

CSE 564

VISUALIZATION & VISUAL ANALYTICS

APPLICATIONS AND BASIC TASKS

**KLAUS MUELLER**

COMPUTER SCIENCE DEPARTMENT  
STONY BROOK UNIVERSITY

# COURSE WEBSITE

<http://www.cs.stonybrook.edu/~mueller/teaching/cse564/>

Everything you need is there:

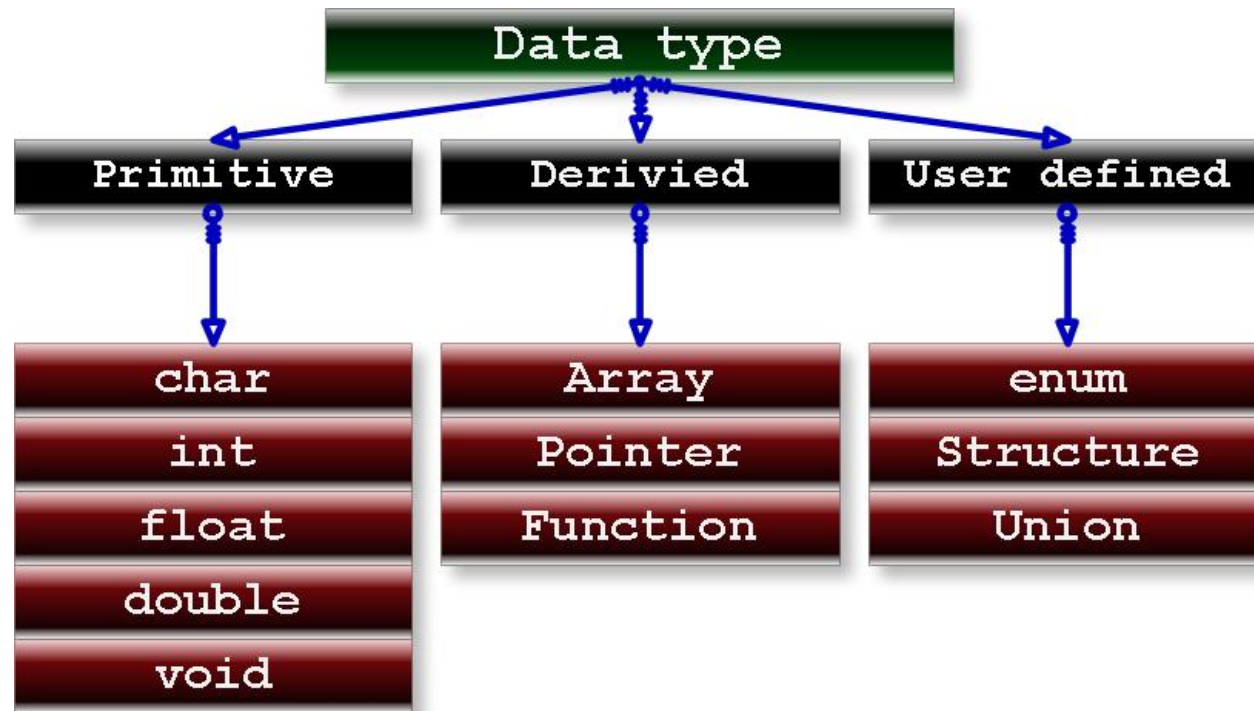
- syllabus
- course notes (slides) posted shortly after the lecture
- lab assignments
- course policy

There will also be (soon to be announced)

- a server for lab assignments
- piazza for online support

Lecture	Topic	Projects
1	Intro, schedule, and logistics	
2	Applications of visual analytics, basic tasks, data types	
3	Introduction to D3, basic vis techniques for non-spatial data	Project #1 out
4	Visual perception and cognition	
5	Visual design and aesthetics	
6	Data types, notion of similarity and distance	
7	Data preparation and reduction	Project #1 due
8	Introduction to R, statistics foundations	Project #2 out
9	Data mining techniques: clusters, text, patterns, classifiers	
10	Data mining techniques: clusters, text, patterns, classifiers	
11	Computer graphics and volume rendering	
12	Techniques to visualize spatial (3D) data	Project #2 due
13	Scientific and medical visualization	Project #3 out
14	Scientific and medical visualization	
15	Midterm #1	
16	High-dimensional data, dimensionality reduction	Project #3 due
17	Big data: data reduction, summarization	
18	Correlation and causal modeling	
19	Principles of interaction	
20	Visual analytics and the visual sense making process	Final project proposal due
21	Evaluation and user studies	
22	Visualization of time-varying and time-series data	
23	Visualization of streaming data	
24	Visualization of graph data	Final Project preliminary report due
25	Visualization of text data	
26	Midterm #2	
27	Data journalism	
	Final project presentations	Final Project slides and final report due

# DATA TYPES EVERY CS PERSON KNOWS



# DATA TYPES IN VISUAL ANALYTICS

Numeric

Categorical

Text

Time series

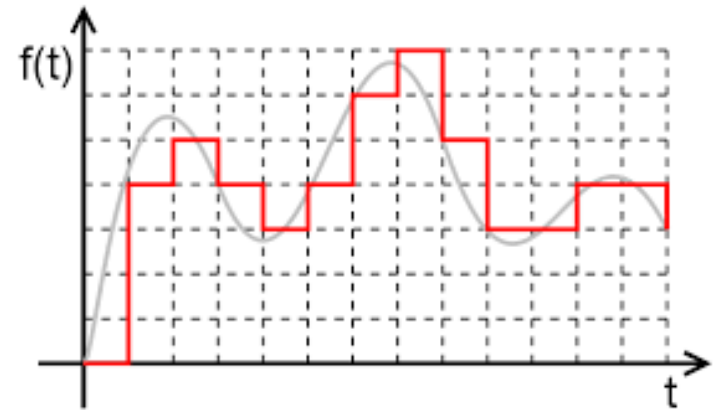
Graphs and networks

Hierarchies

# VARIABLES IN STATISTICS

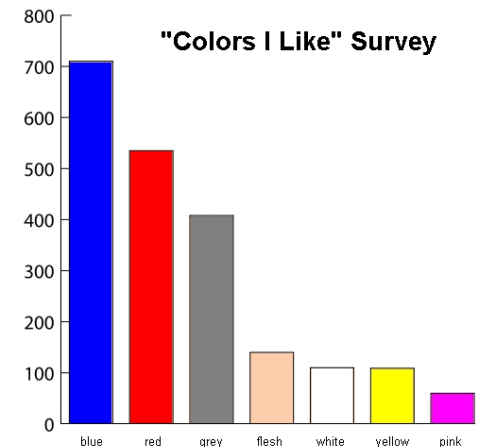
## Numeric variables

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



## Categorical variables

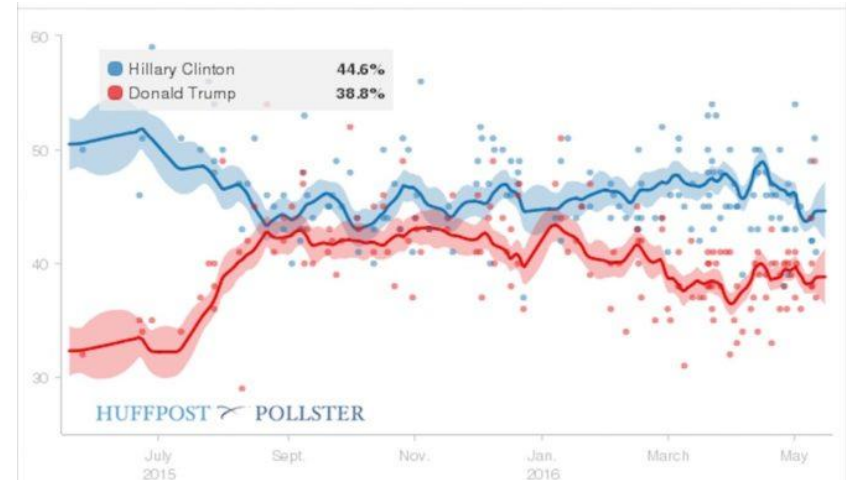
- describe a **quality** or characteristic
- like: 'what type' or 'which category'



# NUMERIC VARIABLES

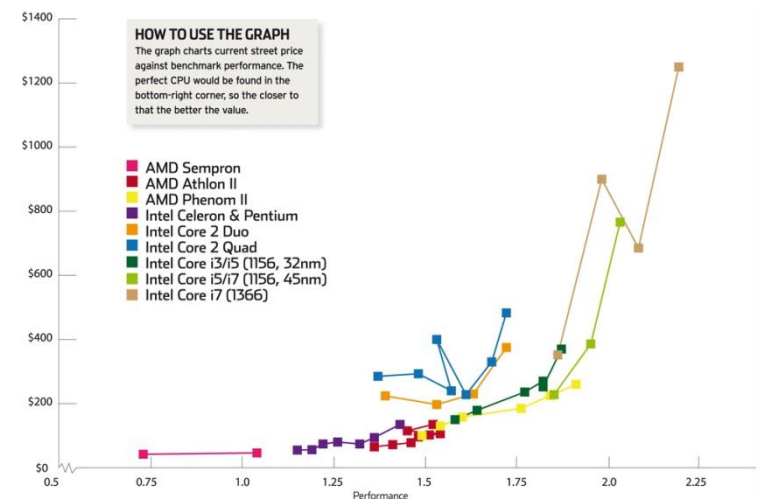
Most often the x-axis is 'time'

- 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

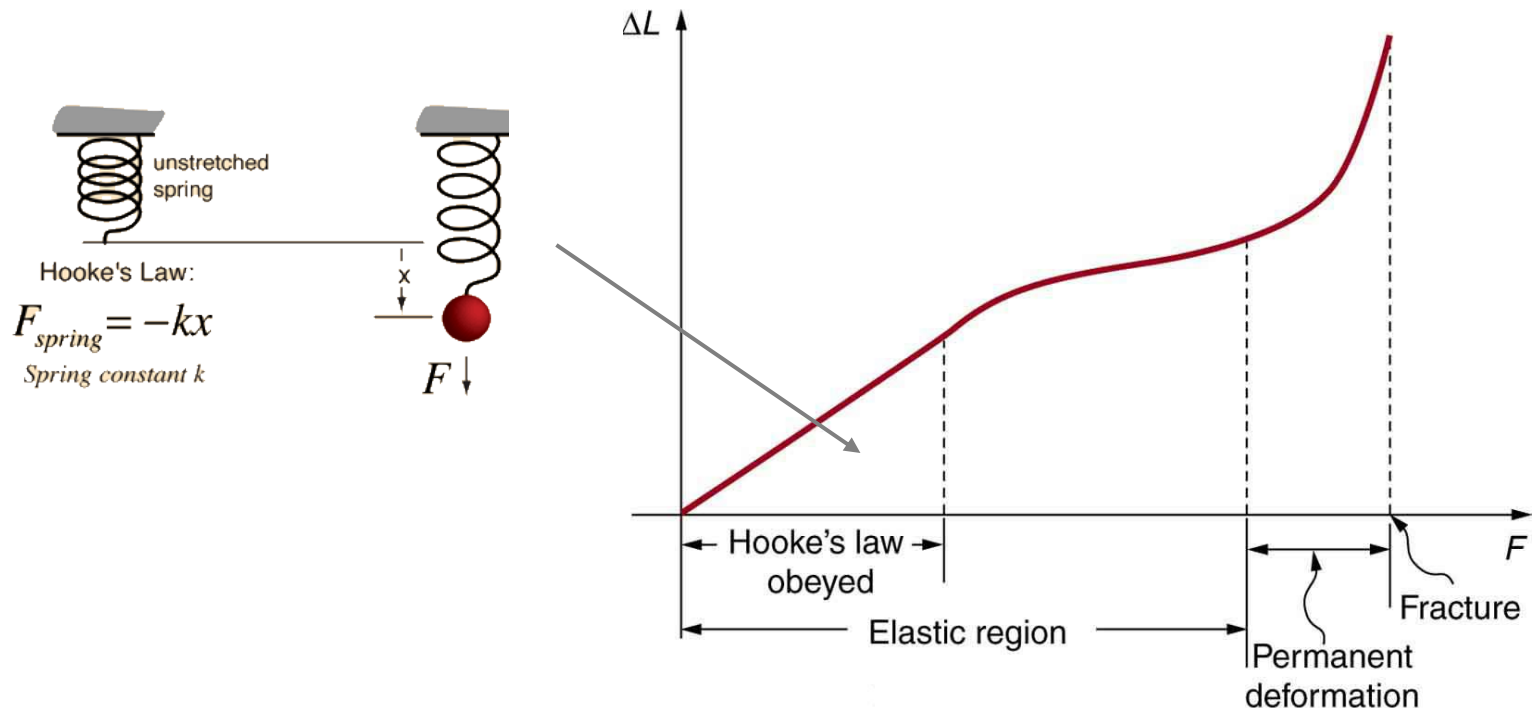
- engineers, statisticians, etc. will be receptive to this idea
- can you think of an example?



