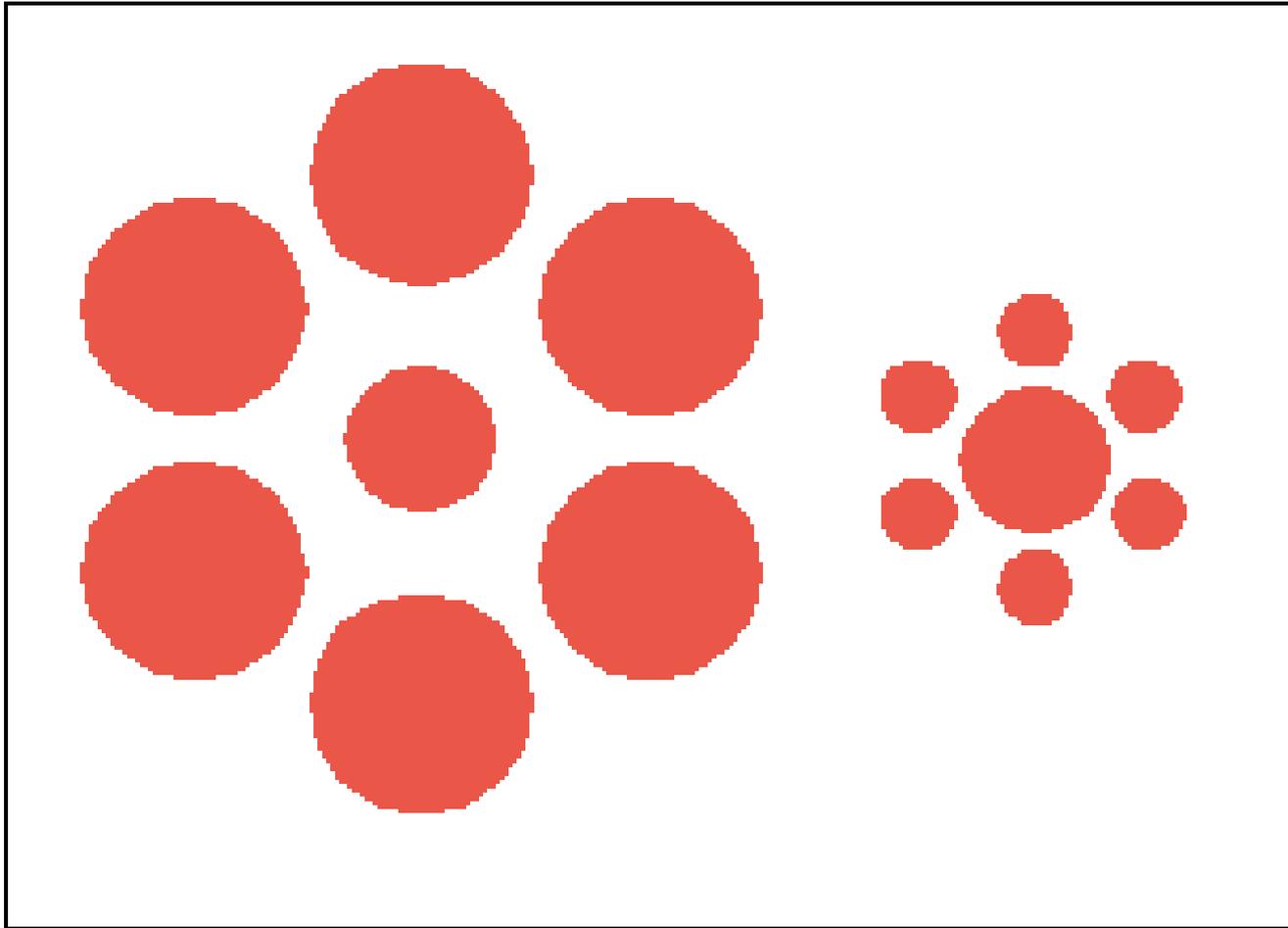


How many legs does this elephant have?

