

GETTING TO KNOW YOUR DATA

Hari Sundaram

hs1@illinois.edu

<http://sundaram.cs.illinois.edu>

adapted from slides by Jiawei Han and Kevin Chang

**thank you for
responding to the
survey**

research

learning new ideas

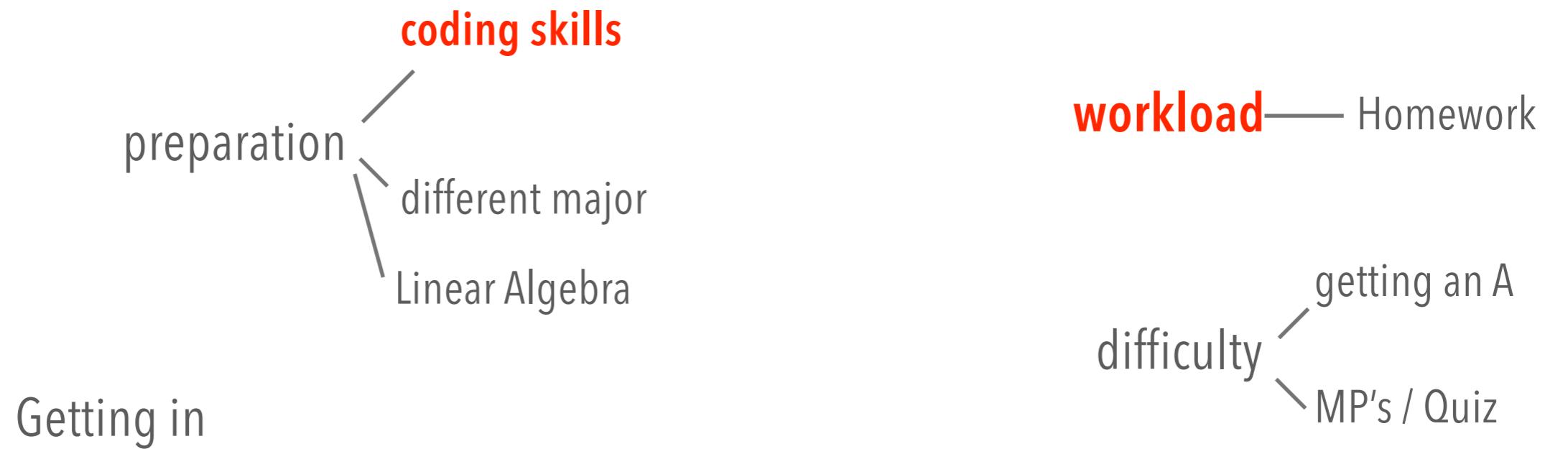
big data

AI

Excited

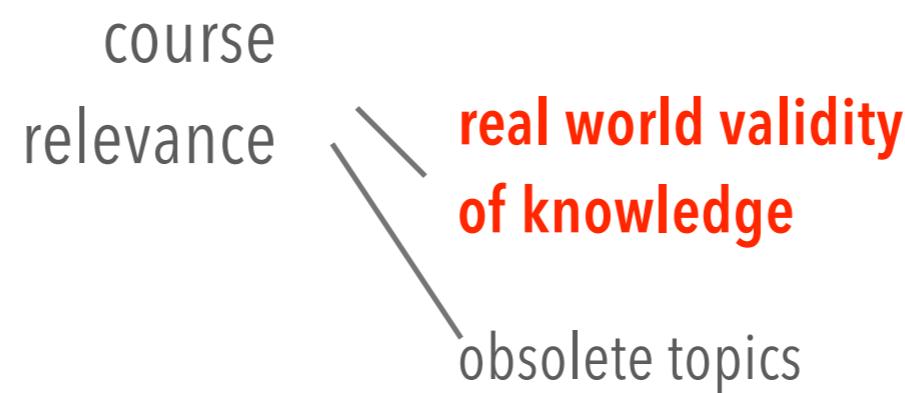
Career

hot topic



Concerned

No concerns



CS 205: data discoveries

More applied



CS 412: data mining

More theory



CS 446 / 498: ML / Applied ML



Eric & Wendy Schmidt
Apply by Jan 31st! **Data Science For Social Good**
Summer Fellowship



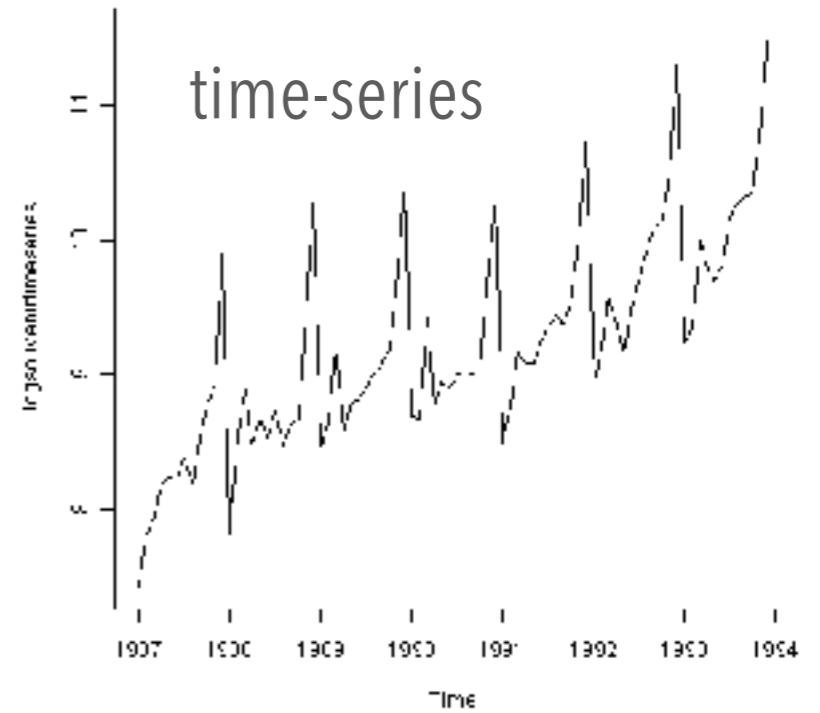
We're training data scientists to tackle problems that really matter.

Applications for Summer 2016 are now open. Click on the links below to learn more and apply.

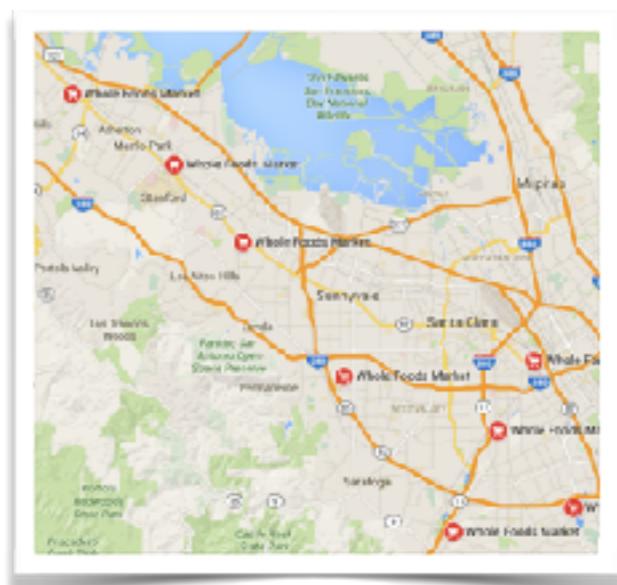
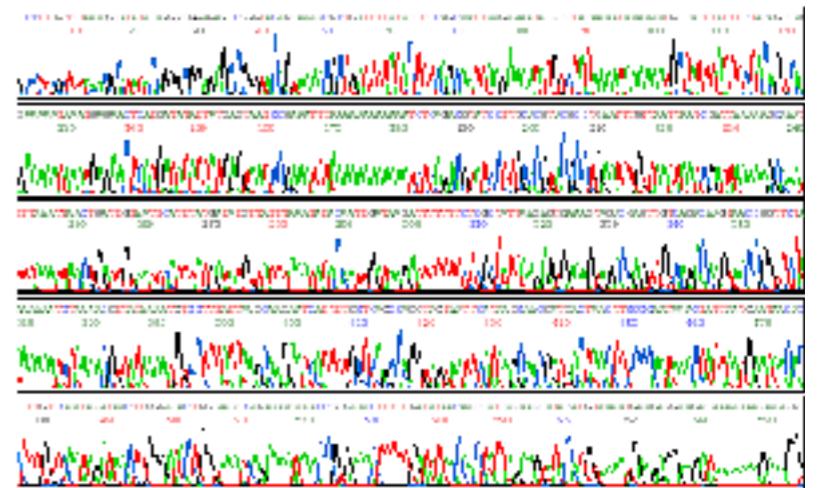
<http://dssg.uchicago.edu/#>

<http://www.datakind.org>

TYPES OF DATASETS



Genetic sequence data



Types of Datasets

Record

Data matrix, e.g., numerical matrix, crosstabs

Relational records

team	coach	play	ball	score	game	n	wi	lost	timeout	season
Document 1	3	0	5	0	2	6	0	2	0	2
Document 2	0	7	0	2	1	0	0	3	0	0
Document 3	0	1	0	0	1	2	2	0	3	0

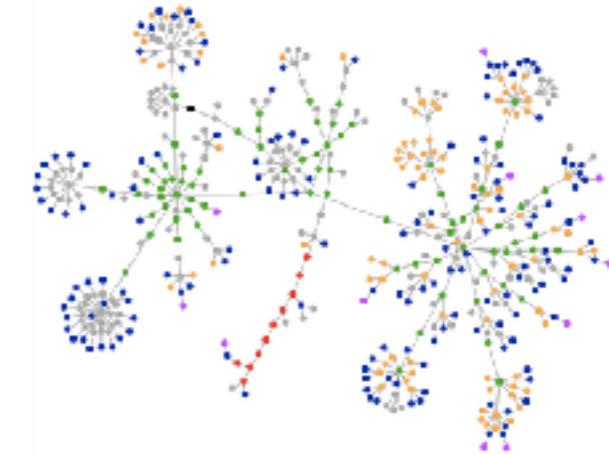
Document data: text documents:
term-frequency vector

Transaction data

TID	Items
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

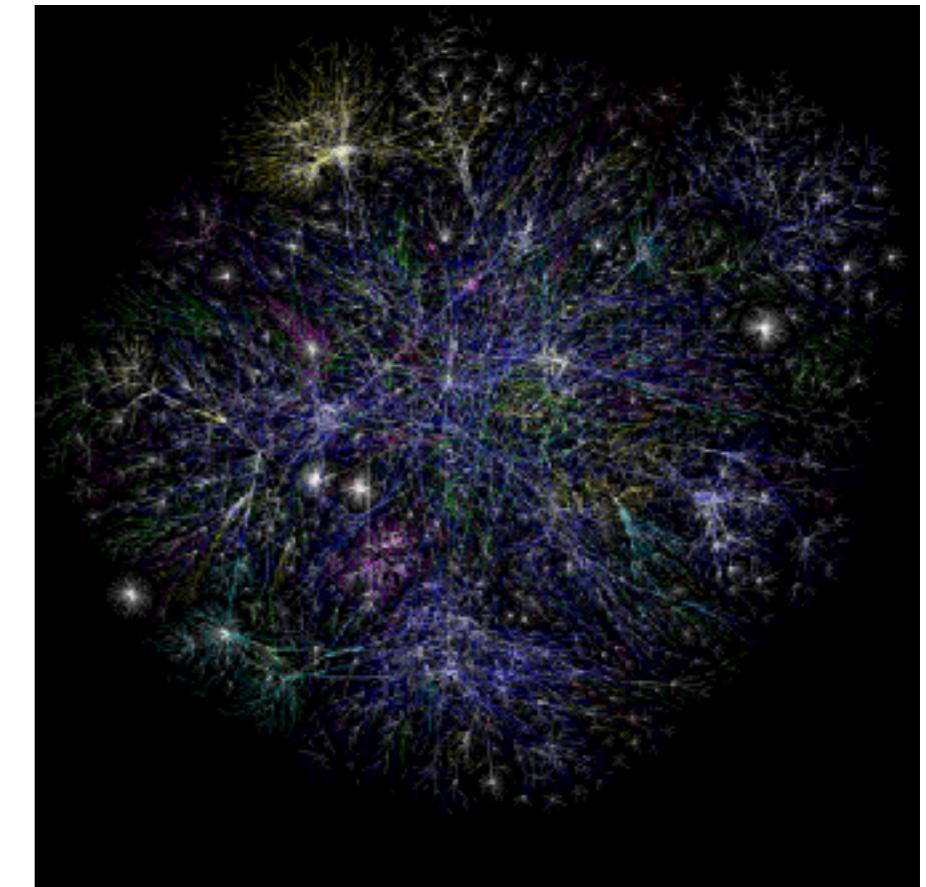
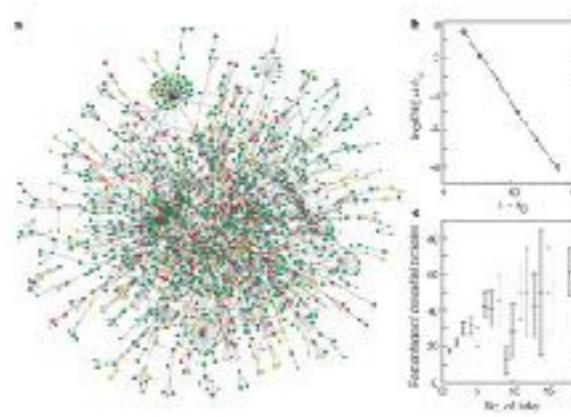
Types of Datasets

Social or
information
networks



Graphs and networks

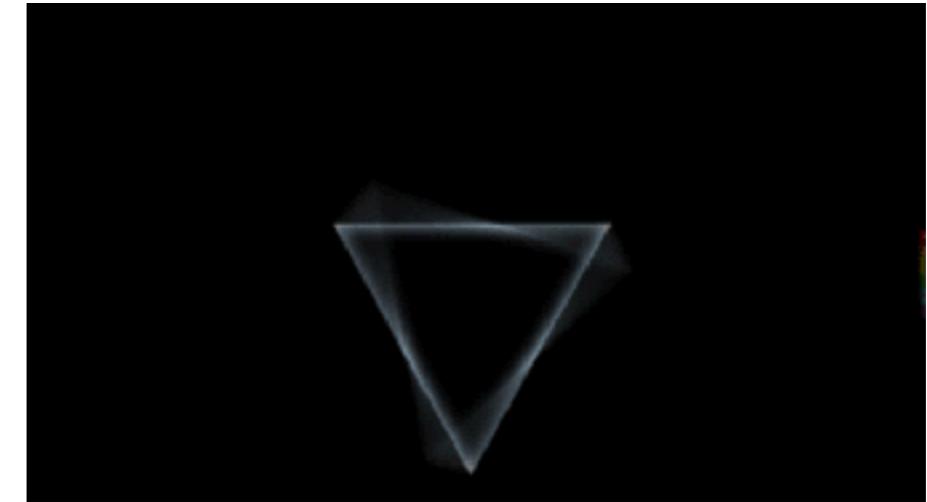
Molecular Structures



World Wide Web

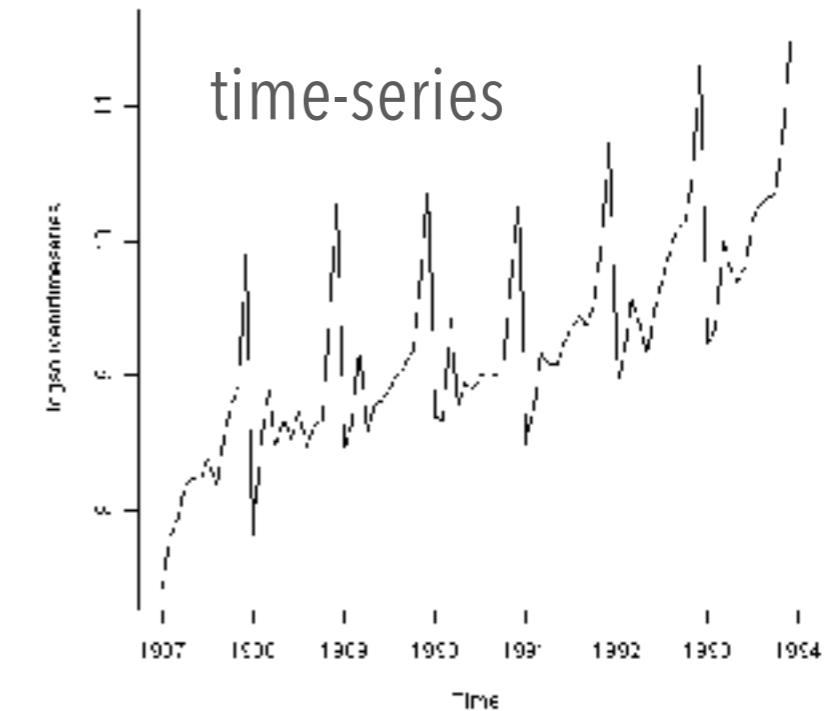
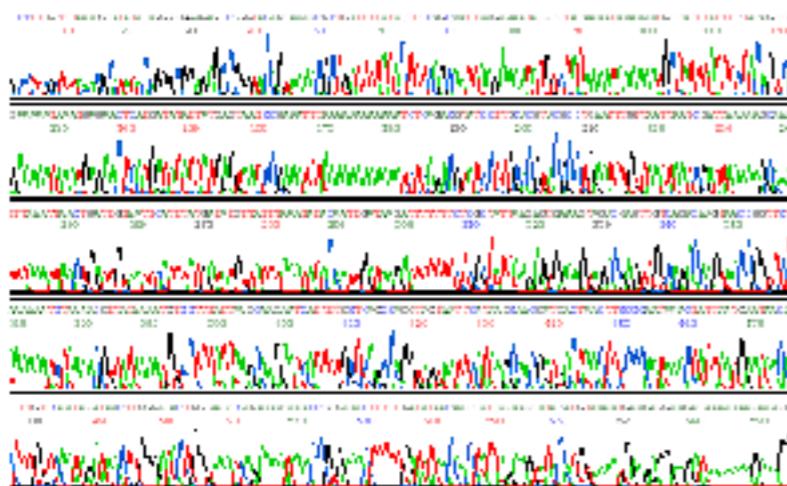
Types of Datasets

Video data:
sequence of
images



Ordered

Genetic sequence data



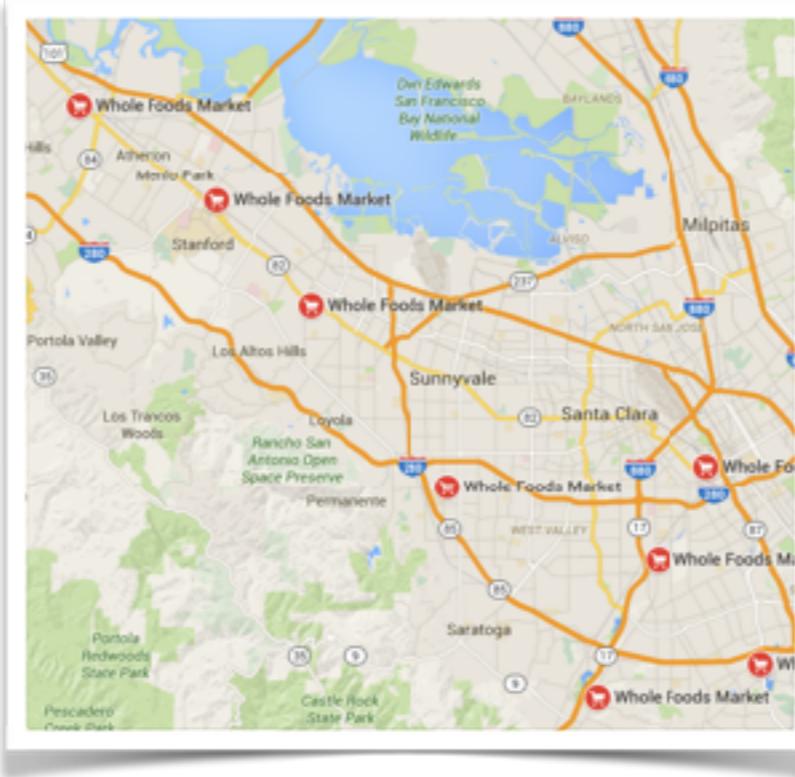
Sequential Data: transaction sequences

Types of Datasets

Spatial,
image and
multimedia

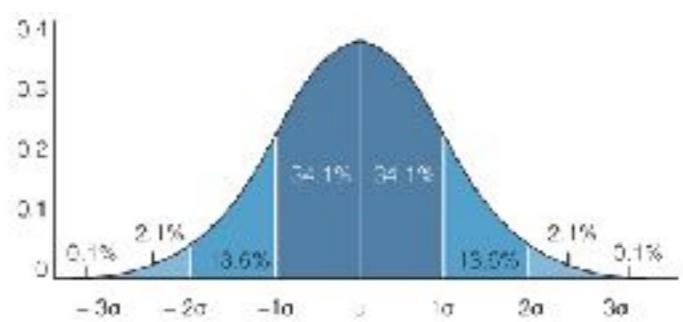
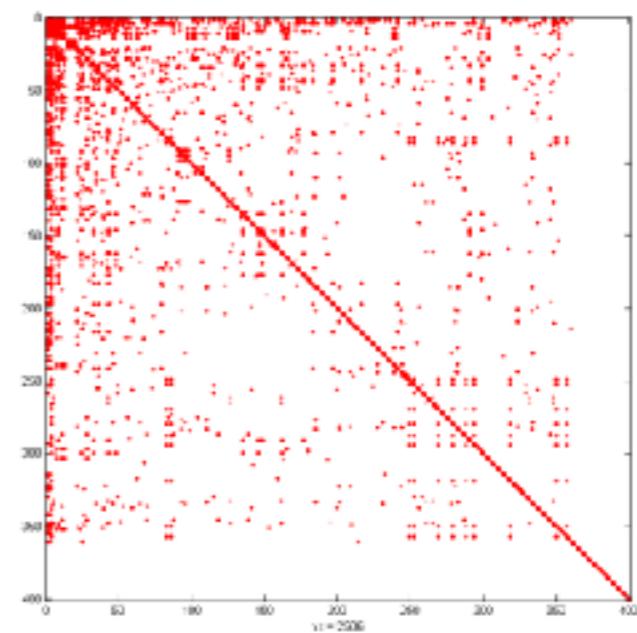


multiple sources



IMPORTANT CHARACTERISTICS OF STRUCTURED DATA

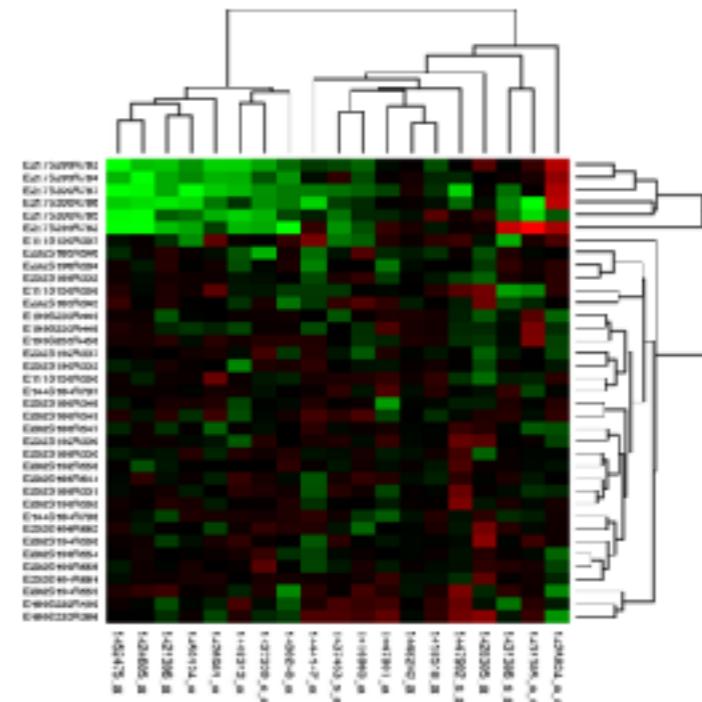
.....



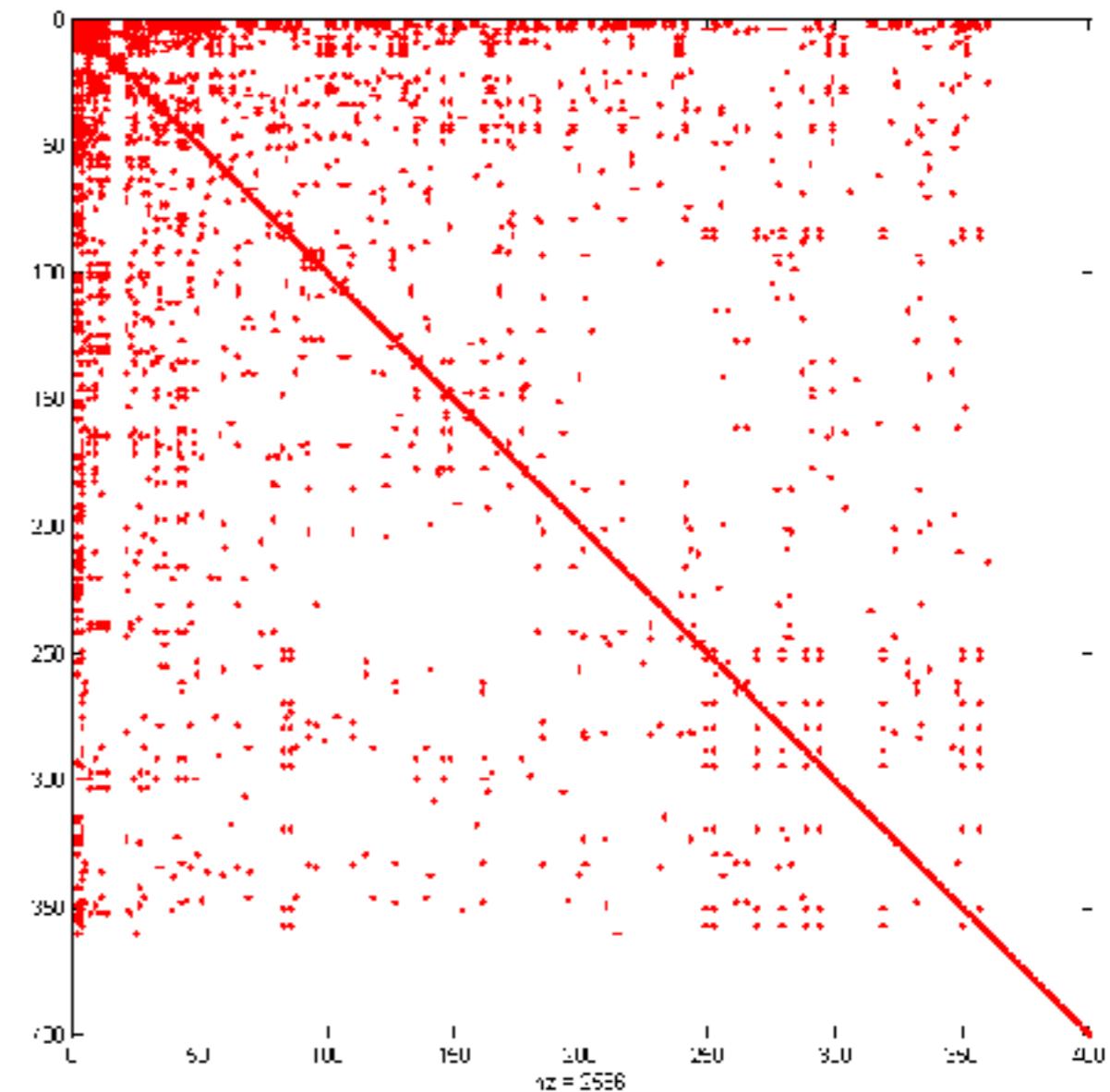
curse



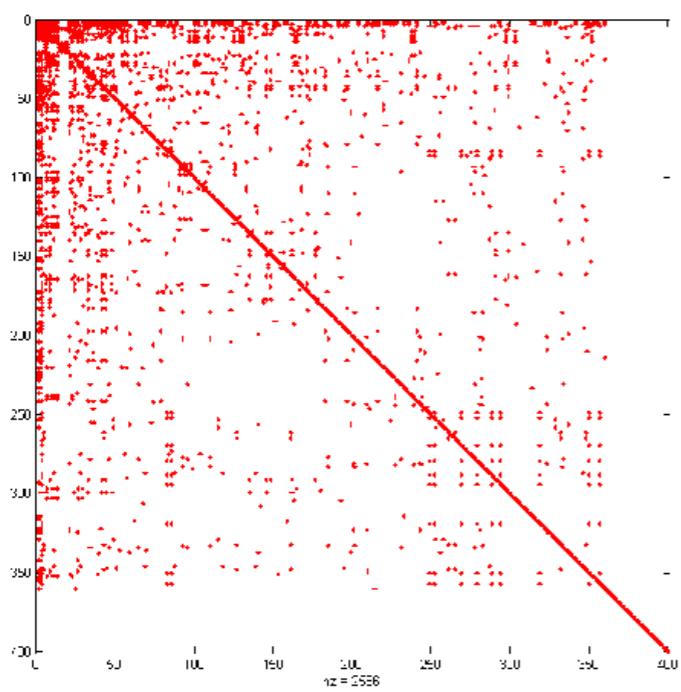
Dimensionality



Sparsity



sparsity pattern of a co-occurrence matrix



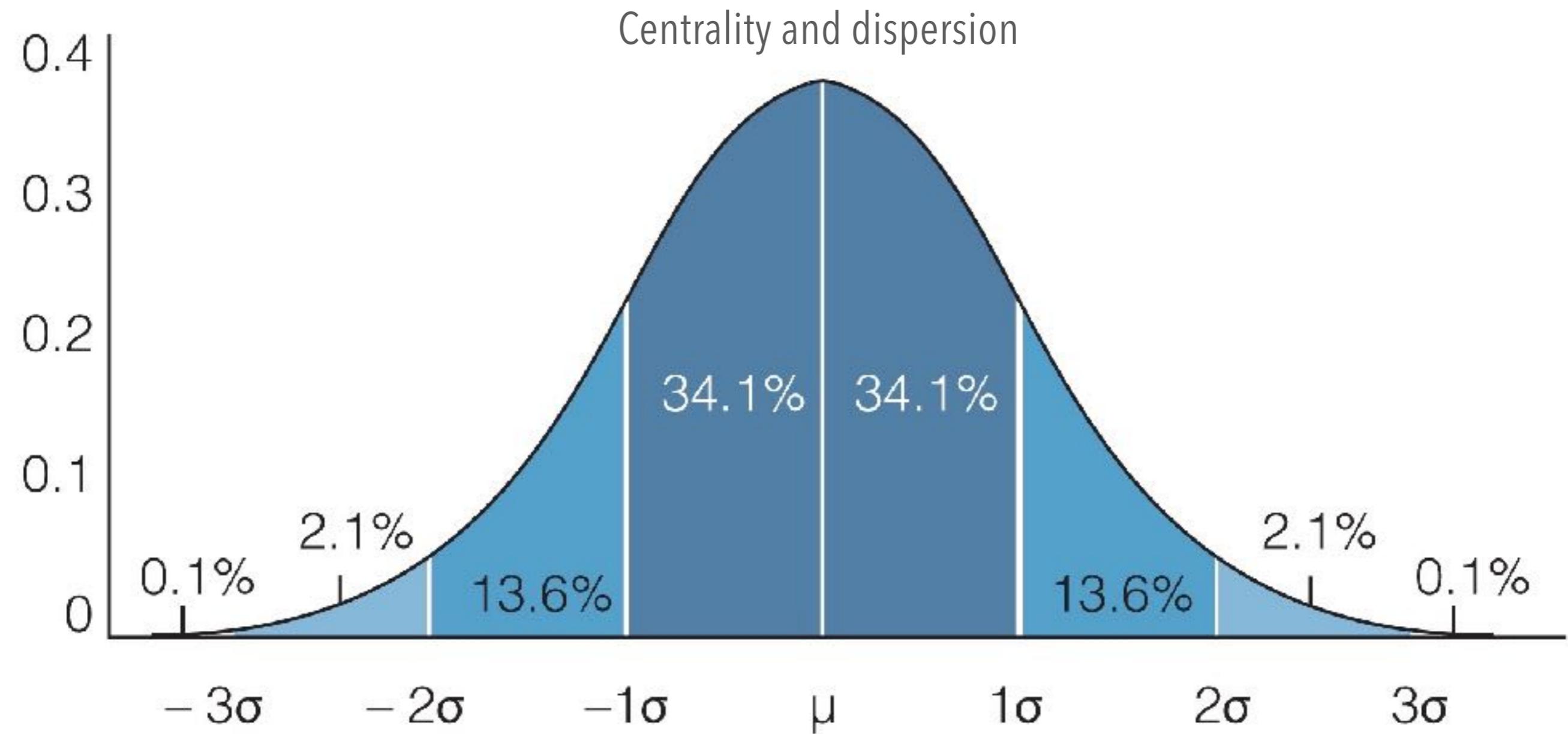


pattern depends on scale

Resolution



Distribution



DATA OBJECTS, ATTRIBUTES



a data object
represents an
entity



sales database:
customers, store
items, sales

medical database:
patients, treatments

university database:
students,
professors, courses

Also called samples ,
examples, instances, data
points, objects, tuples.