# NUMERIC VARIABLES

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

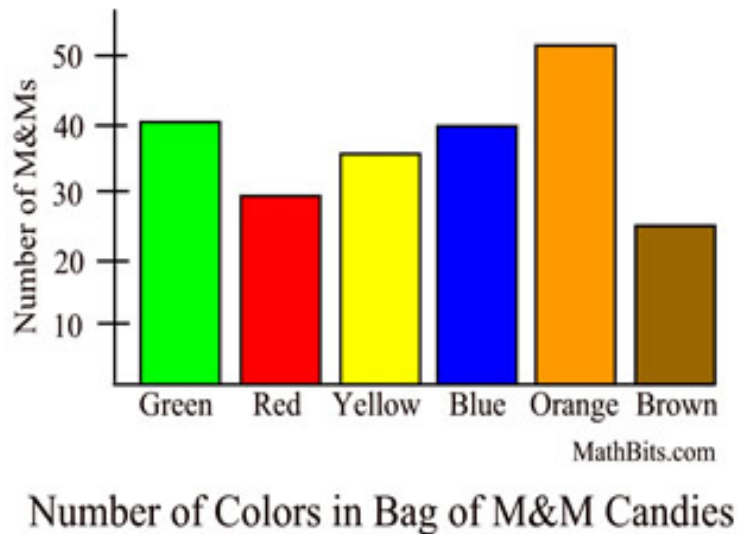
- from the engineering / physics domain
- in some sense, it tells a story





# CATEGORICAL VARIABLES

Usually plotted as bar charts or pie charts



??

nominal  
ordinal



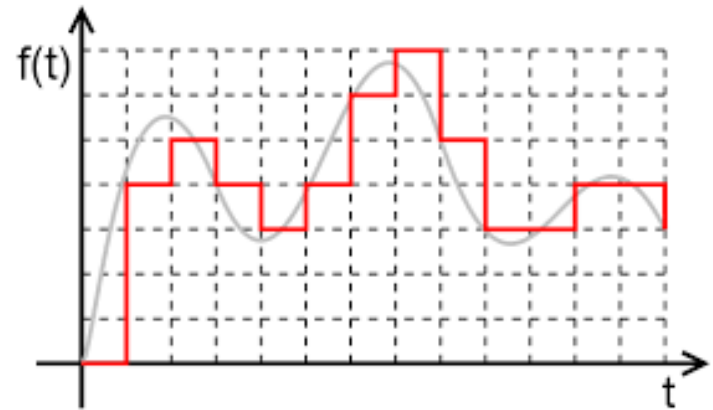
??

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

# 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

# NUMBERS ARE GOOD

But not everything is expressed in numbers

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



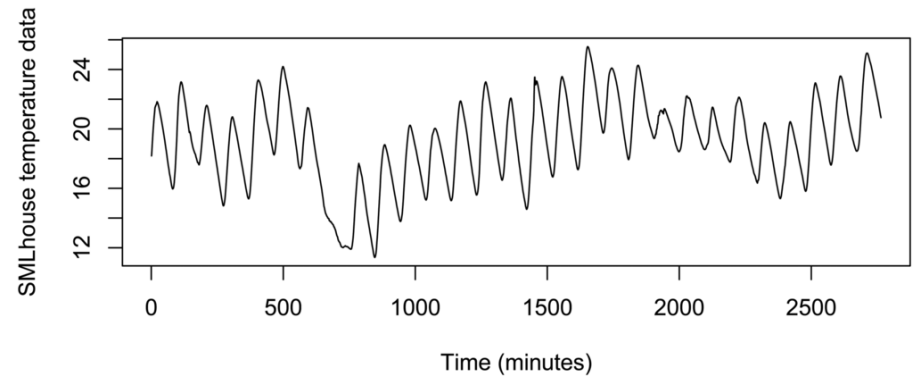
Do **feature analysis** to turn these abstract things into numbers