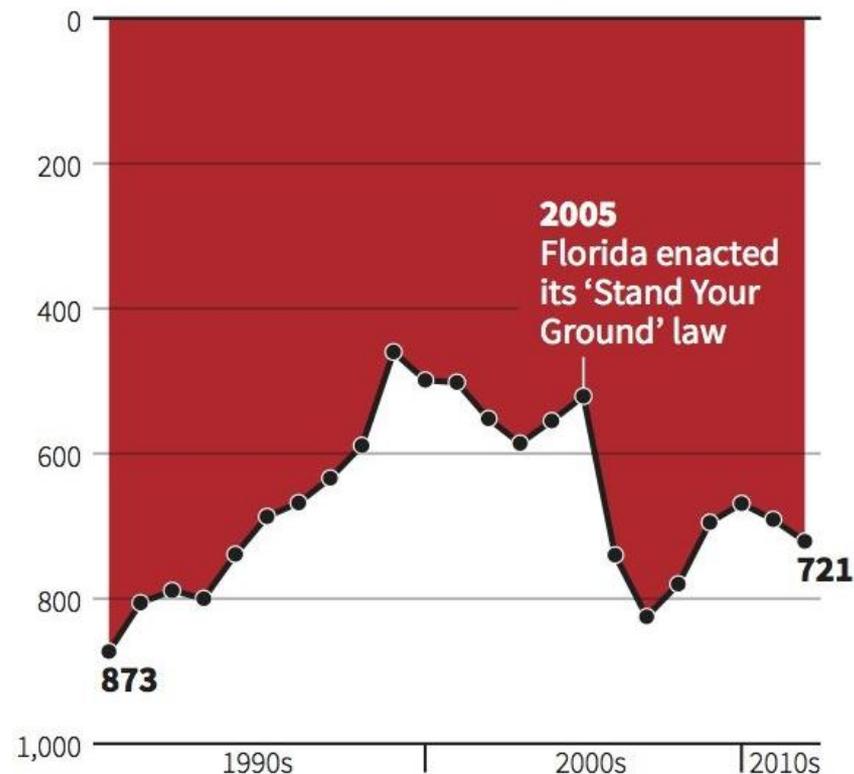# Visualization Can be Deceptive



Julian Beever

# Visualization Can be Deceptive



**Which circle in the middle is bigger?**
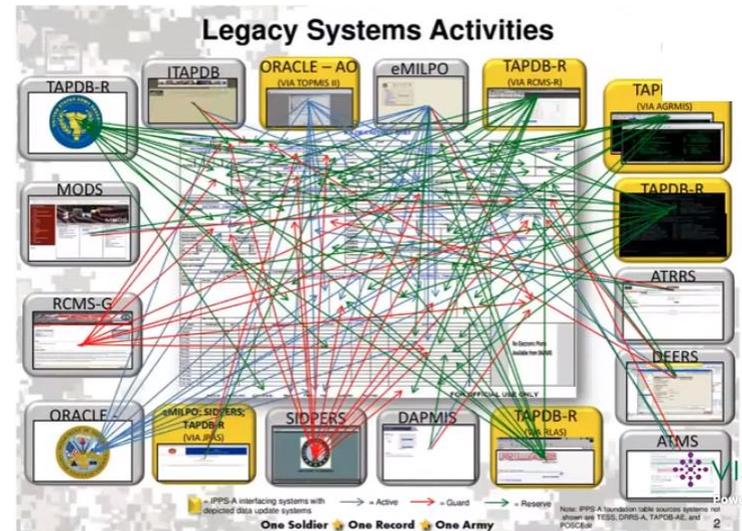
# VISUALIZATION CAN BE DECEPTIVE

# VISUALIZATION CAN BE BS*

From Michael Correll's alt.vis 2021 talk ([link](#))
- Don't relate to the real world
- Don't really help people understand their data
- Don't even have the decency to *lie* to you