datasets

data objects make up datasets

objects

ID	Last Name	First Name	City	State	Gender	Student Status	Major	Country	Age	SAT	Average Score	Honor Roll	(1) Nationality	(class)
10001	JAMES	John	Los Angeles	California	Male	Graduate	Mathematics	US	30	1200	67	SL	1	
10002	JAMES	John	Phoenix	Arizona	Female	Undergraduate	Math	US	19	1000	60	SL	7	
10003	JAMES	John	New York	New York	Male	Graduate	Math	US	24	1222	78	TS	5	
10004	JAMES	John	Los Angeles	California	Male	Graduate	Math	US	23	1001	63	SL	3	
10005	JAMES	John	Chicago	Illinois	Male	Graduate	Math	US	27	1001	68	TS	3	
10006	JAMES	John	Tel Aviv	Israel	Male	Graduate	Math	Israel	27	1200	80	TS	2	
10007	JAMES	John	Cairo	North Carolina	Male	Graduate	Math	US	26	1001	76	TS	1	
10008	JAMES	John	Lagos	Nigeria	Female	Undergraduate	Math	US	21	1002	87	SL	3	
10009	JAMES	John	Montreal	Canada	Female	Undergraduate	Math	Canada	18	1003	84	SL	3	
10010	JAMES	John	New York	New York	Female	Graduate	Math	US	22	1004	71	TS	1	
10011	JAMES	John	Paris	Mississippi	Male	Undergraduate	Math	US	18	1001	89	TS	1	
10012	JAMES	John	London	England	Female	Graduate	Math	US	20	1200	74	TS	1	
10013	JAMES	John	Paris	England	Male	Graduate	Math	US	20	1001	73	TS	1	
10014	JAMES	John	Varna	Bulgaria	Male	Graduate	Math	US	17	1001	68	SL	4	
10015	JAMES	John	Moscow	Russia	Male	Graduate	Math	Russia	18	1202	70	TS	2	
10016	JAMES	John	Orlando	Florida	Female	Undergraduate	Math	US	22	1002	82	SL	3	
10017	JAMES	John	Provo	Utah	Female	Undergraduate	Math	US	18	1001	80	TS	1	
10018	JAMES	John	Amsterdam	Holland	Female	Undergraduate	Math	US	19	1204	75	SL	3	
10019	JAMES	John	Perito	Bolivia	Female	Graduate	Math	Bolivia	21	1204	88	SL	4	
10020	JAMES	John	Caracas	Venezuela	Female	Undergraduate	Math	Venezuela	18	1202	80	TS	1	
10021	JAMES	John	San Jose	Puerto Rico	Male	Graduate	Math	US	22	1003	86	SL	7	
10022	JAMES	John	Portland	Oregon	Female	Undergraduate	Math	US	18	1001	87	SL	7	
10023	JAMES	John	New York	New York	Male	Undergraduate	Math	US	21	1002	88	TS	4	
10024	JAMES	John	The Bronx	Kansas	Female	Graduate	Math	US	23	1203	83	SL	3	
10025	JAMES	John	Bethesda	China	Female	Undergraduate	Math	China	18	1003	79	TS	2	
10026	JAMES	John	Stockholm	Sweden	Male	Undergraduate	Math	Sweden	19	1001	88	SL	4	
10027	JAMES	John	Minneapolis	Minnesota	Male	Graduate	Math	US	20	1204	90	TS	1	
10028	JAMES	John	Interlachen	Ferryland	Male	Undergraduate	Math	US	21	1205	88	TS	1	
10029	JAMES	John	Luca	Ukrania	Female	Undergraduate	Math	US	20	1005	86	SL	3	
10030	JAMES	John	Buenos Aires	Argentina	Male	Graduate	Math	Argentina	19	1206	85	TS	2	
10031	JAMES	John	Asuncion	Uruguay	Male	Undergraduate	Math	Uruguay	18	1007	79	TS	1	

attributes

E.g., customer_ID, name, address

Attributes

(or dimensions, features, variables): a data field, representing a characteristic or feature of a data object.



Nominal (or categorical)

{green, red, blue}

Binary {true, false}

Types

Numeric

quantitative

Interval-scaled

Ordinal

{small, medium, large}

Ratio-scaled

marital status,
occupation, ID numbers,
zip codes

Nominal

(aka Categorical)

Hair color = {auburn, black,
blond, brown, grey, red,
white}

**WHAT IS THE
AVERAGE
COLOR OF
THIS IMAGE?**

.....





Asymmetric binary: outcomes not equally important. e.g., medical test (positive vs. negative) Convention: assign 1 to most important outcome (e.g., HIV positive)

Binary

Nominal attribute with only two states

Symmetric binary: both outcomes equally important

Size = {small, medium,
large}, course grades,
army rankings

Ordinal

Values have a meaningful
order (ranking) but magnitude
between successive values is
not known.



quantity: 7.8, 9,
101.89 etc.

Interval

Measured on a
scale of equal-
sized units

e.g., temperature
in C° or F°,
calendar dates

Values have order

No true zero-point

Numeric



Inherent zero-point

e.g., temperature in
Kelvin, length, counts,
monetary quantities

Ratio

We can speak of values as
being an order of magnitude
larger than the unit of
measurement (10 K° is twice
as high as 5 K°).



It is 0°F in
Chicago;
when will it
be twice as
hot?

E.g., zip codes,
profession, or the set of
words in a collection of
documents

Sometimes
represented as
integer variables

Discrete

Has only a finite or
countably infinite set
of values

Note: Binary
attributes are a special
case of discrete
attributes

Has real
numbers as
attribute values

In practice, attribute
values can only be
measured and
represented using a
finite number of digits

Continuous

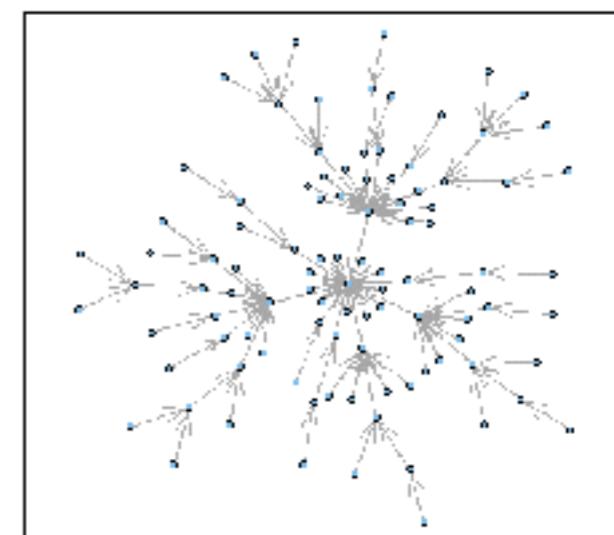
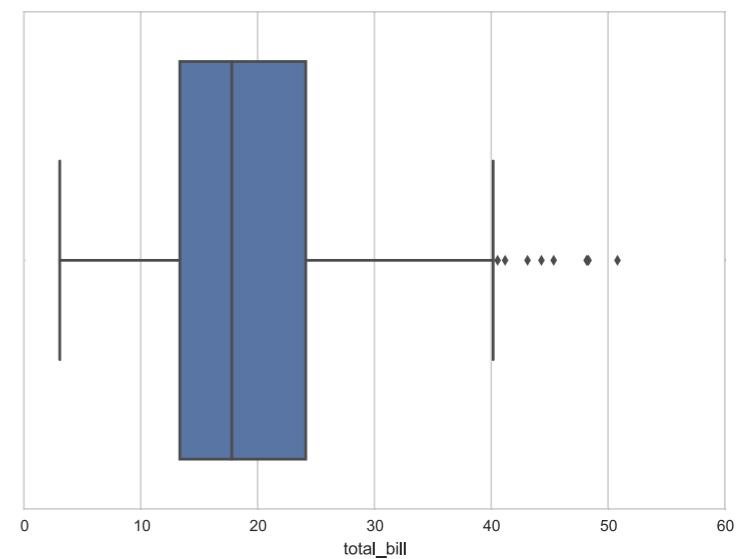
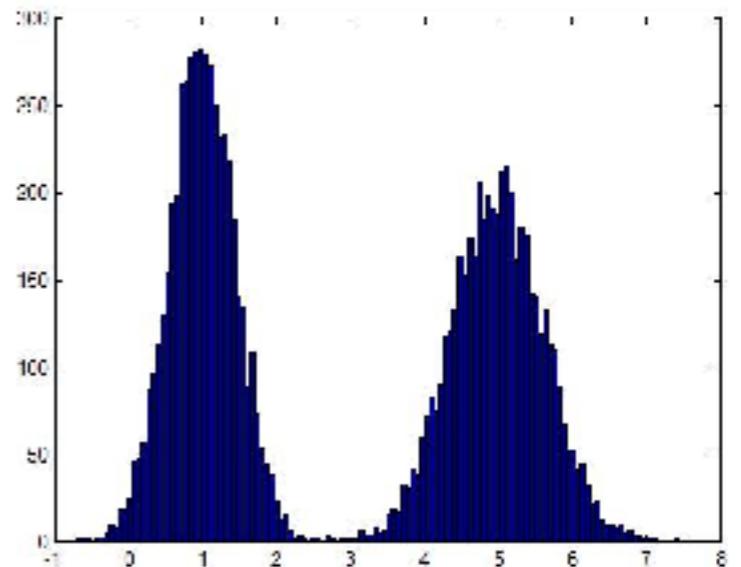
Continuous attributes
are typically
represented as
floating-point variables

e.g., temperature,
height, or weight;
weight=125.6 lb.



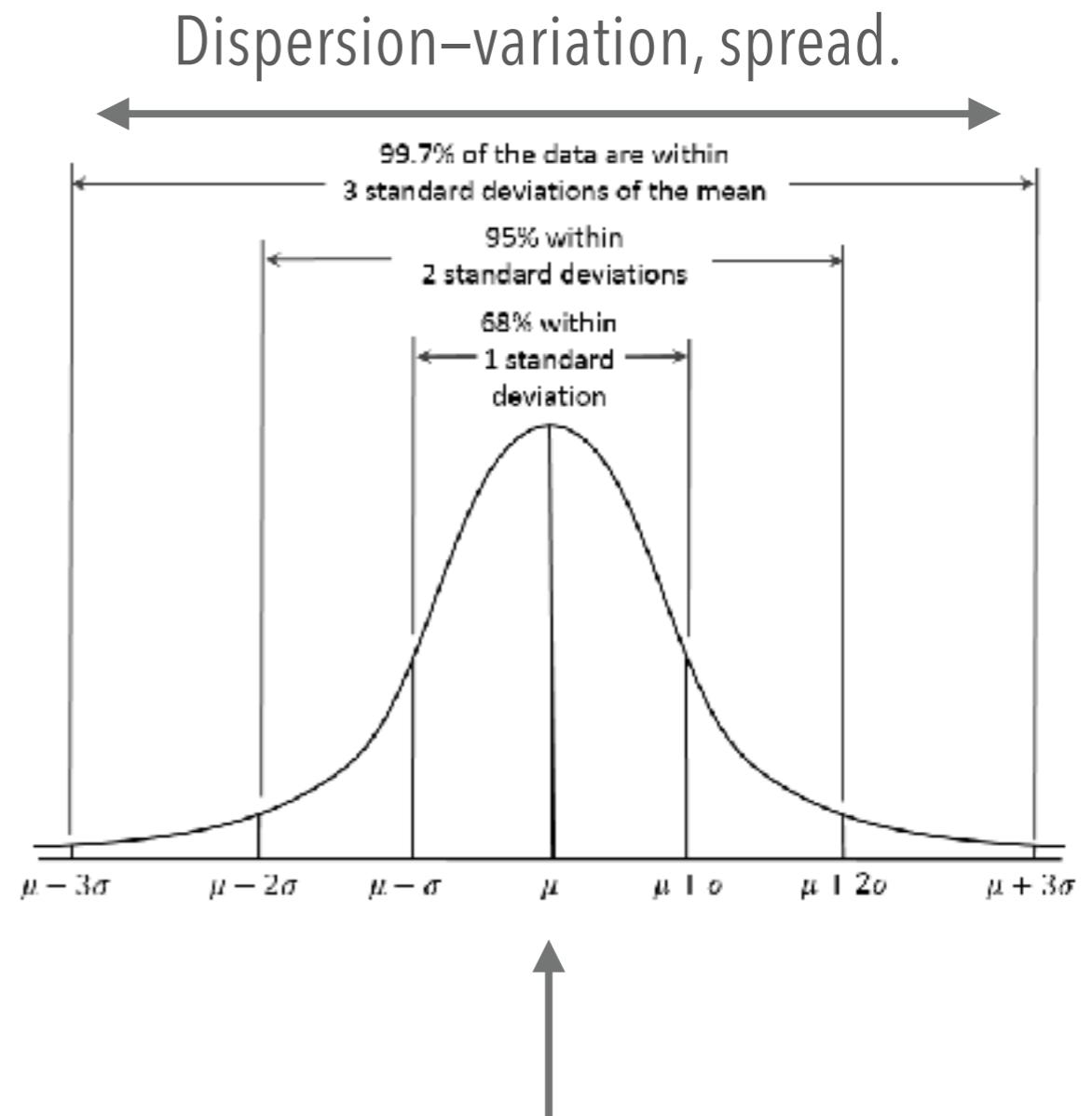
BASIC STATISTICAL DESCRIPTIONS OF DATA

.....



Motivation

To better understand the data.



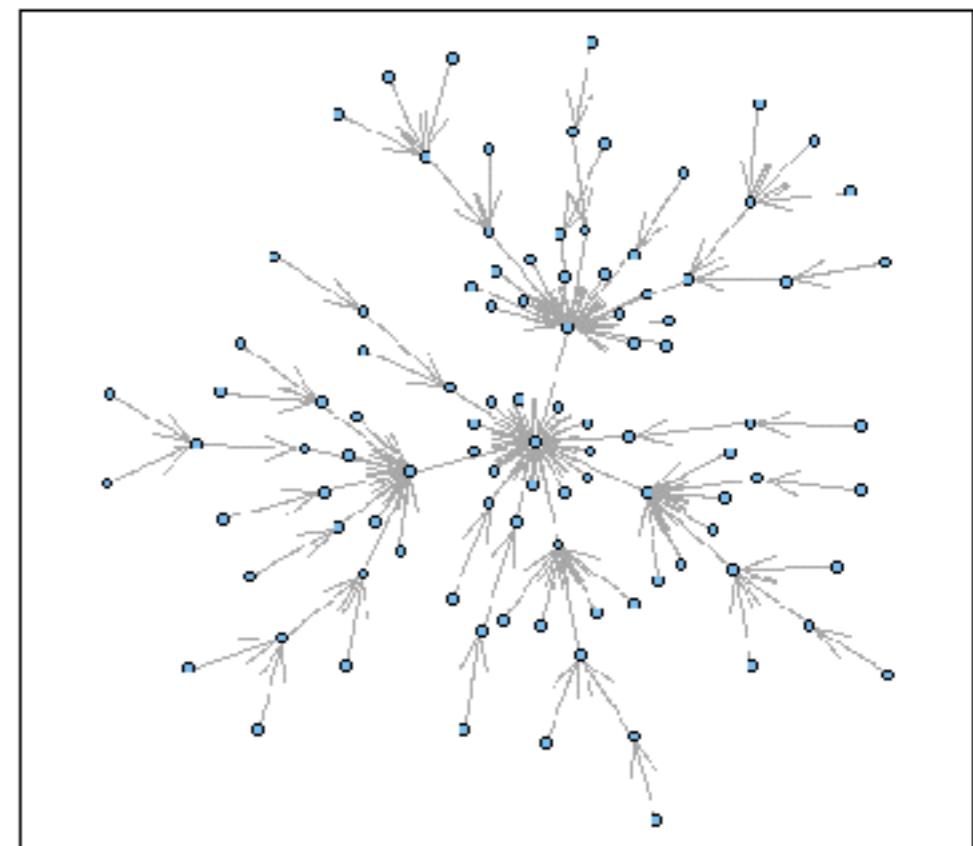
Central tendency—the center, the representative.

perspective
object group surface point map structure orientation
area size sequence proximity adjacency network units
line pattern overlay connectivity symmetry distribution
distance separation gradient center shape boundary density time
field contrast place scale force path neighborhood

TF-IDF

Is that all?

non-numeric data



statistics from graphs

The mean

Trimmed mean: chopping extreme values

weighted
sample
mean

$$\bar{x} = \frac{\sum_i^n \omega_i x_i}{\sum_i^n \omega_i}$$

different notation

population mean

$$\mu = \frac{1}{N} \sum_i^N x_i$$

population size

sample mean

$$\bar{x} = \frac{1}{n} \sum_i^n x_i$$

sample size

Estimated by interpolation (for grouped data):

Middle value if odd number of values, or average of the middle two values otherwise

age	frequency	
1–5	200	
6–15	450	
16–20	300	median is
21–50	1500 ←	to be
51–80	700	found here
81–110	44	

The median

approximate
median
calculation

{

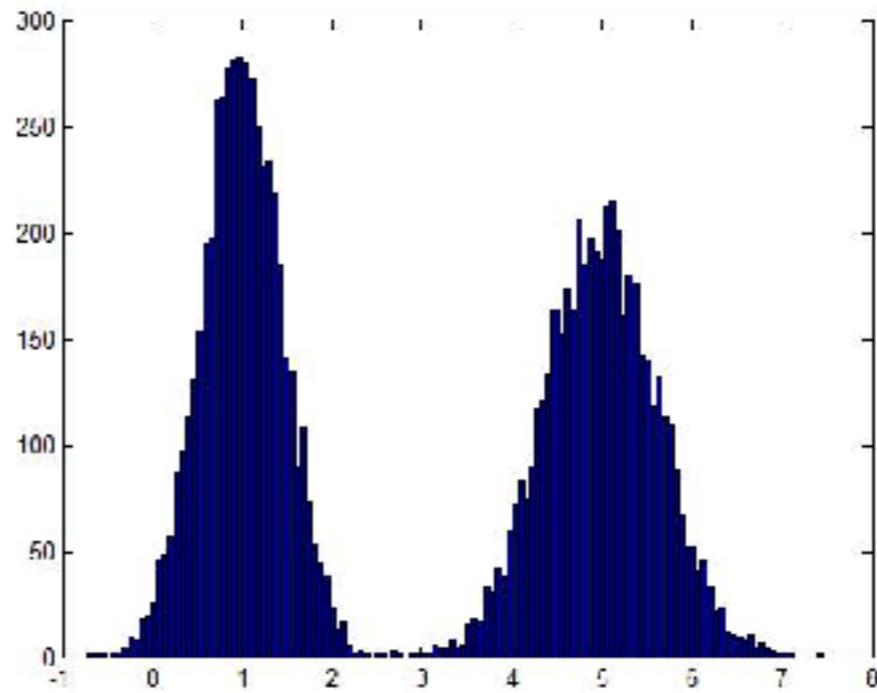
$$L_1 + \left(\frac{n/2 - \sum_l f_l}{f_{median}} \right) \times (L_2 - L_1)$$

lower interval limit

sum before the median interval

interval width

Unimodal,
bimodal, trimodal

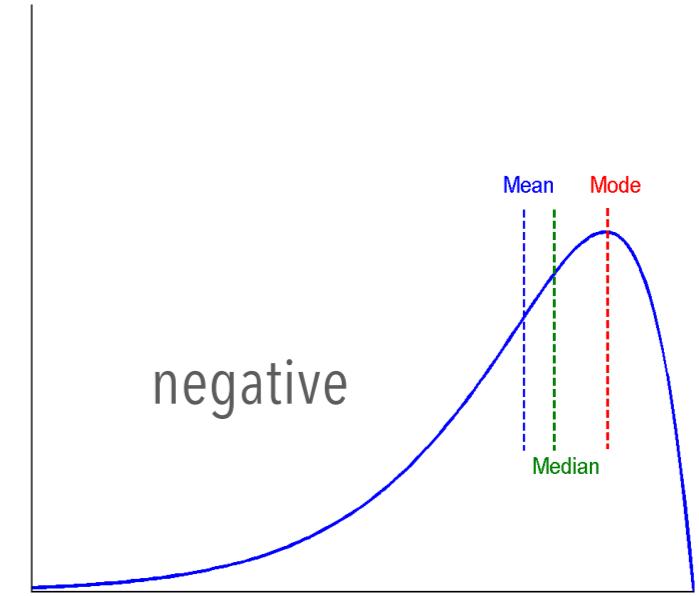
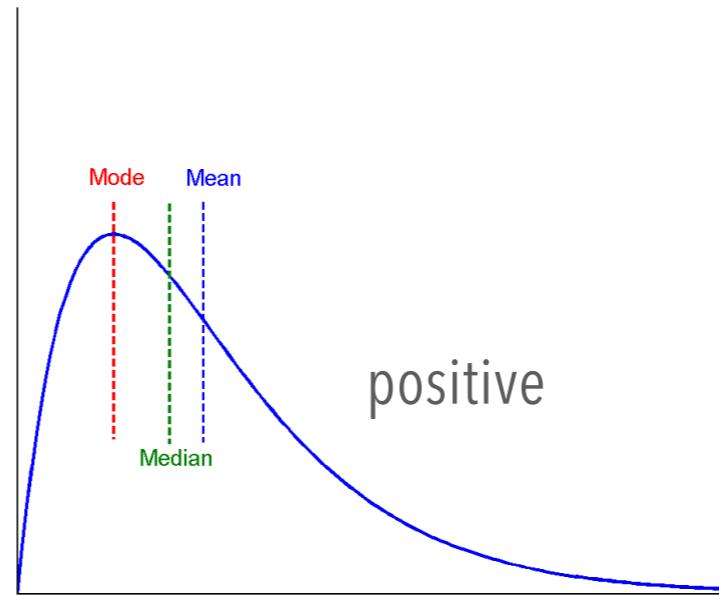
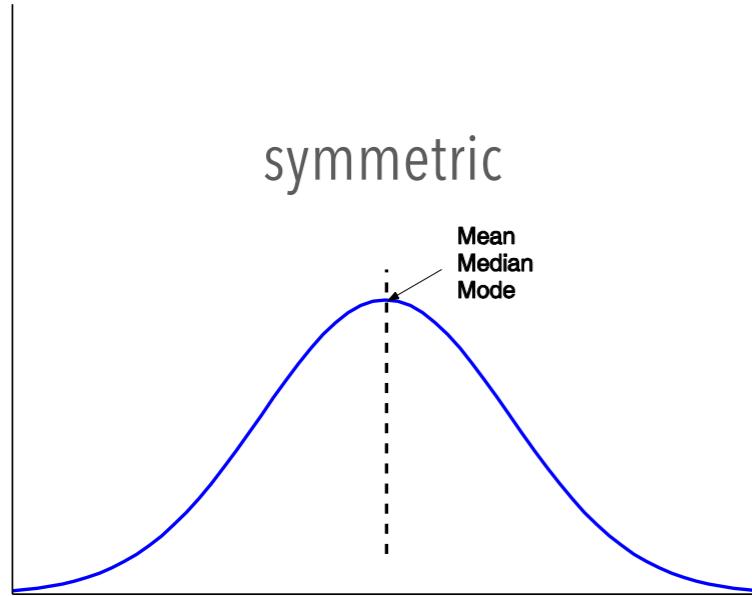


The mode

Value that occurs most frequently in the data

approximate formula

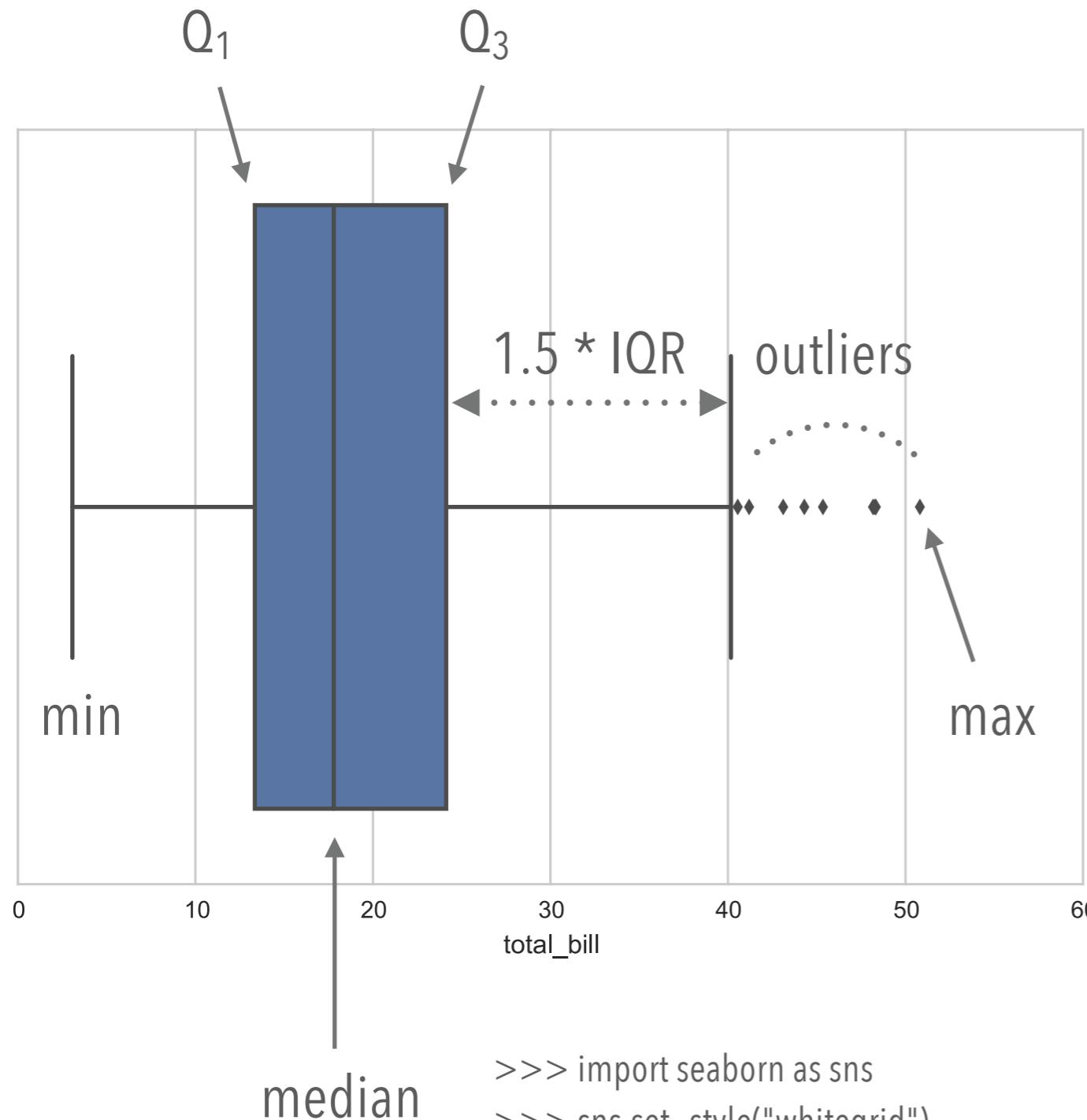
$$\text{mean} - \text{mode} = 3 \times (\text{mean} - \text{median})$$



data skew

Median, mean and mode of symmetric, positively and negatively skewed data

QUARTILES, OUTLIERS AND BOX PLOTS



```
>>> import seaborn as sns  
>>> sns.set_style("whitegrid")  
>>> tips = sns.load_dataset("tips")  
>>> ax = sns.boxplot(x=tips["total_bill"])
```