- 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

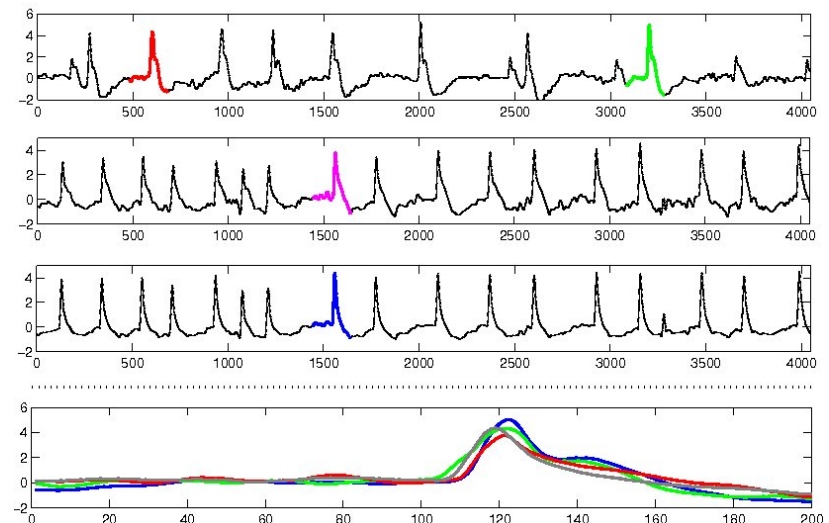
## Characteristics

- often large scale
- time series

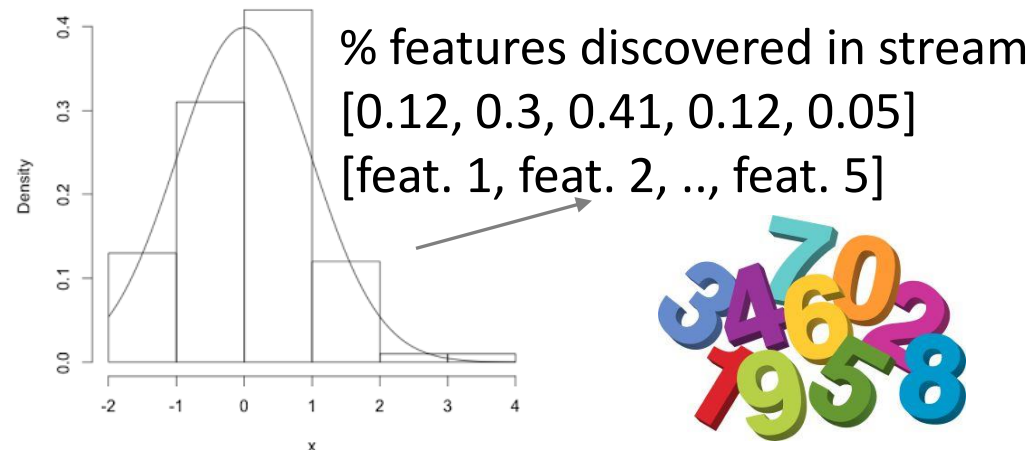


## Feature Analysis

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



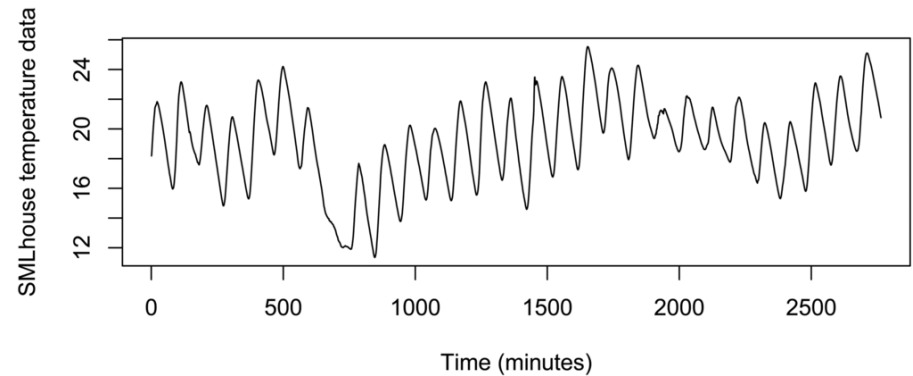
Motif discovery



# SENSOR DATA

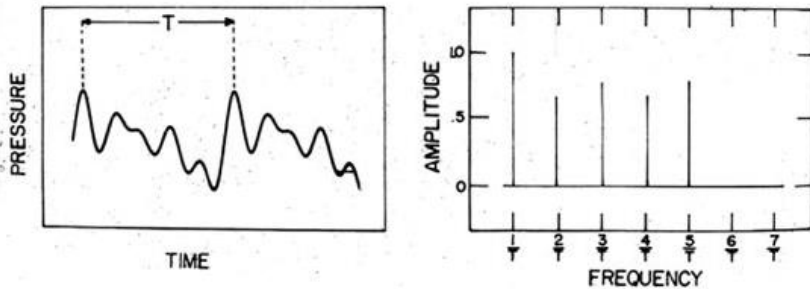
## Characteristics

- often large scale
- time series

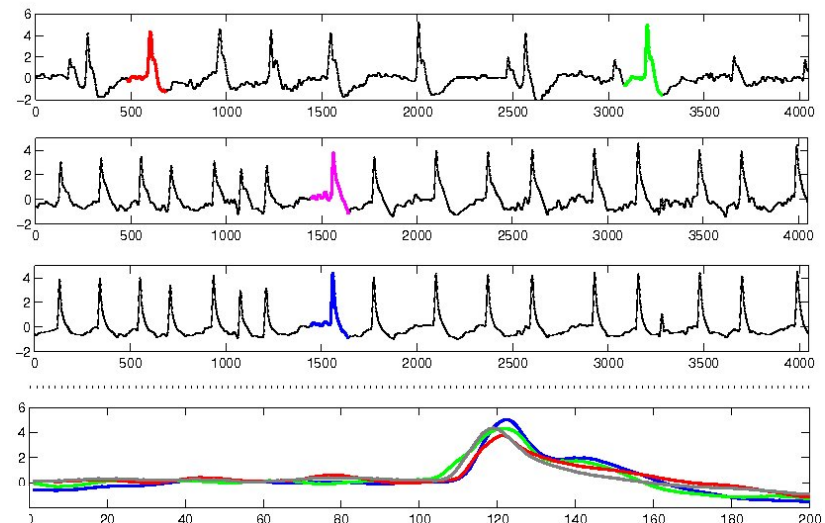


## Feature Analysis

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



Fourier transform



# IMAGE DATA

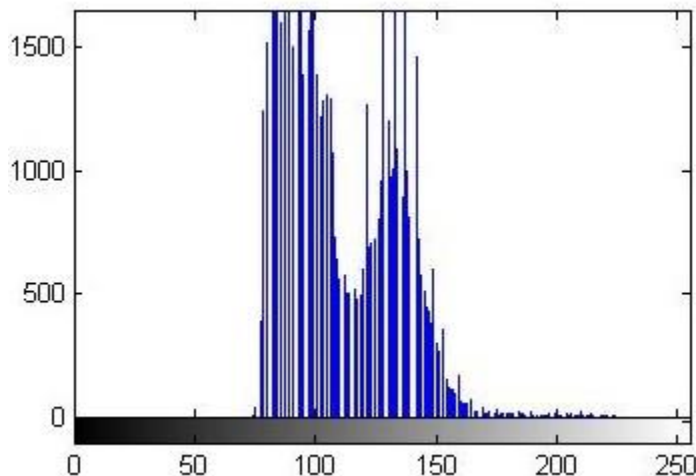
## Characteristics

- array of pixels

## Feature Analysis

- value histograms
- encode into a 256-D vector

histograms



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



# IMAGE DATA

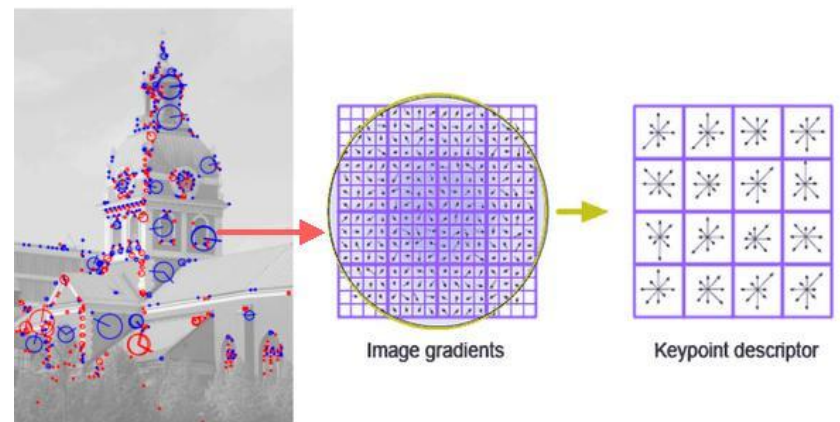
## Characteristics

- array of pixels

## Feature Analysis

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

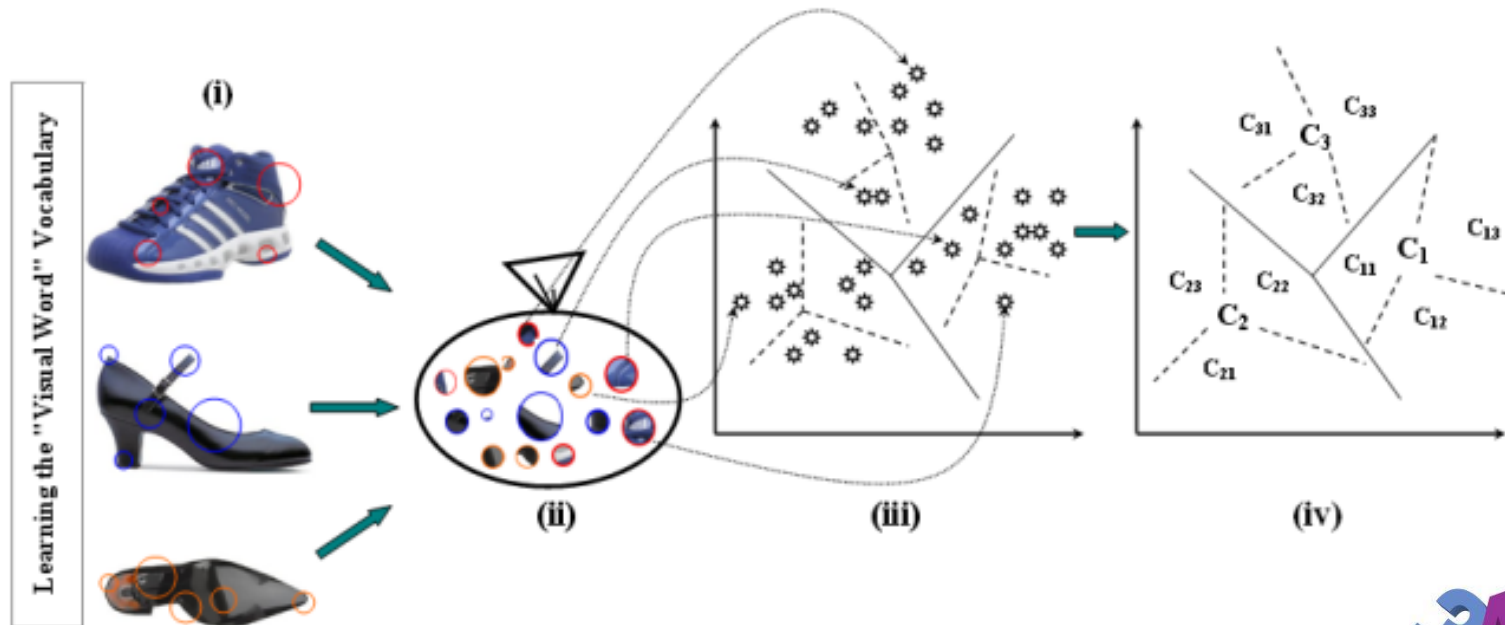
histograms



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]<sub>Person</sub> bought [300 shares]<sub>Quantity</sub> of [Acme Corp.]<sub>Organiz.</sub> in [2006]<sub>Time</sub>
- distinguish between
  - application of grammar rules (old style, need experienced linguists)
  - statistical models (Google etc., need big data to build)

