

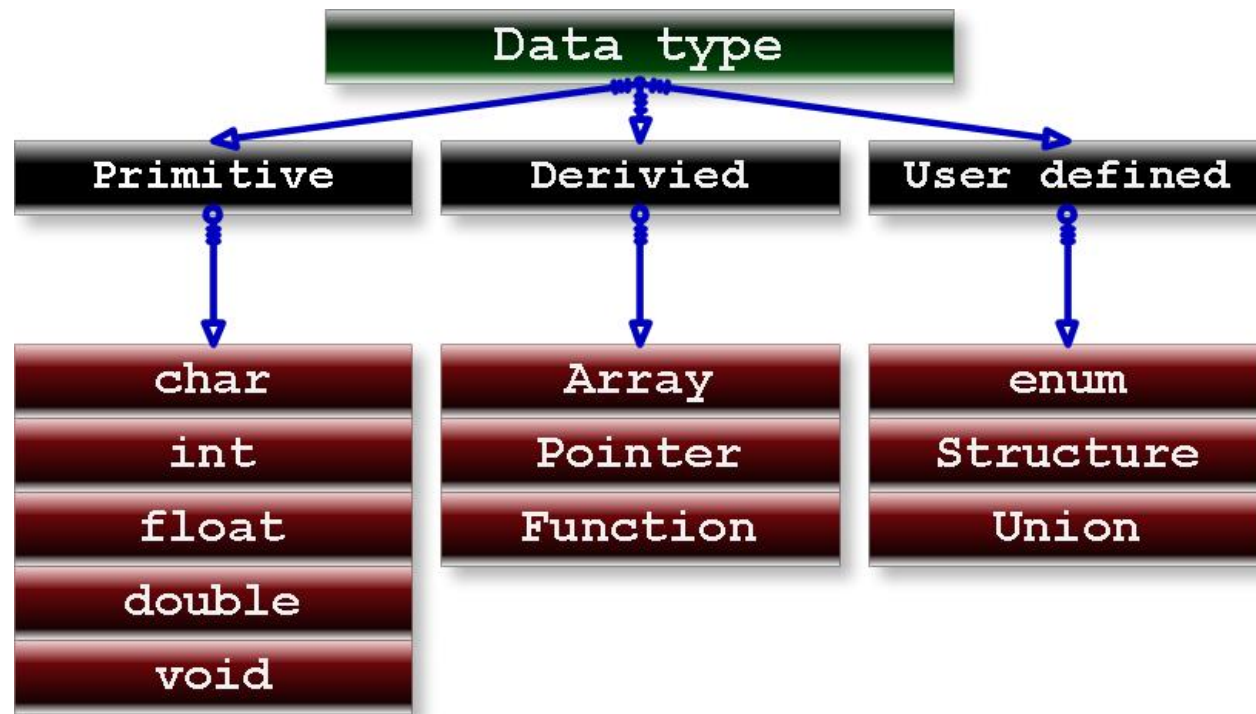
CSE 564
VISUALIZATION & VISUAL ANALYTICS
APPLICATIONS AND BASIC TASKS

KLAUS MUELLER

COMPUTER SCIENCE DEPARTMENT
STONY BROOK UNIVERSITY

Lecture	Topic	Projects
1	Intro, schedule, and logistics	
2	Applications of visual analytics	
3	Basic tasks, data types	Project #1 out
4	Data assimilation and preparation	
5	Introduction to D3	
6	Bias in visualization	
7	Data reduction and dimension reduction	
8	Visual perception	Project #2(a) out
9	Visual cognition	
10	Visual design and aesthetics	
11	Cluster analysis: numerical data	
12	Cluster analysis: categorical data	Project #2(b) out
13	High-dimensional data visualization	
14	Dimensionality reduction and embedding methods	
15	Principles of interaction	
16	Midterm #1	
17	Visual analytics	Final project proposal call out
18	The visual sense making process	
19	Maps	
20	Visualization of hierarchies	Final project proposal due
21	Visualization of time-varying and time-series data	
22	Foundations of scientific and medical visualization	
23	Volume rendering	Project 3 out
24	Scientific and medical visualization	Final Project preliminary report due
25	Visual analytics system design and evaluation	
26	Memorable visualization and embellishments	
27	Infographics design	
28	Midterm #2	

DATA TYPES EVERY CS PERSON KNOWS



DATA TYPES IN VISUAL ANALYTICS

Numeric

Categorical

Text

Time series

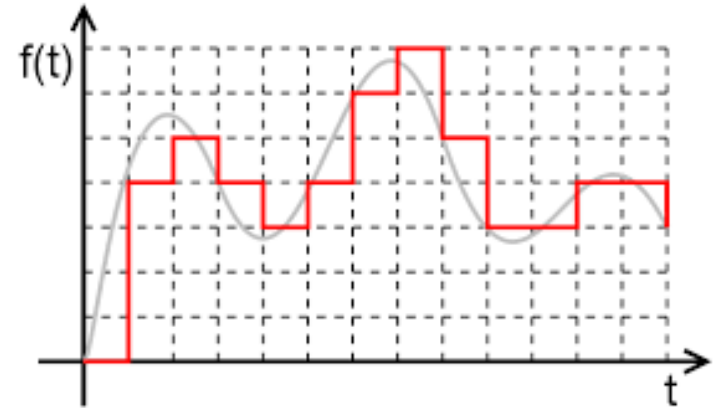
Graphs and networks

Hierarchies

VARIABLES IN STATISTICS

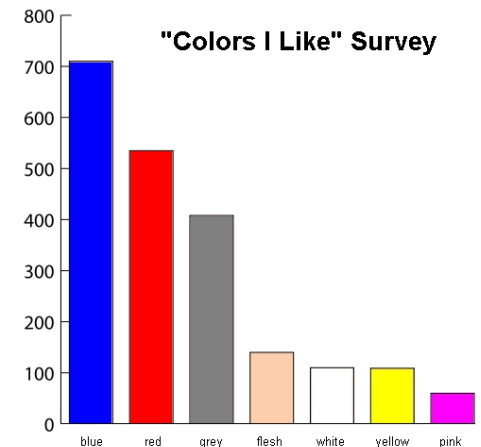
Numeric variables

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



Categorical variables

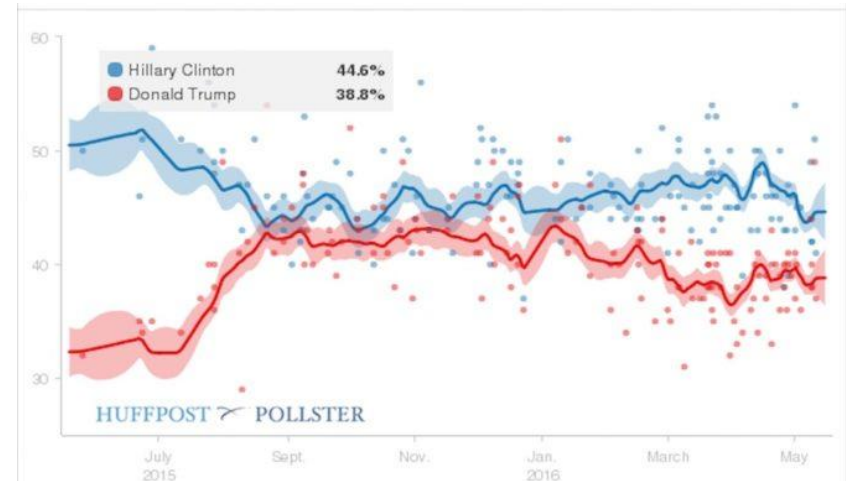
- describe a **quality** or characteristic
- like: 'what type' or 'which category'



NUMERIC VARIABLES

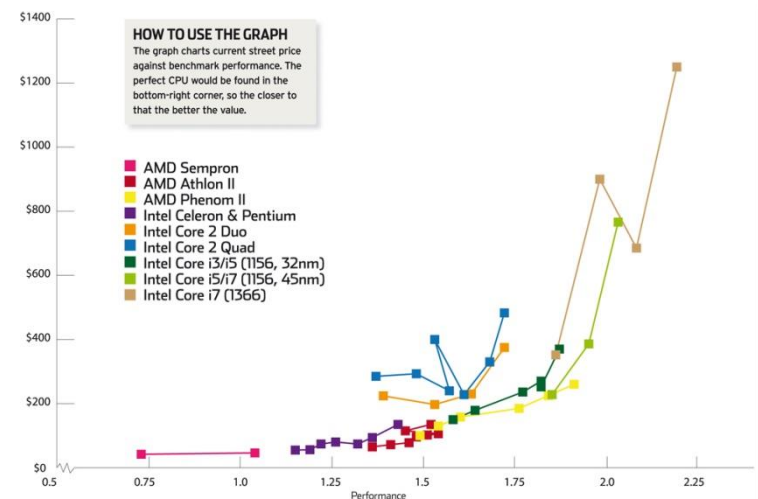
Most often the x-axis is 'time'

- 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

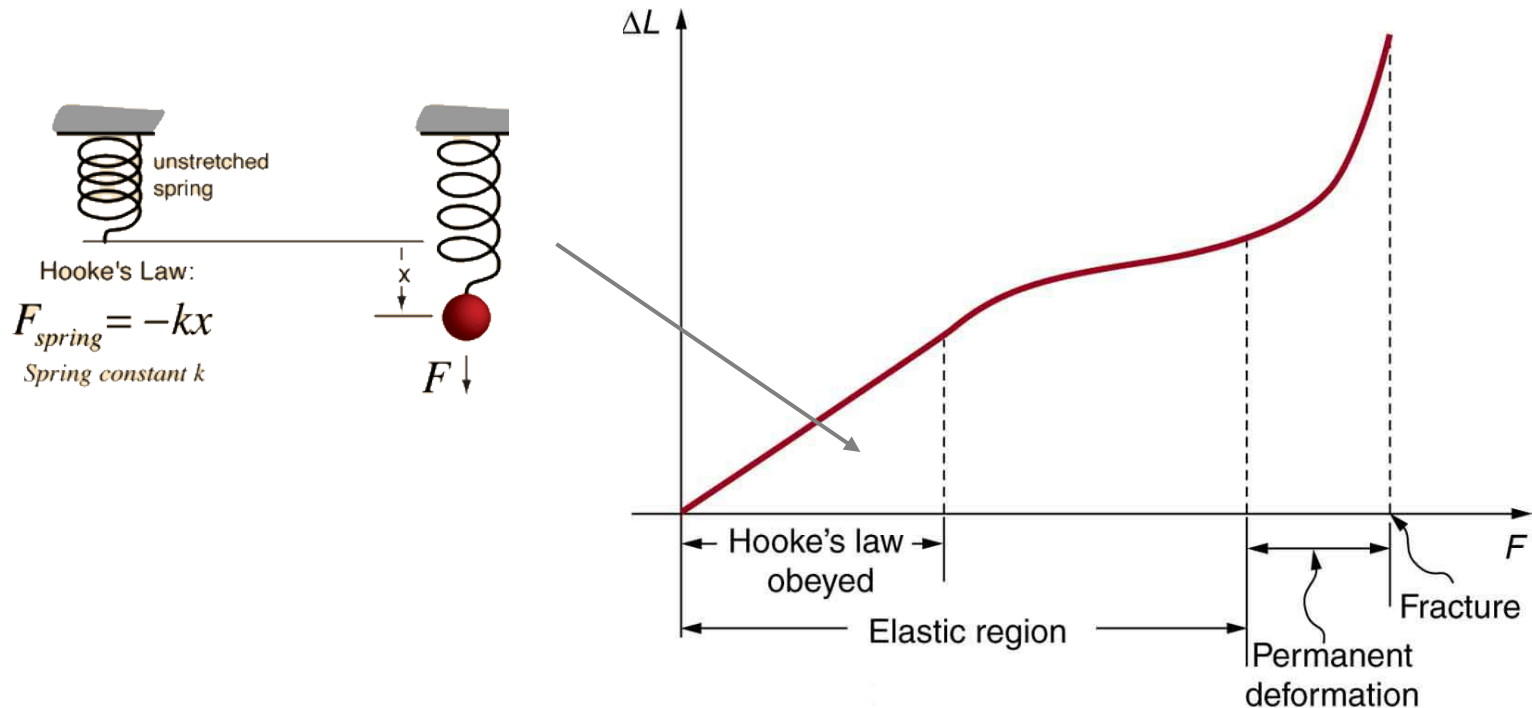
- engineers, statisticians, etc. will be receptive to this idea
- can you think of an example?