<http://bit.ly/2BeUfFg>

Quartiles: Q_1 (25^{th} percentile), Q_3 (75^{th} percentile)

Inter-quartile range: $IQR = Q_3 - Q_1$

Five number summary: min, Q_1 , median, Q_3 , max

Boxplot: ends of the box are the quartiles; median is marked; add whiskers, and plot outliers individually

Outlier: usually, a value higher/lower than $1.5 \times IQR$

Variance

sample variance

$$s^2 = \frac{1}{n-1} \sum_i^n (x_i - \bar{x})^2$$

$$\sigma^2 = \frac{1}{N} \sum_i^N (x_i - \mu)^2$$

population variance

that obvious
formula is an
underestimate

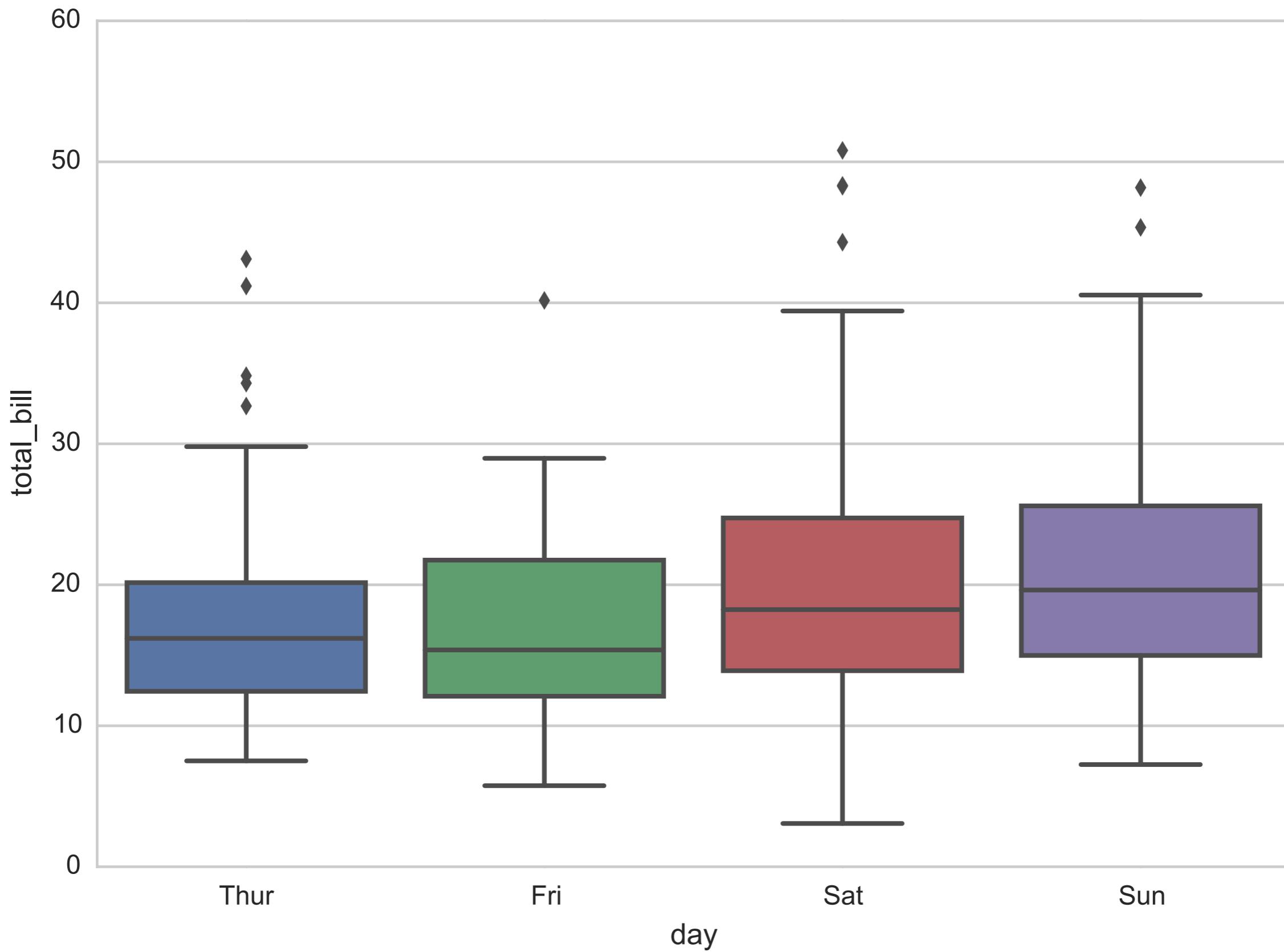
$$s^2 = \frac{1}{n} \sum_i^n (x_i - \bar{x})^2$$

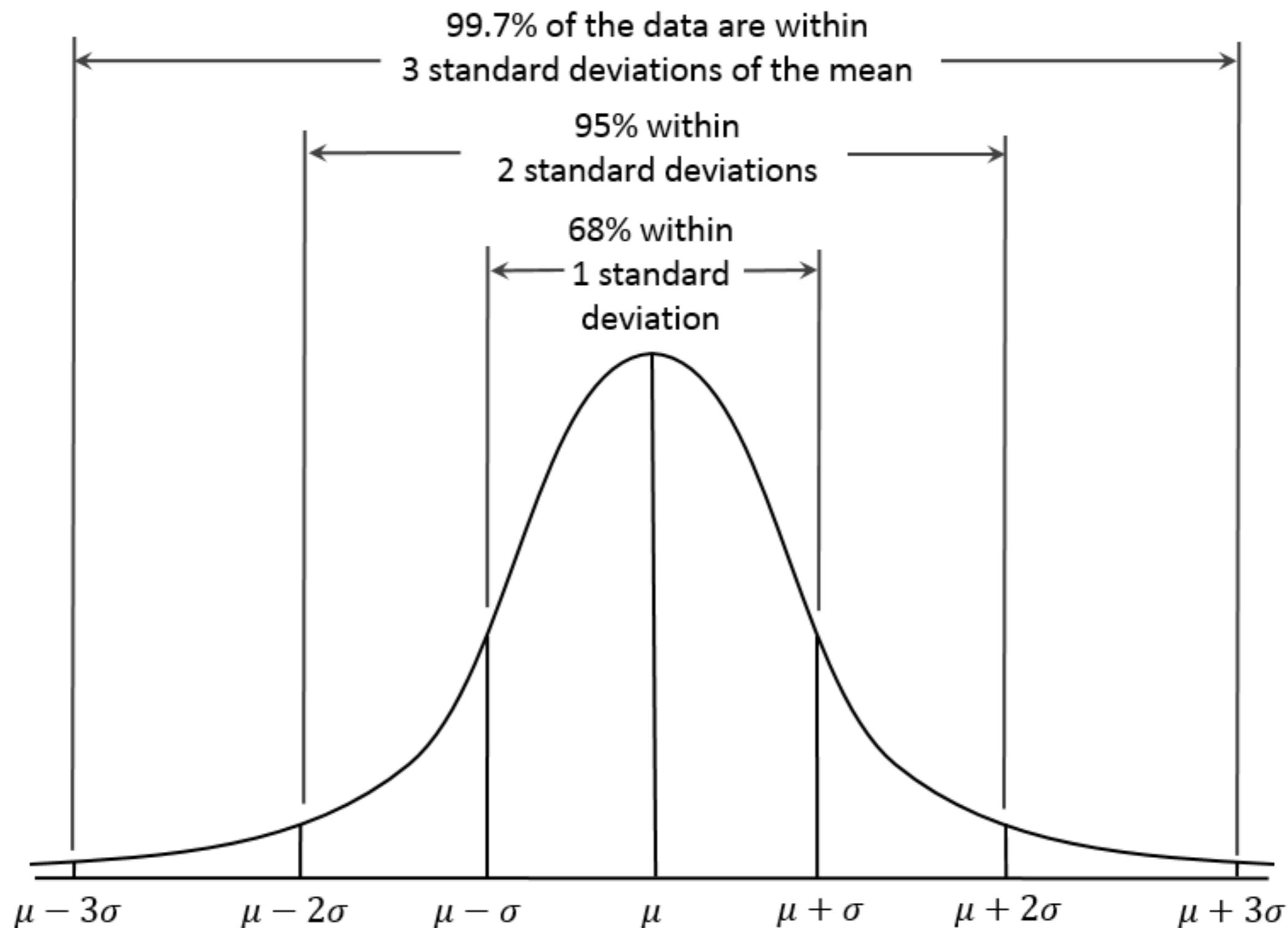
explaining the puzzle

$$s^2 = \frac{1}{n} \left(\sum_i x_i^2 + \sum_i \left(\frac{\sum_j x_j}{n} \right)^2 - 2 \sum_i x_i \left(\frac{\sum_j x_j}{n} \right) \right) \quad s^2 = \frac{1}{n} \left(\sum_i^n x_i^2 - \frac{\sum_i x_i^2 + 2 \sum_{i \neq j} x_i x_j}{n} \right)$$

$$s^2 = \frac{n-1}{n} \sigma^2$$

$$\begin{aligned} E(x_i) &= \mu \\ \text{using } E(x_i^2) &= \sigma^2 + \mu^2 \end{aligned}$$

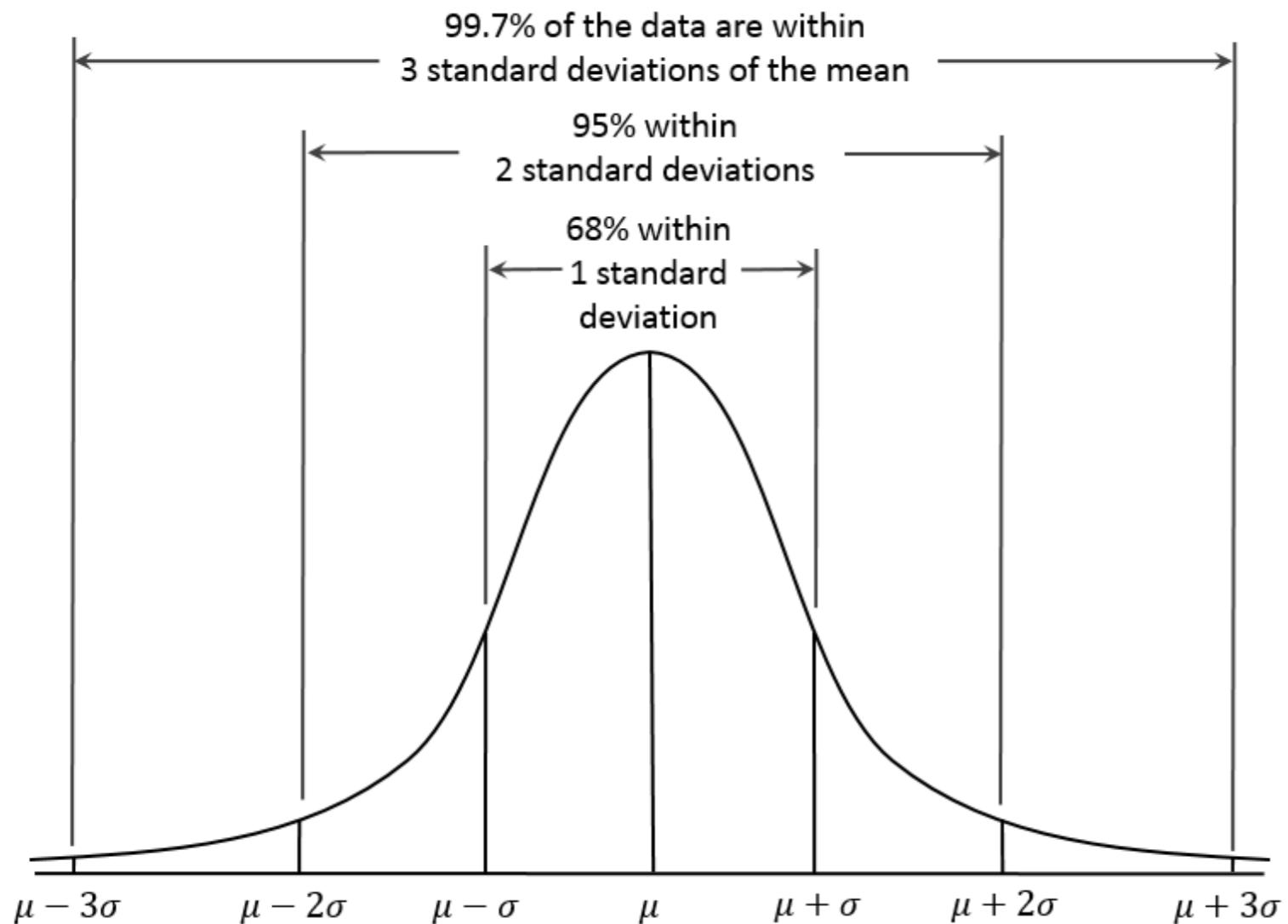




the normal distribution

Why is the normal distribution so important?

Assume that we have n
observations; $n \geq 30$



$$\bar{x} = \frac{1}{n} \sum_i x_i$$

$$E(\bar{x}) = \mu$$

$$Var(\bar{x}) = \frac{\sigma^2}{n}$$

Central Limit Theorem

graphical displays of statistical properties

quantile plots

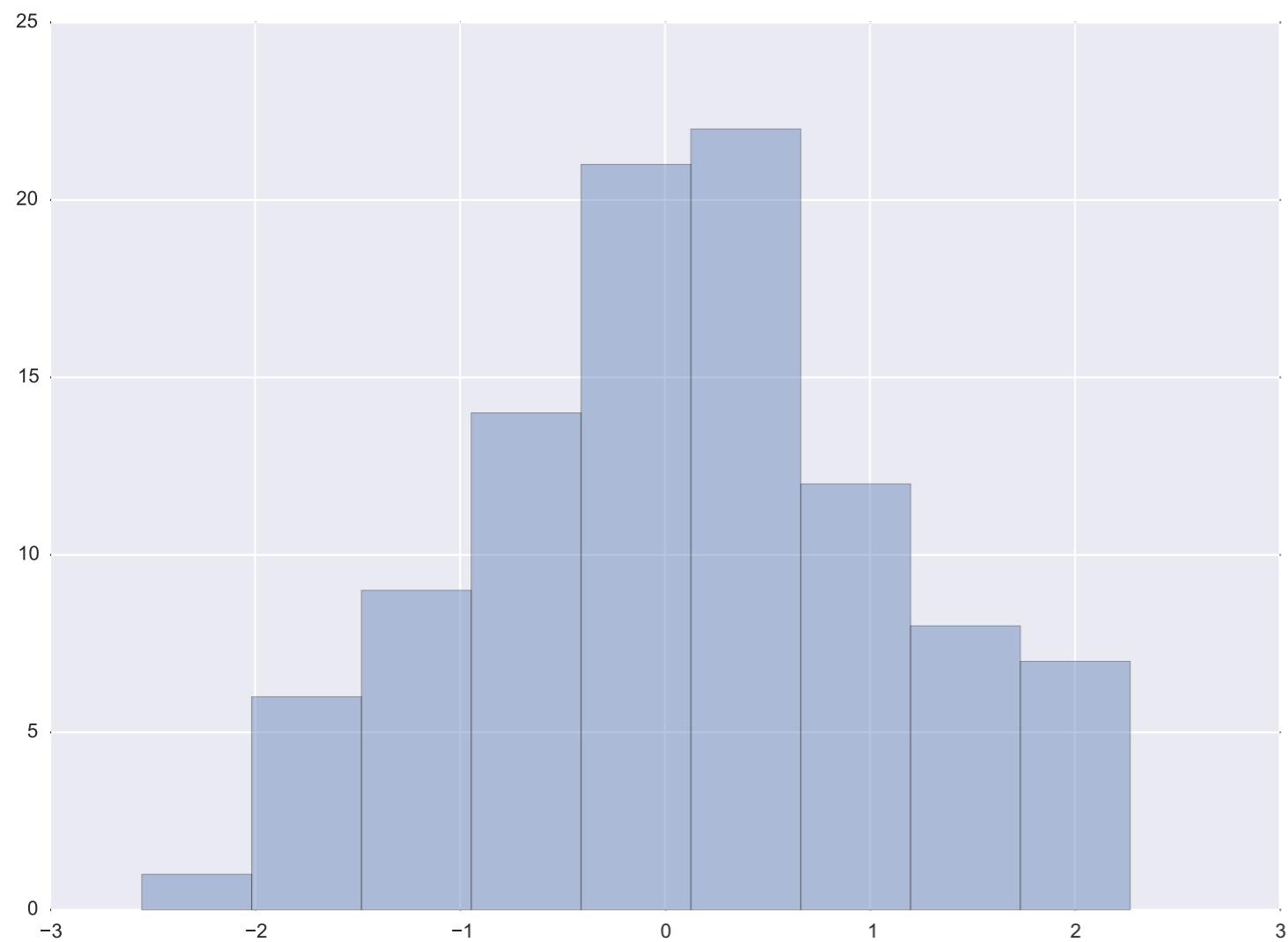
quantile-quantile
plots

box plots

histograms

scatter

HISTOGRAMS



```
>>> import seaborn as sns, numpy as np  
>>> sns.set(rc={"figure.figsize": (8, 4)}); np.random.seed(0)  
>>> x = np.random.randn(100)  
>>> ax = sns.distplot(x,kde=False)
```

<http://stanford.io/1U7Uki3>

Histogram: Graph display of tabulated frequencies, shown as bars

Differences between histograms and bar charts:

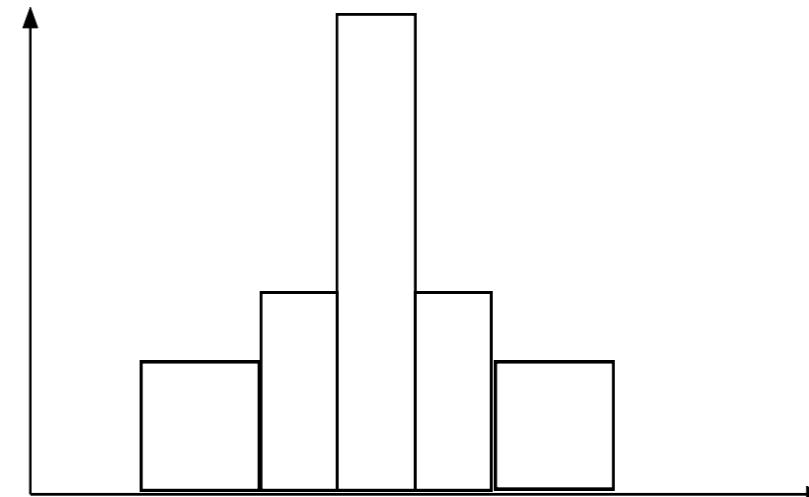
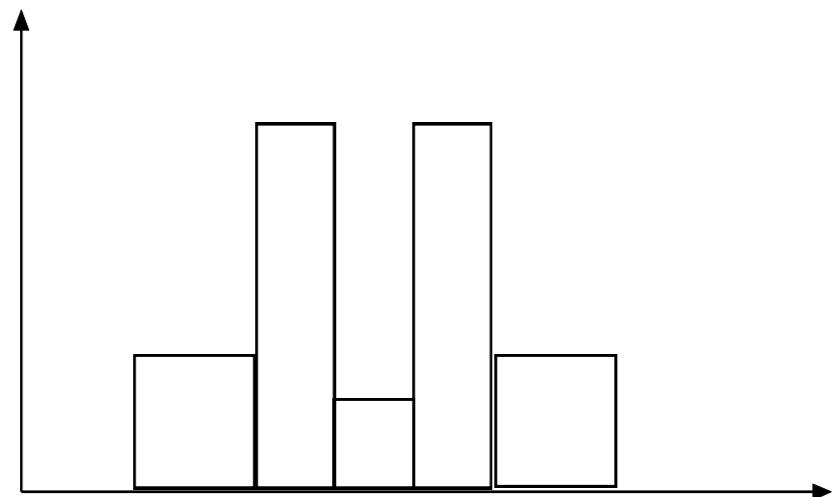
Histograms are used to show distributions of variables while bar charts are used to compare variables.

Histograms plot binned quantitative data while bar charts plot categorical data.

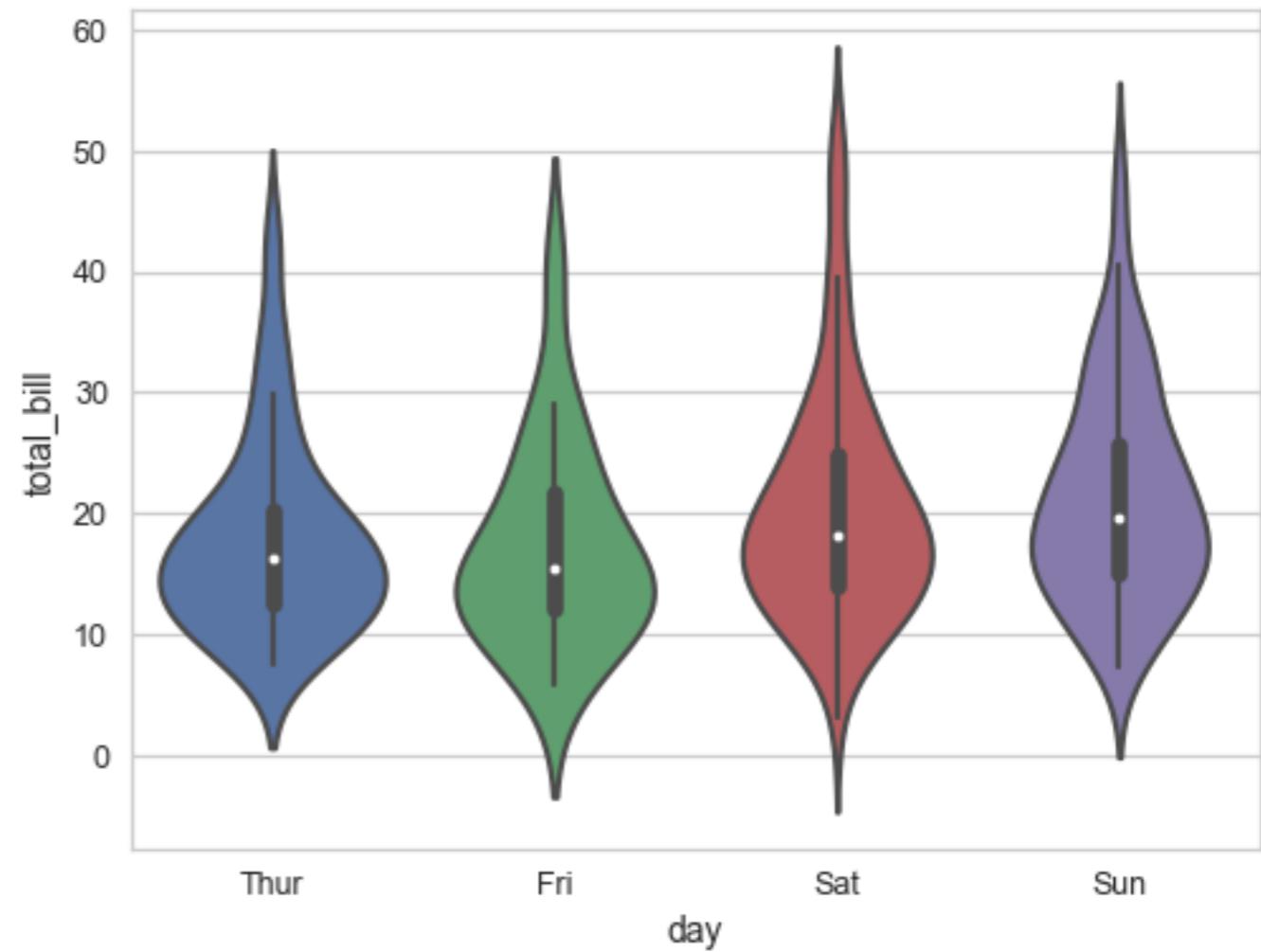
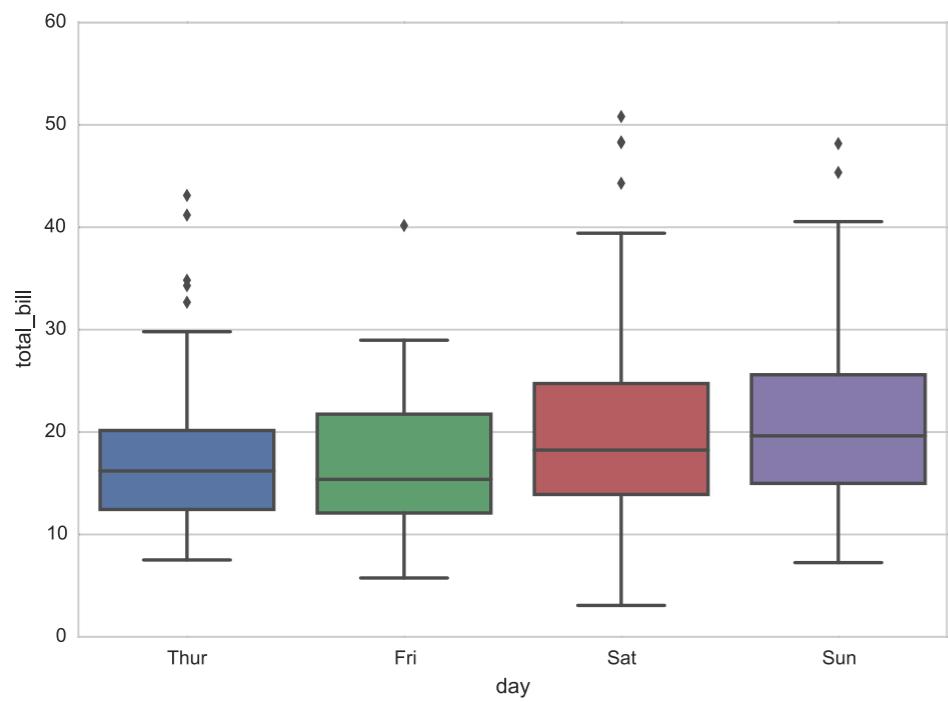
Bars can be reordered in bar charts but not in histograms.