"Stock footage chart"
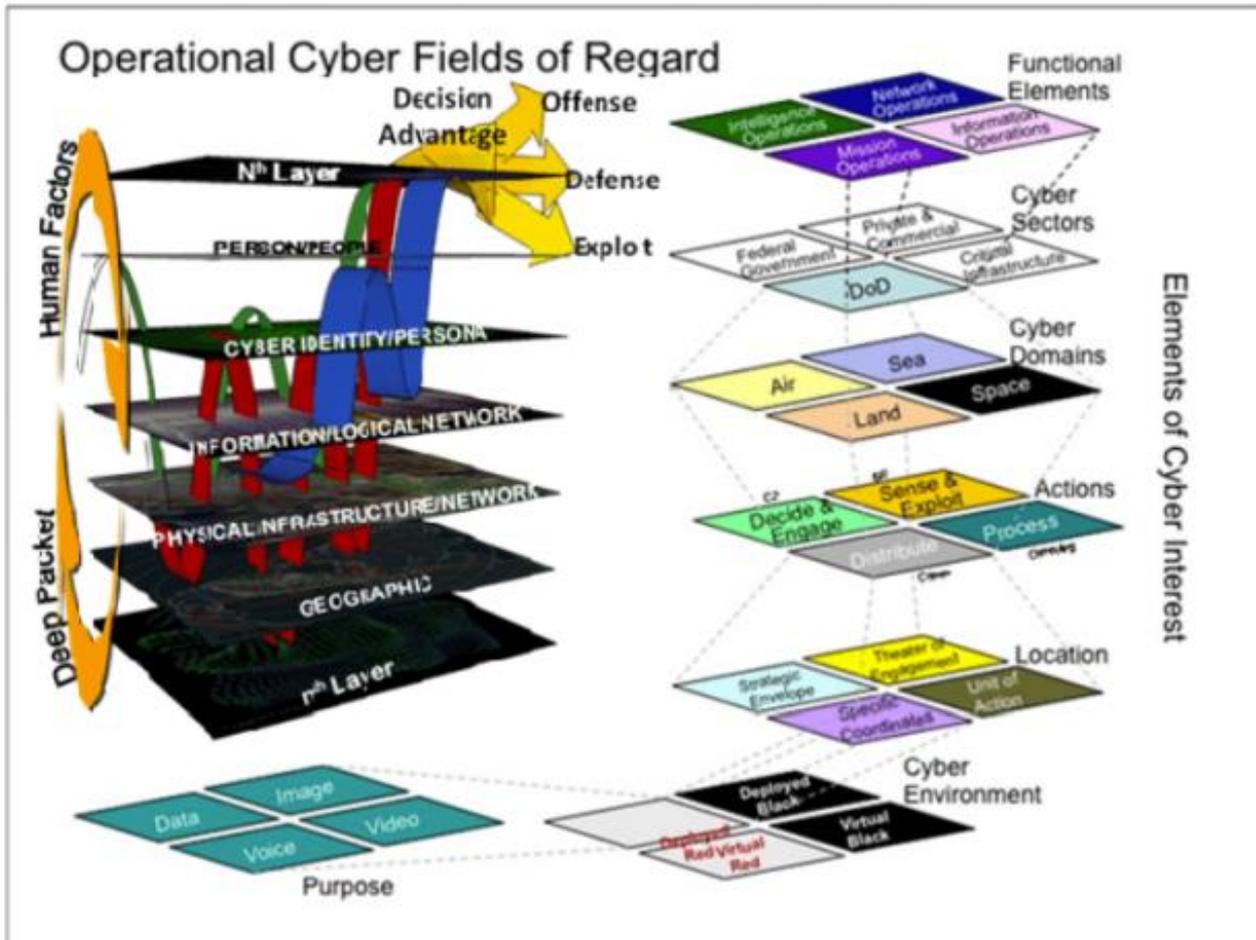


"Novocaine chart"

This stuff is way too complex for you to understand. Aren't you glad there's somebody smart like me taking care of it?
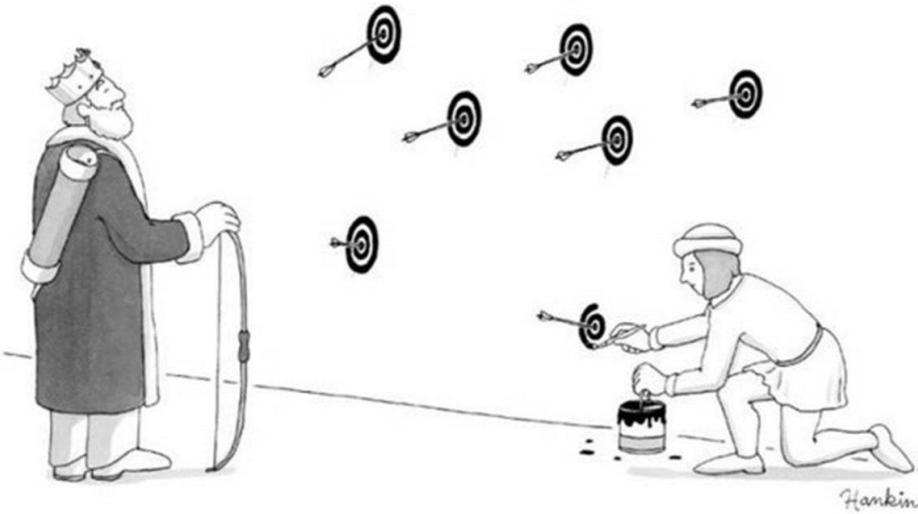
* Bullshit