NUMERIC VARIABLES

Another plot where 'time' is not the x-axis

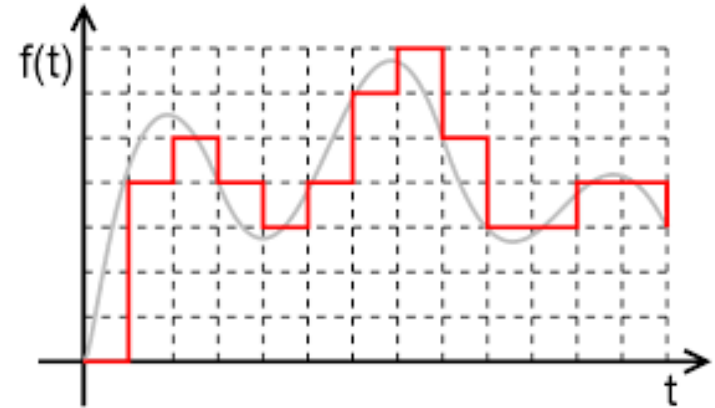
- from the engineering / physics domain
- in some sense, it tells a story



VARIABLES IN STATISTICS

Numeric variables

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

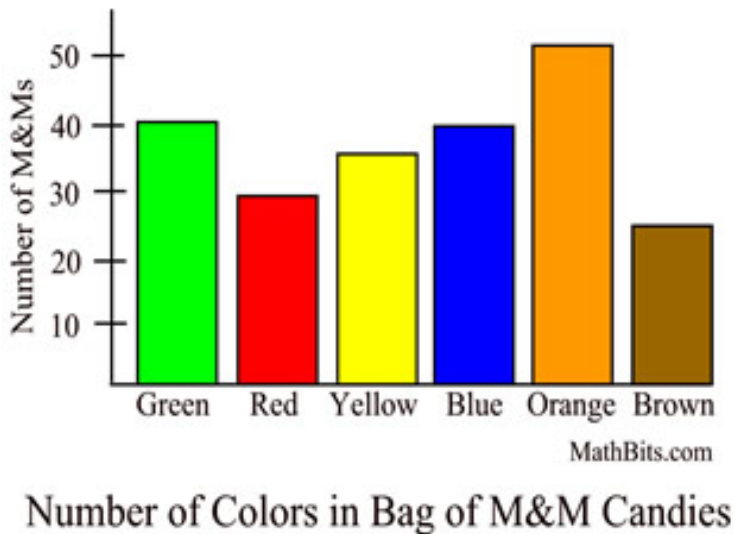


Categorical variables

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



??

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
- ...



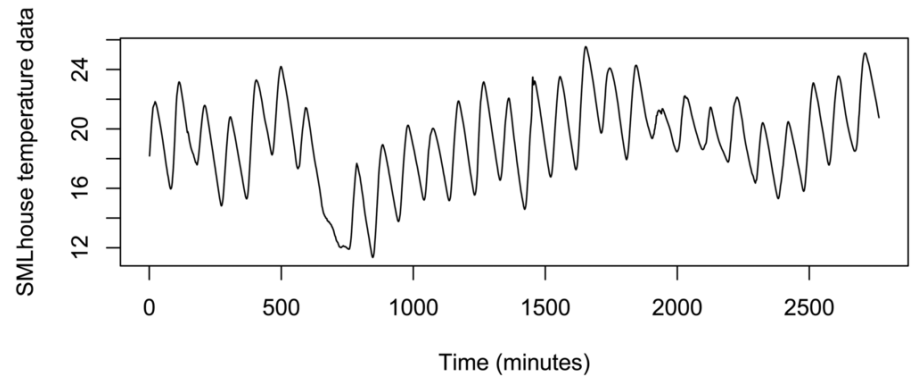
Do **feature analysis** to turn these abstract things into numbers

- 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

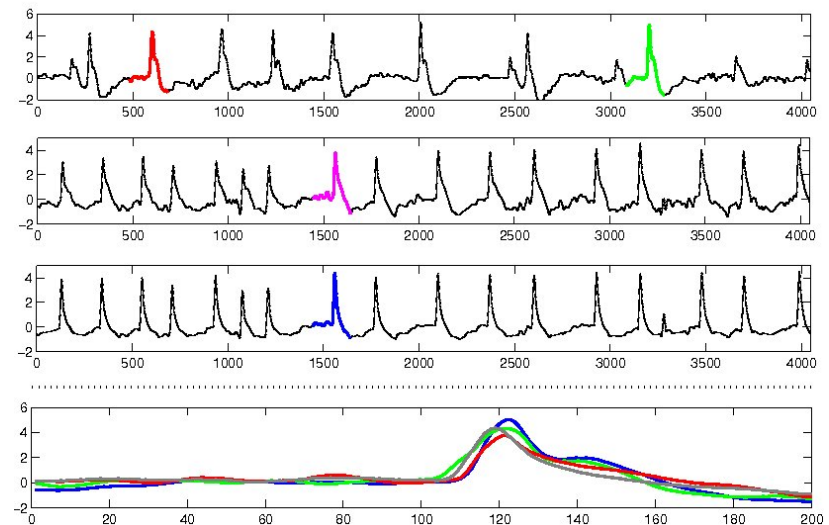
Characteristics

- often large scale
- time series

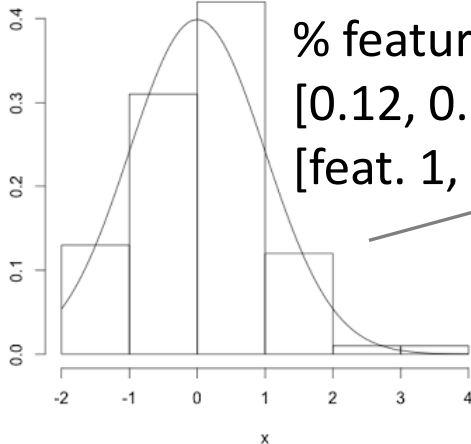


Feature Analysis

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



Motif discovery



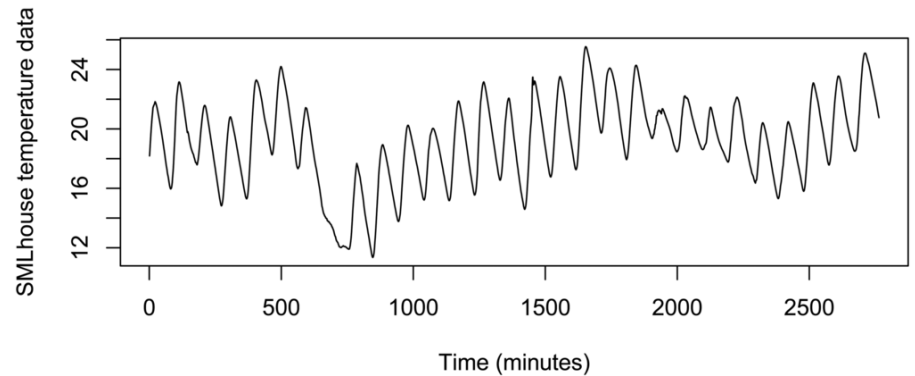
% features discovered in stream
[0.12, 0.3, 0.41, 0.12, 0.05]
[feat. 1, feat. 2, .., feat. 5]