Differs from a bar chart in that it is the **area** of the bar that denotes the value, not the height as in bar charts, a crucial distinction when the categories are not of uniform width

histograms tell more than box plots

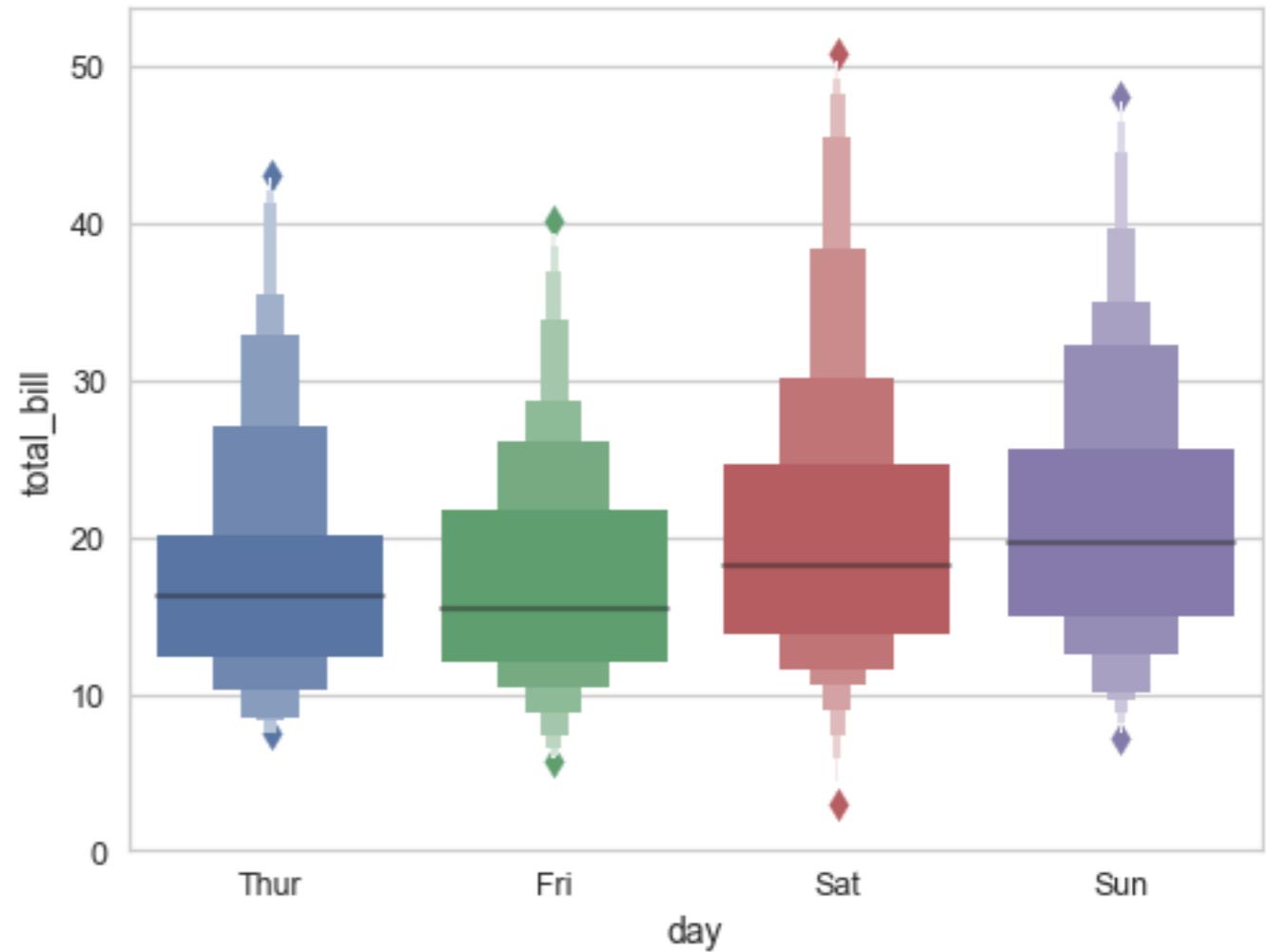
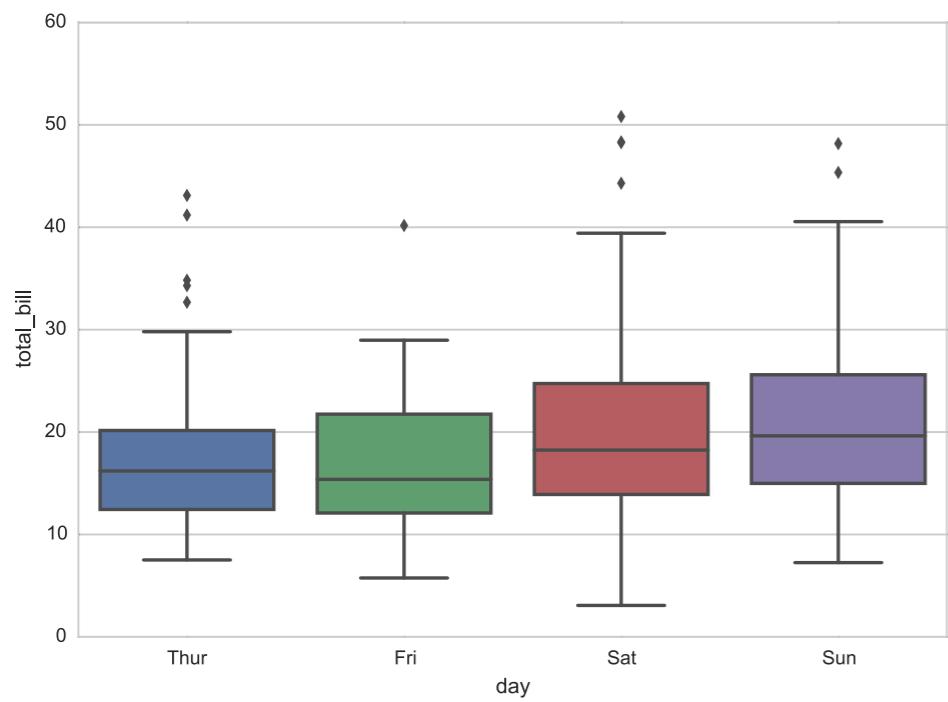


different distributions, with identical box plots



Violin Plots

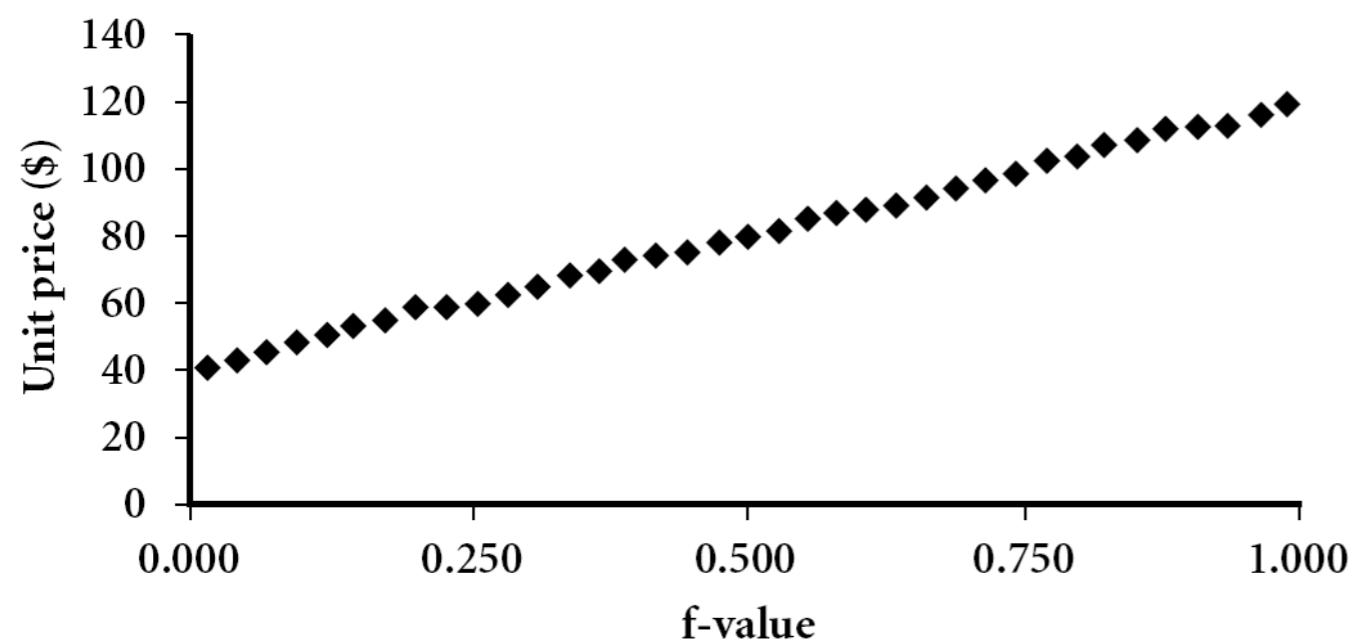
<http://bit.ly/2Dv5UFb>



Letter Value Plots

<http://bit.ly/2DtPRY3>

QUANTILES

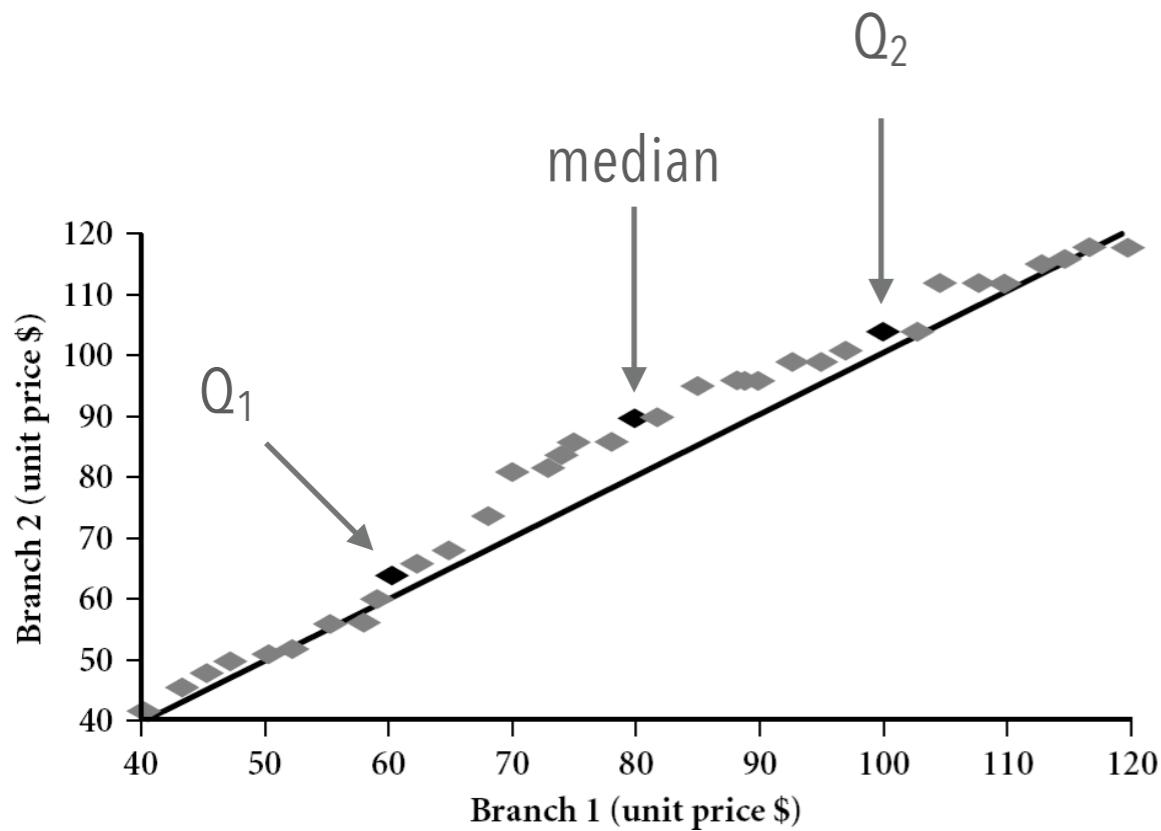


Displays all of the data
(allowing the user to assess
both the overall behavior and
unusual occurrences)

Plots quantile information

For a data x_i data sorted in
increasing order, f_i indicates
that approximately $100f_i\%$
of the data are **below or**
equal to the value x_i

Q-Q PLOT

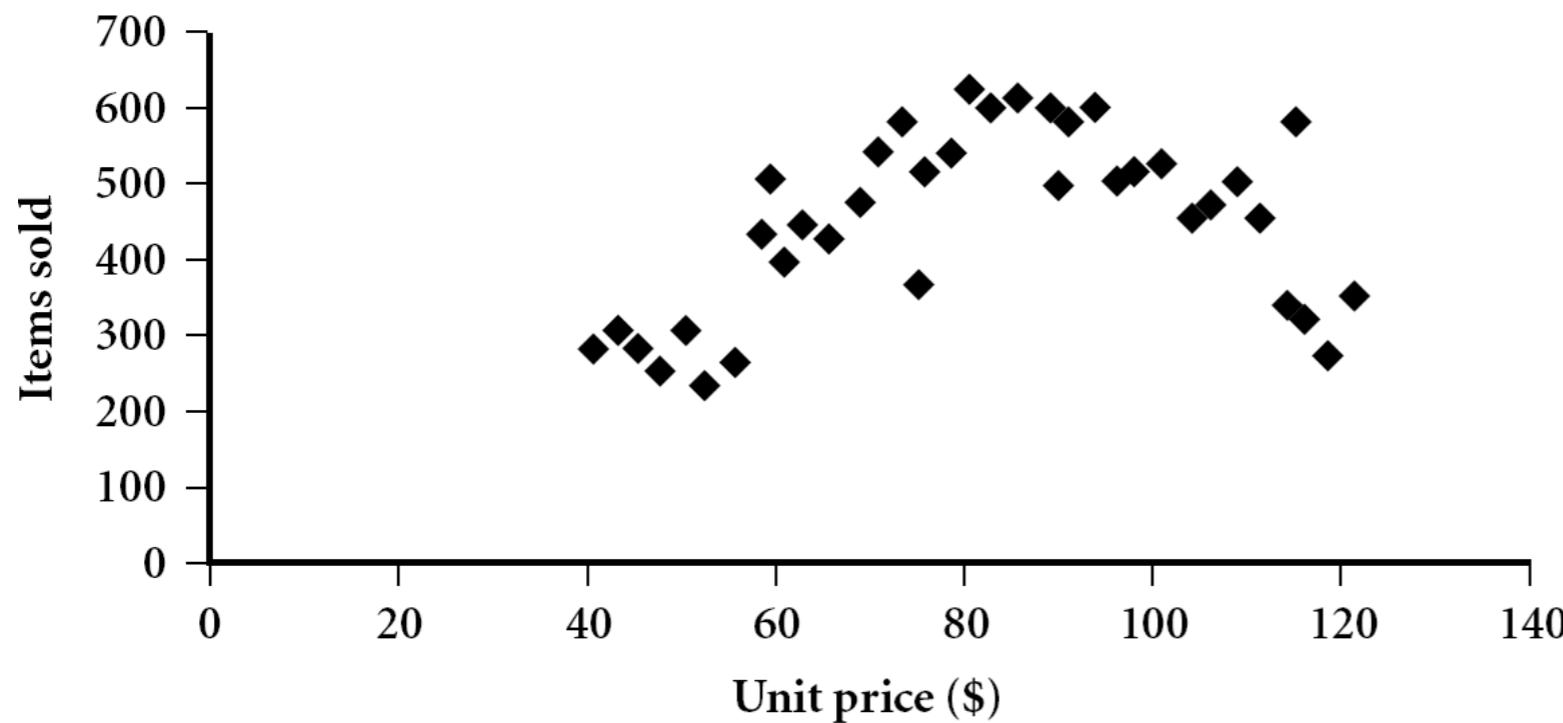


Graphs the quantiles of one univariate distribution against the corresponding quantiles of another

Is there is a shift in going from one distribution to another?

Example shows unit price of items sold at **Branch 1** vs. **Branch 2** for each quantile. Unit prices of items sold at **Branch 1** tend to be lower than those at **Branch 2**.

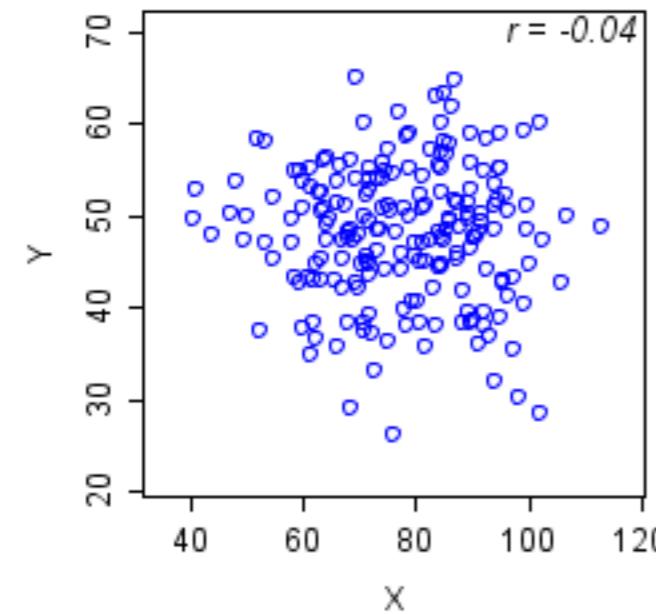
Each pair of values is treated as a pair of coordinates and plotted as points in the plane



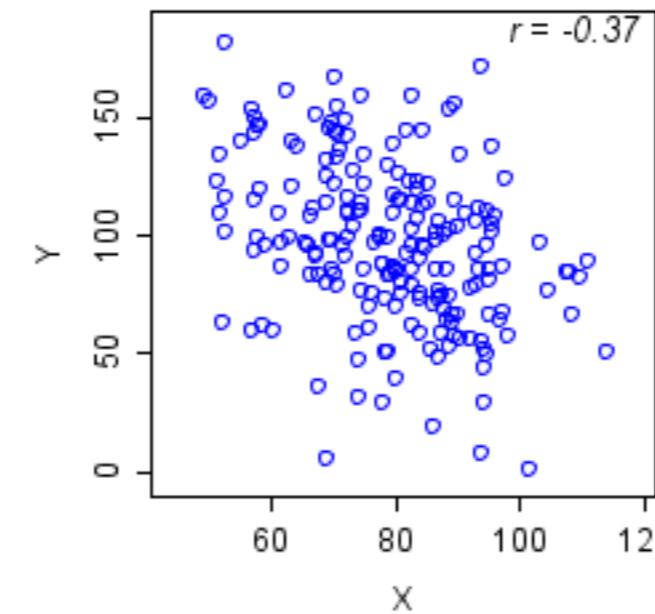
Scatter Plot

Provides a first look at bivariate data to see clusters of points, outliers, etc.

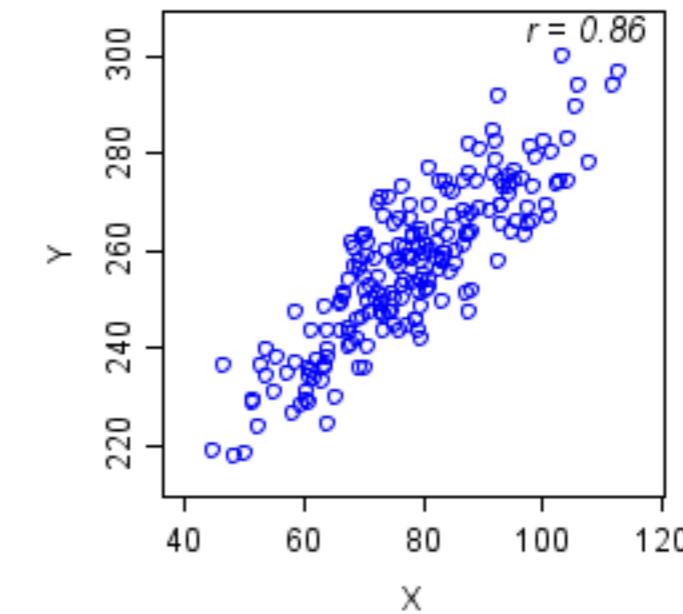
types of correlation



almost uncorrelated



negatively correlated



positively correlated

On writing reference letters

Extra Credit on Piazza

Students get points for
making substantial
contributions to answers
endorsed by TA

TA will look at questions **8** hours after being posted

Person with most
endorsements gets
5 extra credit points

Everyone else is **proportionally** scaled down

Arjun gets **5** points

Arjun has 10
endorsements

John gets **3** points

John has 6 endorsements

Off topic questions

**Gaming answers
will get you a 0**

TA will keep an excel sheet

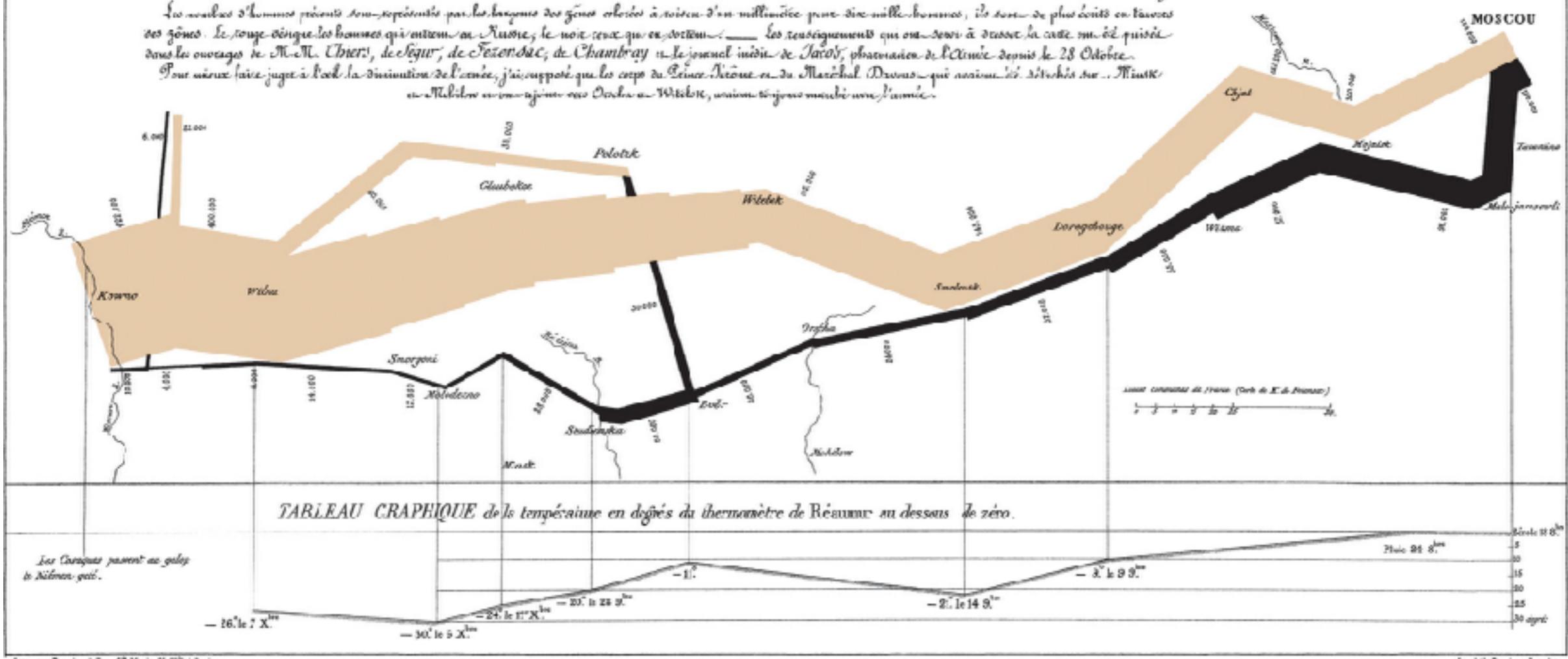
Carte Figurative des pertes successives en hommes de l'Armée Française dans la Campagne de Russie 1812-1813.

Dessiné par N. Minard, Inspecteur Général des Ponts et Chaussées en retraite.

Paris, le 20 Novembre 1869.

Si malais d'humain prédire sous-représenter par les langues des gènes colorés à soient d'un millionième pour six mille hommes, il sera de plus douteux de faire ces gènes. Le rouge évoque les hommes qui entrent en Russie; le noir ceux qui en sortent. — Les renseignements qui ont servi à dresser la carte n'ont été pris que dans les ouvrages de M. L. Ubier, de Segur, de Péreire, de Chambray et le journal indien de Jacob, pharmacien de l'Armée depuis le 28 Octobre.

Pour mieux faire juger à l'œil la diminution de l'armée, j'ai supposé que les corps du Prince Napoléon et du Maréchal Soult qui avaient été détruits sur la Moscova, étaient arrivés jusqu'à Orel et Witebsk, un peu moins qu'une moitié de l'armée.

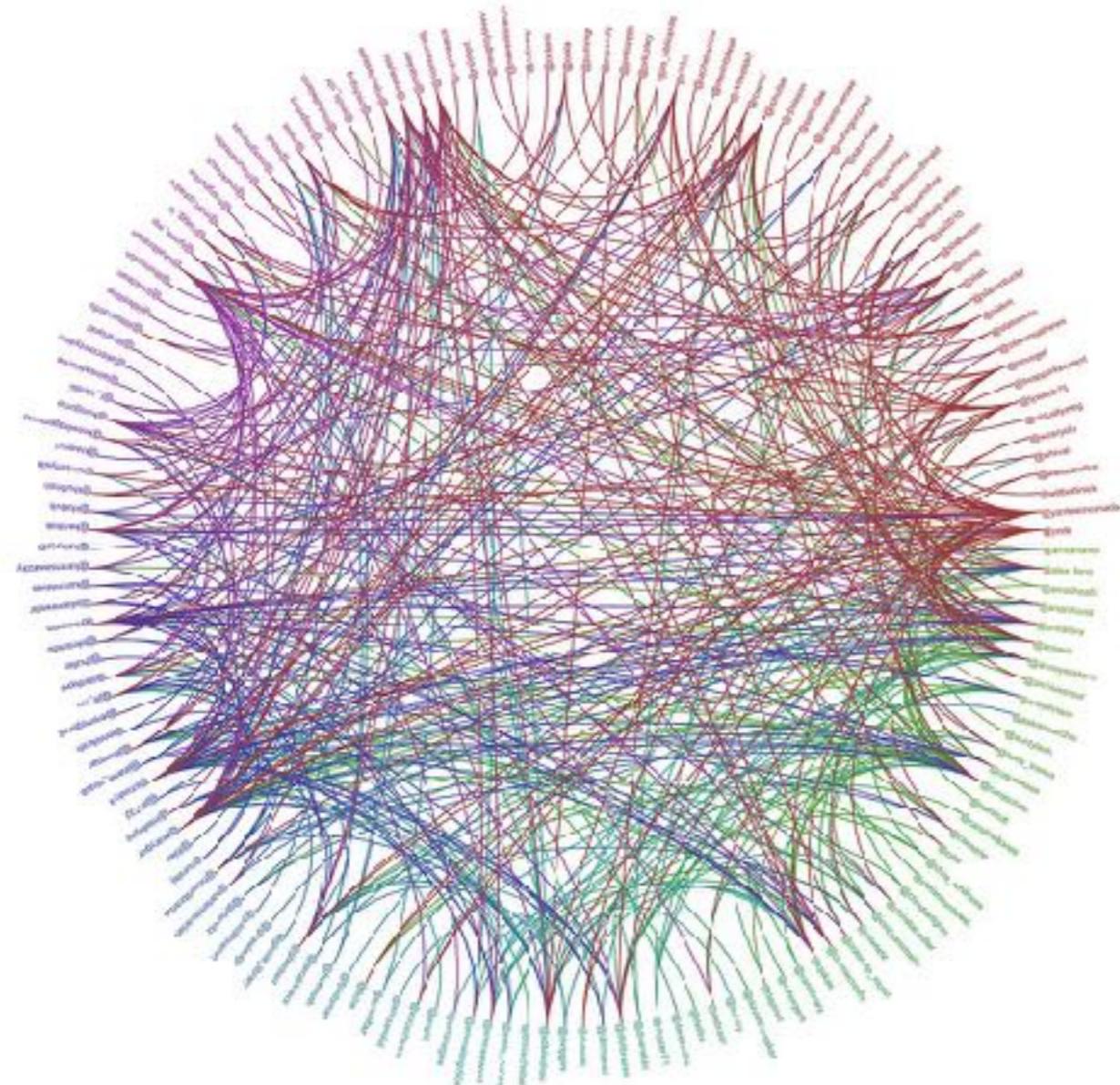


DATA VISUALIZATION

Minard's plot of Napoleon's march

WHY?

.....



Enron email corpus

Gain **insight** into an information space by mapping data onto graphical primitives

Provide qualitative **overview** of large data sets

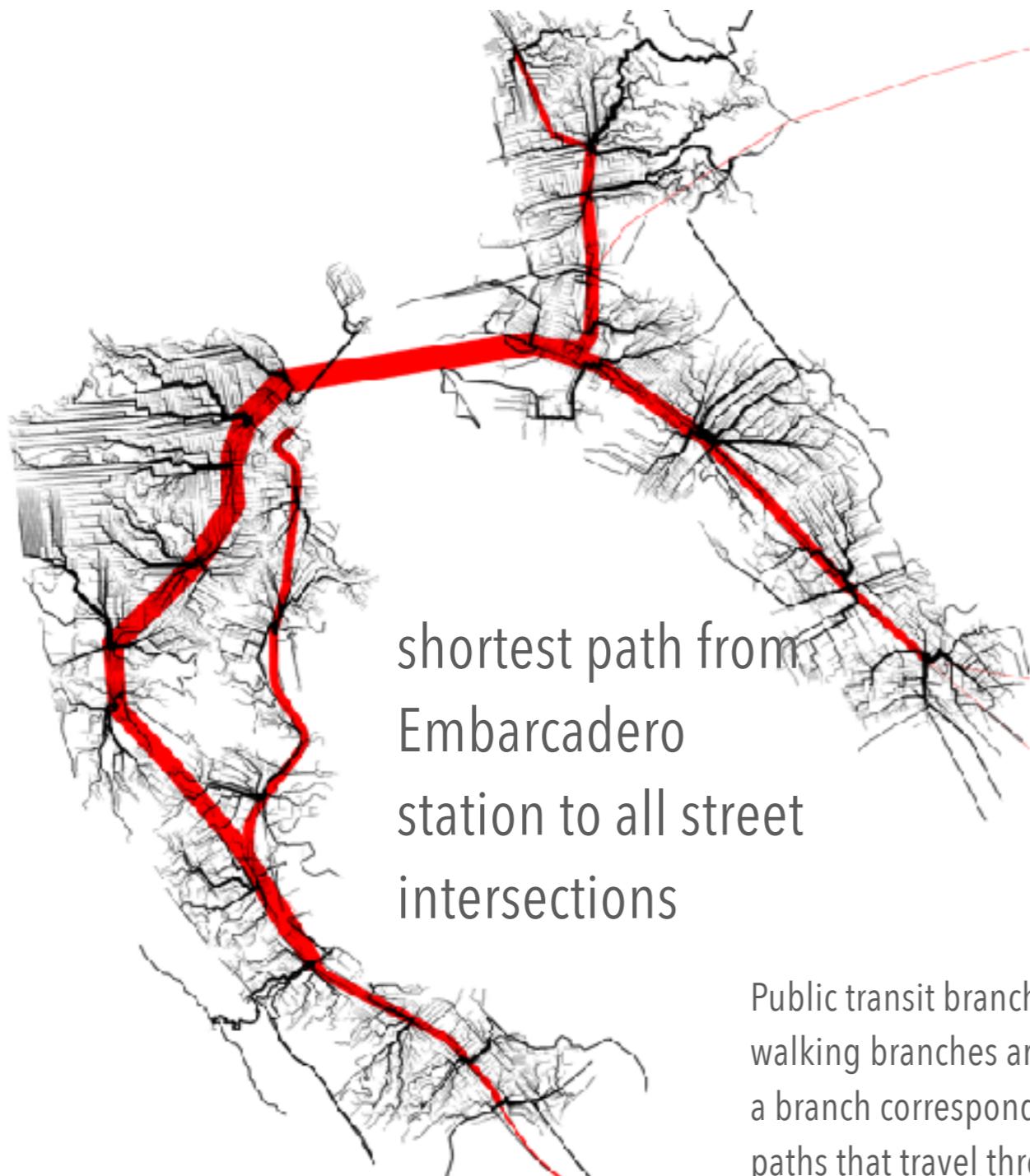
Search for patterns, trends, structure, irregularities, relationships among data

Help find **interesting regions** and suitable parameters for further quantitative analysis

Provide a **visual proof** of computer representations derived

WHY?