# TEXT TO NUMERIC DATA

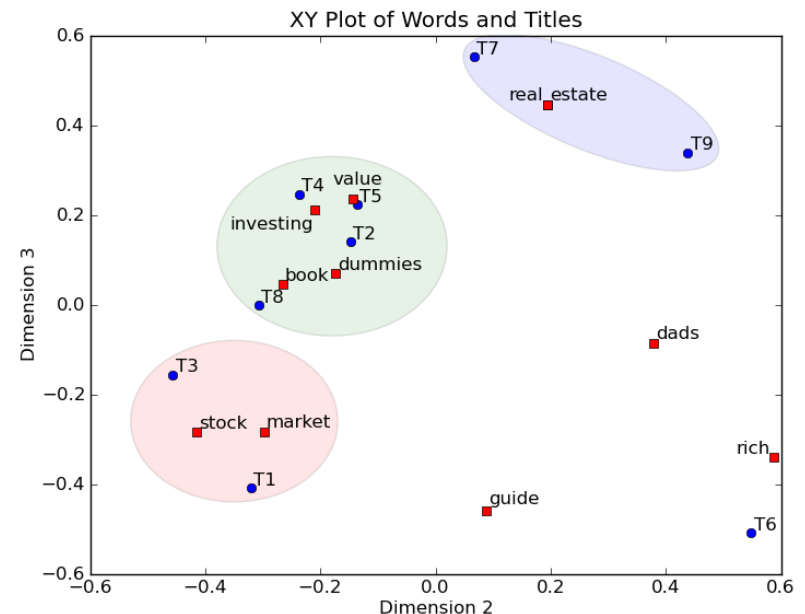
Create a term-document matrix

- 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

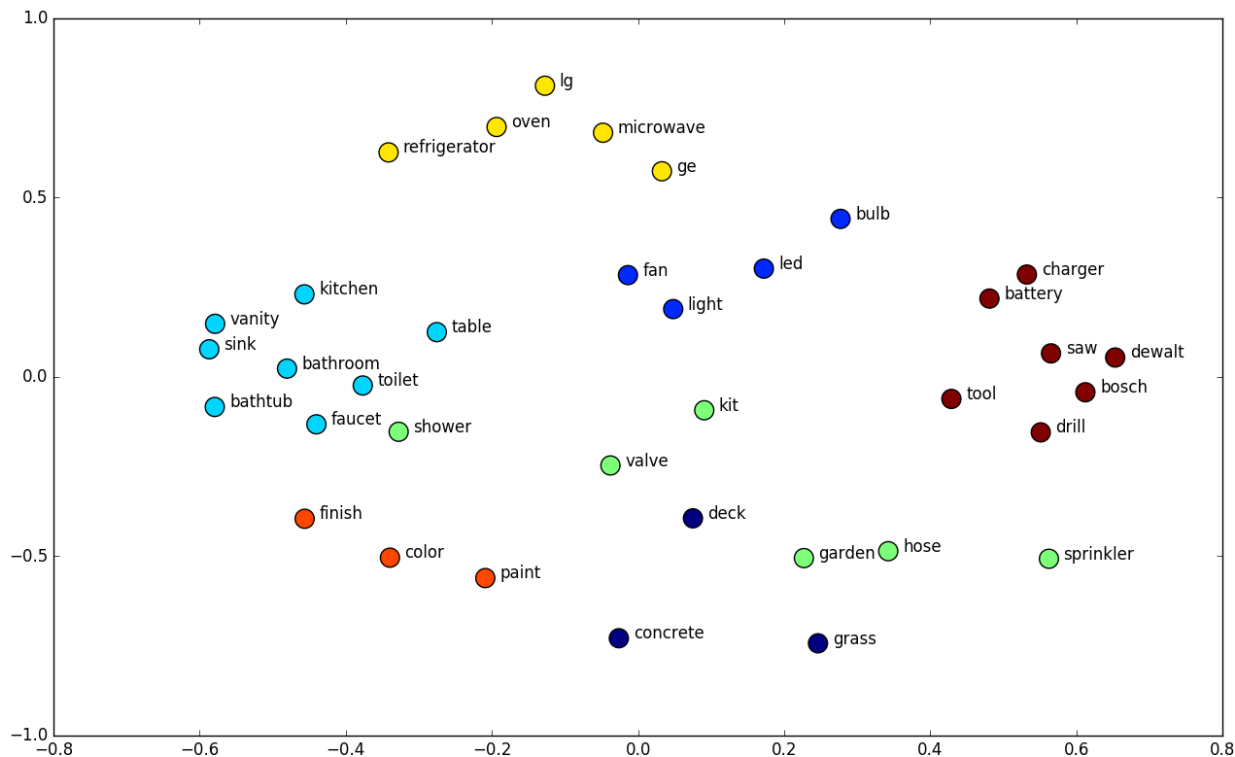


Word/document cluster

# WORD EMBEDDING

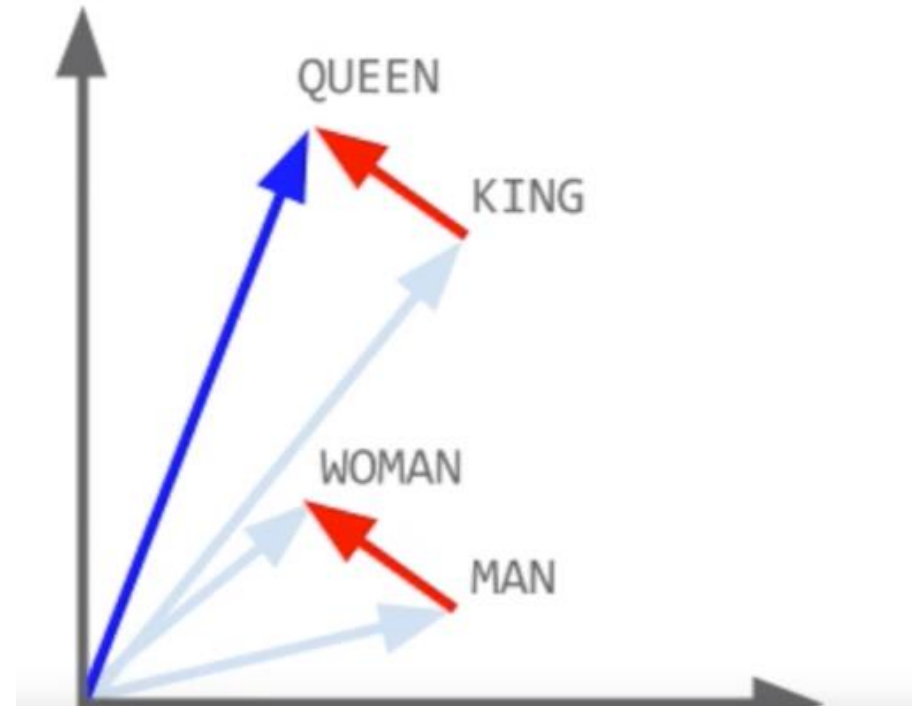
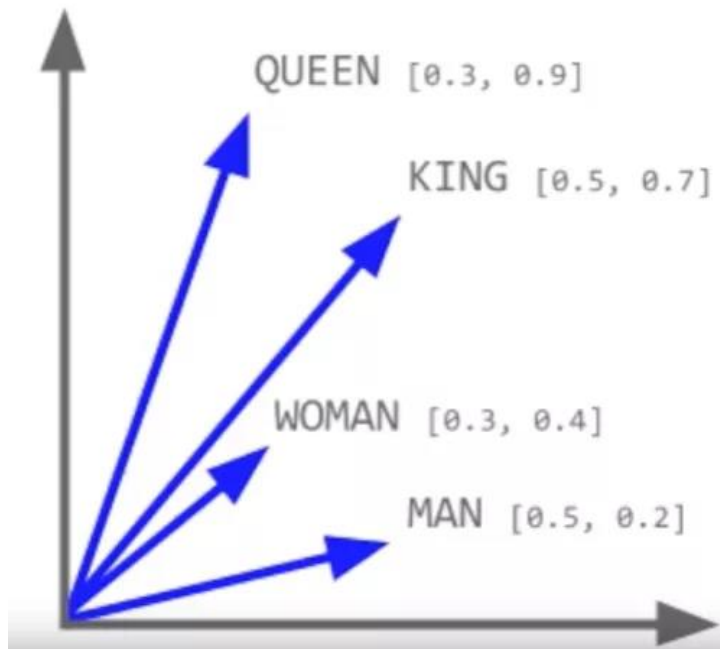
Train a shallow neural network (NN) on a corpus of text