SENSOR DATA

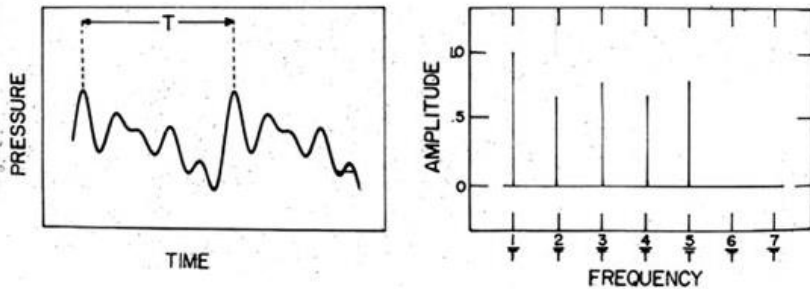
Characteristics

- often large scale
- time series



Feature Analysis

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



Fourier transform

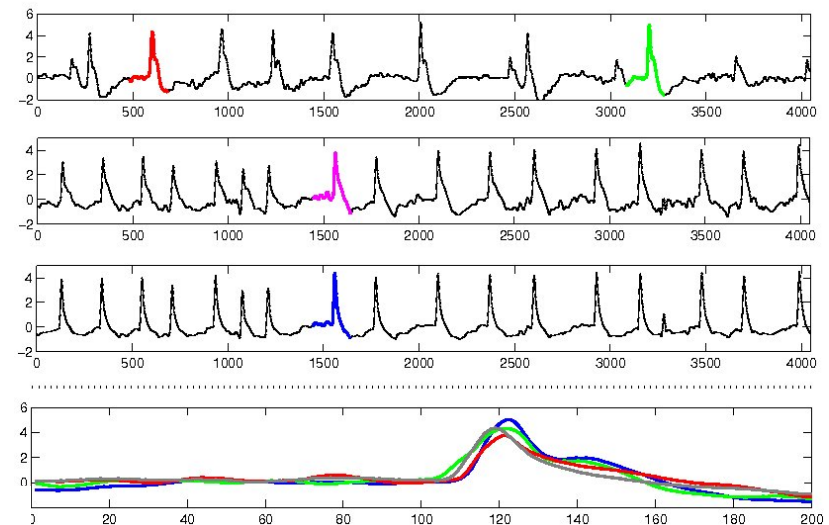


IMAGE DATA

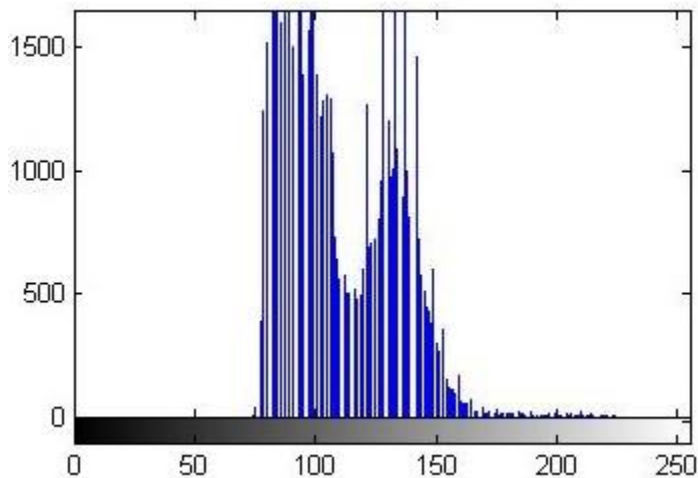
Characteristics

- array of pixels

Feature Analysis

- value histograms
- encode into a 256-D vector

histograms



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



IMAGE DATA

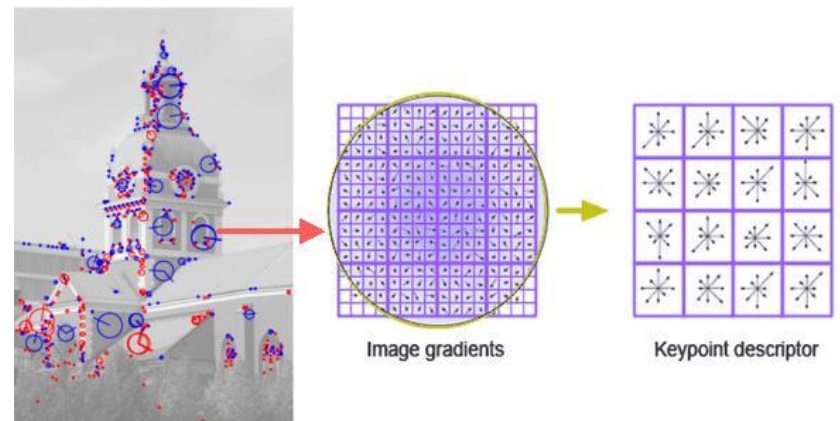
Characteristics

- array of pixels

Feature Analysis

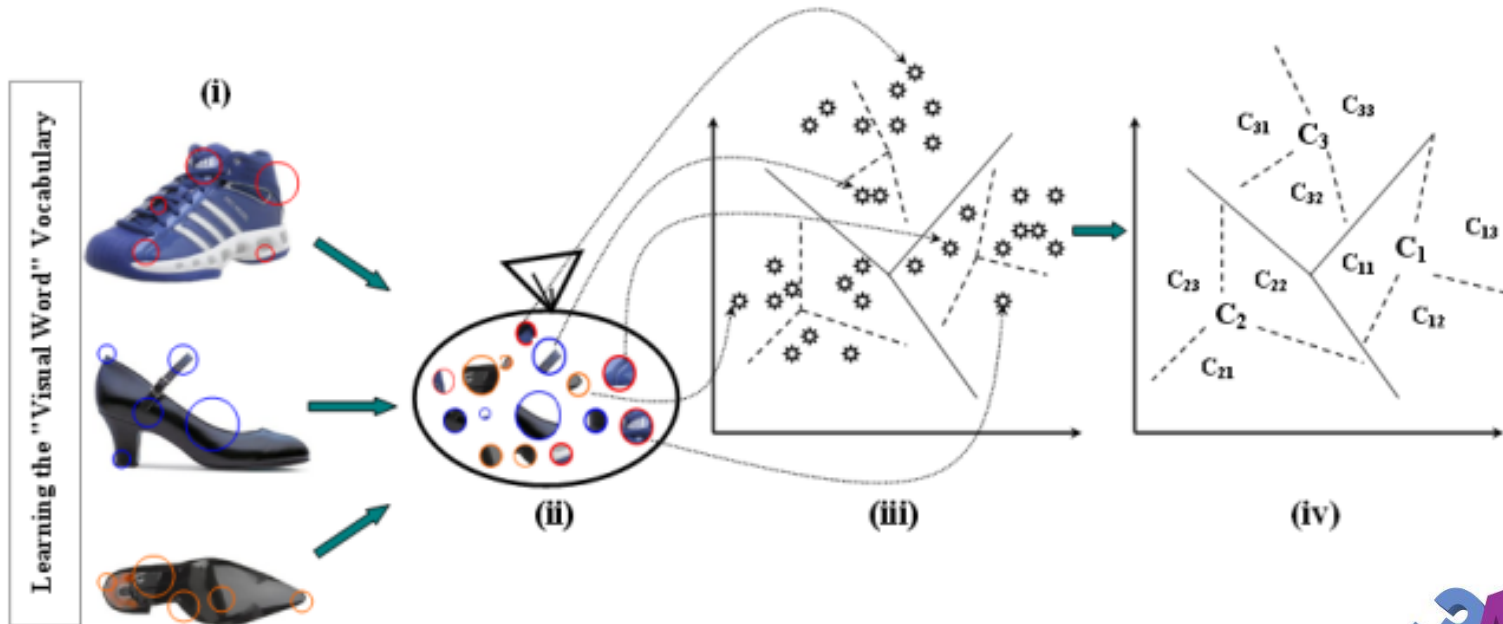
- value histograms
- gradient histograms
- FFT, FWT
- Scale Invariant Feature Transform (SIFT)
- Bag of Features (BoF)
- visual words

histograms



SIFT

BAG OF FEATURES (BoF)



BAG OF FEATURES (BoF)

1. Obtain the set of bags of features

- (i) Select a large set of images
- (ii) Extract the SIFT feature points of all the images in the set and obtain the SIFT descriptor for each feature point extracted from each image
- (iii) Cluster the set of feature descriptors for the amount of bags we defined and train the bags with clustered feature descriptors
- (iv) Obtain the visual vocabulary

2. Obtain the BoF descriptor for a given image/video frame

- (v) Extract SIFT feature points of the given image
- (vi) Obtain SIFT descriptor for each feature point
- (vii) Match the feature descriptors with the vocabulary we created in the first step
- (viii) Build the histogram

[More information](#)

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

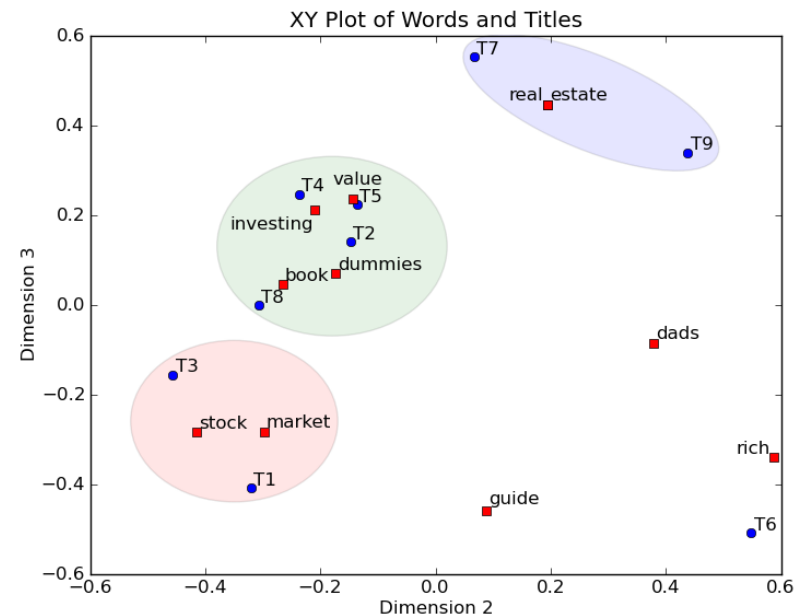
Create a term-document matrix

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

Index Words	Titles								
	T1	T2	T3	T4	T5	T6	T7	T8	T9
book			1	1					
dads						1			1
dummies		1						1	
estate							1		1
guide	1					1			
investing	1	1	1	1	1	1	1	1	1
market	1		1						
real							1		1
rich						2			1
stock	1		1					1	
value				1	1				