.....



This image shows the **shortest path** from Embarcadero Station on August 1st, 2008 11:43:33 AM local time to all street intersections in the San Francisco Bay area, assuming travel by public transit (namely, BART and Caltrain) and walking.

Public transit branches are colored red, while walking branches are colored black. The width of a branch corresponds to the number of shortest paths that travel through that branch. For example, the shortest route from Embarcadero to almost all destinations involves taking the BART. As a result, the BART lines are very thick.

Pixel-oriented

Hierarchical

Icon-based

Geometric

Visualization Categories

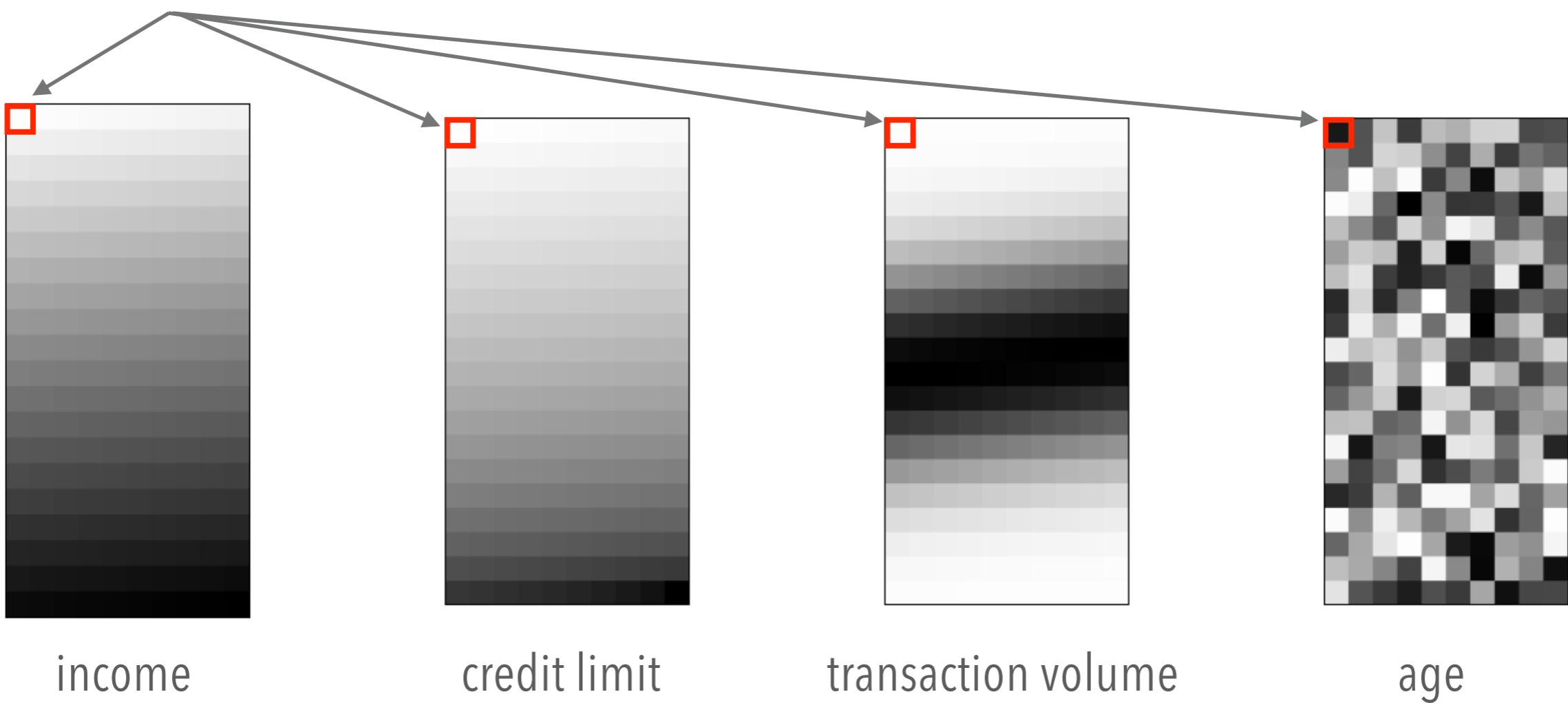
Visualizing complex data & relationships

pixel oriented visualization techniques

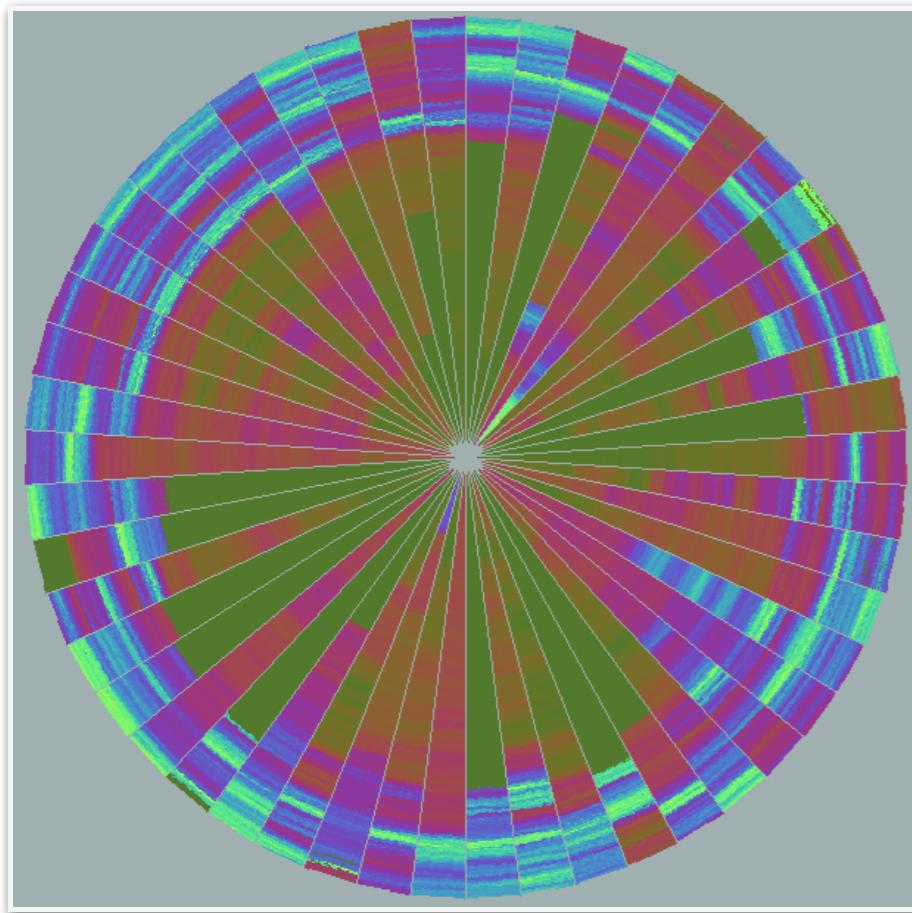
For a data set of m dimensions, create m windows on the screen, one for each dimension

The m dimension values of a record are mapped to m pixels at the corresponding positions in the windows

The colors of the pixels reflect the corresponding values

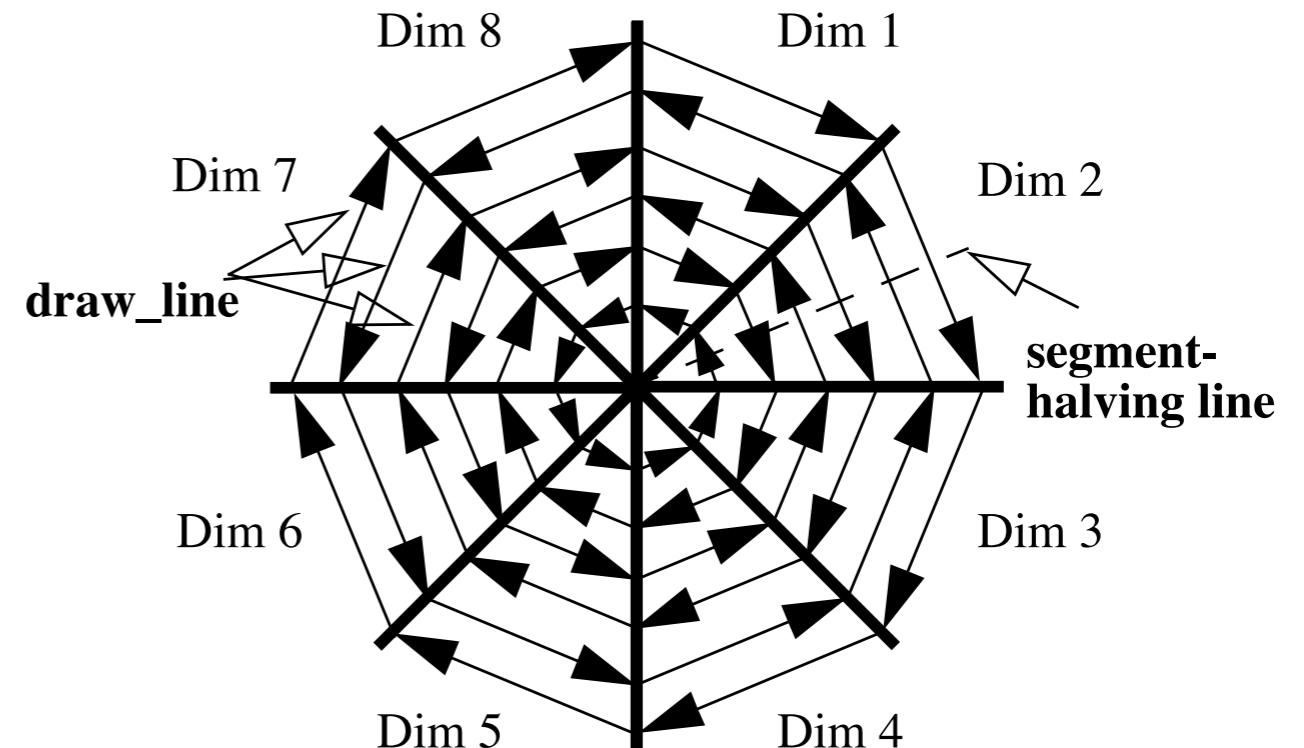


Space Filling Curves



Representing about 265,000 50-dimensional Data Items with the 'Circle Segments' Technique

M. Ankerst, D. A. Keim, and H.-P. Kriegel. **Circle segments: A technique for visually exploring large multidimensional data sets.** In Visualization, 1996.



To save space and show the connections among multiple dimensions, space filling is often done in a circle segment

Circular Segments

Direct visualization

Scatterplot matrices

Projection views

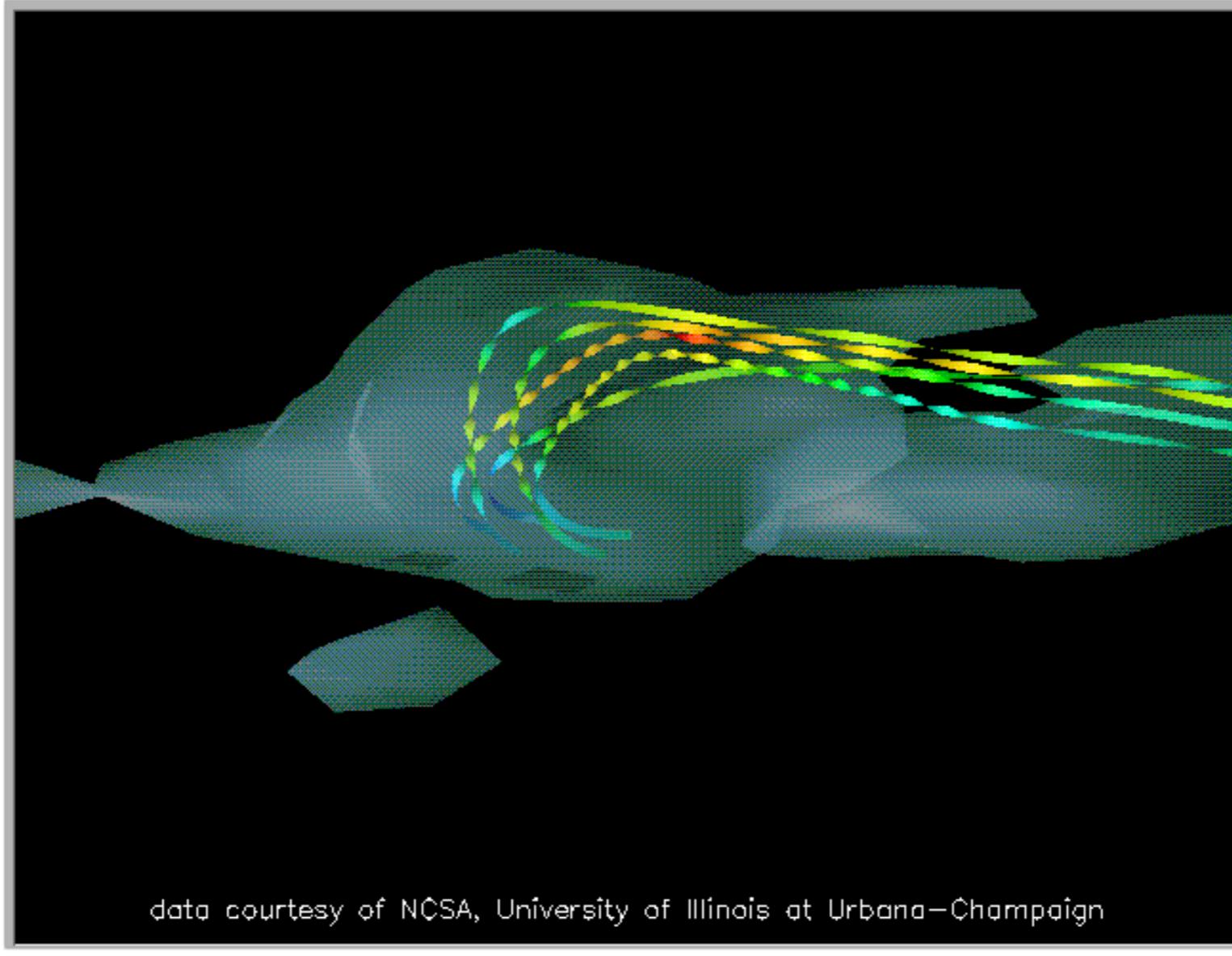
Landscapes

Geometric Projection Visualization

Hyperslice

Projection pursuit technique:
Help users find meaningful
projections of
multidimensional data

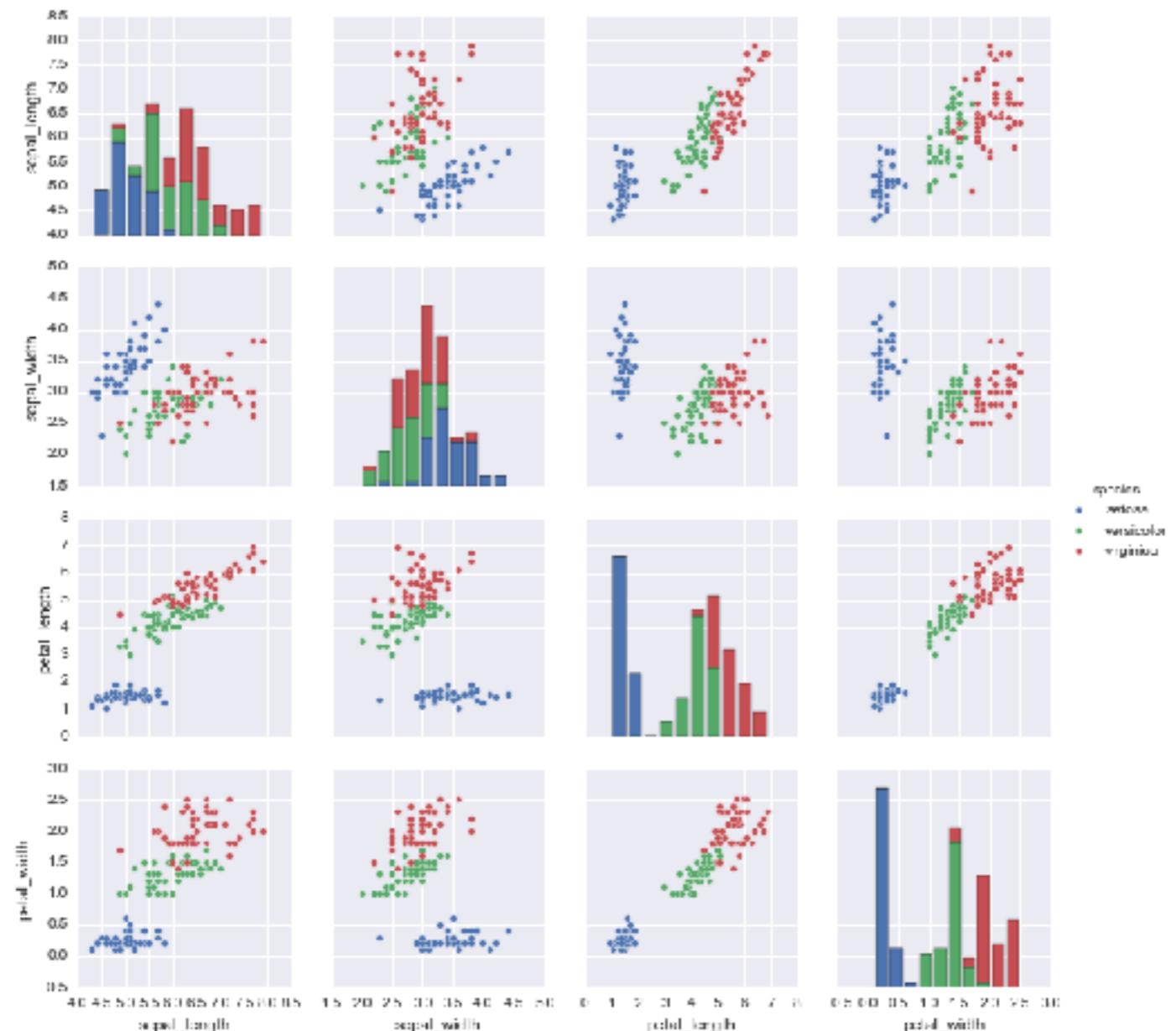
Parallel coordinates



data courtesy of NCSA, University of Illinois at Urbana-Champaign

direct data visualization

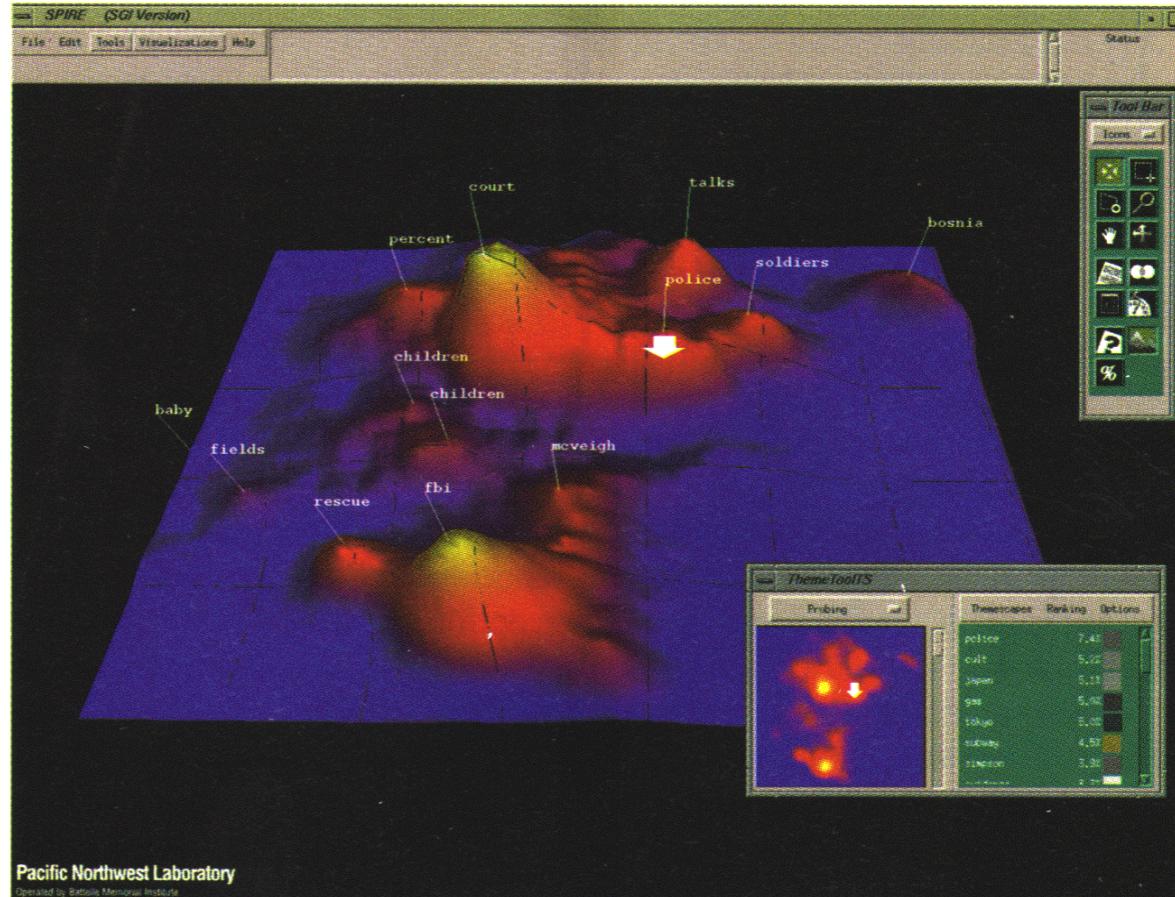
Matrix of scatterplots
(x-y-diagrams) of the k-
dim. data [total of
 $k*(k-1)/2$ scatterplots]



scatter plot matrices

Visualization of the data as perspective landscape

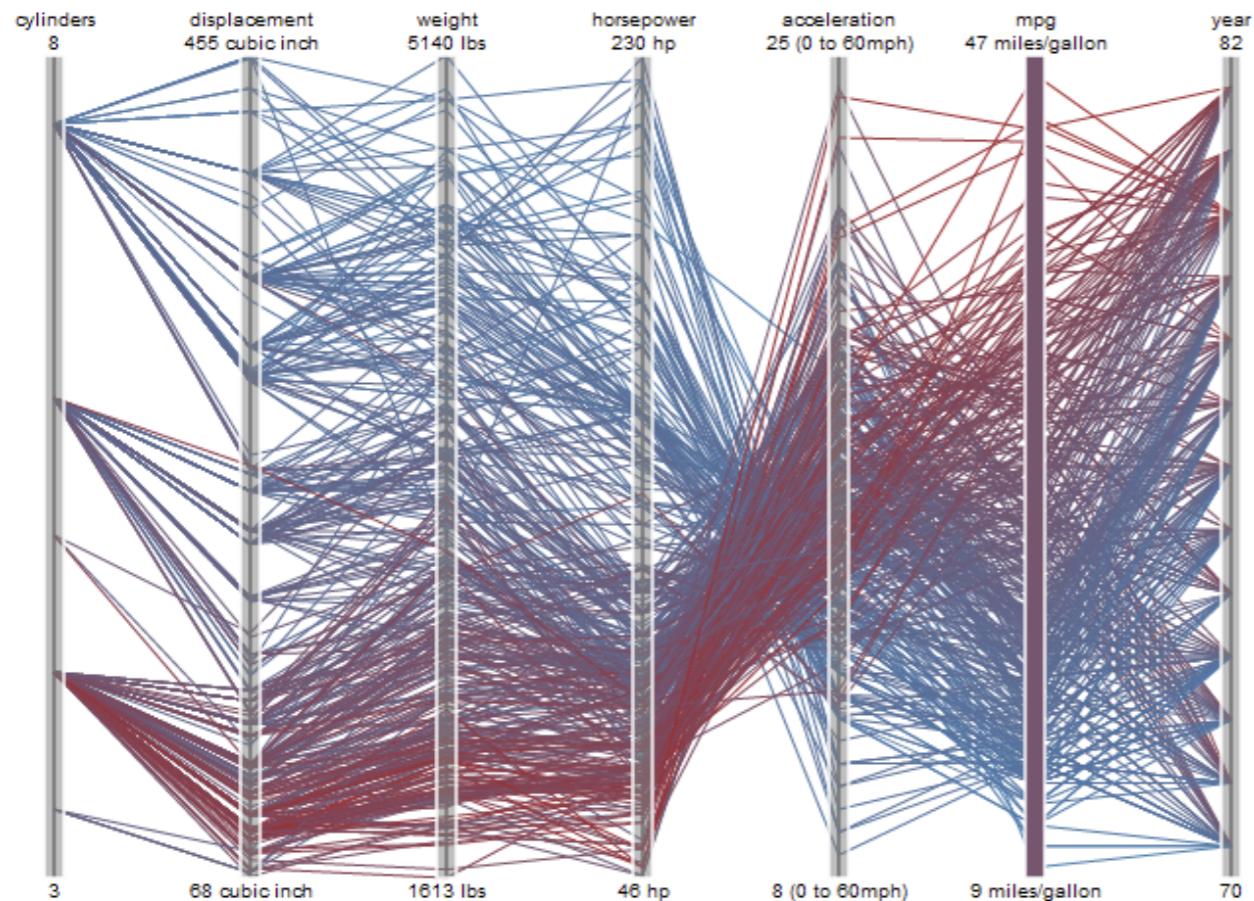
The data needs to be transformed into a (possibly artificial) 2D spatial representation which preserves the characteristics of the data



news articles visualized as a landscape

landscapes

PARALLEL COORDINATES



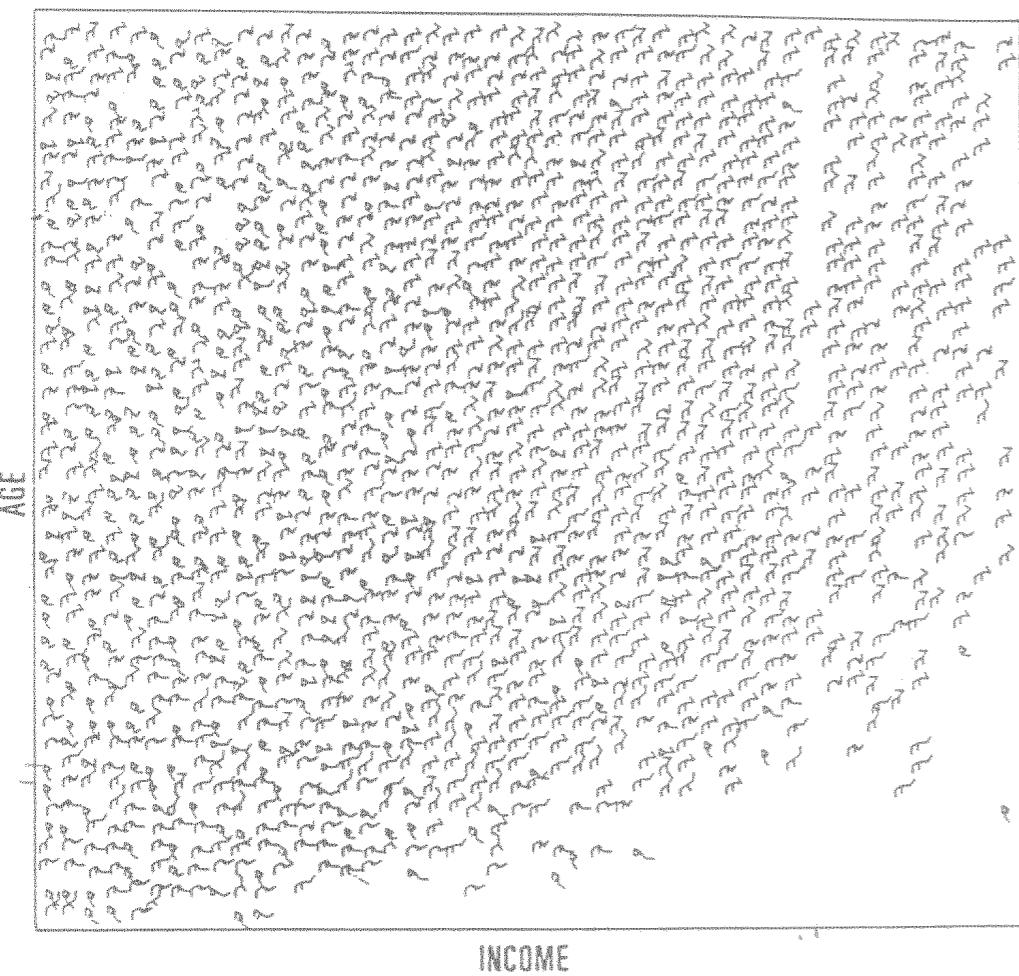
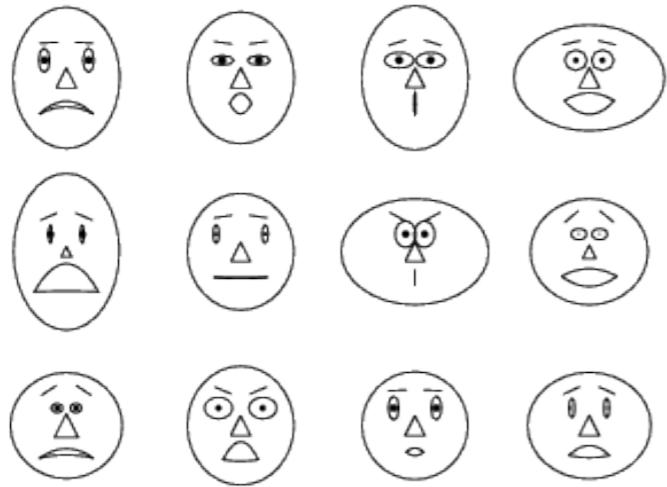
n equidistant axes which are parallel to one of the screen axes and correspond to the attributes

The axes are scaled to the [minimum, maximum]: range of the corresponding attribute

Every data item corresponds to a polygonal line which intersects each of the axes at the point which corresponds to the value for the attribute

ICONS

.....