- 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



$\text{gender} = \text{WOMAN} - \text{MAN}$

$\text{QUEEN} = \text{KING} + \text{gender}$

$\text{QUEEN} = \text{KING} - \text{MAN} + \text{WOMAN}$

# 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

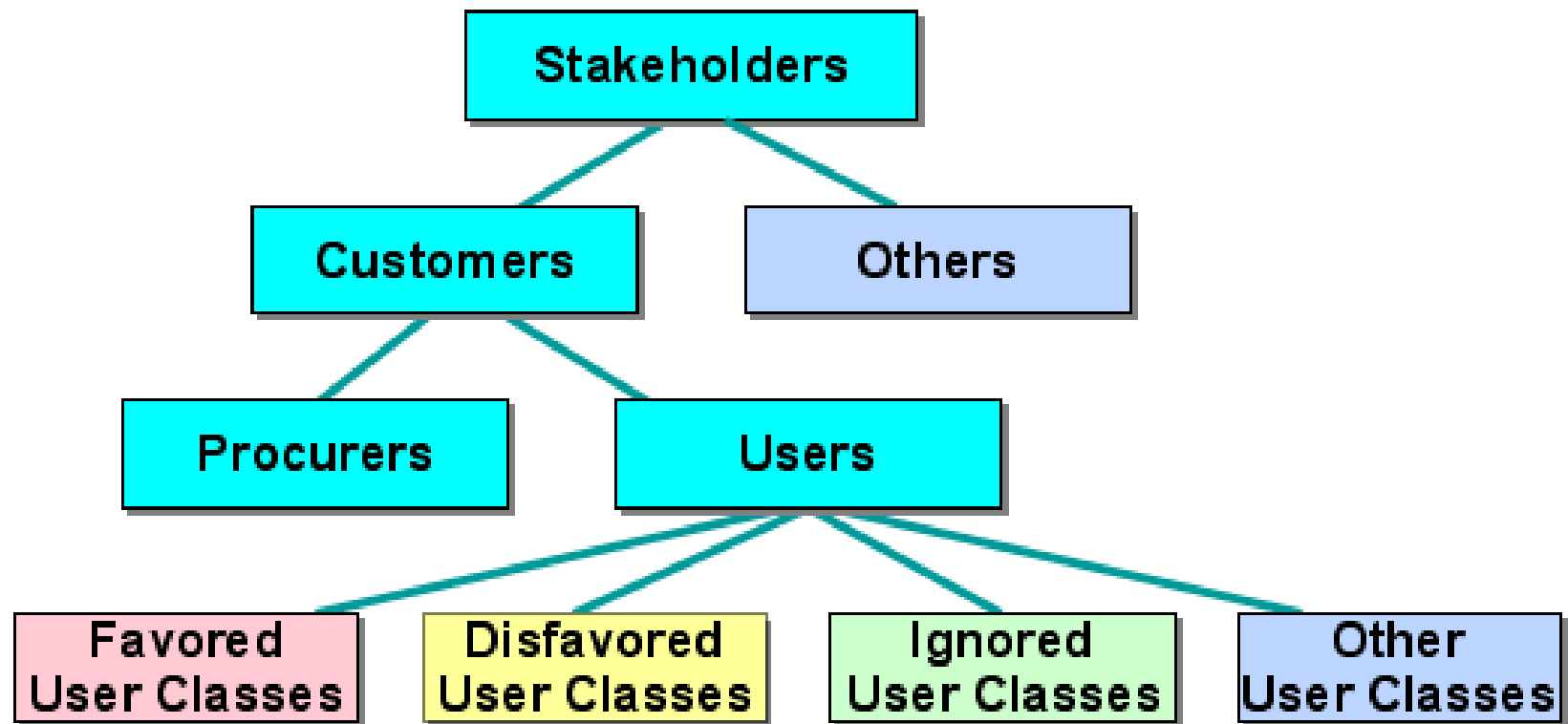




LET'S LOOK AT SOME ESSENTIAL  
GRAPHICAL REPRESENTATIONS

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

<http://bl.ocks.org/kerryrodden/7090426>

# BUBBLE CHARTS

No hierarchy information, just quantities

<http://bl.ocks.org/mbostock/4063269>

# 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

<http://bl.ocks.org/mbostock/4062006>

# 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

<http://bl.ocks.org/mbostock/4060366>

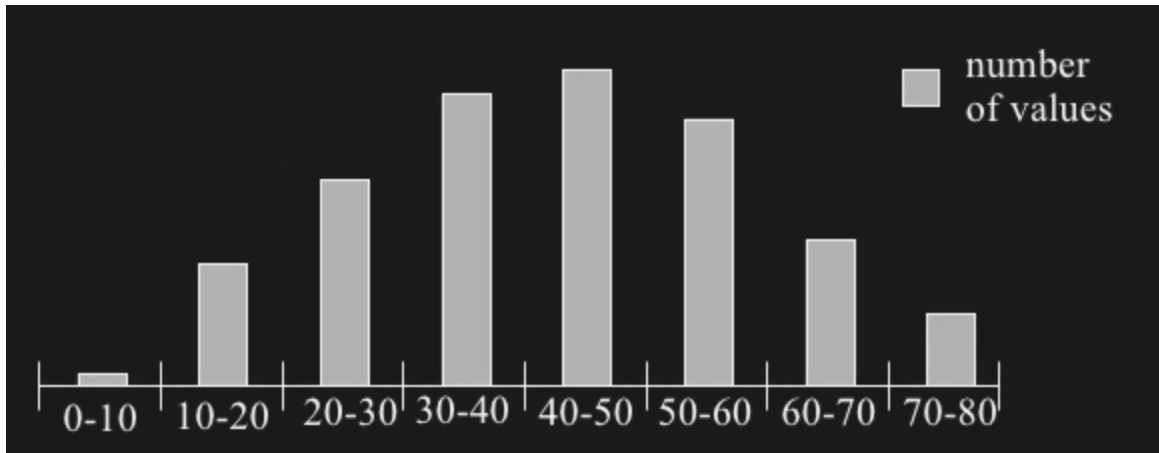


# DATA TYPE CONVERSIONS AND TRANSFORMATION

# NUMERIC TO CATEGORICAL DATA: DISCRETIZATION (1)

## Solution 1:

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

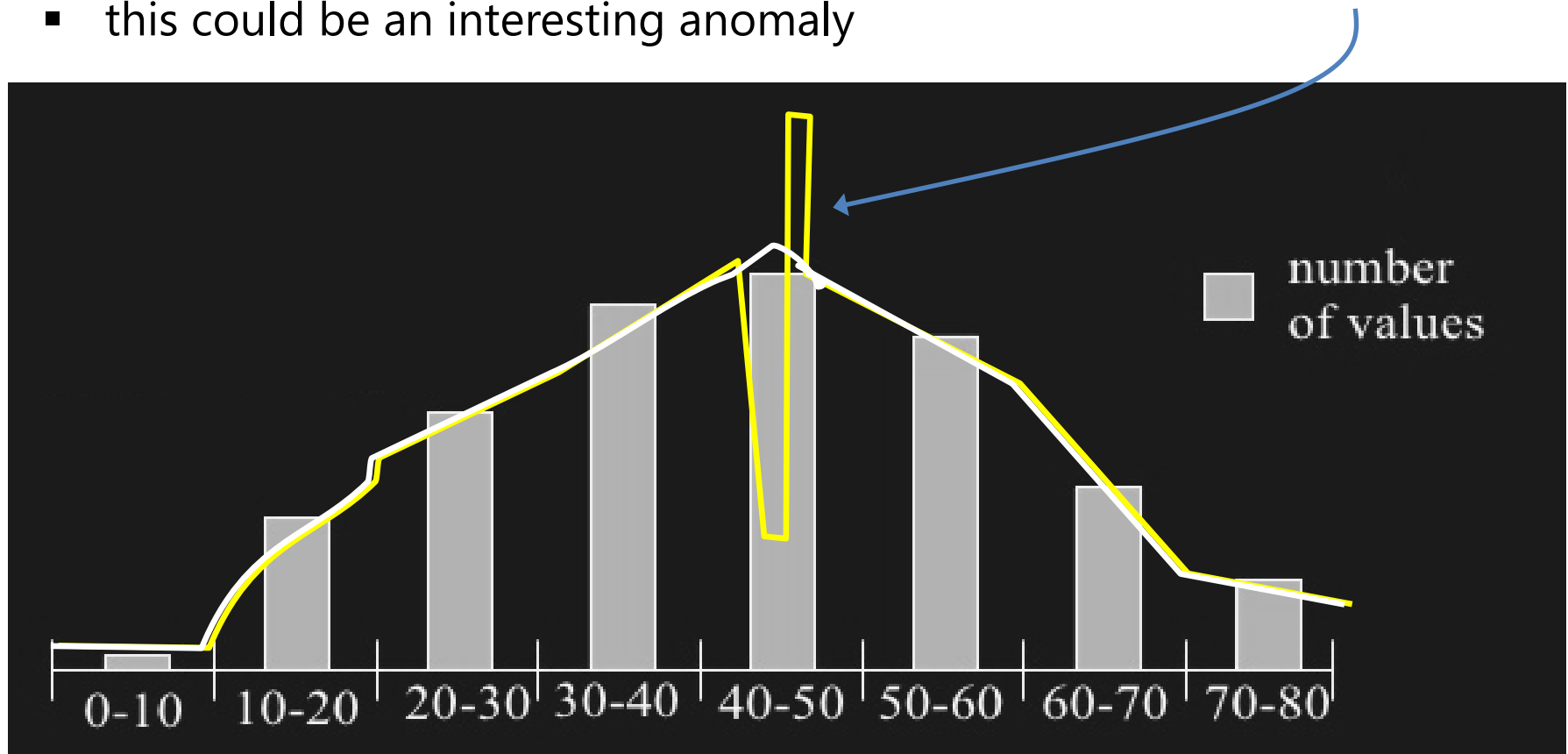


- what is lost here?

# PROBLEM WITH EQUI-WIDTH HISTOGRAM

Age ranges of customers could be unevenly distributed within a bin

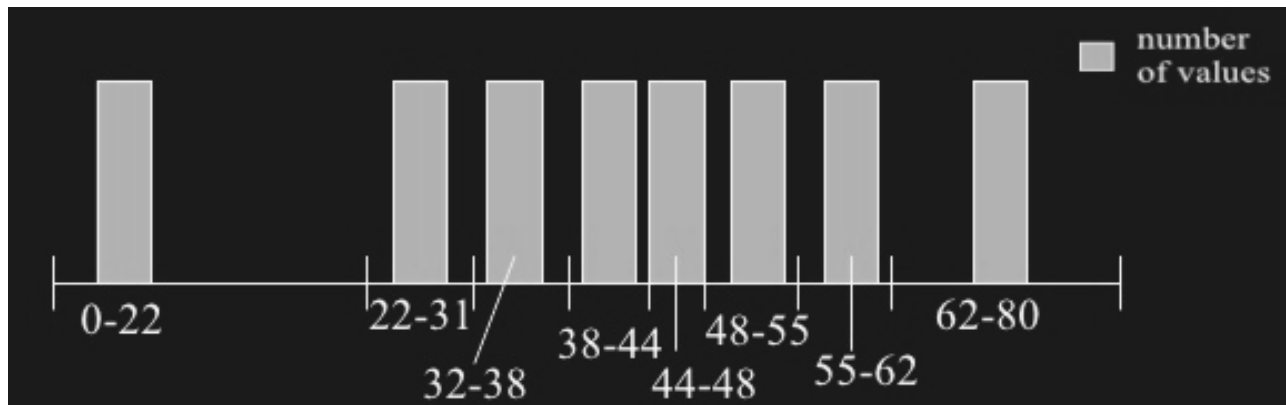
- this could be an interesting anomaly



# NUMERIC TO CATEGORICAL DATA: DISCRETIZATION (2)

## Solution 2:

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

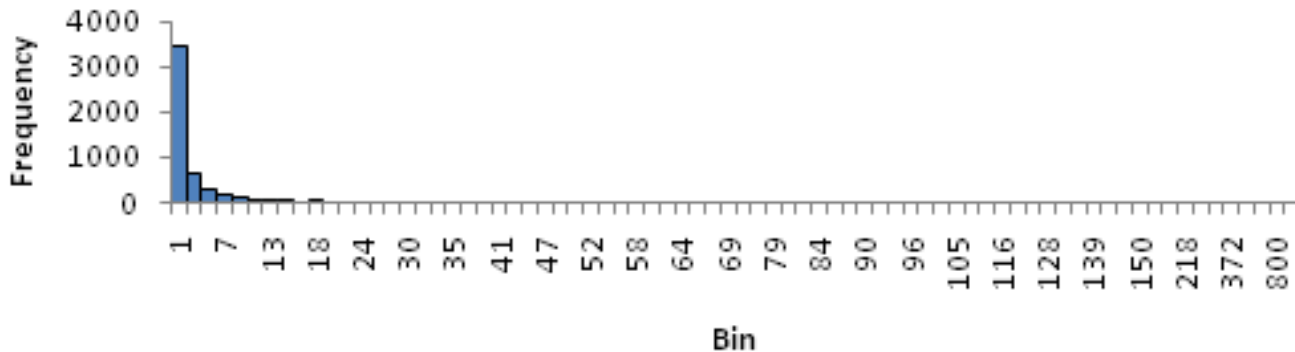


- what is the disadvantage here?
- extra storage needed: must store the start/end value for each bin

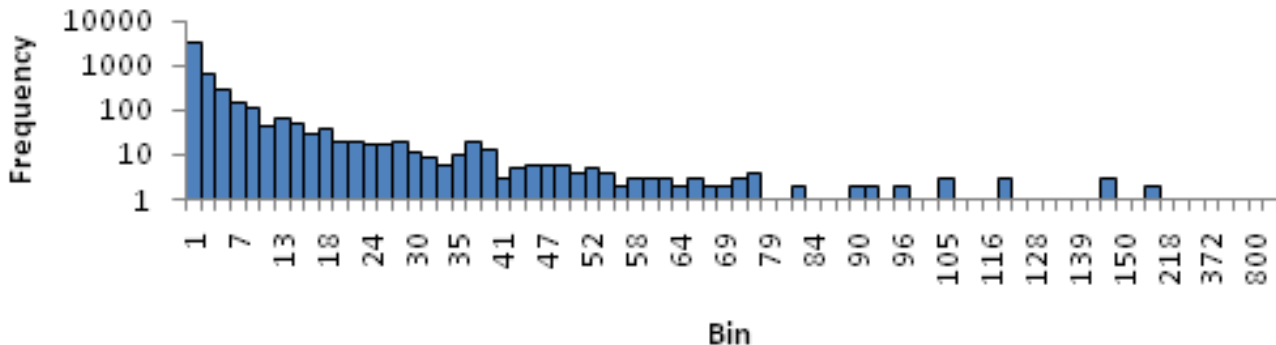
# NUMERIC TO CATEGORICAL DATA: DISCRETIZATION (3)

## Solution 3:

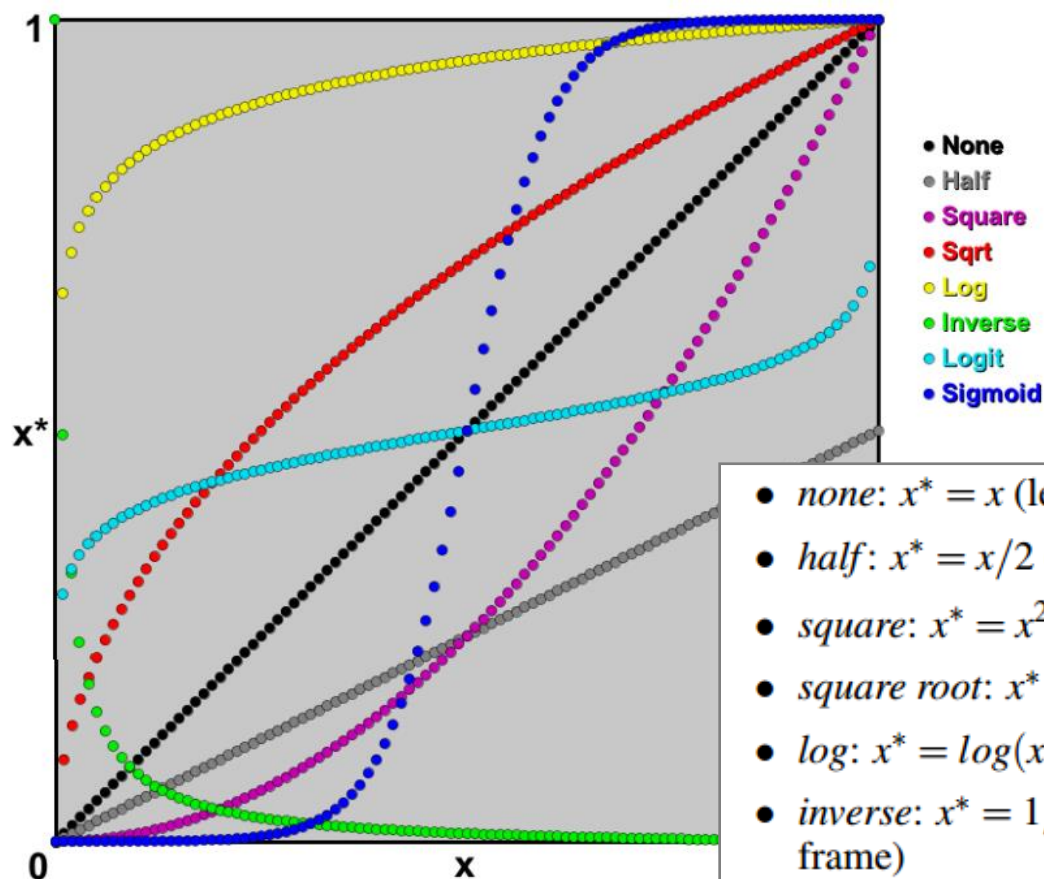
- what if all the bars have seemingly height
- or are dominated by one large peak