Term-Document Matrix

LSA
→

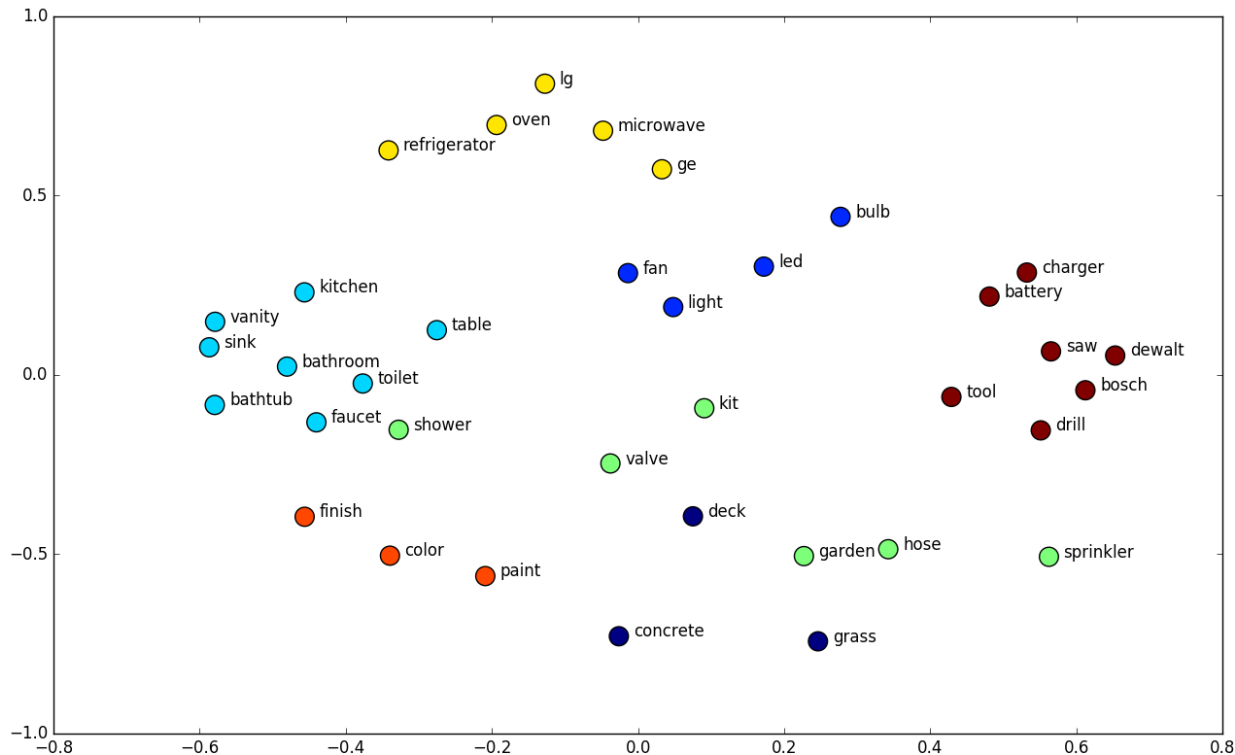


Word/document cluster

WORD EMBEDDING

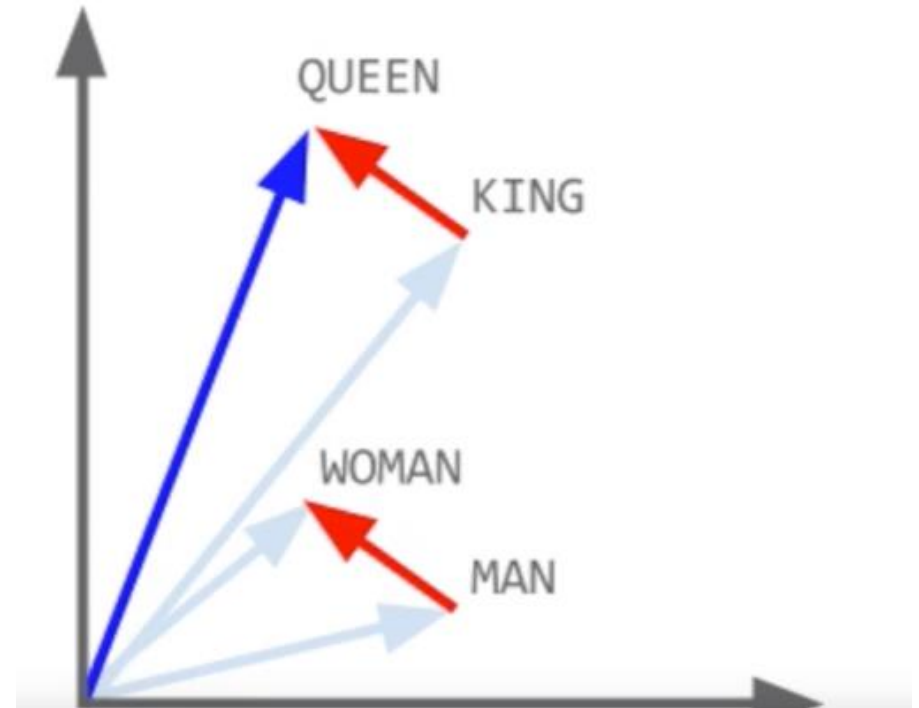
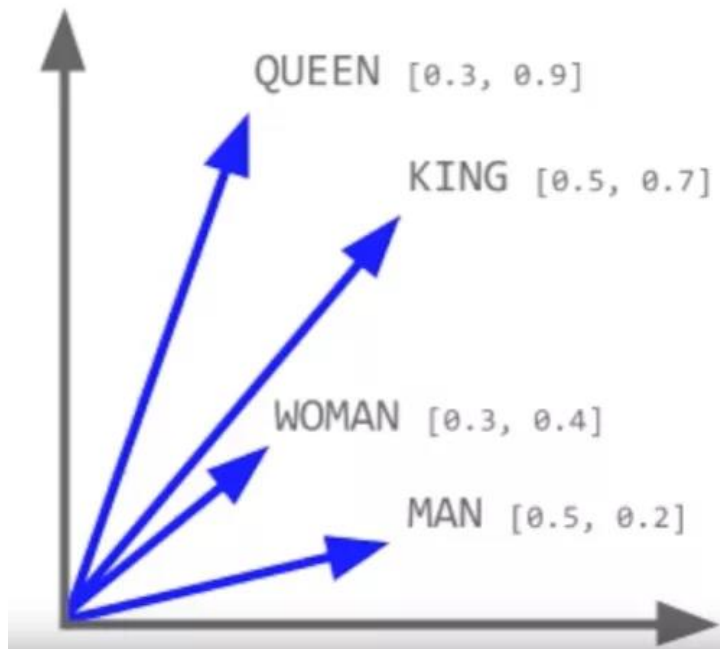
Train a shallow neural network (NN) on a corpus of text

- the NN weight vectors encode word similarity as a high-D vector
- use a 2D embedding technique to display



WORD EMBEDDING ALGEBRA

Load up the word vectors



gender = WOMAN - MAN

QUEEN = KING + **gender**

QUEEN = KING - MAN + WOMAN

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

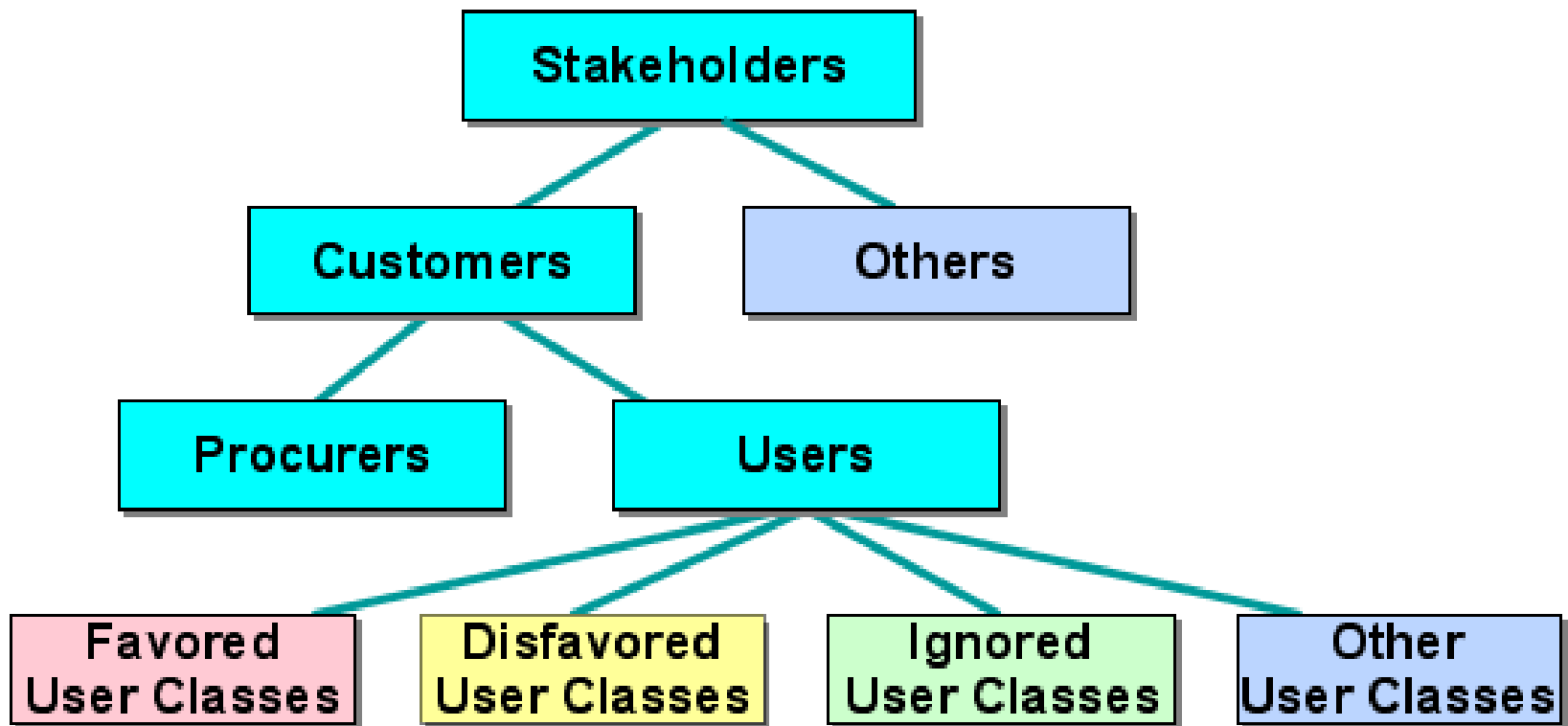
- characteristics of the network packets are used to analyze intrusions or other interesting activity
- a variety of features may be extracted from these packets
 - the number of bytes transferred
 - the network protocol used
 - IP ports used



LET'S LOOK AT SOME ESSENTIAL
GRAPHICAL REPRESENTATIONS

AND DO SOME ADVERTISING FOR D3

STAKEHOLDER HIERARCHY



FUNCTION CALL TREE



MORE COMPLEX STAKEHOLDER HIERARCHY



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>

VORONOI TESSELLATION

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

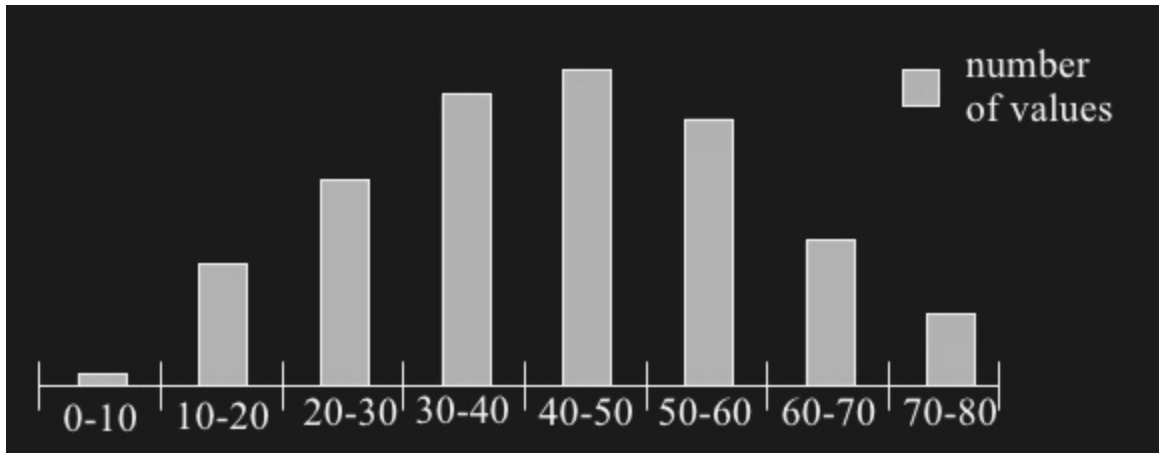
<https://observablehq.com/collection/@d3/d3-delaunay>

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

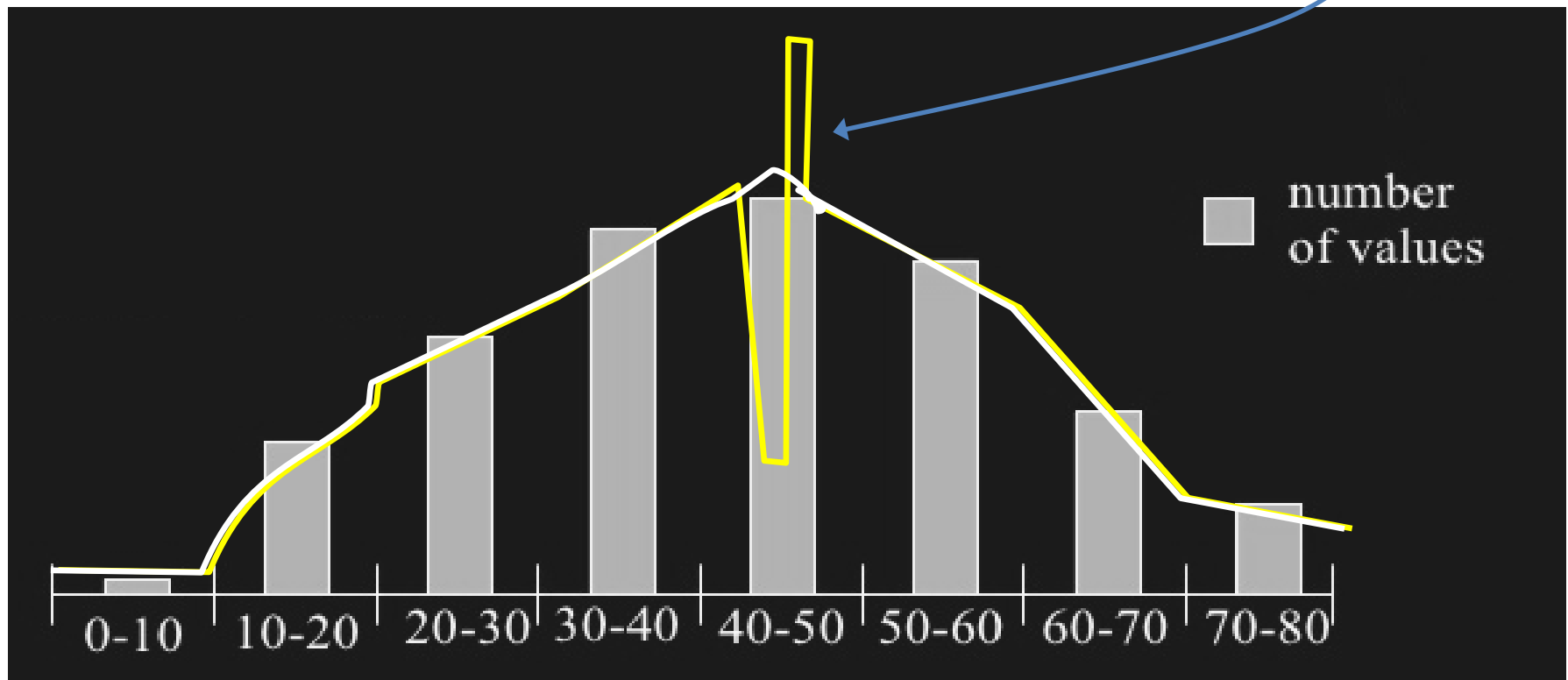


- what is lost here?