Visualization of the data values as features of icons

Typical visualization methods

Chernoff Faces

Stick Figures

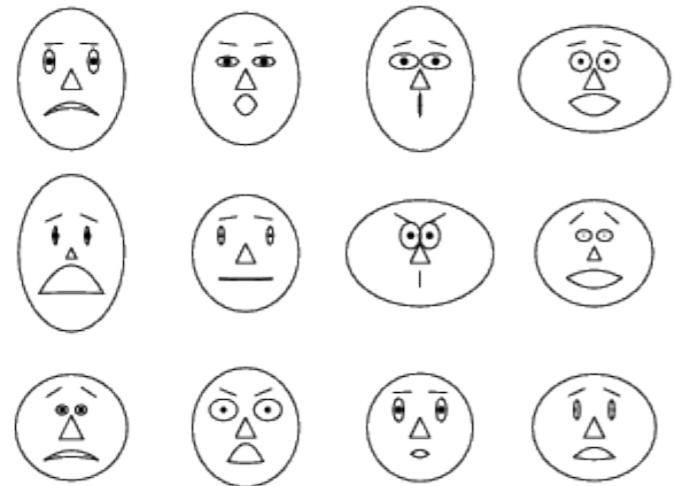
General techniques

Shape coding: Use shape to represent certain information encoding

Color icons: Use color icons to encode more information

Tile bars: Use small icons to represent the relevant feature vectors in document retrieval

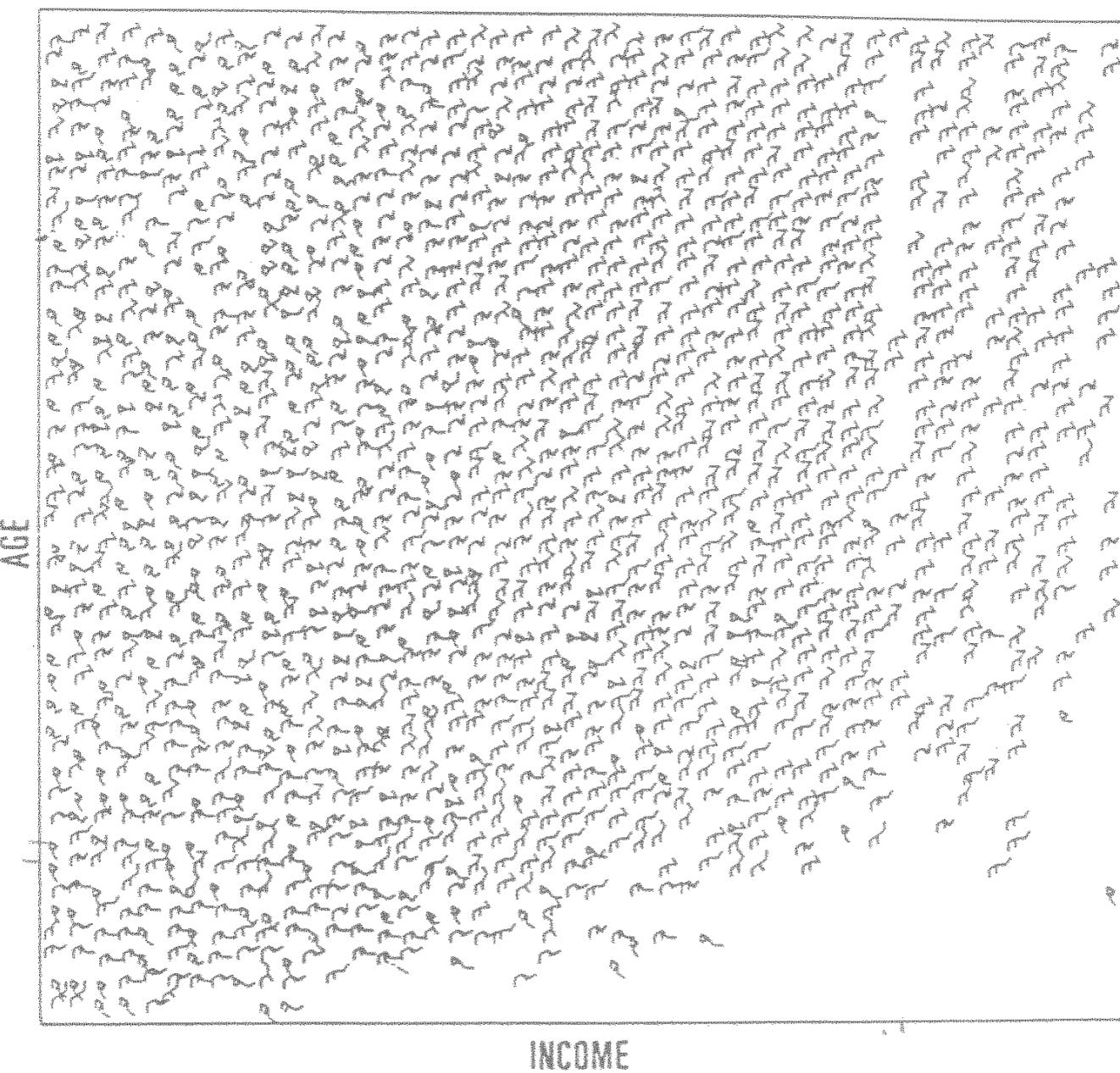
CHERNOFF FACES



A way to display variables on a two-dimensional surface, e.g., let **x** be eyebrow slant, **y** be eye size, **z** be nose length, etc.

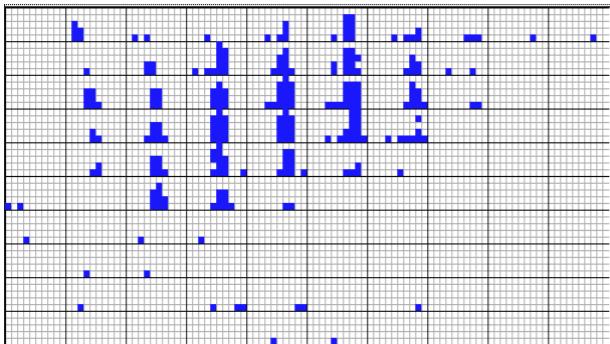
The figure shows faces produced using 10 characteristics—head eccentricity, eye size, eye spacing, eye eccentricity, pupil size, eyebrow slant, nose size, mouth shape, mouth size, and mouth opening; Each assigned one of 10 possible values

stick figures



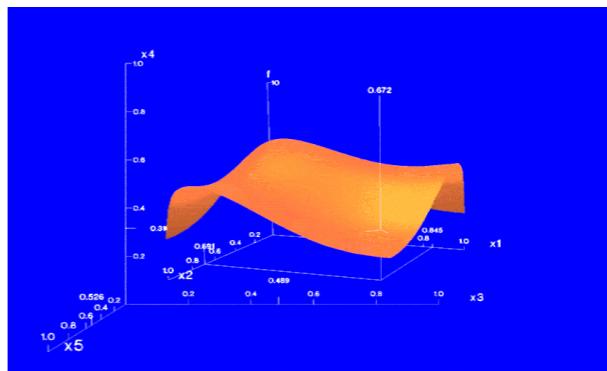
A census data figure
showing age, income,
gender, education, etc.

A 5-piece stick
figure (1 body and
4 limbs w. different
angle/length)



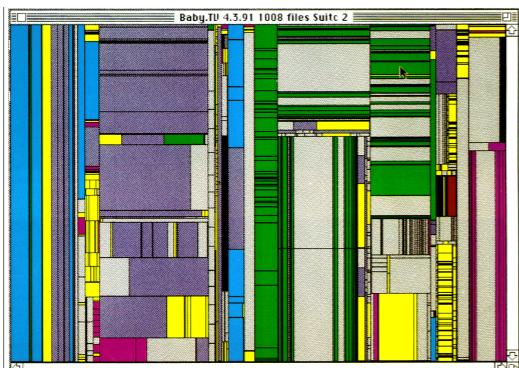
Dimensional Stacking

Worlds-within-Worlds

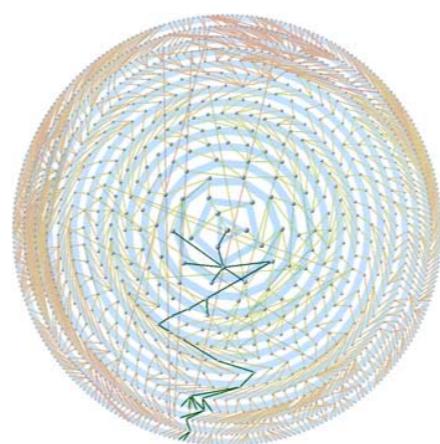


hierarchical visualization

Tree-Map



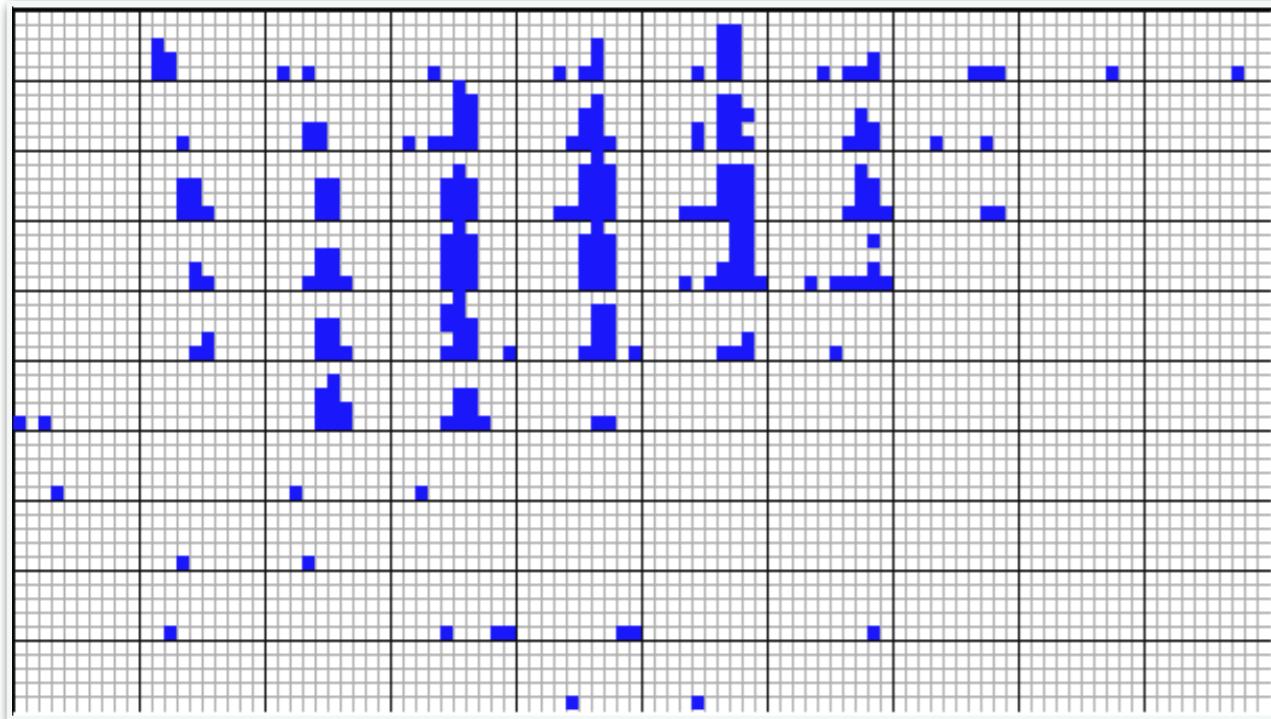
Cone Trees



InfoCube

DIMENSIONAL STACKING

.....



Visualization of oil mining data with longitude and latitude mapped to the outer **x**, **y**-axes and ore grade and depth mapped to the inner **x**, **y**-axes

Partitioning of the **n**-dimensional attribute space in 2-D subspaces, which are ‘stacked’ into each other

Partitioning of the attribute value ranges into classes. The important attributes should be used on the outer levels.

Adequate for data with ordinal attributes of low cardinality

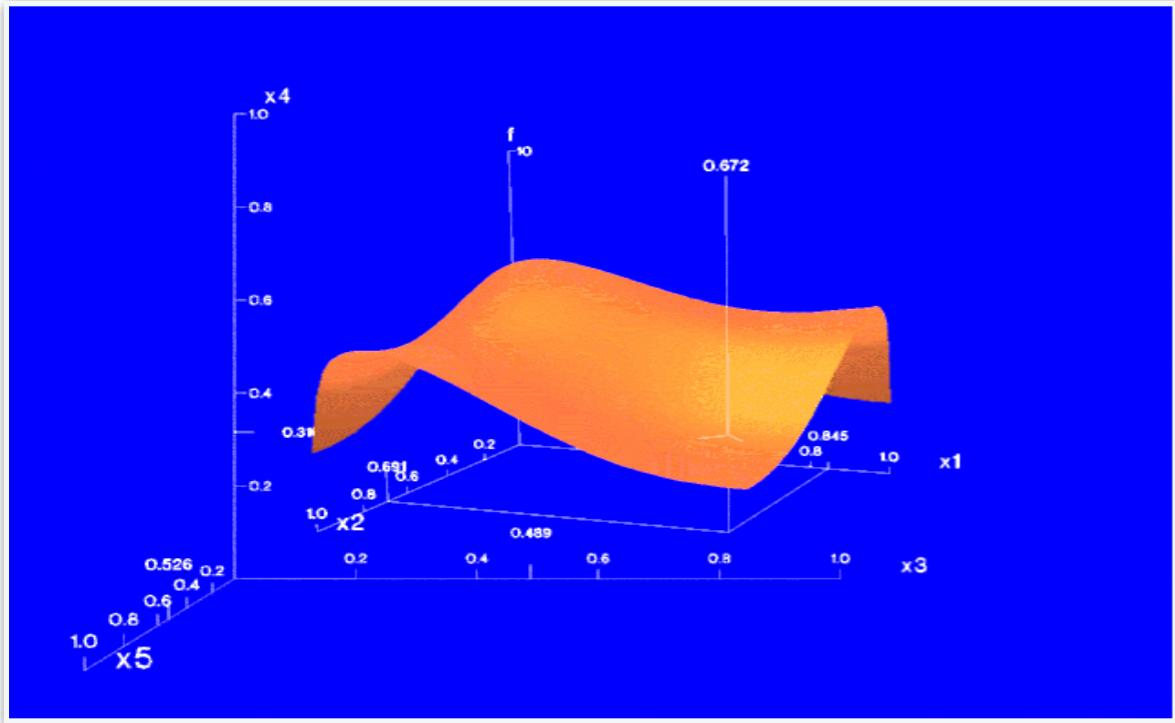
But, difficult to display more than nine dimensions

Important to map dimensions appropriately

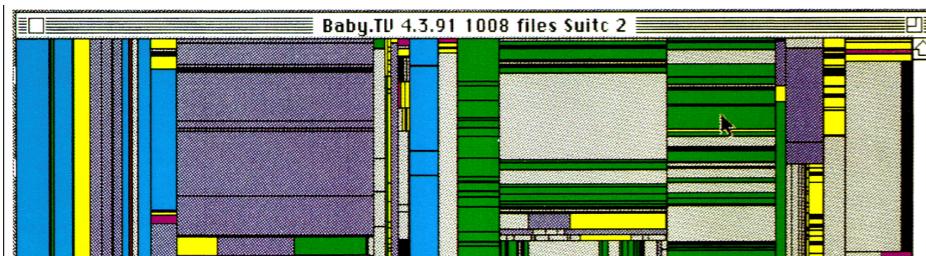
WORLDS WITHIN WORLDS

Assign the function and two most important parameters to innermost world

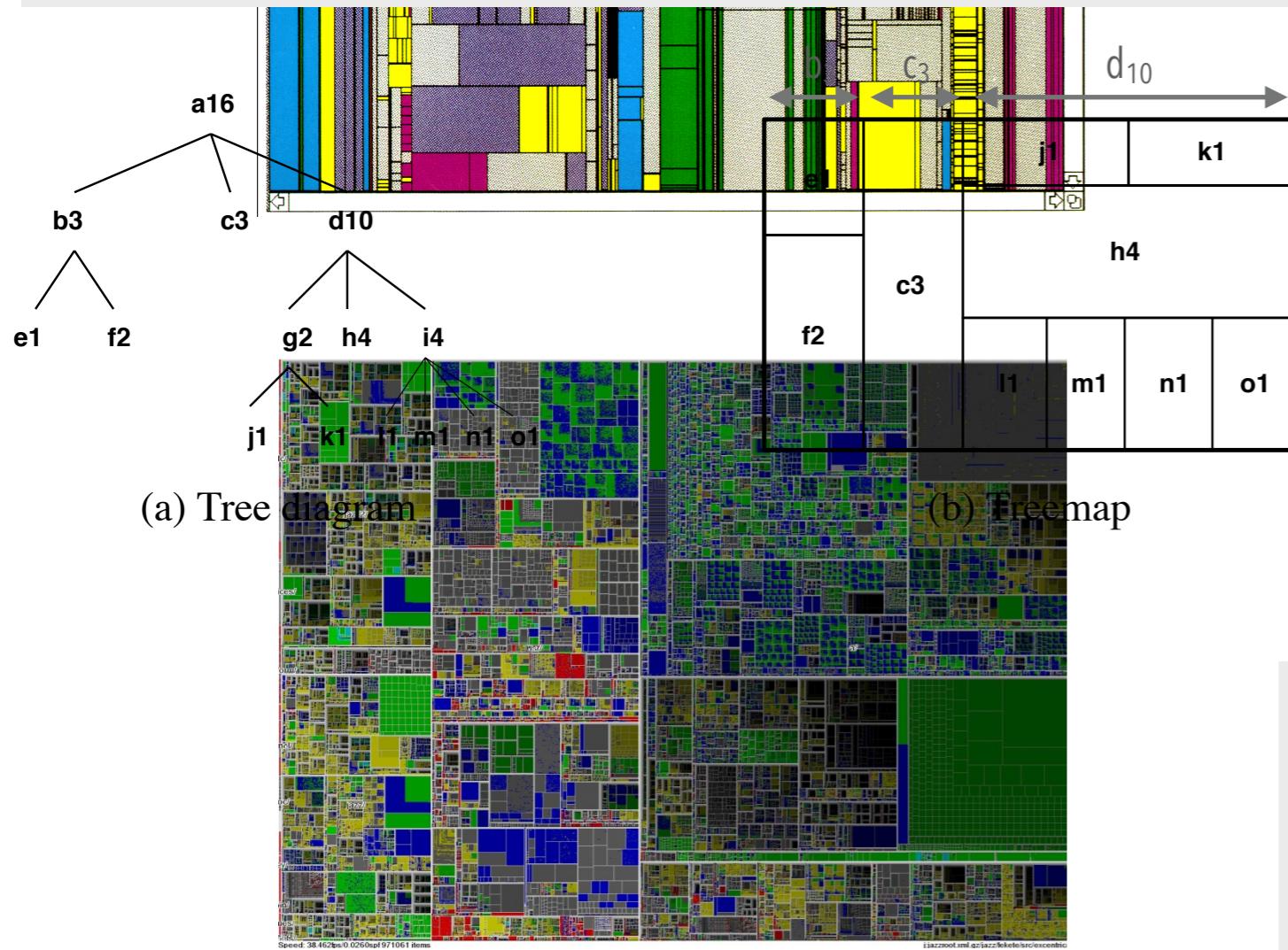
Fix all other parameters at constant values—draw other (1 or 2 or 3 dimensional worlds choosing these as the axes)



treemap of a file system



Bruls, M., Huizing, K., & Van Wijk, J. J. (2000). **Squareified treemaps**. In Data Visualization 2000 (pp. 33-42). Springer Vienna.



TREEMAP

.....

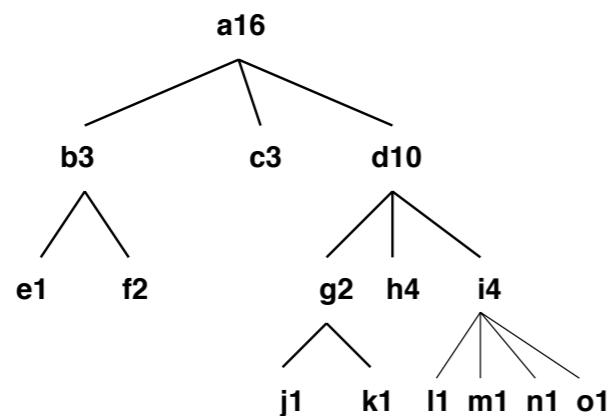
Screen-filling method which uses a hierarchical partitioning of the screen into regions depending on the attribute values

The **x**- and **y**-dimension of the screen are partitioned alternately according to the attribute values (classes)

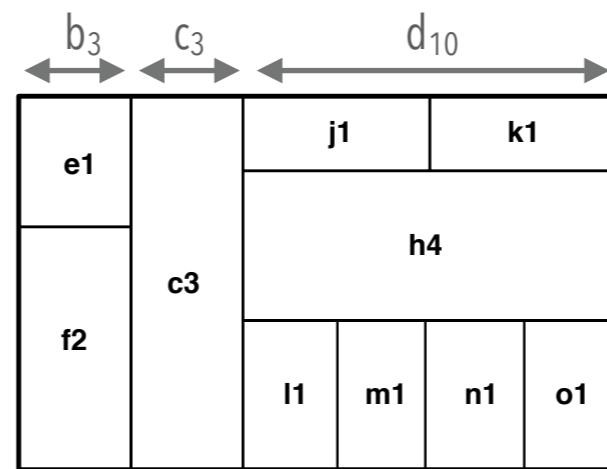
Martin Wattenberg. **Visualizing the stock market**. In CHI'99 extended abstracts on Human factors in computing systems, pages 188-189. ACM, 1999.

<http://newsmap.jp>

Bruls, M., Huizing, K., & Van Wijk, J. J. (2000). **Squareified treemaps**. In Data Visualization 2000 (pp. 33-42). Springer Vienna.

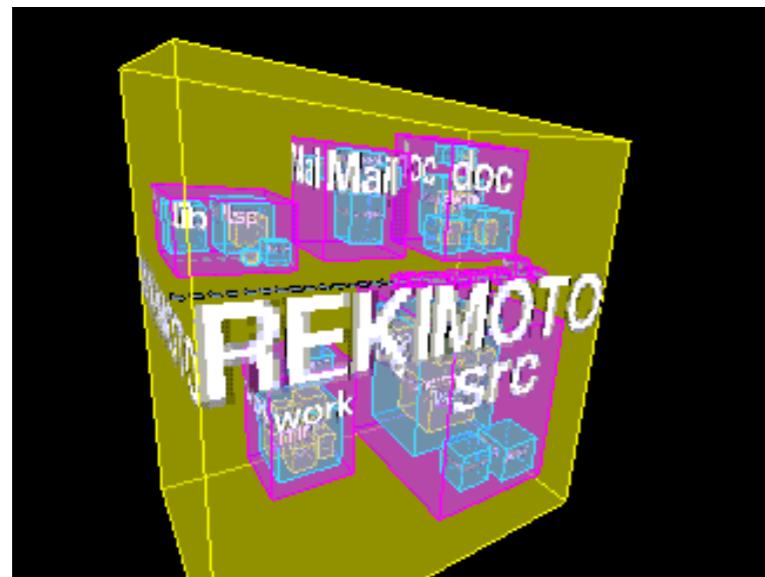


(a) Tree diagram



(b) Treemap

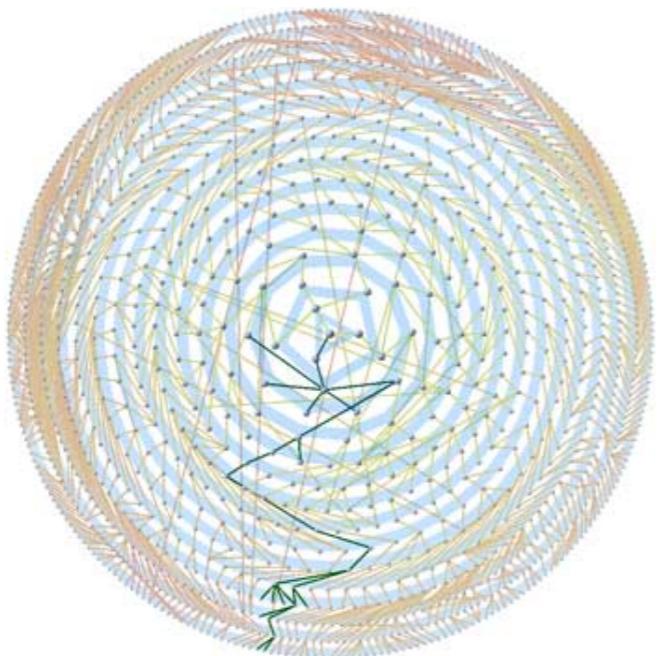
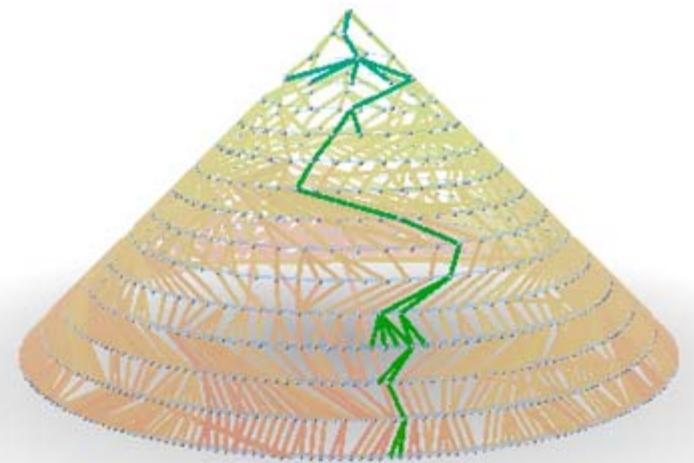
INFOCUBE



A 3-D visualization technique where hierarchical information is displayed as nested semi-transparent cubes

The outermost cubes correspond to the top level data, while the sub nodes or the lower level data are represented as smaller cubes inside the outermost cubes, and so on

3D CONE TREES



Visualize a social network data set that models the way an infection spreads from one person to the next

3D cone tree visualization technique works well for data with **hierarchical structure** for up to a thousand nodes or so

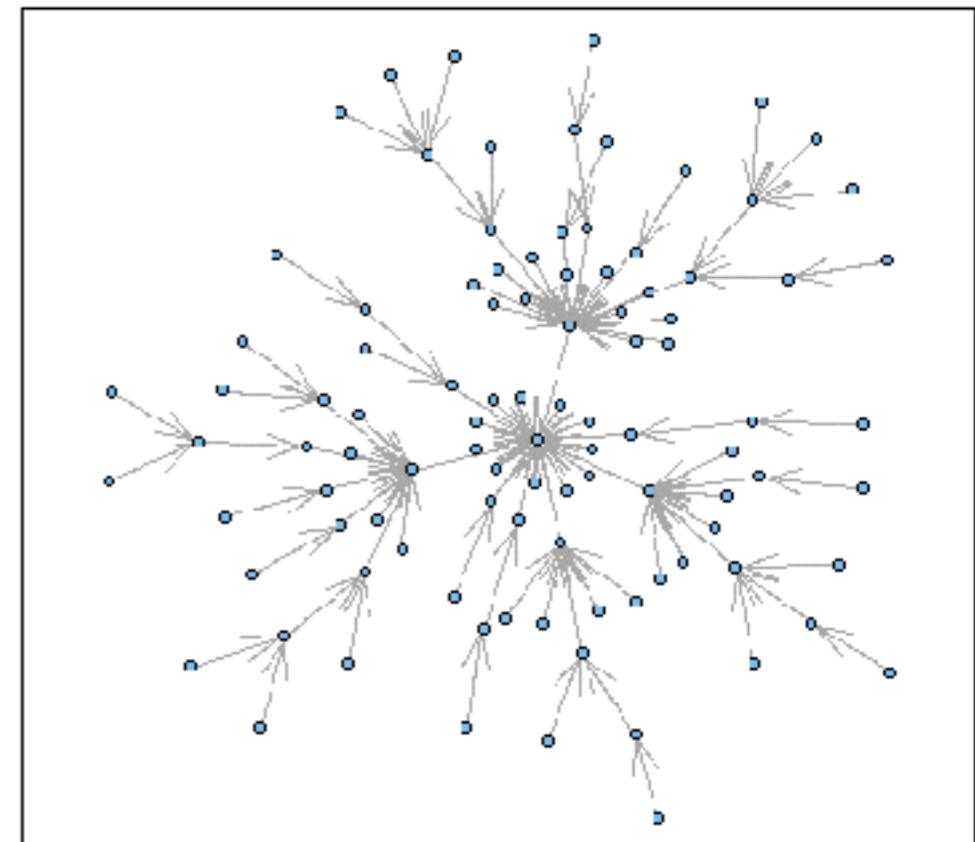
First build a 2D circle tree that arranges its nodes in concentric circles centered on the root node