- switch to log scaling of the y-value



# OTHER TRANSFORMATIONS

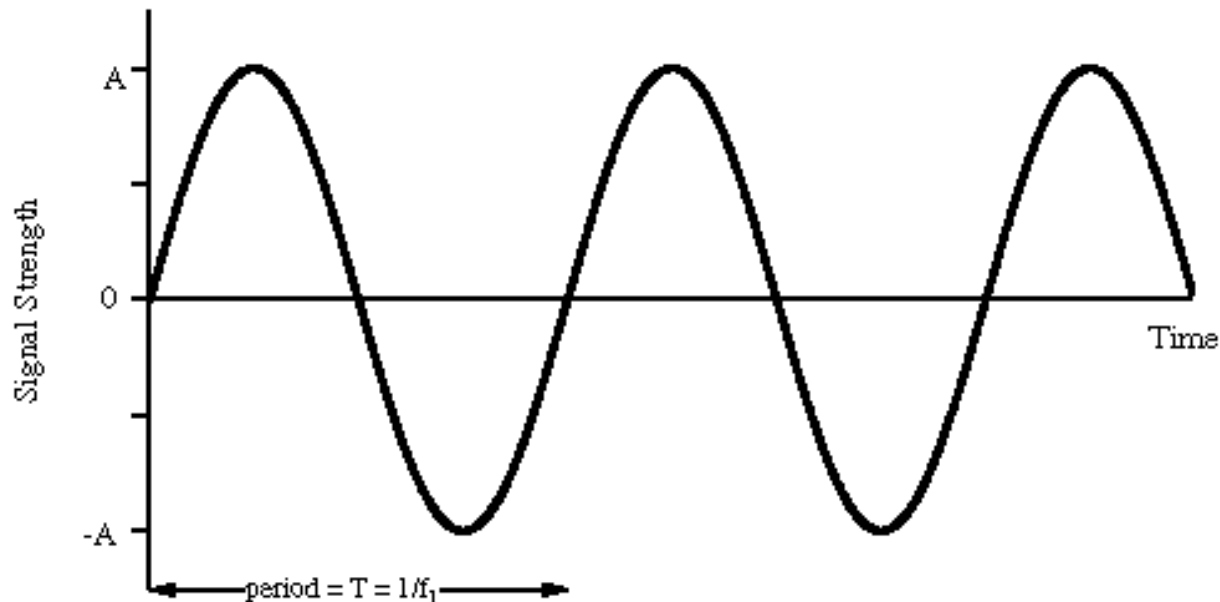


- *none*:  $x^* = x$  (leaves points unchanged)
- *half*:  $x^* = x/2$  (squeezes all points together)
- *square*:  $x^* = x^2$  (pulls points toward left of frame)
- *square root*:  $x^* = \sqrt{x}$  (mildly pulls points toward right of frame)
- *log*:  $x^* = \log(x)$  (strongly pulls points toward 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)

# CONTINUOUS TO DISCRETE

## Why discrete?

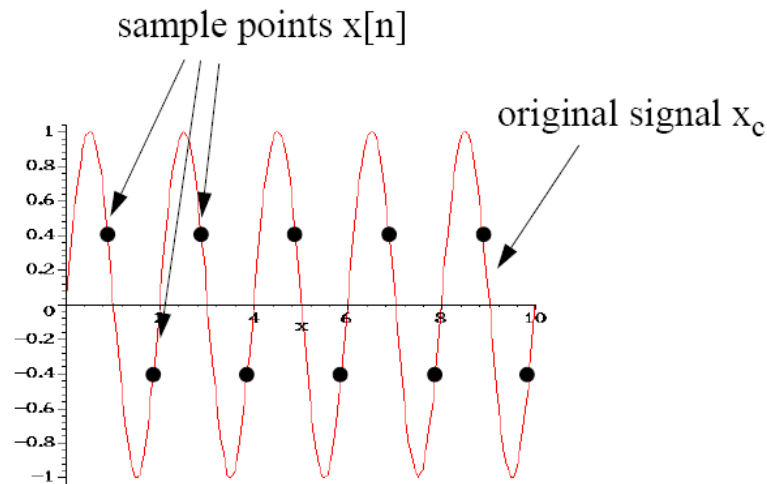
- because we can't store continuous data
- we can only store samples of the continuous data
- how many samples do we need?
- also keep this in mind for data reduction



# CONTINUOUS TO DISCRETE

## Why discrete?

- because we can't store continuous data
- we can only store samples of the continuous data
- how many samples do we need?

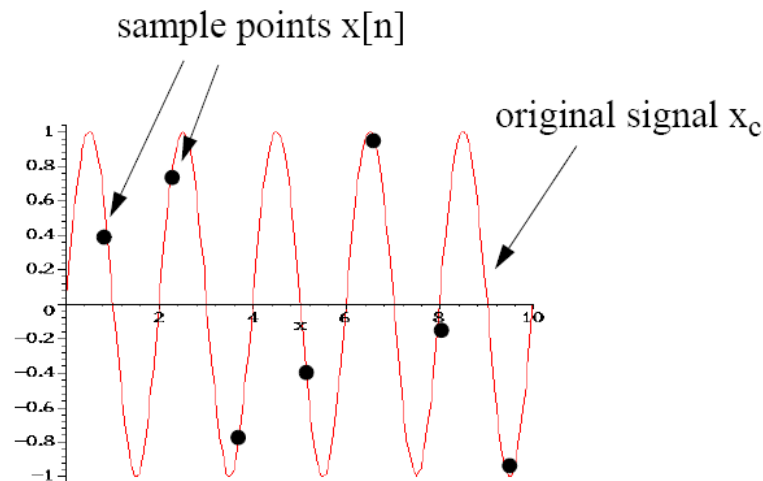




# CONTINUOUS TO DISCRETE

## Why discrete?

- because we can't store continuous data
- we can only store samples of the continuous data
- how many samples do we need?



# CONTINUOUS TO DISCRETE

Why discrete?

- because we can't store continuous data
- we can only store samples of the continuous data
- how many samples do we need?

We need a certain number of samples to represent a continuous phenomenon

- twice as many samples as the highest frequency in the signal
- called the *Nyquist frequency*
- else we get *aliasing*

# PRACTICAL IMPLICATIONS

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



This is aliasing

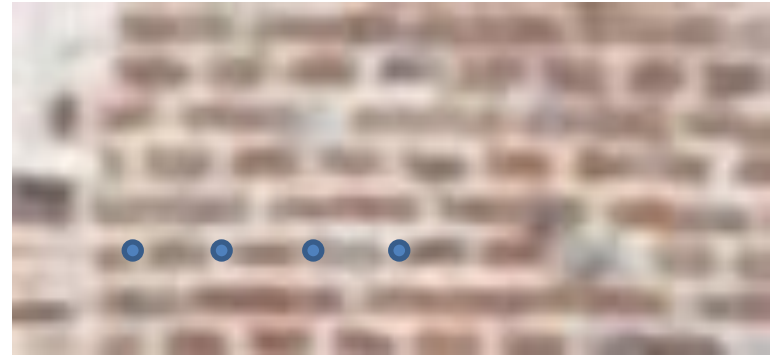
# PRACTICAL IMPLICATIONS

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

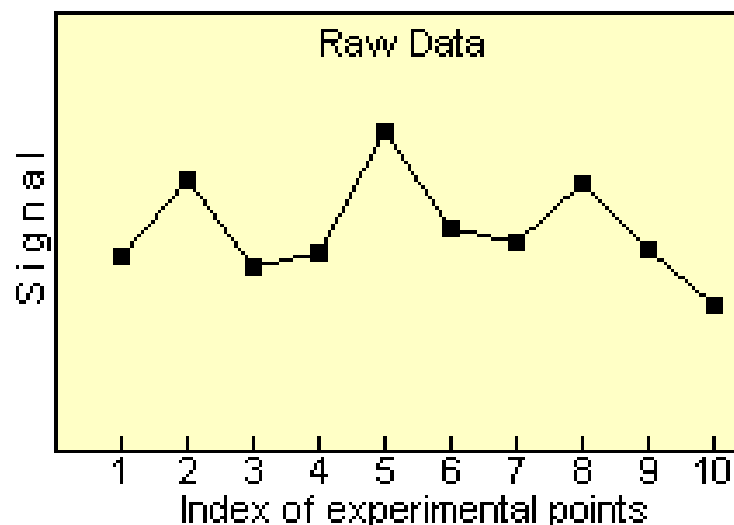




# WHAT IS SMOOTHING?

Slide a window across the signal

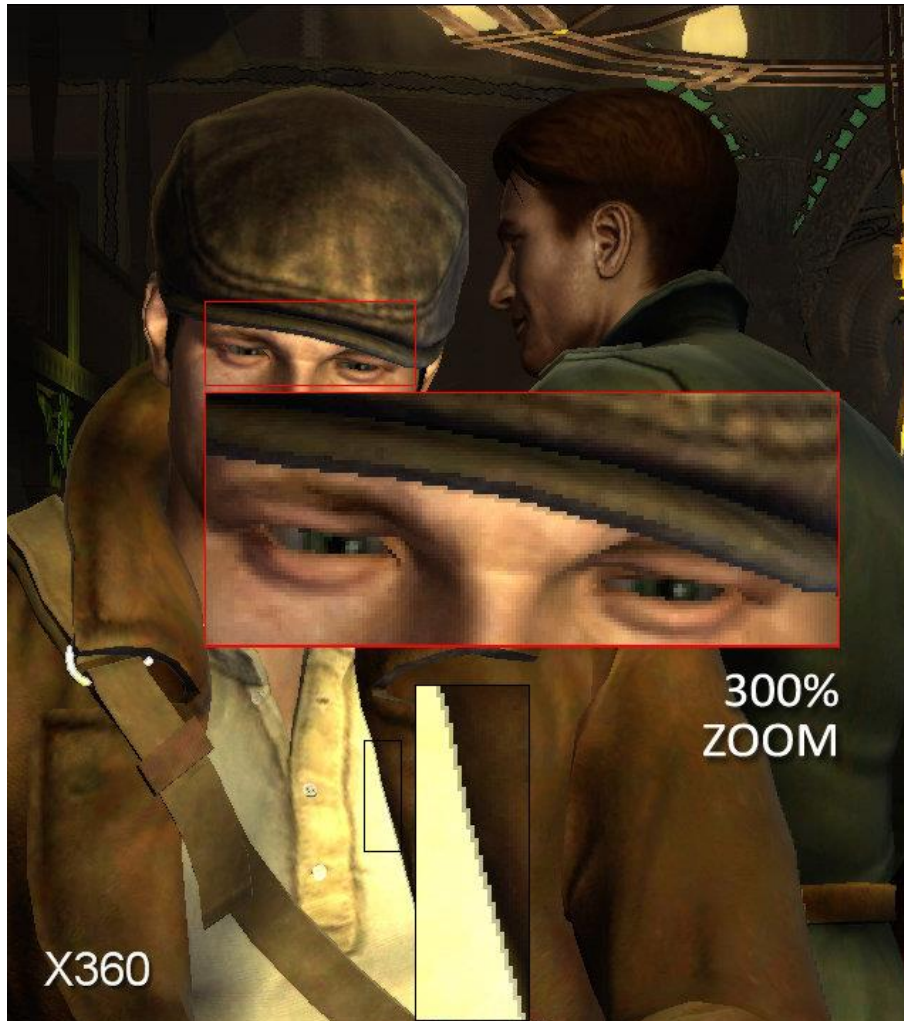
- stop at each discrete sample point
- average the original data points that fall into the window
- store this average value at the sample point
- move the window to the next sample point
- repeat



# ANTI-ALIASING VIA SMOOTHING: TRADEOFFS

looks sharper, but has “jaggies”

a bit blurred, but no more jaggies



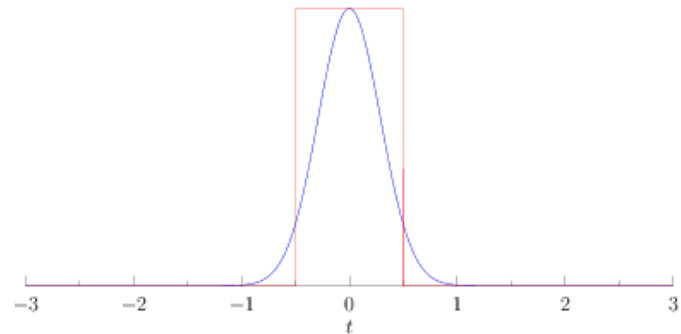
# FILTERS

What is the filter we just used called?