PROBLEM WITH EQUI-WIDTH HISTOGRAM

Age ranges of customers could be unevenly distributed within a bin

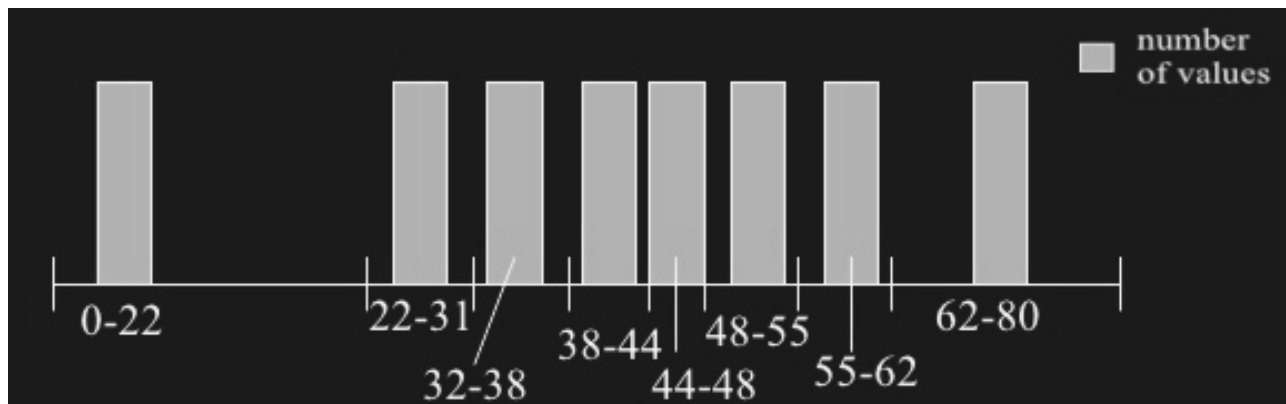
- this could be an interesting anomaly



NUMERIC TO CATEGORICAL DATA: DISCRETIZATION (2)

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

ONE MORE BINNING METHOD

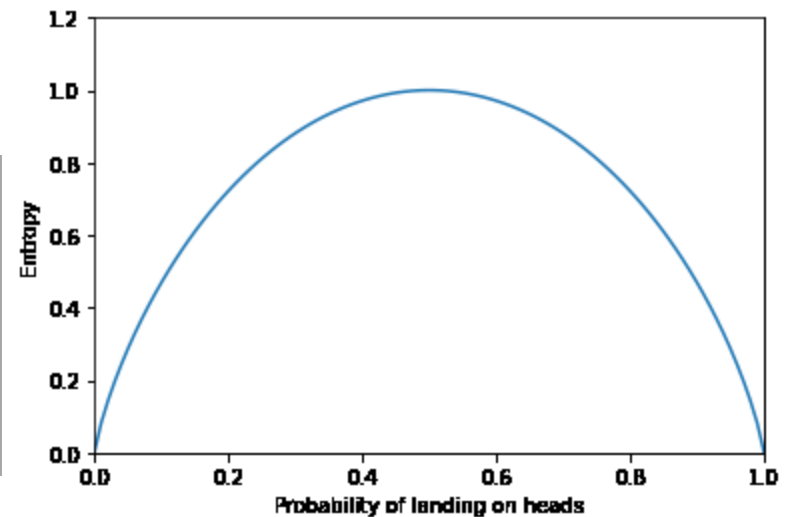
Entropy-based binning

$S(\text{solid}) < S(\text{liquid}) < S(\text{gas})$



Entropy is the amount of surprise to make a certain observation

$$H(X) = - \sum_{i=1}^n P(x_i) \log_b P(x_i)$$



THE DATA

O-Ring Failure	Temperature
Y	53
Y	56
Y	57
N	63
N	66
N	67
N	67
N	67
N	68
N	69
N	70
Y	70
Y	70
Y	70
N	72
N	73
N	75
Y	75
N	76
N	76
N	78
N	79
N	80
N	81

from
https://www.saedsayad.com/supervised_binning.htm

ENTROPY BASED BINNING (EBB)

Aim:

- find the best split so that the bins are as pure as possible that is the majority of the values in a bin correspond to have the same class label
- formally, it is characterized by finding the split with the maximal information gain.

Step 1: Calculate "Entropy" for the target.

$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

O-Ring Failure	
Y	N
7	17

$$E(\text{Failure}) = E(7, 17) = E(0.29, .71) = -0.29 \times \log_2(0.29) - 0.71 \times \log_2(0.71) = \mathbf{0.871}$$

ENTROPY BASED BINNING (EBB)

Step 2: Calculate "Entropy" for the target given a bin.

$$E(S,A) = \sum_{v \in A} \frac{|S_v|}{|S|} E(S_v)$$

		O-Ring Failure	
		Y	N
Temperature	<= 60	3	0
	> 60	4	17

$$E(\text{Failure, Temperature}) = P(<=60) \times E(3,0) + P(>60) \times E(4,17) = 3/24 \times 0 + 21/24 \times 0.7 = \mathbf{0.615}$$

Step 3: Calculate "Information Gain" given a bin.

$$\mathbf{Information\ Gain} = E(S) - E(S,A)$$

$$\mathbf{Information\ Gain\ (Failure,\ Temperature)} = \mathbf{0.256}$$

ENTROPY BASED BINNING (EBB)

[≤ 60 , > 60] turns out to be the best split

Iterate for further splits for bins with highest entropies

Gain = 0.256

		O-Ring Failure	
		Y	N
Temperature	≤ 60	3	0
	> 60	4	17

Gain = 0.101

		O-Ring Failure	
		Y	N
Temperature	≤ 70	6	8
	> 70	1	9

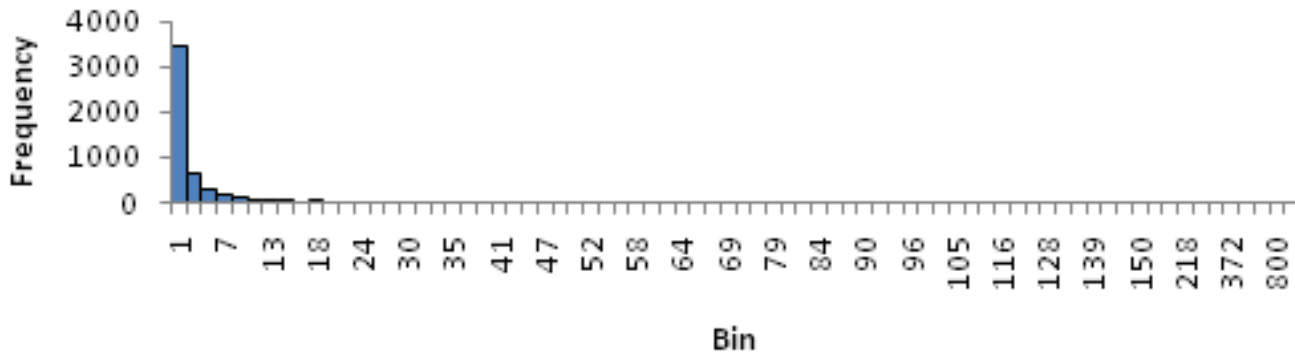
Gain = 0.148

		O-Ring Failure	
		Y	N
Temperature	≤ 75	7	11
	> 75	0	6

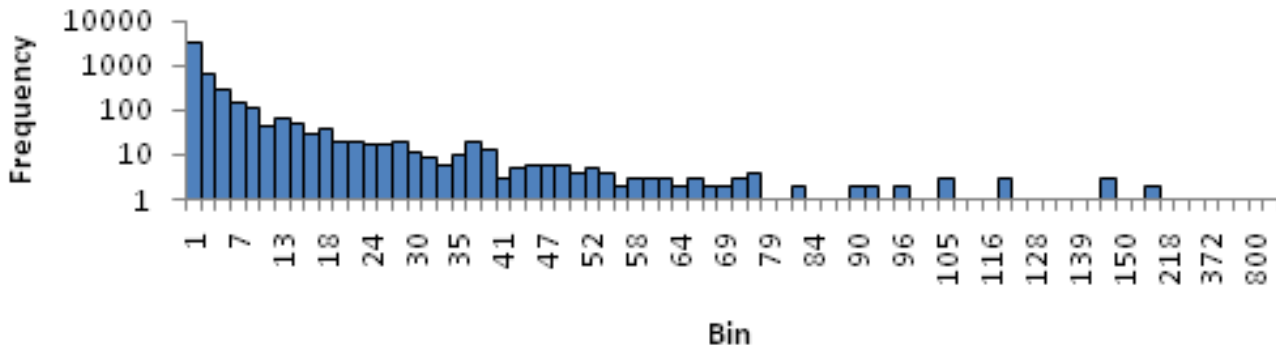
NUMERIC TO CATEGORICAL DATA: DISCRETIZATION (3)

Solution 3:

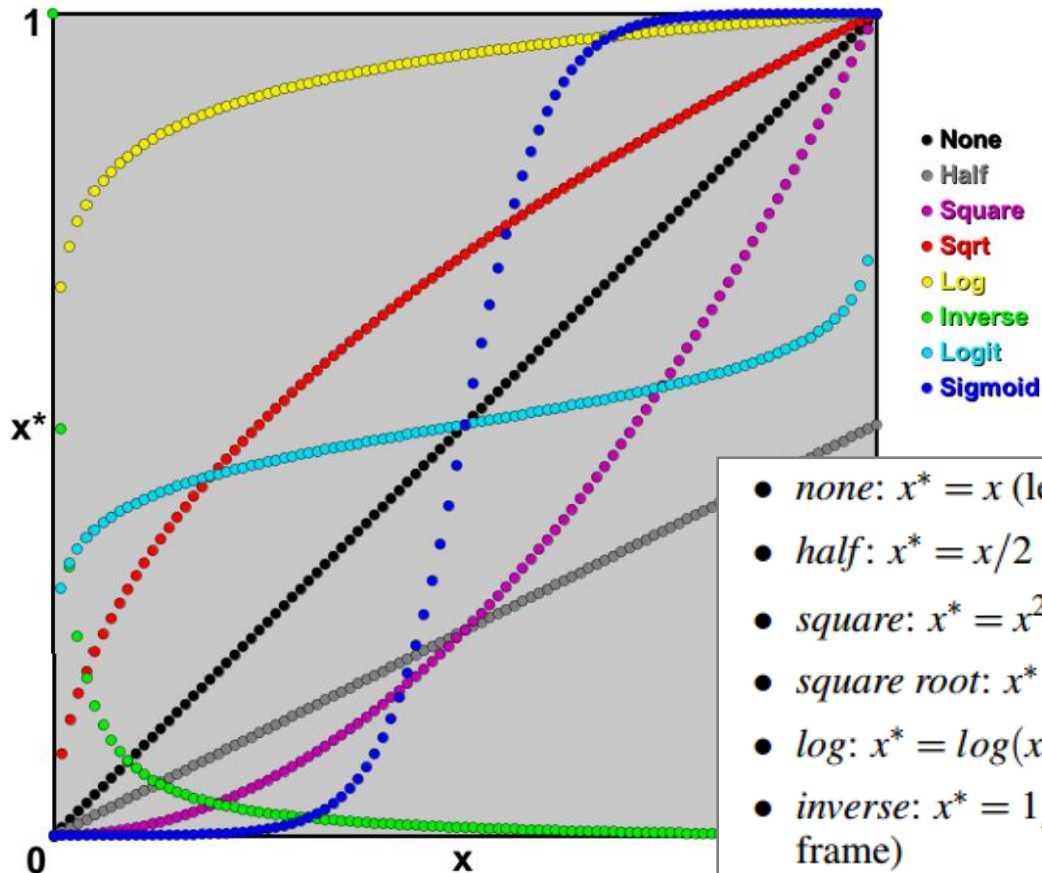
- what if all the bars have seemingly the same height
- or are dominated by one large peak



- 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)

DATA REPRESENTATION

Ever tried to reduce the size of an image and you got this?



This is aliasing

DATA REPRESENTATION

But what you really wanted is this:



This is *anti-aliasing*

WHY IS THIS HAPPENING?



The smaller image resolution cannot represent the image detail captured at the higher resolution

- skipping this small detail leads to these undesired artifacts

WHAT IS ANTI-ALIASING

Procedure

- either sample at a higher rate
- or smooth the signal before sampling it
- the latter is called *filtering*

ANTI-ALIASING VIA SMOOTHING



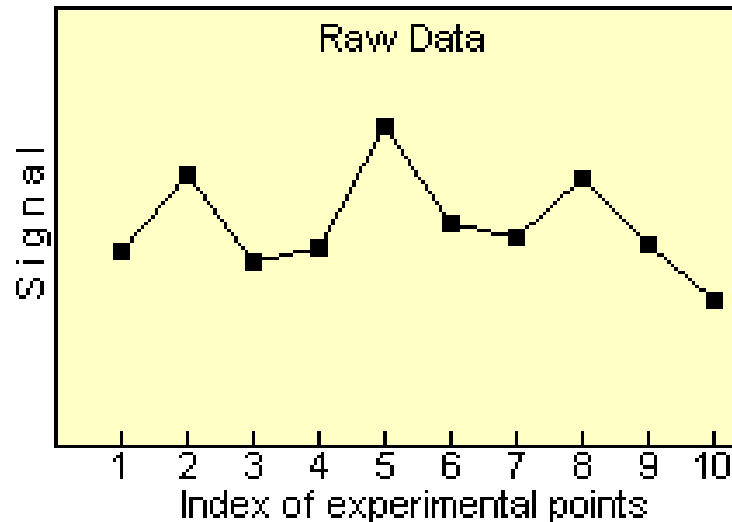
ANTI-ALIASING VIA SMOOTHING



WHAT IS SMOOTHING?

Slide a window across the signal

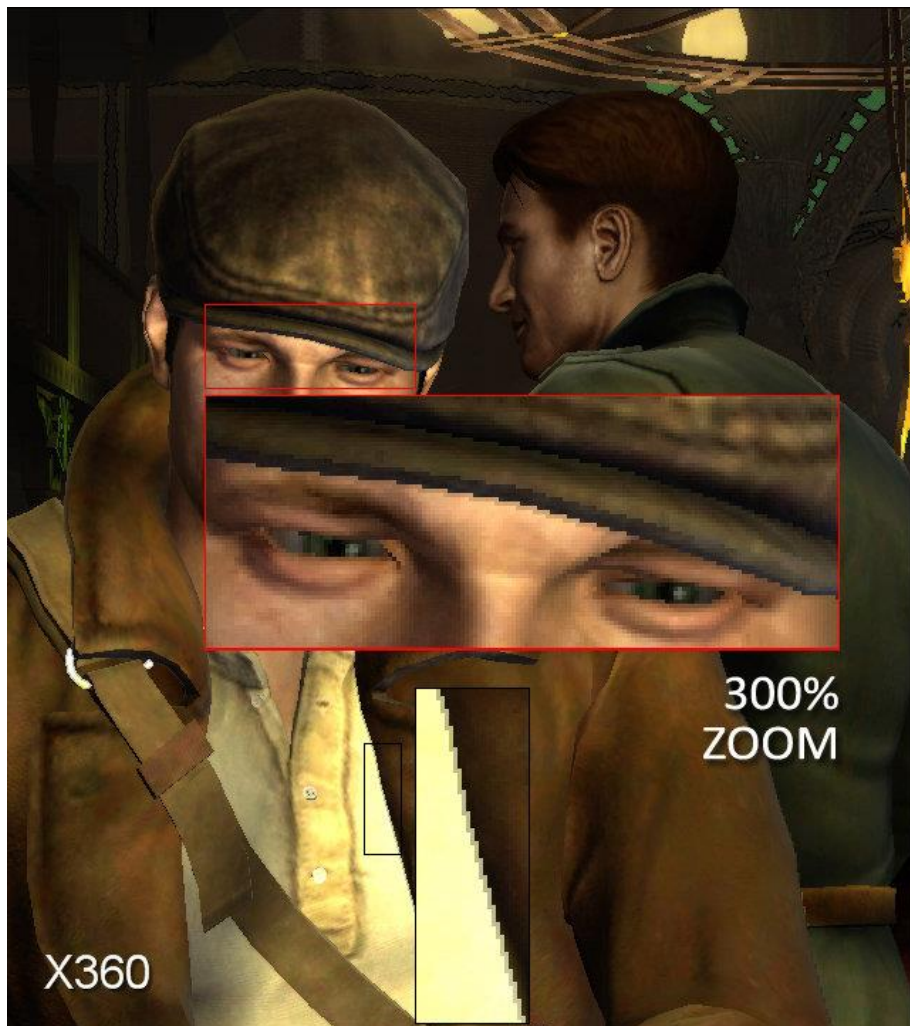
- stop at each discrete sample point
- average the original data points that fall into the window
- store this average value at the sample point
- move the window to the next sample point
- repeat



ANTI-ALIASING VIA SMOOTHING: TRADEOFFS

looks sharper, but has “jaggies”

a bit blurred, but no more jaggies



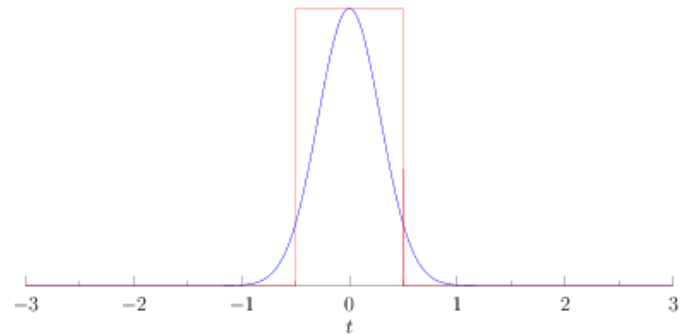
FILTERS

What is the filter we just used called?

- it's called a *box filter*

There are other filters

- for example, Gaussian filter
- yields a smoother result
- box filtering is simplest



BOX FILTER VS. GAUSSIAN FILTER

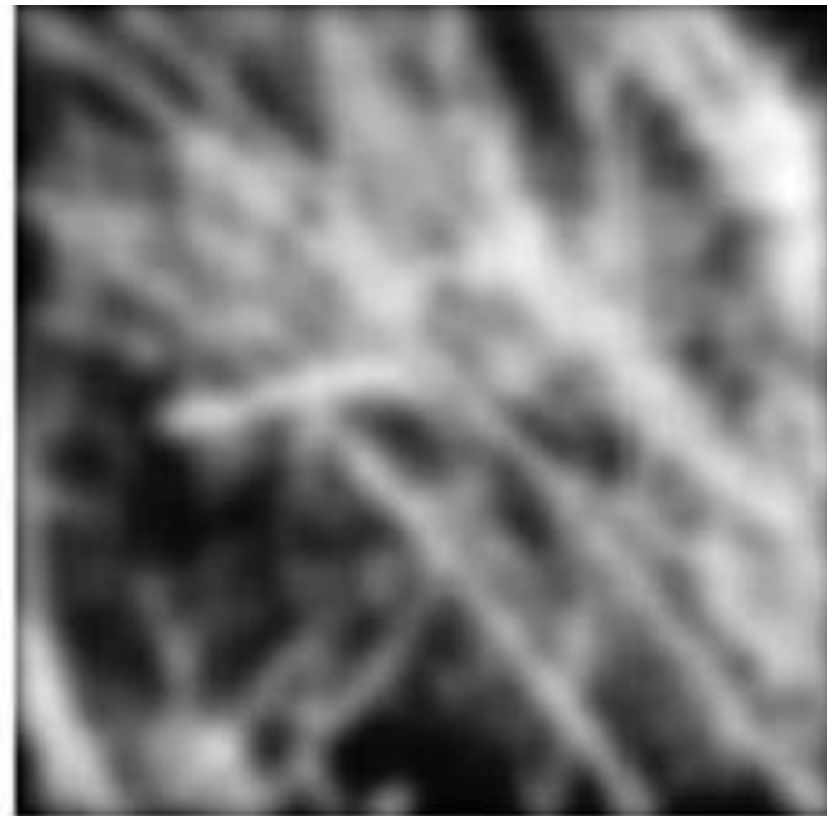


Can you see
some patterns?