# Visualization Can Be BS



Hwang's @DefenseCharts Twitter account, "dedicated to the presentational aesthetics of the defense-industrial complex"
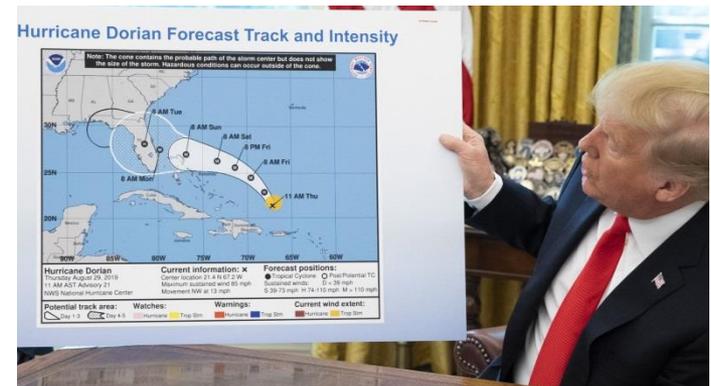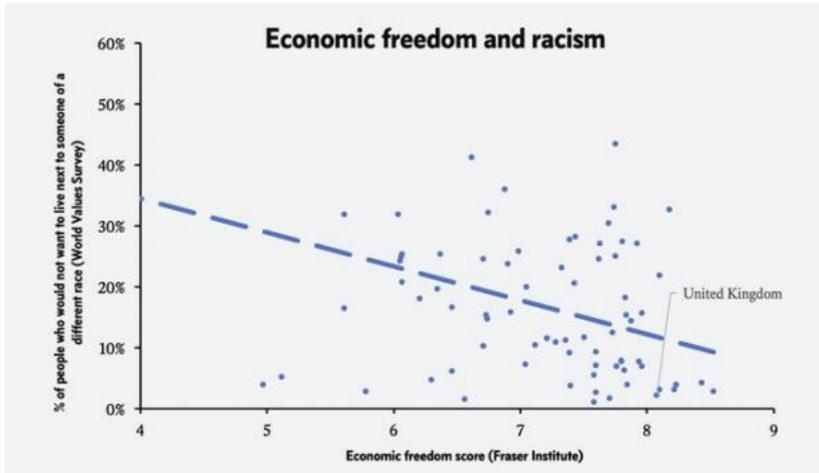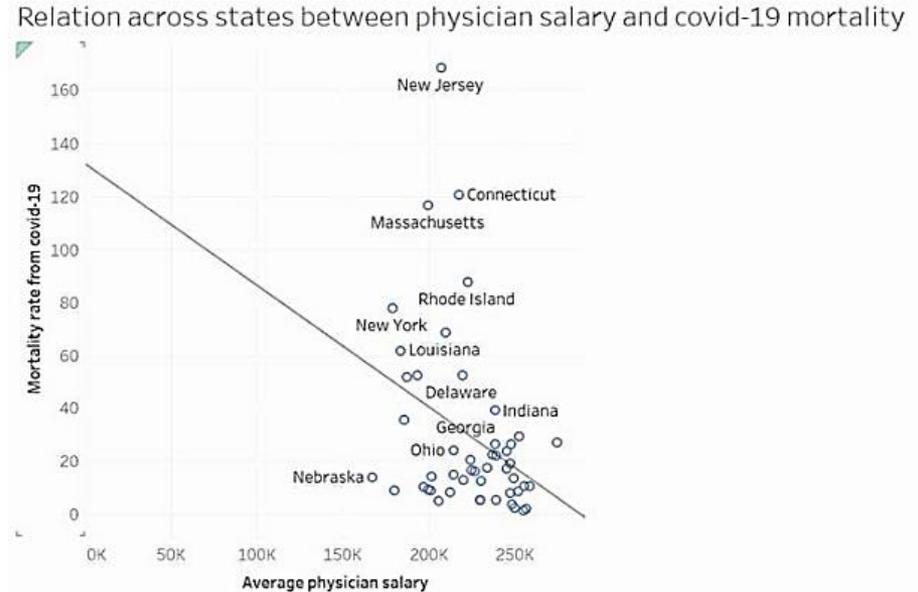
# VISUALIZATION CAN BE BS



Example: Sharpiegate

"Texas sharpshooter chart"

# VISUALIZATION CAN BE BS



Economic freedom and racism

To show: "Countries with more economic freedom have less racist attitudes"



Relation across states between physician salary and covid-19 mortality

To show: "States where physicians are highly paid have lower COVID-19 mortality per capita"



the DECLINE of ARTISTIC STANDARDS

Artificial noise added to make the chart look like there is a complex metric being measured precisely over time (when it is really not)

# DETECT AND FIX MISLEADING VISUALIZATIONS

There is an app for that:

- MisVisFix: An Interactive Dashboard for Detecting, Explaining & Correcting Misleading Visualization (Das, Mueller, IEEE VIS 2025)



youtube