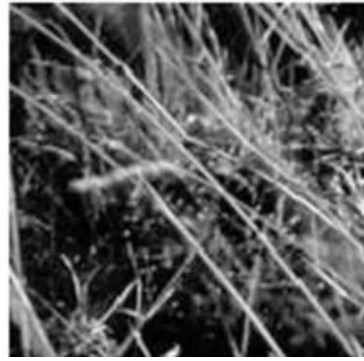
- it's called a *box filter*

There are other filters

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



# BOX FILTER VS. GAUSSIAN FILTER

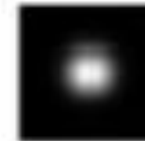
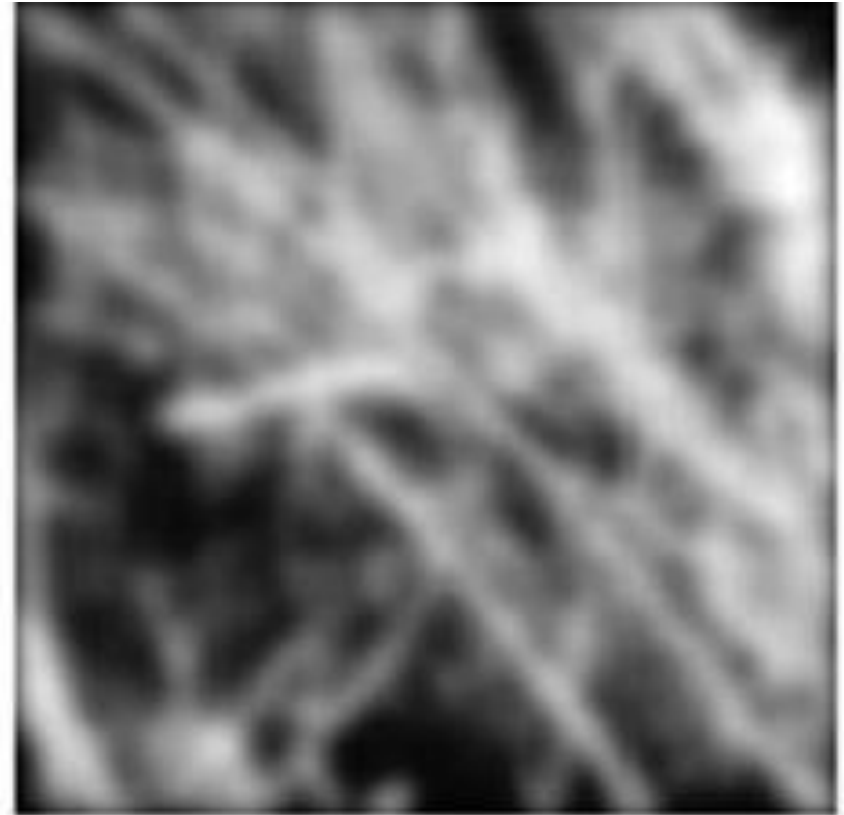


Can you see  
some patterns?

It's another form  
of aliasing



2D box

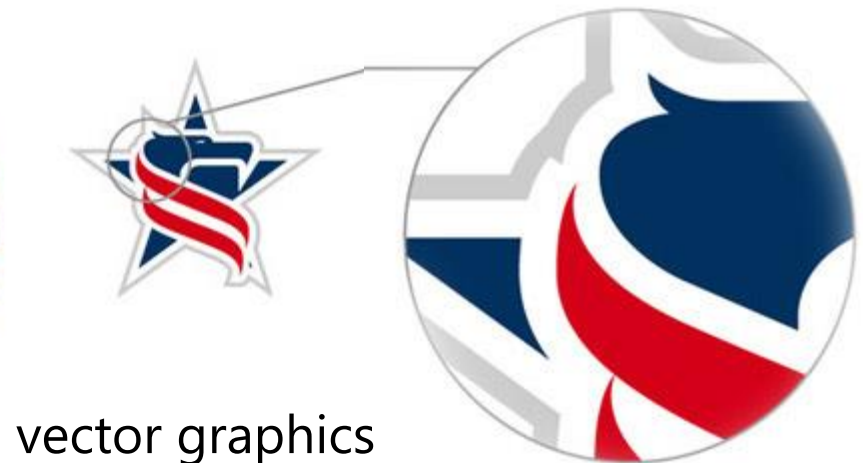
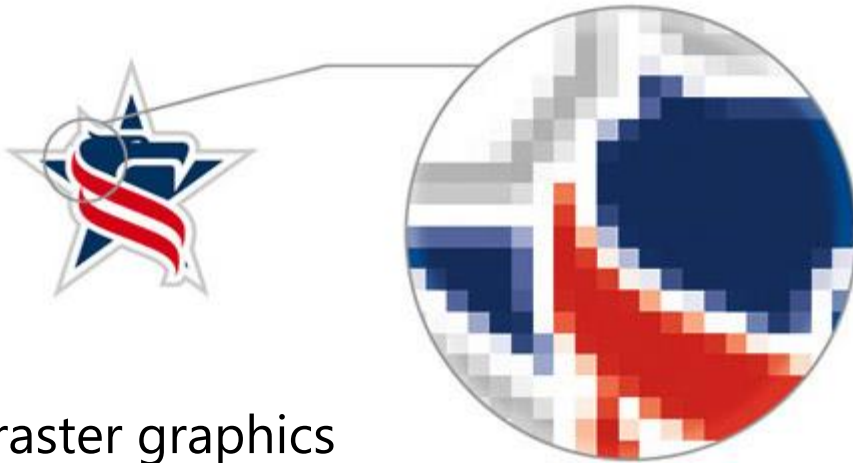


2D Gaussian

# THE SOLUTION

What's the underlying problem?

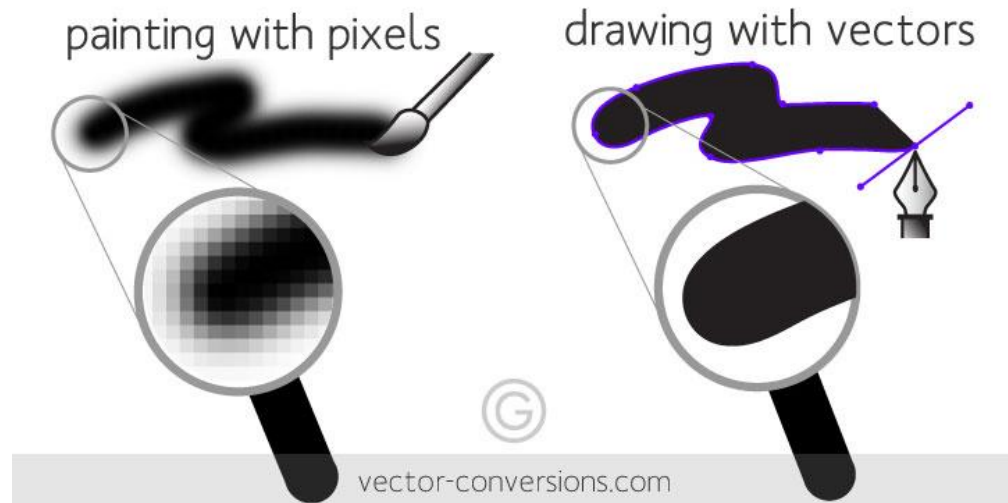
- detail can't be refined upon zoom
- can just be replicated or blurred



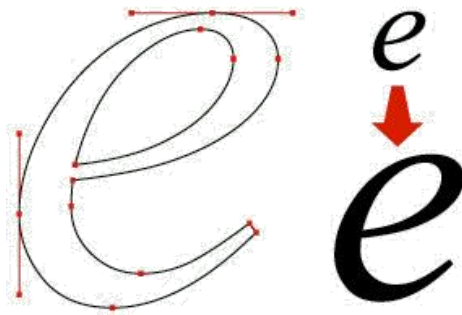
The solution...

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

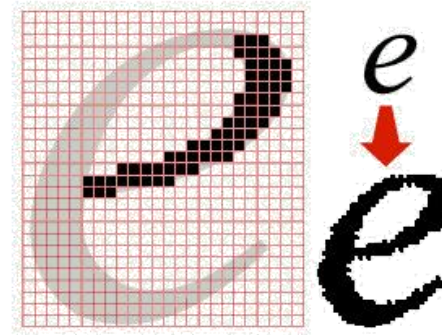
# SCALABLE VECTOR GRAPHICS (SVG)



VECTOR GRAPHICS



BITMAPMED (RASTER) GRAPHICS





# PHOTOGRAPHS AND IMAGES IN SVG

Vector graphics tends to have an “cartoonish” look



raster graphics



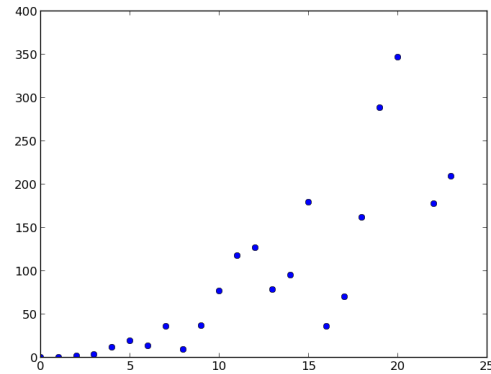
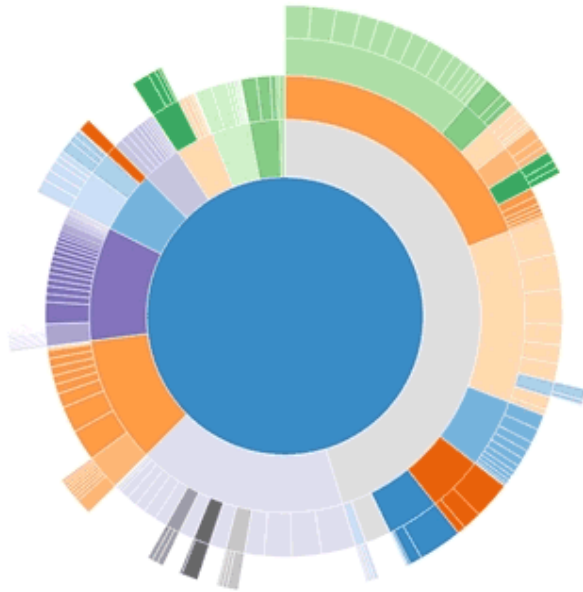
vector graphics