Cannot avoid overlaps when projected to 2D

perspective
object
group
surface
size
spatial area
line pattern
distance
separation gradient
location
size
order
line
center
shape
sequence
overlap connectivity
field
contrast
place
center
space
magnitude
proximity adjacency network units
symmetry distribution
boundary
interpolation
interaction
representation
position
scale force
path neighborhood

tag clouds

non-
numerical
data

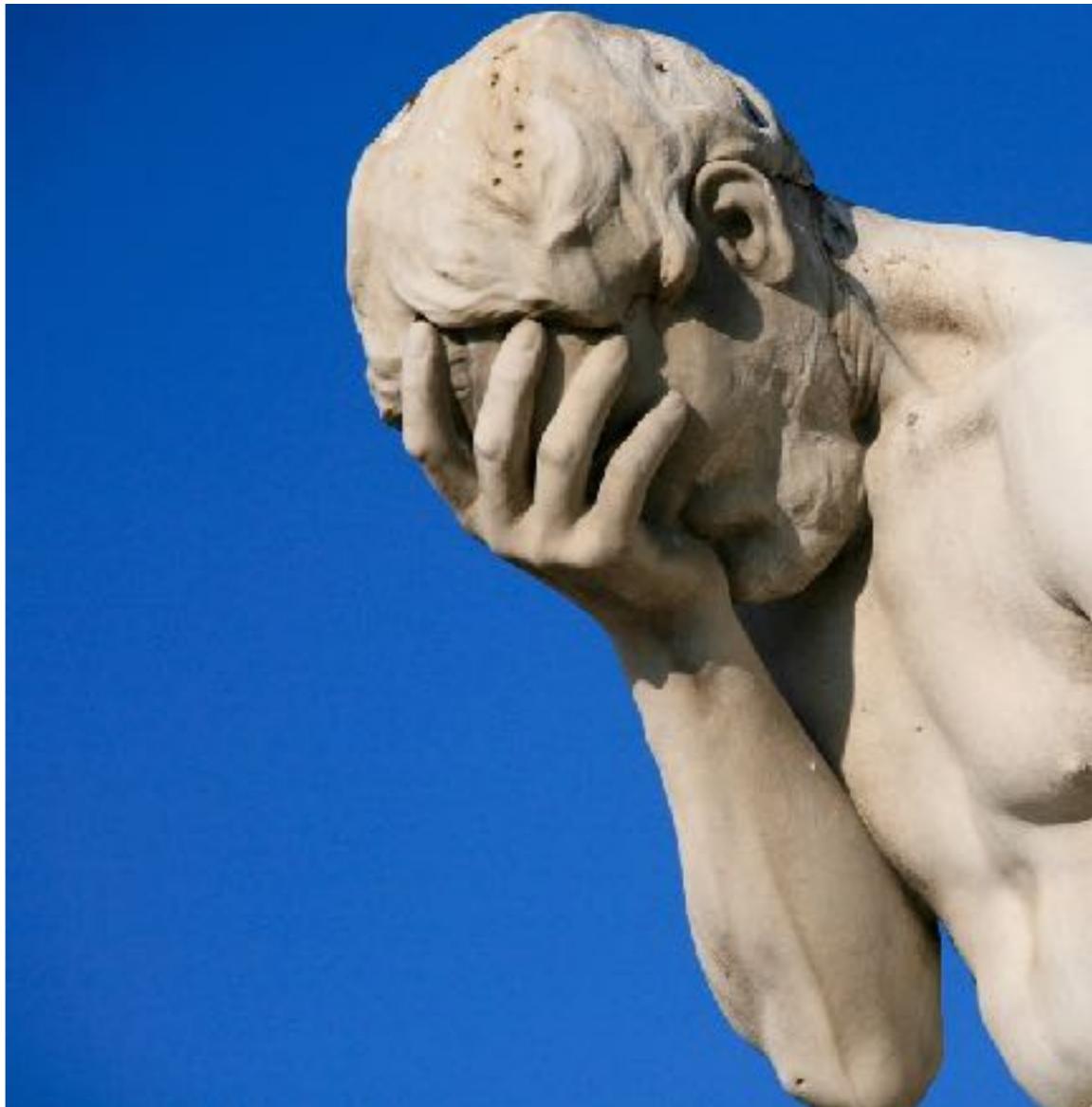


many network visualizations

WHAT CAN YOU DO RIGHT NOW?

.....

tools and code

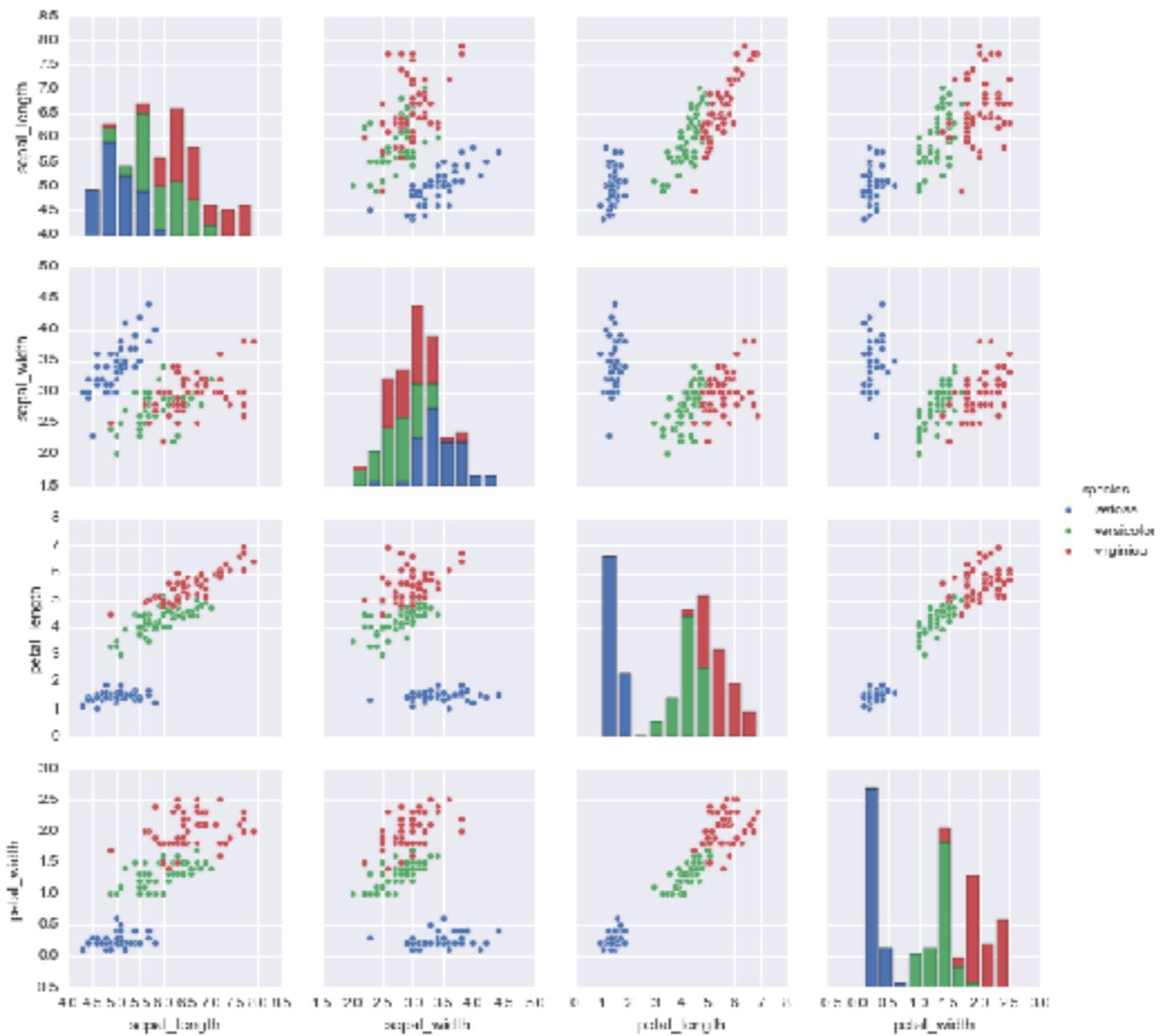




<http://www.tableau.com>

seaborn

<http://stanford.io/1PhonDZ>



Radial Reingold–Tilford Tree

<http://d3js.org>

<http://bokeh.pydata.org/en/latest/>



The `tree` layout implements the Reingold-Tilford algorithm for efficient, tidy arrangement of layered nodes. The depth of nodes is computed by distance from the root, leading to a ragged appearance. Cartesian orientations are also supported. Implementation based on work by [Jeff Heer](#) and [Jason Davies](#) using [Buchheim et al.](#)'s linear-time variant of the Reingold-Tilford algorithm. Data shows the [Flare class hierarchy](#), also courtesy [Jeff Heer](#).

[Open in a new window](#)

[Compare to this Cartesian layout.](#)

SeeDB: Visualizing Database Queries Efficiently

[Vision Paper]

Aditya Parameswaran
Stanford University
adityagp@cs.stanford.edu

Neoklis Polyzotis
Google & UCSC
alkis@cs.ucsc.edu

Hector Garcia-Molina
Stanford University
hector@cs.stanford.edu

ABSTRACT

Data scientists rely on visualizations to interpret the data returned by queries, but finding the right visualization remains a manual task that is often laborious. We propose a DBMS that partially automates the task of finding the right visualizations for a query. In a nutshell, given an input query Q , the new DBMS optimizer will explore not only the space of physical plans for Q , but also the space of possible visualizations for the results of Q . The output will comprise a recommendation of potentially “interesting” or “useful” visualizations, where each visualization is coupled with a suitable query execution plan. We discuss the technical challenges in building this system and outline an agenda for future research.

1. INTRODUCTION

... today's researchers must consume ever higher volumes of numbers that gush, as if from a fire hose ...

— R.M. Friedhoff and T. Kiely

Data analysts must sift through huge volumes of data looking for valuable data-specific insights, trends, or anomalies. This process involves selecting the “right” subset of the data, and the “right” way to view it, so that the “insights” become apparent. Moreover, the process is often ad-hoc and consumes a lot of the analyst’s time. Our vision is that *some especially cumbersome aspects* of this search for interesting insights can be automated.

To illustrate, consider the following interactive exploration workflow, which we believe is often used in practice.

Step (1): First, the analyst poses a relational query to extract some subset of data they are interested in exploring. For example, the analyst may select all records associated with “stapler” products.

Step (2): Then, the analyst considers several candidate views over this subset of data, formed by, say, aggregation and grouping; the analyst must study all of these views one by one. For example, one view may be total stapler sales by year, while another view may be the quantity in stock by sales region. Since these views have two-attributes each, we can view them as 2-dimensional graphs. For example, Figure 1(a) may be the stapler sales (y-axis) by year (x-axis), while Figure 1(c) may be the quantity (y-axis) by region (x-axis). (Figures 1(b) and 1(d) are discussed below.)

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Articles from this volume were invited to present their results at The 40th International Conference on Very Large Data Bases, September 1st - 5th 2014, Hangzhou, China.

Proceedings of the VLDB Endowment, Vol. 7, No. X
Copyright 2013 VLDB Endowment 2150-8097/13/09... \$ 10.00.

Step (3): Next, the analyst steps through each view, and decides which views are “interesting.” This of course is the critical and time-consuming step. What makes a view like Figure 1(a) interesting or not? Well, it all depends on the application semantics and what we are comparing against. For example, Figure 1(a) shows decreasing sales over time. If we are in a recession and all product sales are down then this observation is not very interesting. However, say that Figure 1(b) shows the aggregate (all) product sales over the same time periods. Then the stapler sales view goes against the general trend: overall sales are up, but stapler sales are going down. In this case, the view is *potentially* “interesting” because it depicts a trend in the subset of data that the analyst is interested in (i.e., stapler-related data) that *deviates* from the trend in the overall data. Of course, the analyst must decide if this deviation is truly an insight for this application. Even so, our key insight is that we may be able to *identify and highlight to the analyst potentially interesting views using automated mechanisms based on deviation*. By doing so, we eliminate the laborious process of stepping through all possible views that the analyst currently performs. Once we recommend potentially interesting views, we can let the analyst make the final decision.

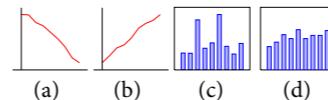


Figure 1: Views (a), (b): Sales over Time. (c), (d): Quantity by Region.

Figures 1(c) and 1(d) illustrate a different type of deviation. The first figure shows the distribution of staplers across regions, while the second figure shows the overall product distribution. Again, the stapler view does not follow the general trend: the regions that have the most staplers are not the larger regions that have most product in stock. The analysis must decide if this observation is interesting: perhaps the region that has many staplers is near the world-famous stapler-gun-wrestling contest, in which case the observation is expected. But perhaps there is a problem with the product shipping strategy, in which case the deviation is very important.

In this vision paper, we sketch our design for a new data base management system (DBMS), SeeDB, that automates the especially laborious aspects of the search for useful data insights. Figure 2 depicts SeeDB, together with a conventional DBMS. In the conventional system, the user or analyst submits a query Q and obtains data subsets. Thus, conventional systems do not provide any means for the analyst to get intuitive visual insights directly. In SeeDB, the analyst also submits a query, but instead automatically obtains views (or visualizations) of the query result that are potentially of interest. As illustrated in our examples, these visualizations help the analyst quickly interpret and understand specific “interesting/useful” aspects of the query result. Thus, with SeeDB, we fundamentally modify the query-result paradigm of databases: SeeDB is provided a query Q , and outputs visualizations of interesting aspects of Q .

demo on Friday

THE WORK OF EDWARD TUFTE AND GRAPHICS PRESS

GRAPHICS PRESS LLC P.O. BOX 430 CHESHIRE, CT 06410 800 822-2454

Edward Tufte is a statistician and artist, and Professor Emeritus of Political Science, Statistics, and Computer Science at Yale University. He wrote, designed, and self-published 4 classic books on data visualization. *The New York Times* described ET as the "Leonardo da Vinci of data," and *Business Week* as the "Galileo of graphics." He is now writing a book/film *The Thinking Eye* and constructing a 234-acre tree farm and sculpture park in northwest Connecticut, which will show his artworks and remain open space in perpetuity. He founded Graphics Press, ET Modern gallery/studio, and Hogpen Hill Farms LLC.

<http://www.edwardtufte.com/tufte/>



MEMBERSHIP TUTORIALS GUIDES BOOKS FEATURES

Recent SEE ALL →

Satellite time-lapse of Earth

Charlie Loyd, who works with satellite imagery at Mapbox, put together a 12-second time-lapse of Earth using a day of ...



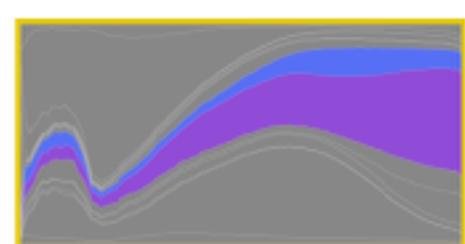
Amanda Cox is new editor of The Upshot

So great and well-deserved.

Features SEE ALL →

How You Will Die

So far we've seen when you will die and how other people tend to die. Now let's put the two together to see how and when you will die, given your sex, race, and age.



Playing with fonts using neural networks

Erik Bernhardsson downloaded 50,000 fonts and then threw them to the neural networks to see what sort of letters ...

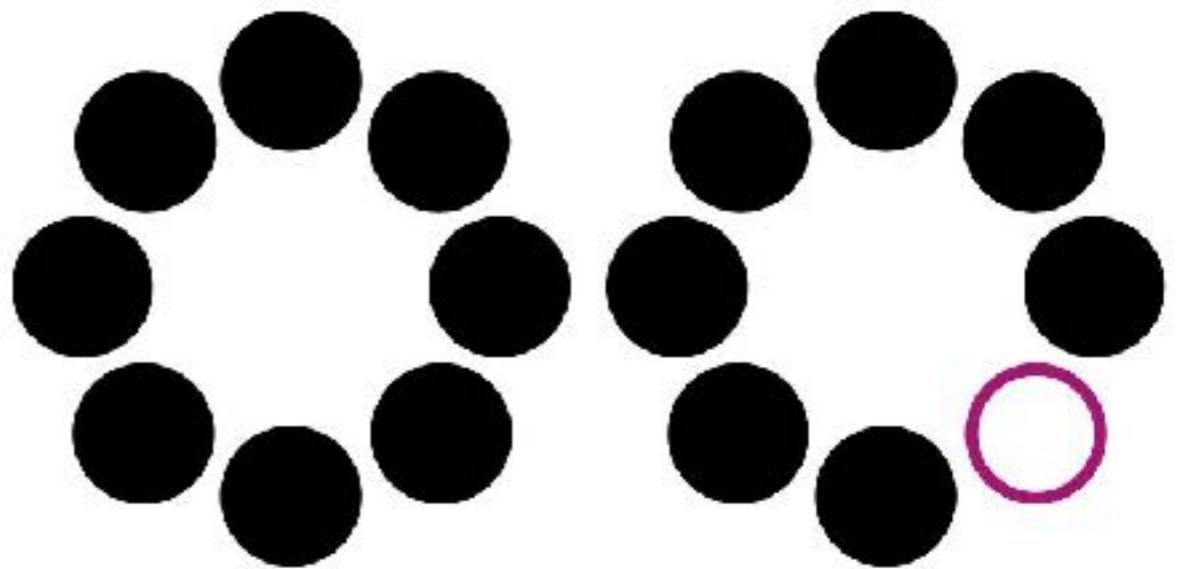


Miccing 11th of the

[View original](#)

etc.

<http://flowingdata.com>



SIMILARITY

Numerical measure of
how **different** two data
objects are

Minimum
dissimilarity is often **0**

Upper limit **varies**

dis-similarity

Lower when objects are more alike

Numerical measure of
how alike two data
objects are

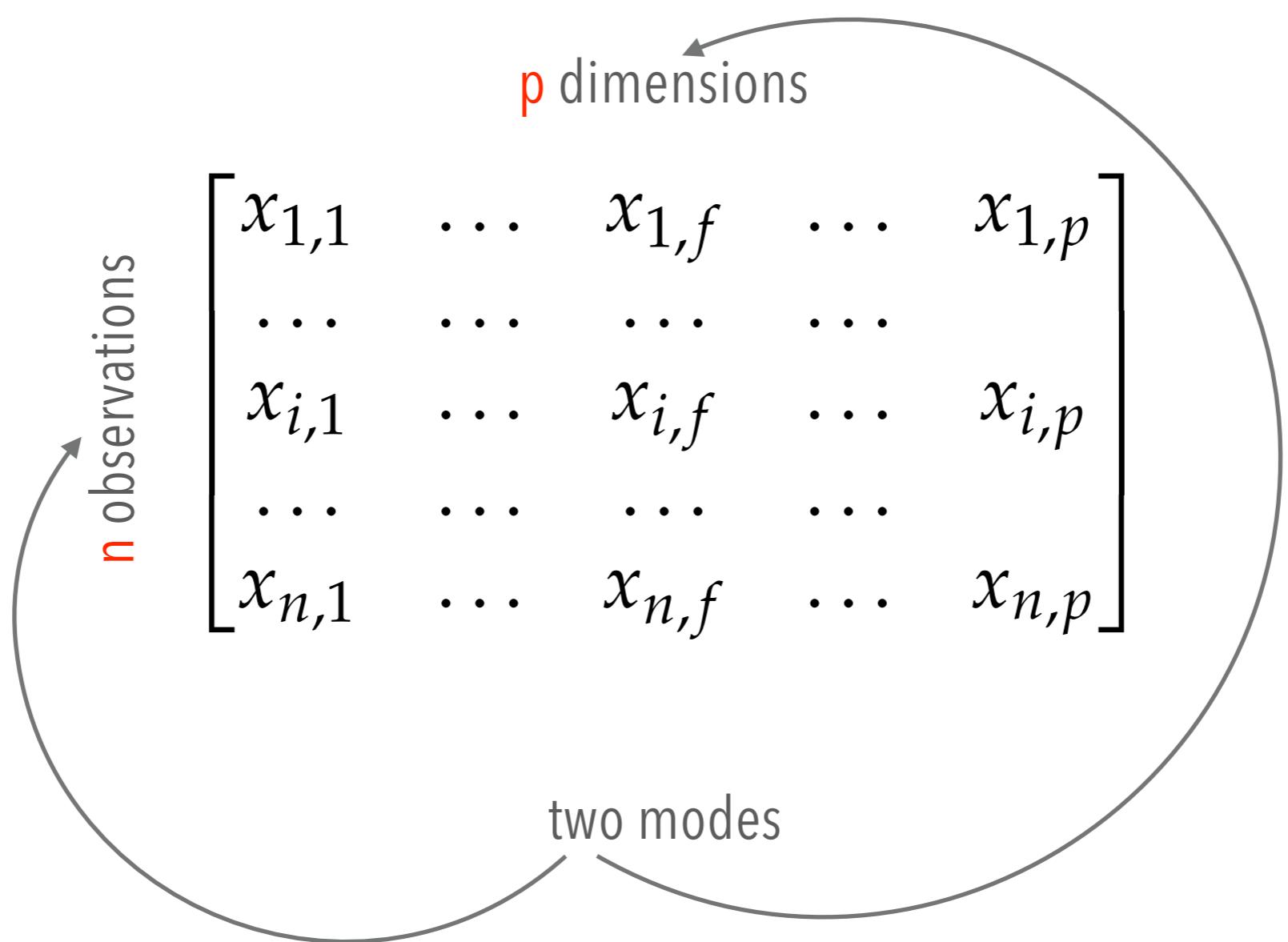
Often falls in the
range [0,1]

similarity

Value is higher when objects are more alike

Proximity refers
to similarity or to
dissimilarity

since similarity and dis-similarity are related



n objects

n objects

$$\begin{bmatrix} 0 & & & & \\ d(2,1) & 0 & & & \\ d(3,1) & d(3,2) & 0 & & \\ \vdots & \vdots & \vdots & \ddots & \\ d(n,1) & d(n,2) & \dots & \dots & 0 \end{bmatrix}$$

a triangular matrix

only registers distance \longrightarrow one mode

dissimilarity matrix

Proximity for Nominal Attributes

1

two objects

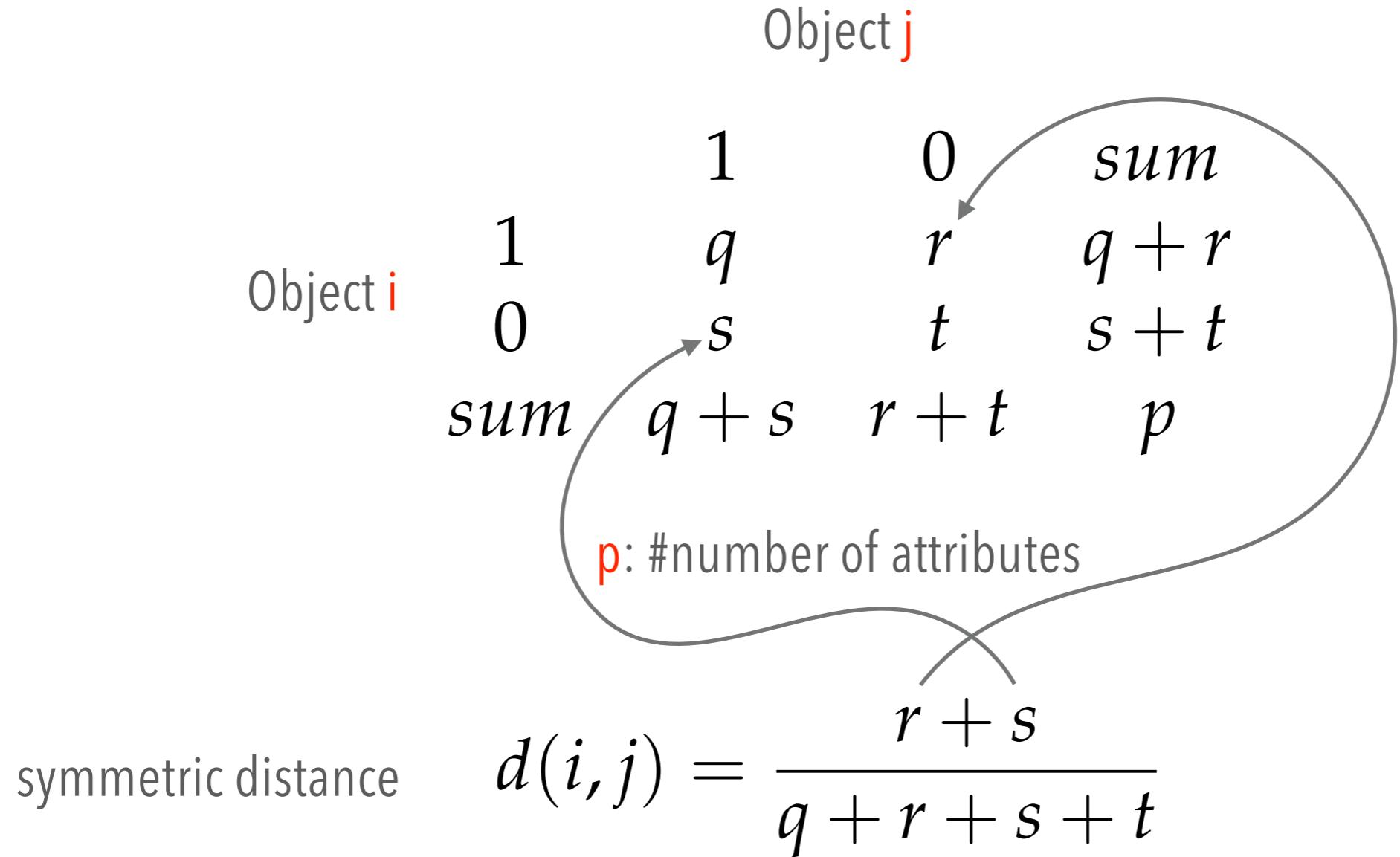
$$d(i, j) = \frac{p - m}{p}$$

m: #matches p: #variables

2

binary





binary attributes

binary attributes

asymmetric distance

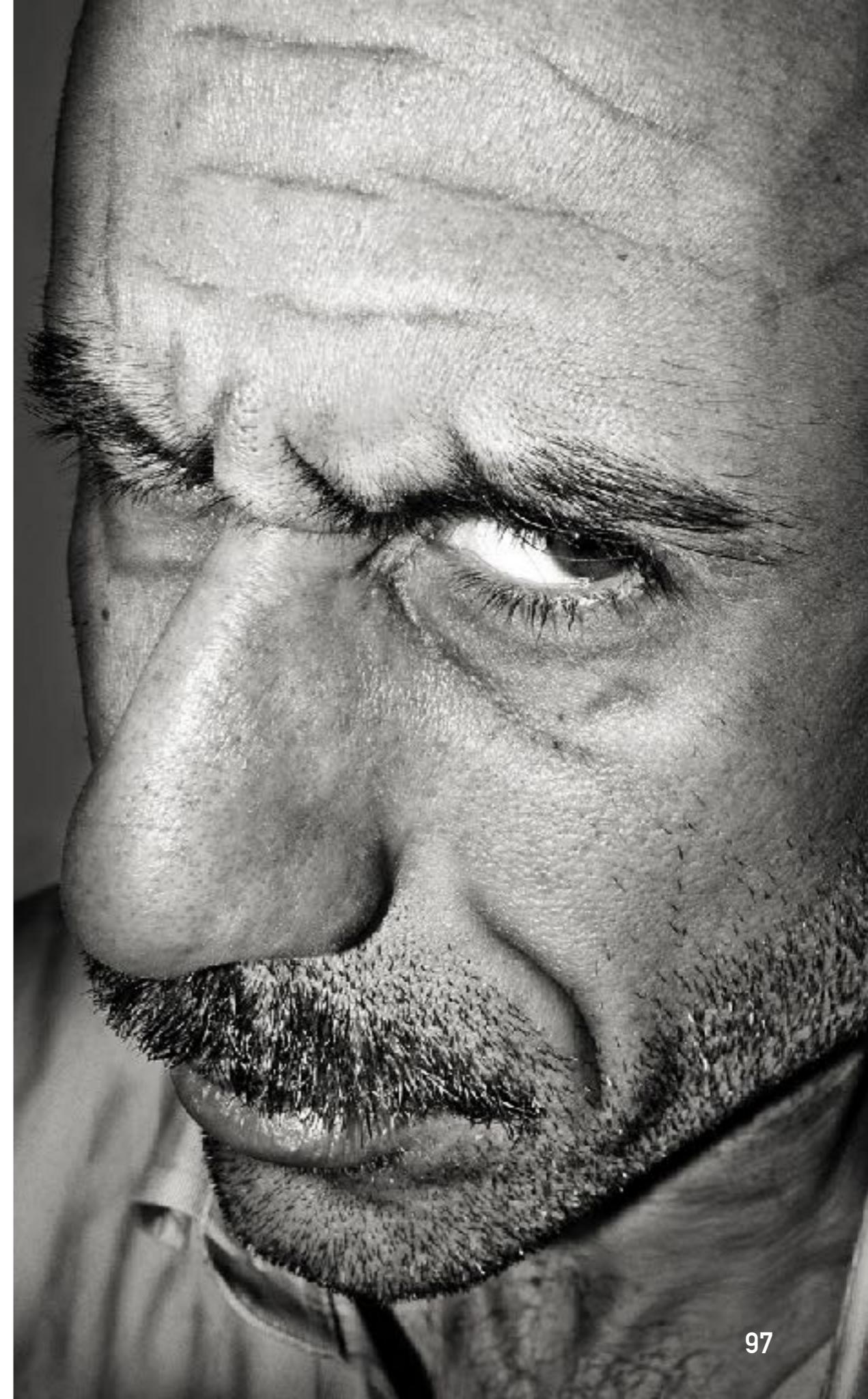
$$d(i, j) = \frac{r + s}{q + r + s}$$

		Object j		
		1	0	sum
Object i	1	q	r	$q + r$
	0	s	t	$s + t$
sum	$q + s$	$r + t$...	p

p : #number of attributes

both objects register 0 for these attributes

Why drop
the case
when both
attributes
register
False?