It's another form
of aliasing



2D box

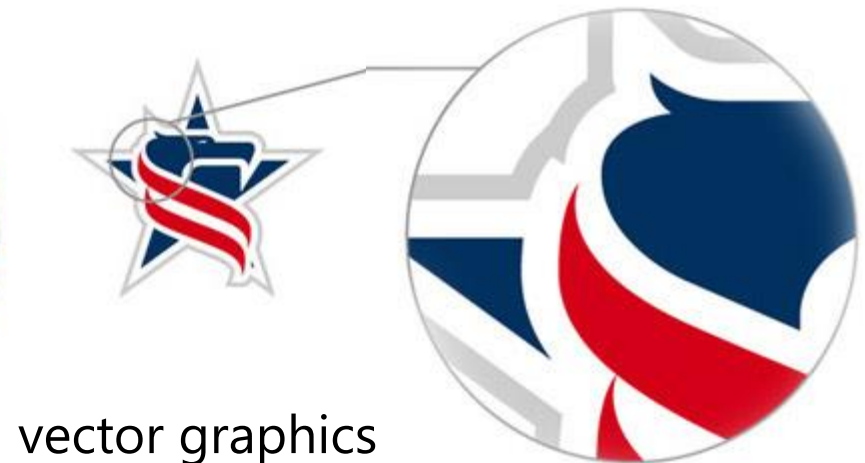
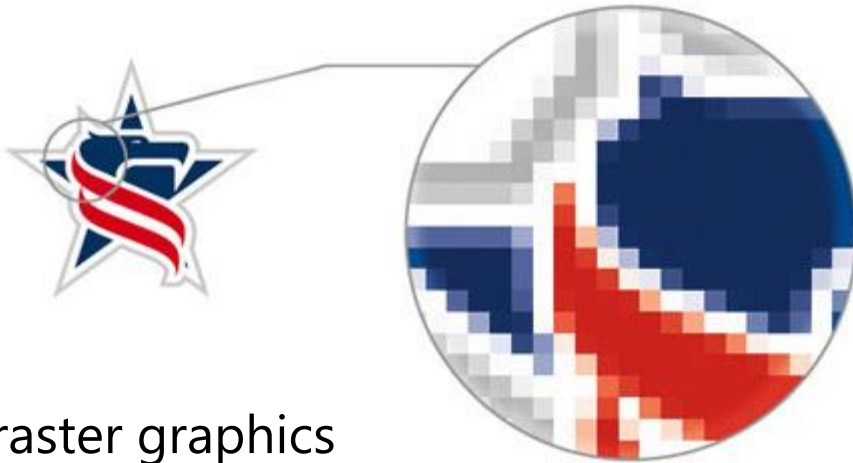


2D Gaussian

THE SOLUTION

What's the underlying problem?

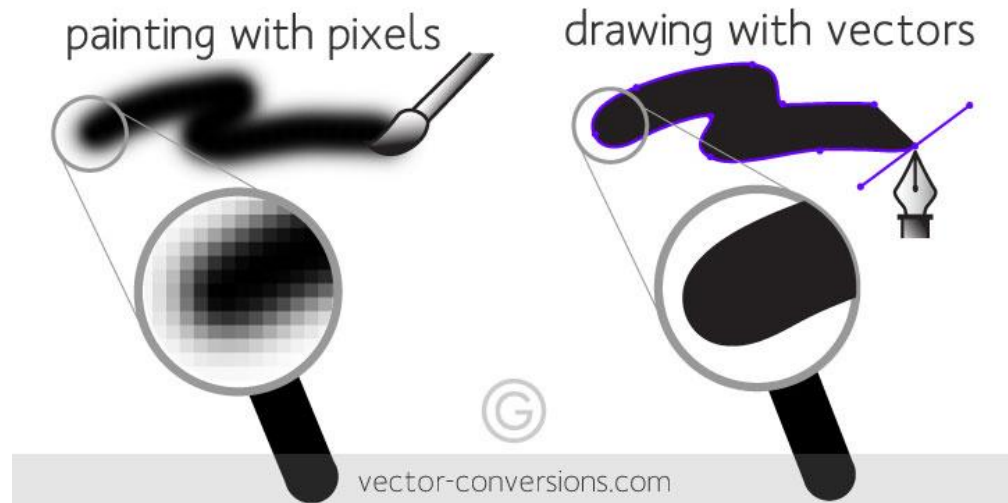
- detail can't be refined upon zoom
- can just be replicated or blurred



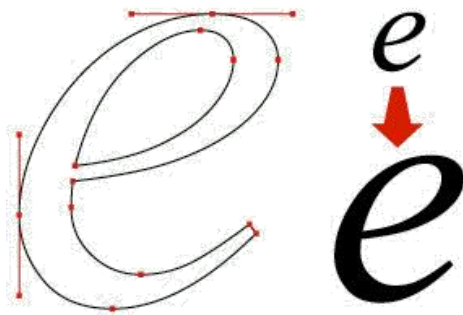
The solution...

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

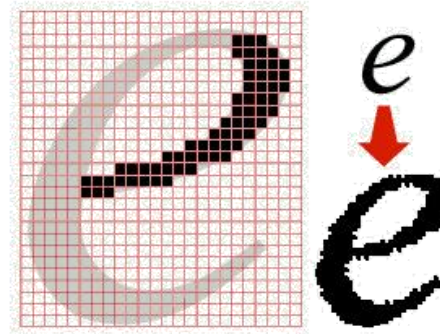
SCALABLE VECTOR GRAPHICS (SVG)



VECTOR GRAPHICS



BITMAPMED (RASTER) GRAPHICS



PHOTOGRAPHS AND IMAGES IN SVG

Vector graphics tends to have an “cartoonish” look



raster graphics

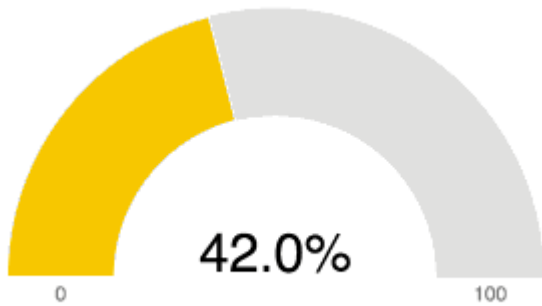
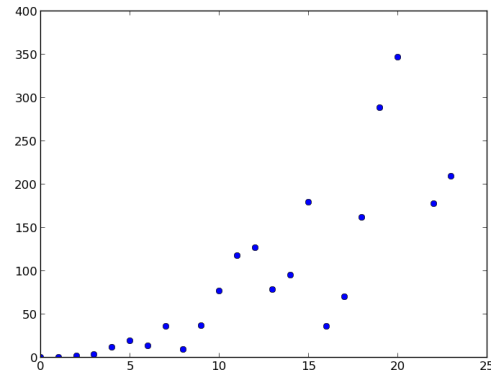
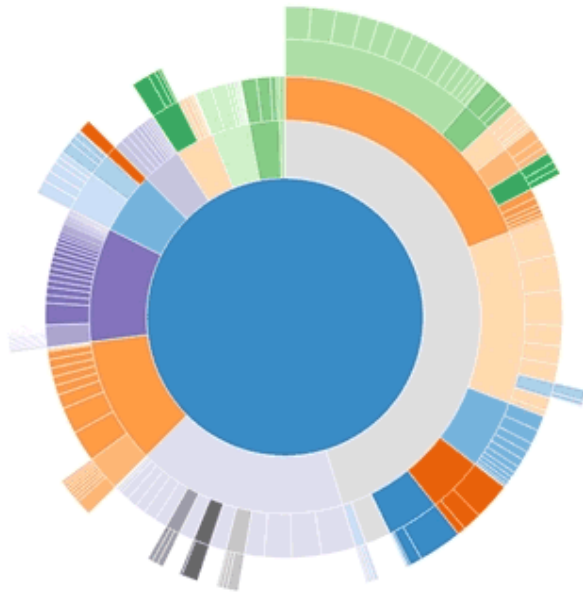