# PHOTOGRAPHS AND IMAGES IN SVG

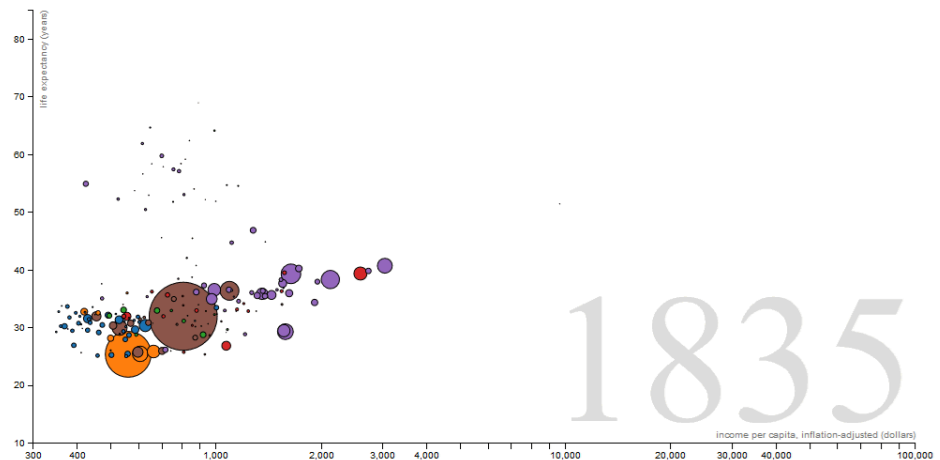
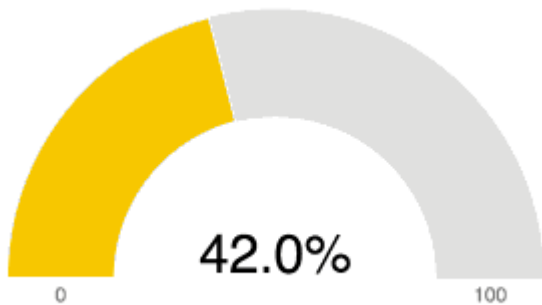




# D3 USES SVG

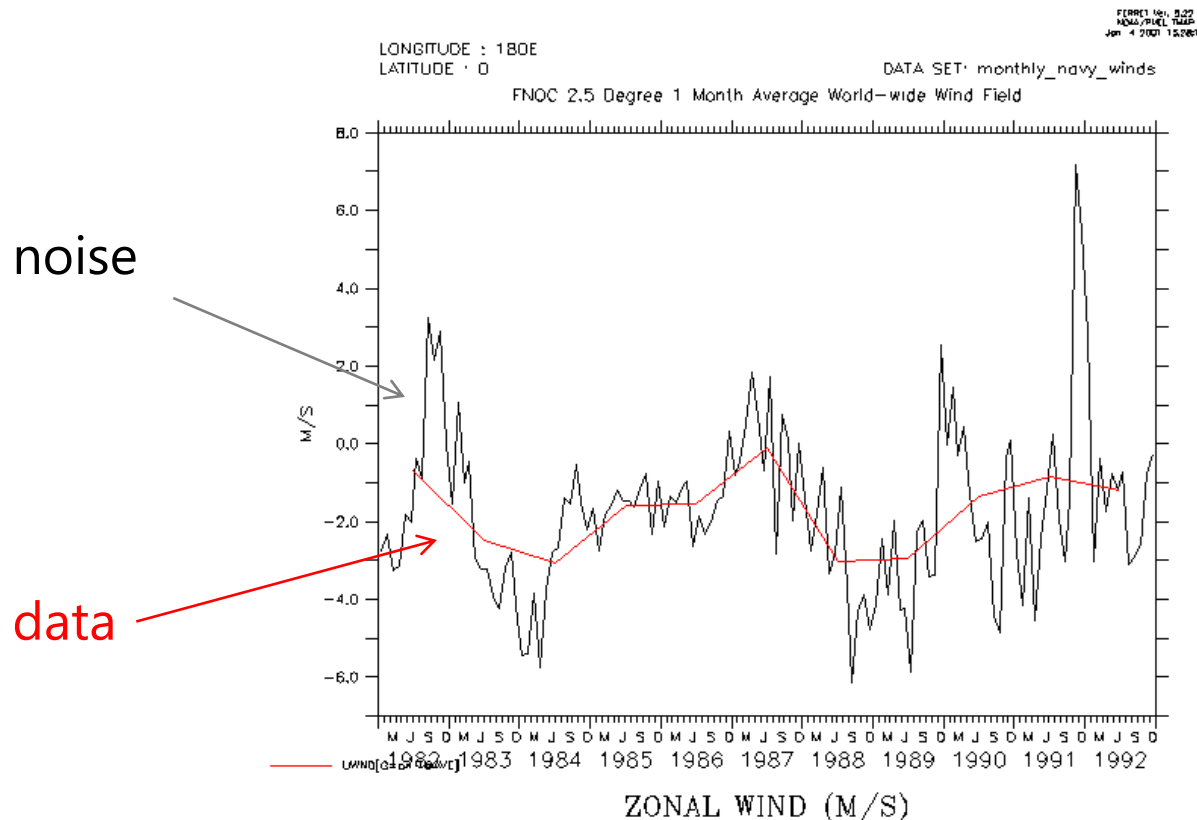


## The Wealth & Health of Nations



# SMOOTHING FOR DE-NOISING

Filtering/smoothing also eliminates noise in the data



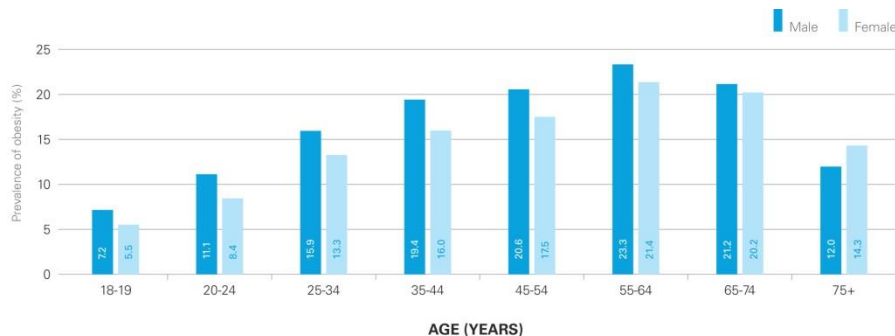
# BACK TO BAR CHARTS

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

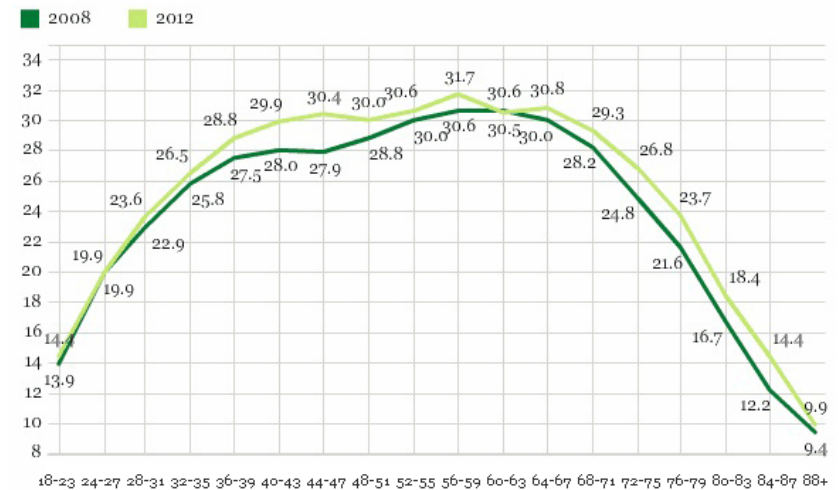
- the bins do the smoothing

Example:

- obesity over age (group)



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



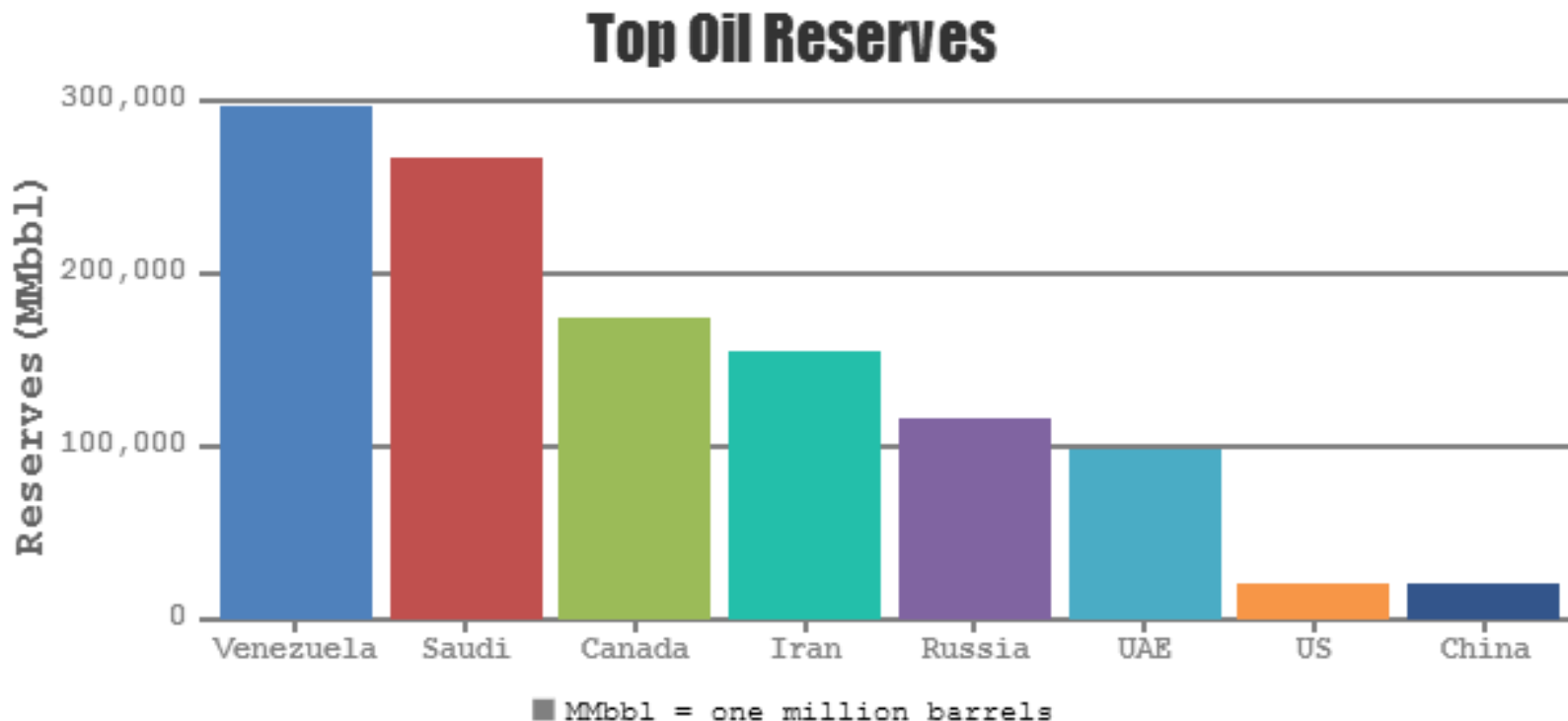
Gallup-Healthways Well-Being Index

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 IN D3

<http://bl.ocks.org/mbostock/3885304>

Working with bar charts will be your job for Lab 1

- the next two slides offer some help with calculations

# BAR CHART CALCULATIONS – BINNING

Determine bin size

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

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

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

- but first initialize bin val array to 0

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

# BAR CHART 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