vector graphics

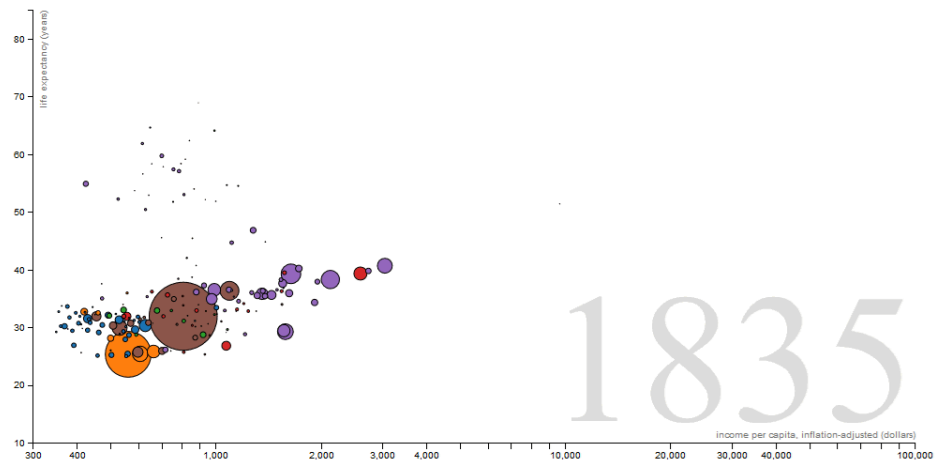
PHOTOGRAPHS AND IMAGES IN SVG



D3 USES SVG

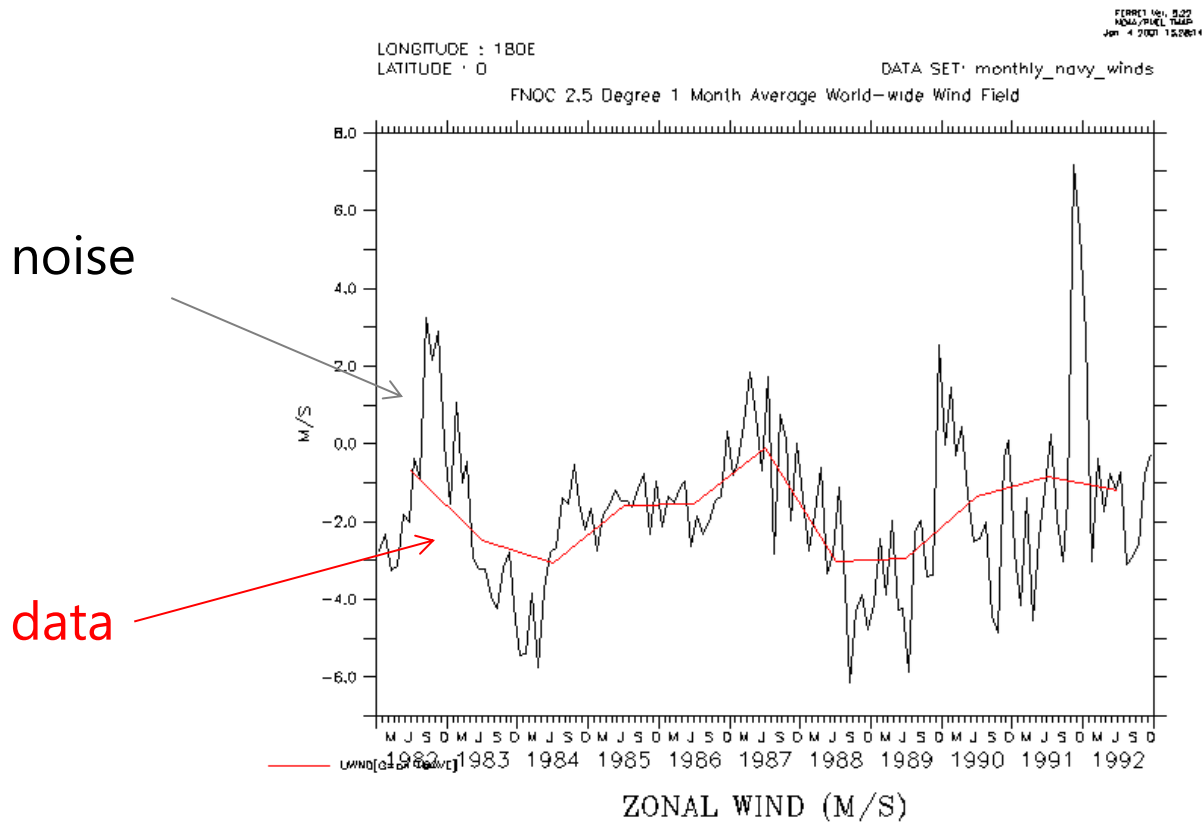


The Wealth & Health of Nations



SMOOTHING FOR DE-NOISING

Filtering/smoothing also eliminates noise in the data



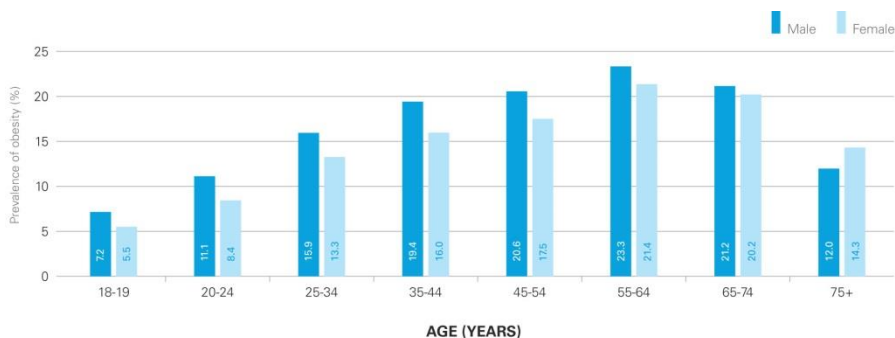
BACK TO BAR CHARTS

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

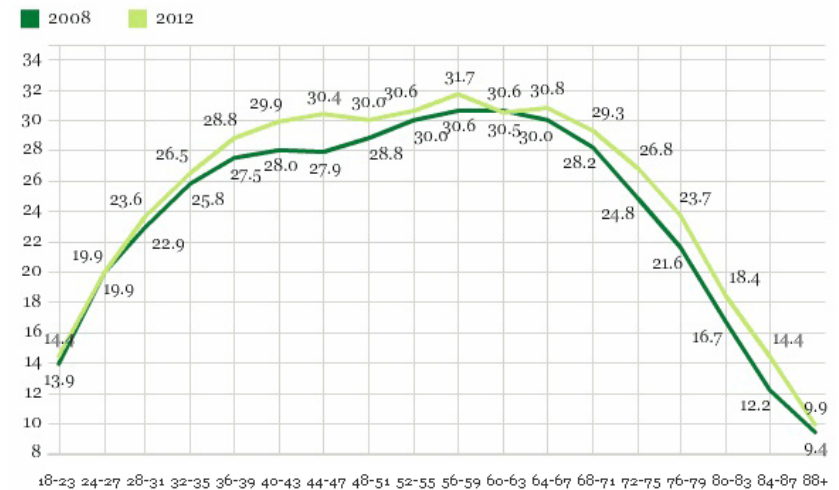
- the bins do the smoothing

Example:

- obesity over age (group)



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



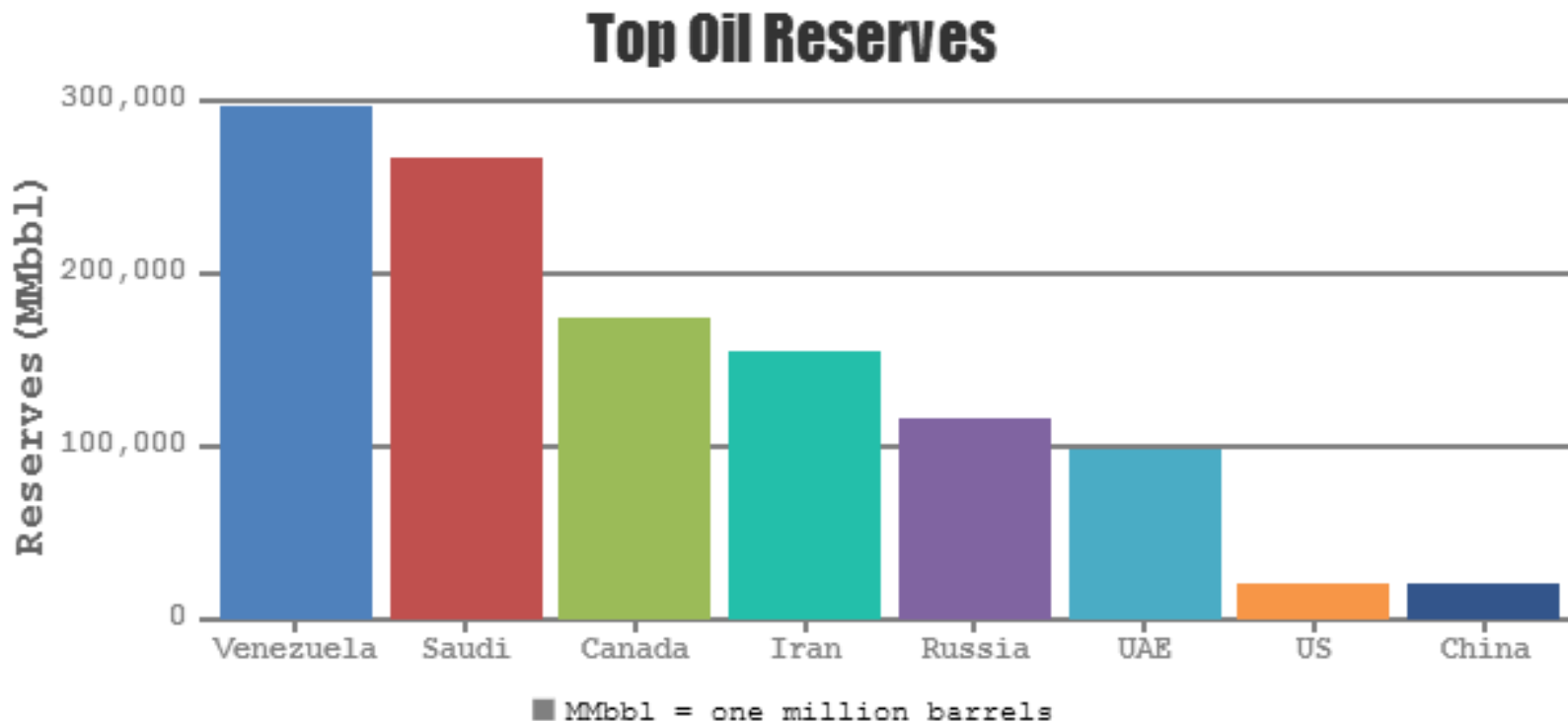
Gallup-Healthways Well-Being Index

GALLUP

BAR CHARTS

Of course, bar charts can also hold categorical data

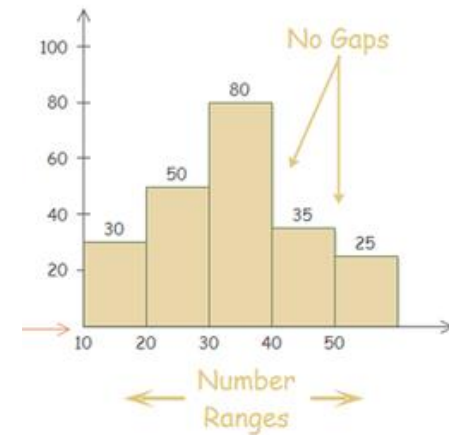
- smoothing by semantic grouping
- for example, Europe vs. {France, Spain, Italy, Germany, ...}



BAR CHARTS VS. HISTOGRAMS

Histograms

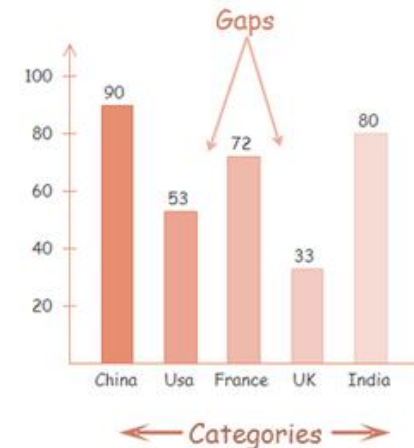
- bars show the frequency of numerical data
- quantitative data
- elements are grouped together, so that they are considered as ranges
- bars cannot be reordered
- width of bars need not be the same



Histogram

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

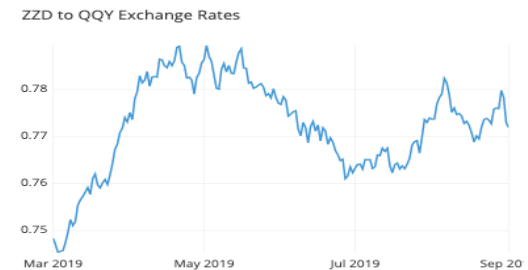
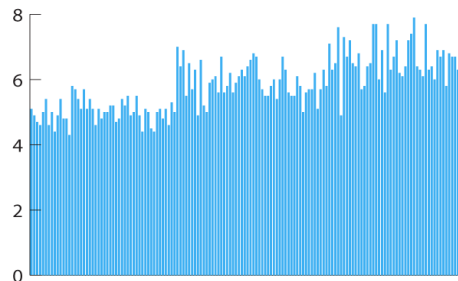
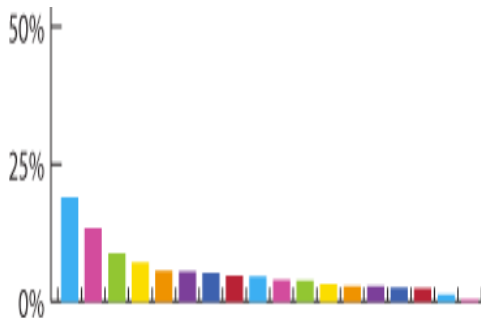


Bar Chart

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](#)

BAR CHARTS IN D3

<https://observablehq.com/@d3/bar-chart>

Working with bar charts and histograms is the topic of Lab 1

- the next two slides offer some help with calculations

HISTOGRAM CALCULATIONS – BINNING

Determine bin size

- $\min(\text{data})$ is optional, can also use 0 or some reasonable value
- $\max(\text{data})$ is optional, can also use some reasonable value

$$\text{bin size} = \frac{\max(\text{data}) - \min(\text{data})}{\text{number of bins}}$$

Given a data value val increment (++) the bin value

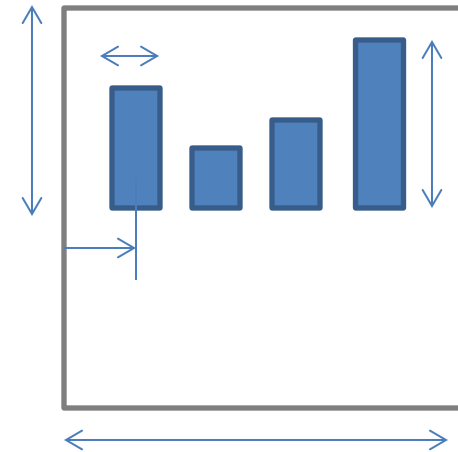
- but first initialize bin val array to 0

$$\text{bin val array} \left[\left[\frac{val - \min(\text{data})}{\text{bin size}} \right] \right] ++$$

HISTOGRAM CALCULATIONS – PLOTTING

Determine bin size on the screen

$$\text{bin size on screen} = \frac{\text{chart width}}{\text{number of bins}}$$



Center of a bar for bin with index *bin index*

$$\text{bar center on screen} = (\text{bin index} \cdot \text{bin size on screen}) + 0.5$$

Height of the bar for a bin with index *bin index*

$$\text{bar height}(\text{bin index}) = \text{bin val array}(\text{bin index}) \cdot \frac{\text{chart height}}{\max(\text{bin val array})}$$

Do not forget that the origin of a web page is the top left corner