		Object j		
		1	0	sum
Object i	1	q	r	$q + r$
	0	s	t	$s + t$
sum	$q + s$	$r + t$	\dots	p

p : #number of attributes

drop the case
when both
objects register 0

asymmetric similarity
Jacquard coefficient

$$\text{sim}(i, j) = 1 - d(i, j) = \frac{q}{q + r + s}$$

binary attributes

compute dissimilarity between Jack and Jim on asymmetric attributes

Y:yes N:no, negative P: positive

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	M	Y	N	P	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

$$d(Jack, Jim) = \frac{r + s}{q + r + s} = \frac{2}{3}$$

on how many attributes do they disagree?

what are the attributes that we care about?

z is zero
mean, unit
variance

$$z = \frac{x - \bar{x}}{s}$$

raw score
mean
standard deviation

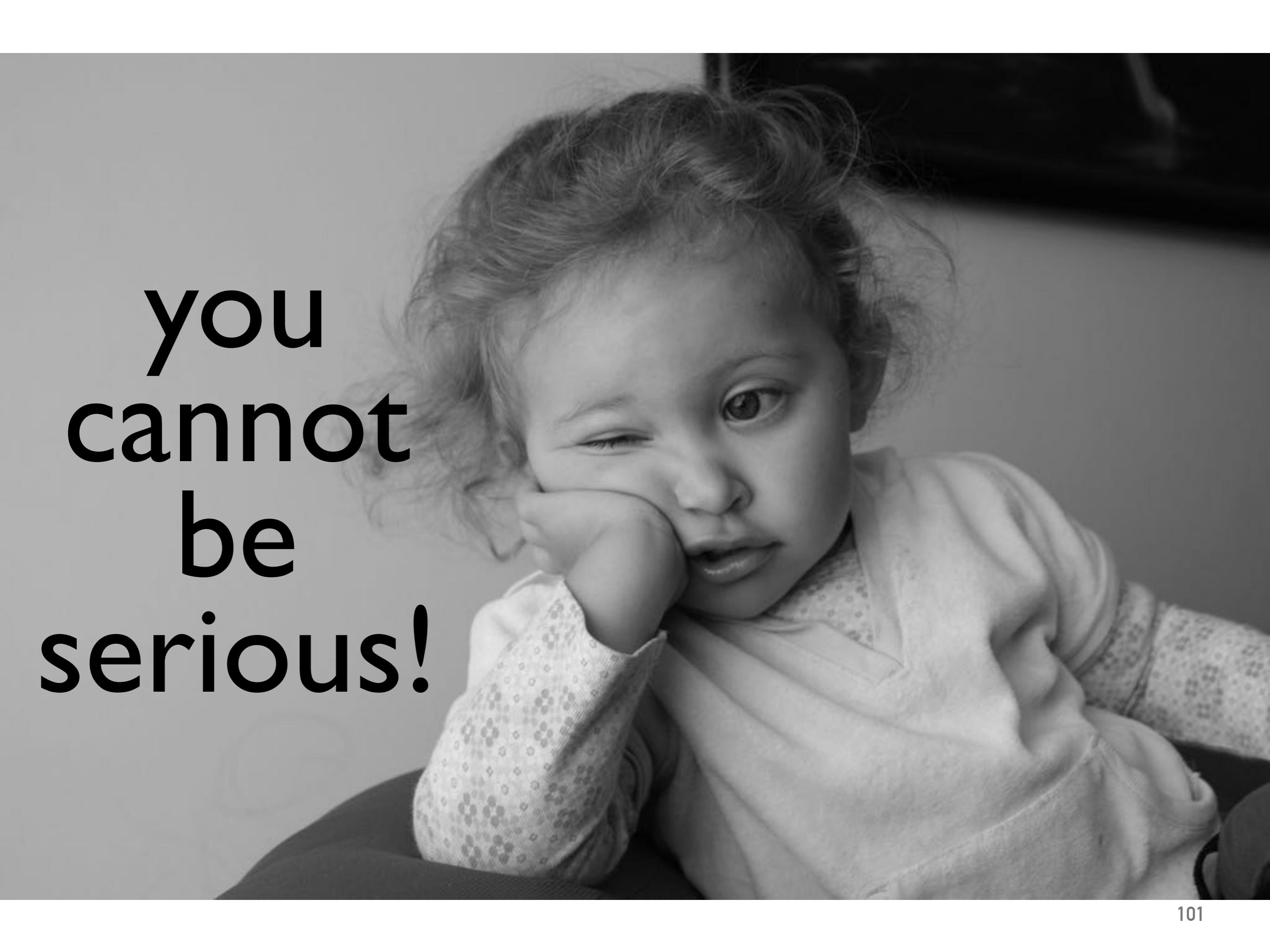
z-score

$$z = \frac{x - \bar{x}}{s_f}$$

$$s_f = \frac{1}{n} \sum_i^n |x_i - \bar{x}|$$

Apply z-score normalization **prior** to
computing distance, when **combining**
attributes with different ranges

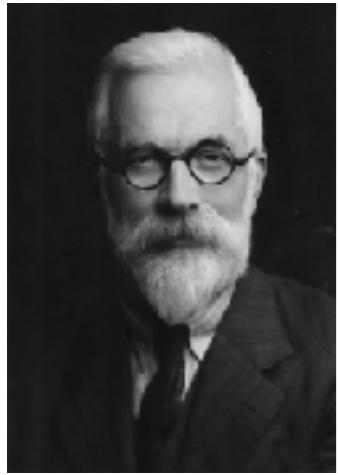
mean absolute deviation



**you
cannot
be
serious!**

consistent sufficient **efficient**

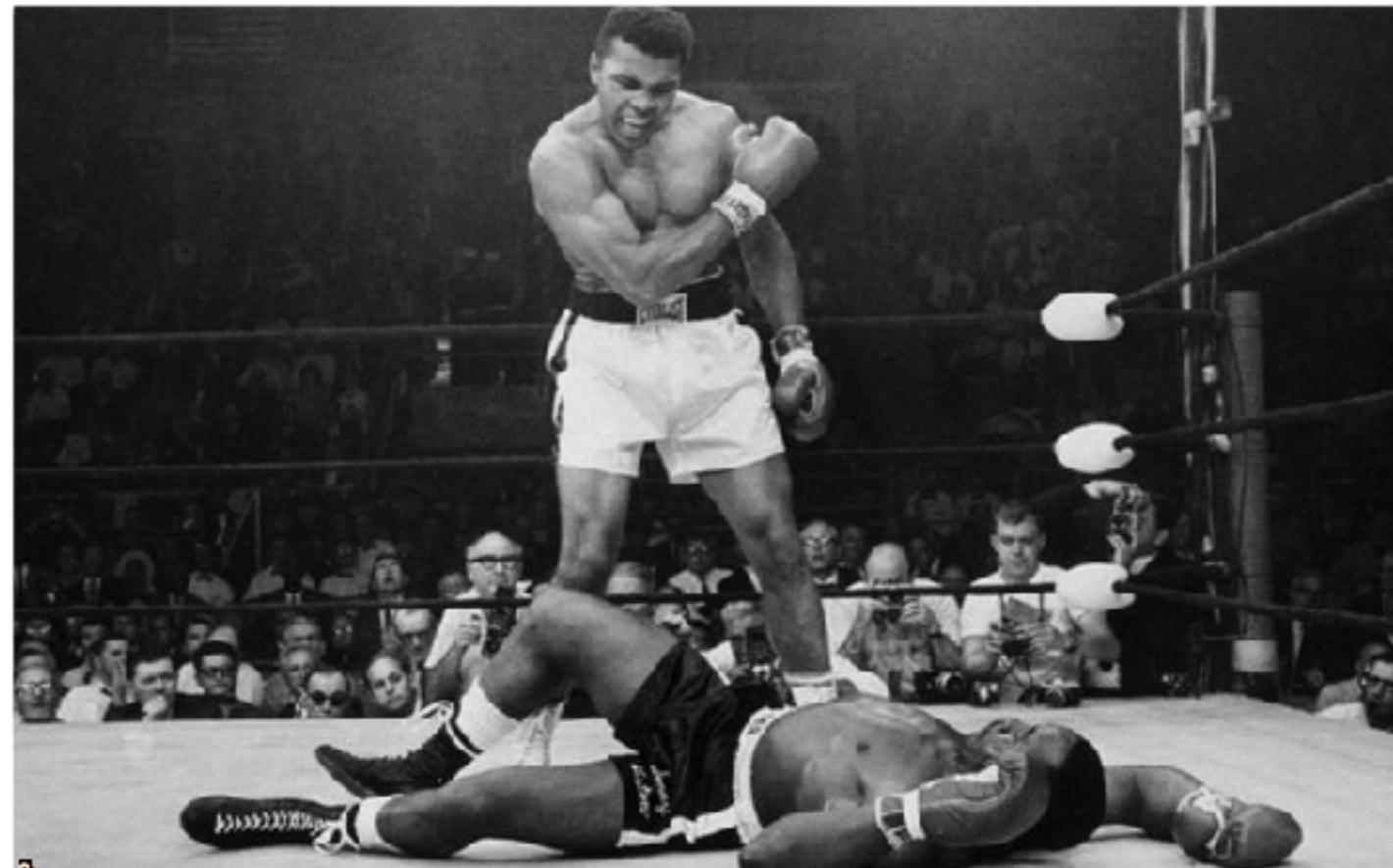
standard deviation is more efficient for **normally** distributed data than MAD



RA Fisher



Eddington

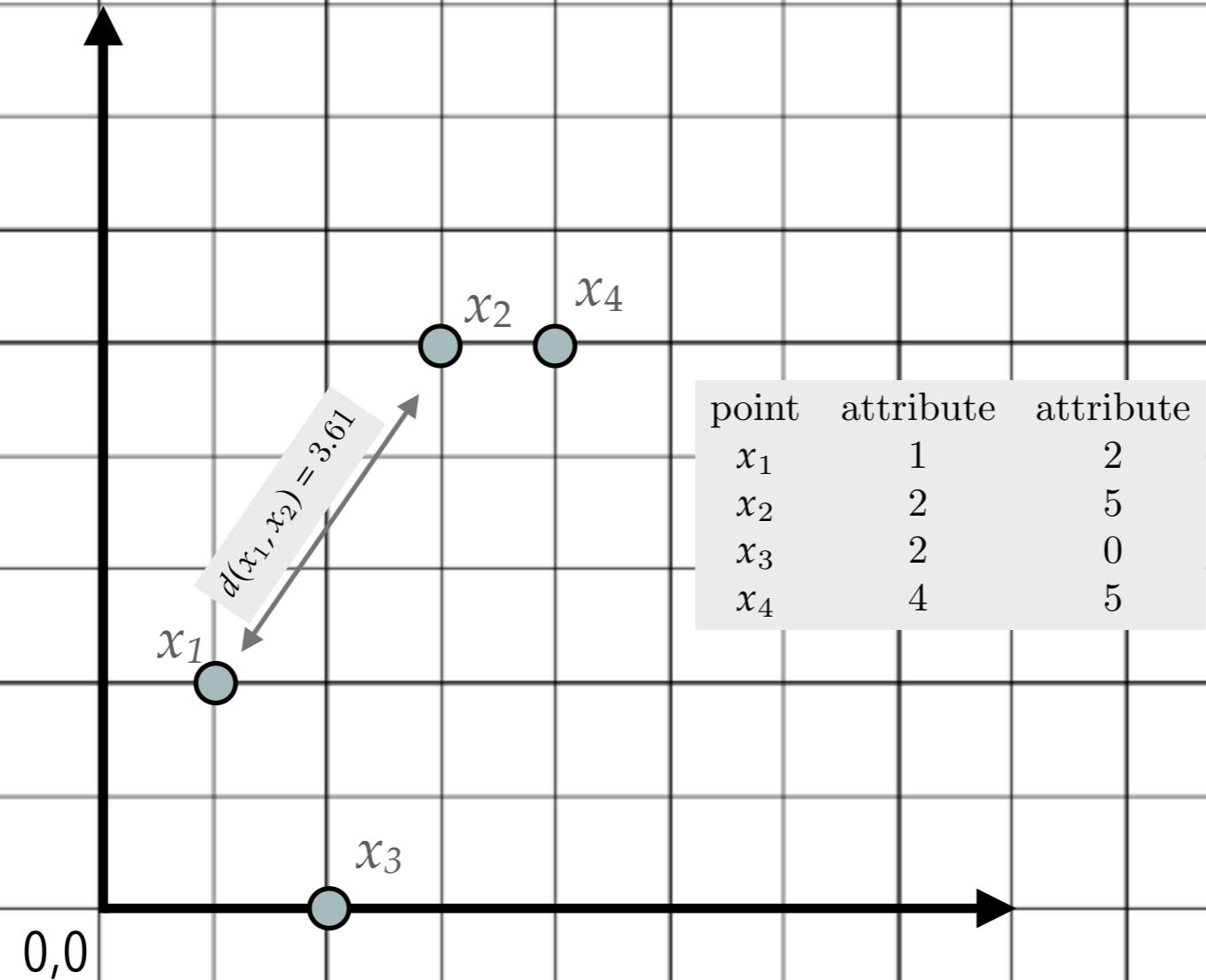


RA Fisher's work settled the debate in favor of σ

Gorard, S. (2005). Revisiting a 90-Year-Old Debate: The Advantages of the Mean Deviation. *British Journal of Educational Studies*, 53(4), 417–430.

noise, non-normal data

euclidean dissimilarity



$$d(x, y) = \left(\sum_i^p |x_i - y_i|^h \right)^{\frac{1}{h}} \quad L_h \text{ norm}$$

Minkowski distance

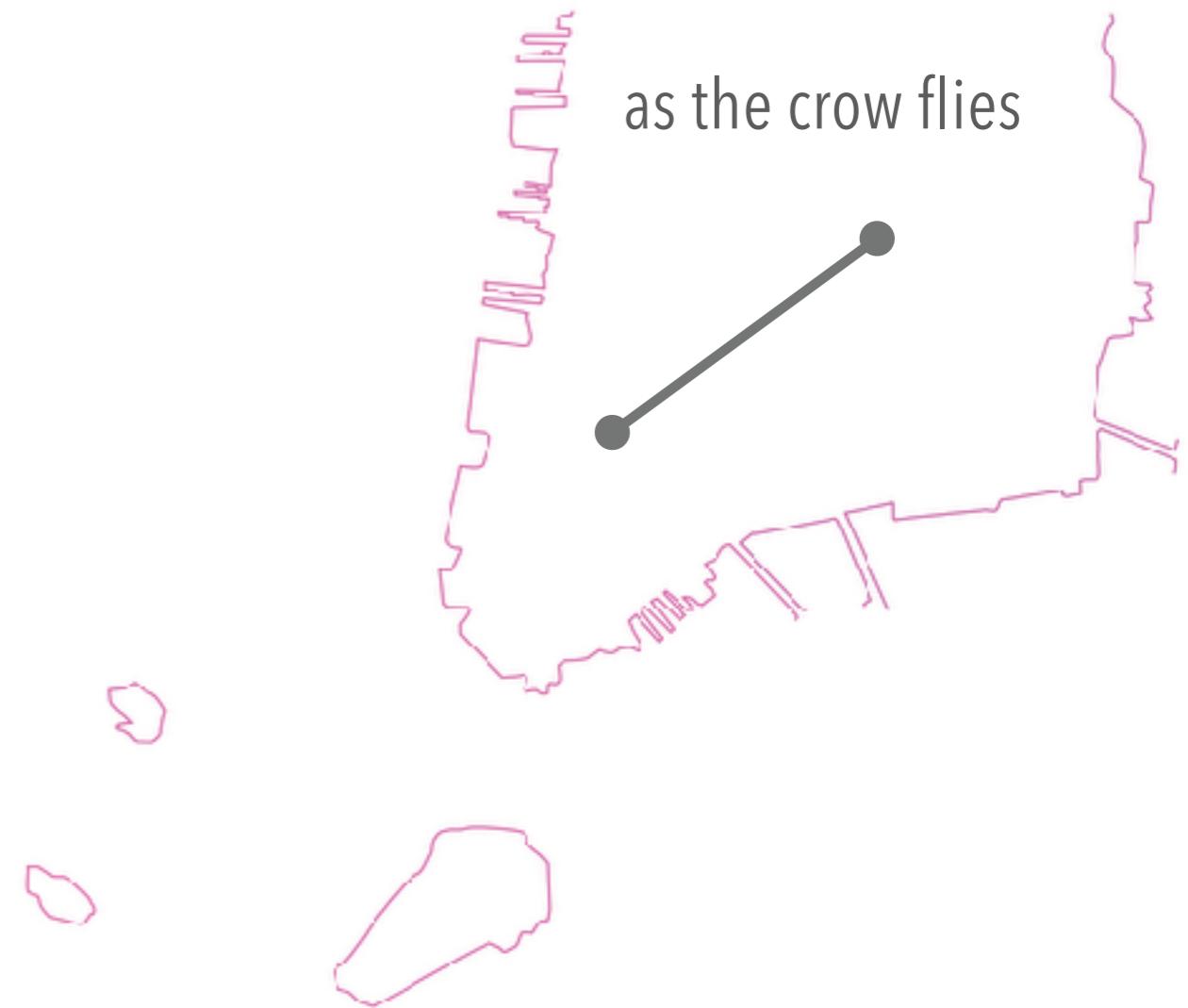
positive $d(x, y) \geq 0$

symmetry $d(x, y) = d(y, x)$

triangle inequality $d(x, y) \leq d(x, z) + d(z, y)$

$$d(x, y) = \left(\sum_i^p |x_i - y_i|^2 \right)^{\frac{1}{2}}$$

Special cases:
Euclidean distance



$$d(x, y) = \sum_i^p |x_i - y_i|$$

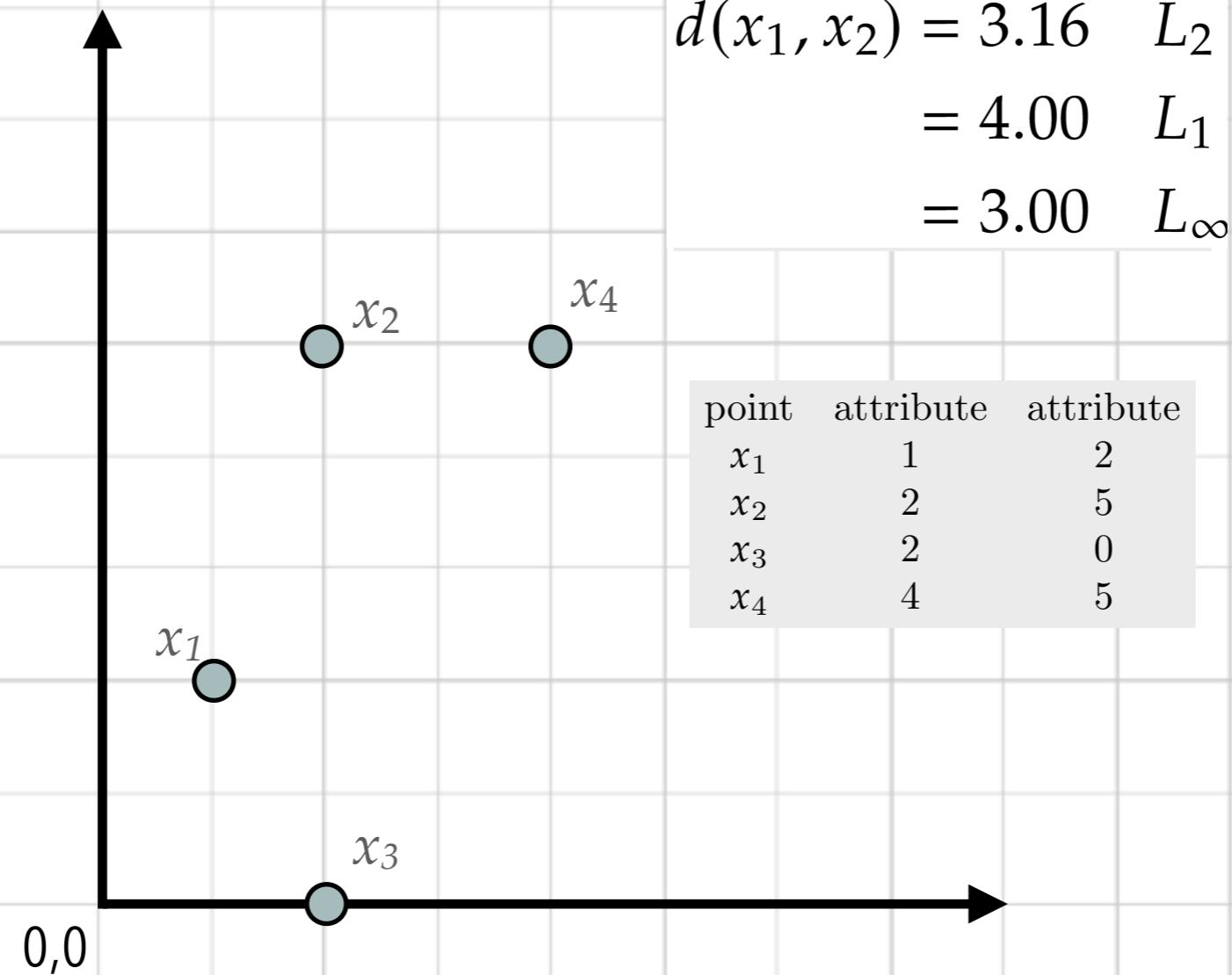


Special cases:
Manhattan distance

$$d(\mathbf{x}, \mathbf{y}) = \lim_{h \rightarrow \infty} \left(\sum_i^p |\mathbf{x}_i - \mathbf{y}_i|^h \right)^{\frac{1}{h}}$$
$$= \max_i |\mathbf{x}_i - \mathbf{y}_i|$$

This is the maximum difference between any component (attribute) of the vectors

Special cases:
Supremum norm



Minkowski comparison

An ordinal variable
can be discrete or
continuous

Order is important,
e.g., rank

Can be treated like
interval-scaled

ordinal variables

$$x_i \rightarrow r_i, r_i \in \{1, \dots, M\}$$

replace x_i by their rank

$$z_i = \frac{r_i - 1}{M - 1}$$

map

compute dissimilarity
using an appropriate
Minkowski distance

use the distance
corresponding to the
attribute type

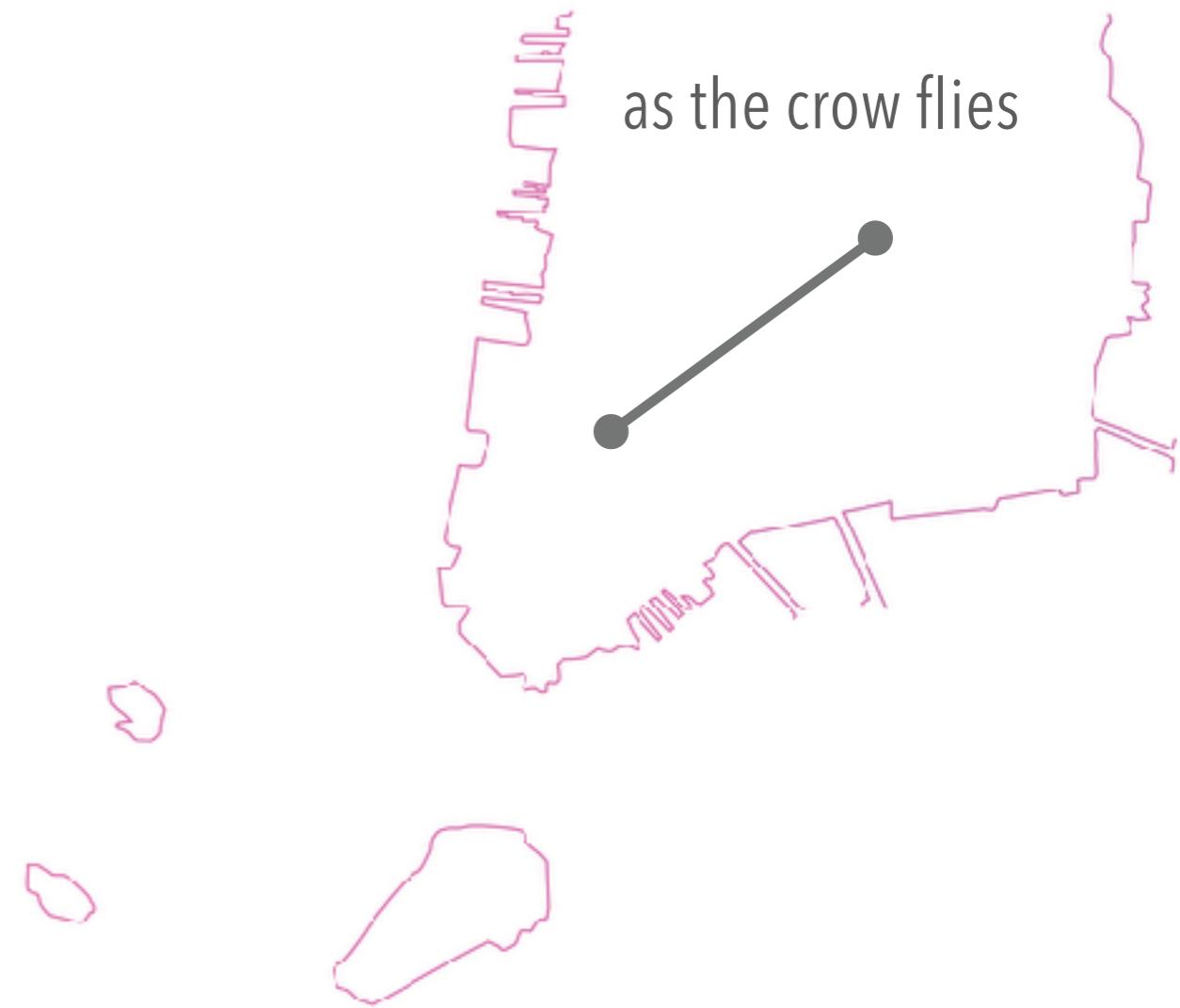
$$d(x, y) = \frac{\sum_{i=1}^p \omega_i d(x, y)^{(i)}}{\sum_{i=1}^p \omega_i}$$

↓
attribute

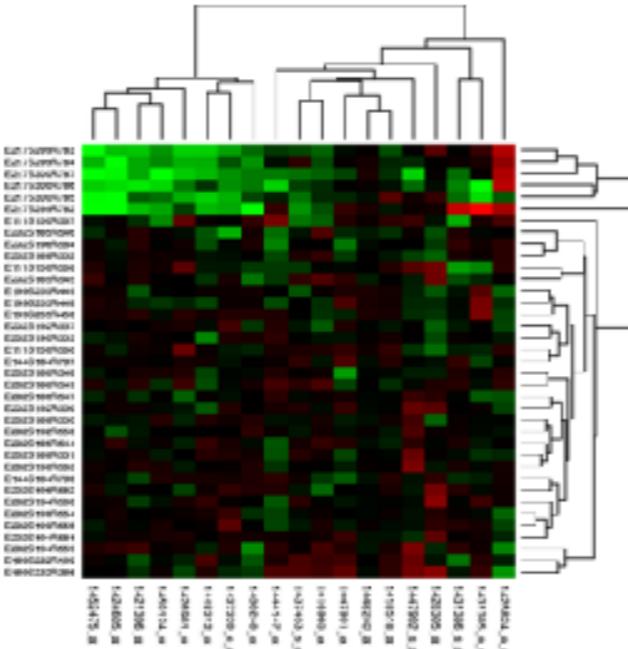
mixed attributes

Nominal, symmetric
binary, asymmetric
binary, numeric, ordinal

$$d(x, y) = \left(\sum_i^p |x_i - y_i|^2 \right)^{\frac{1}{2}}$$



**what happens when the number
of dimensions is very large?**



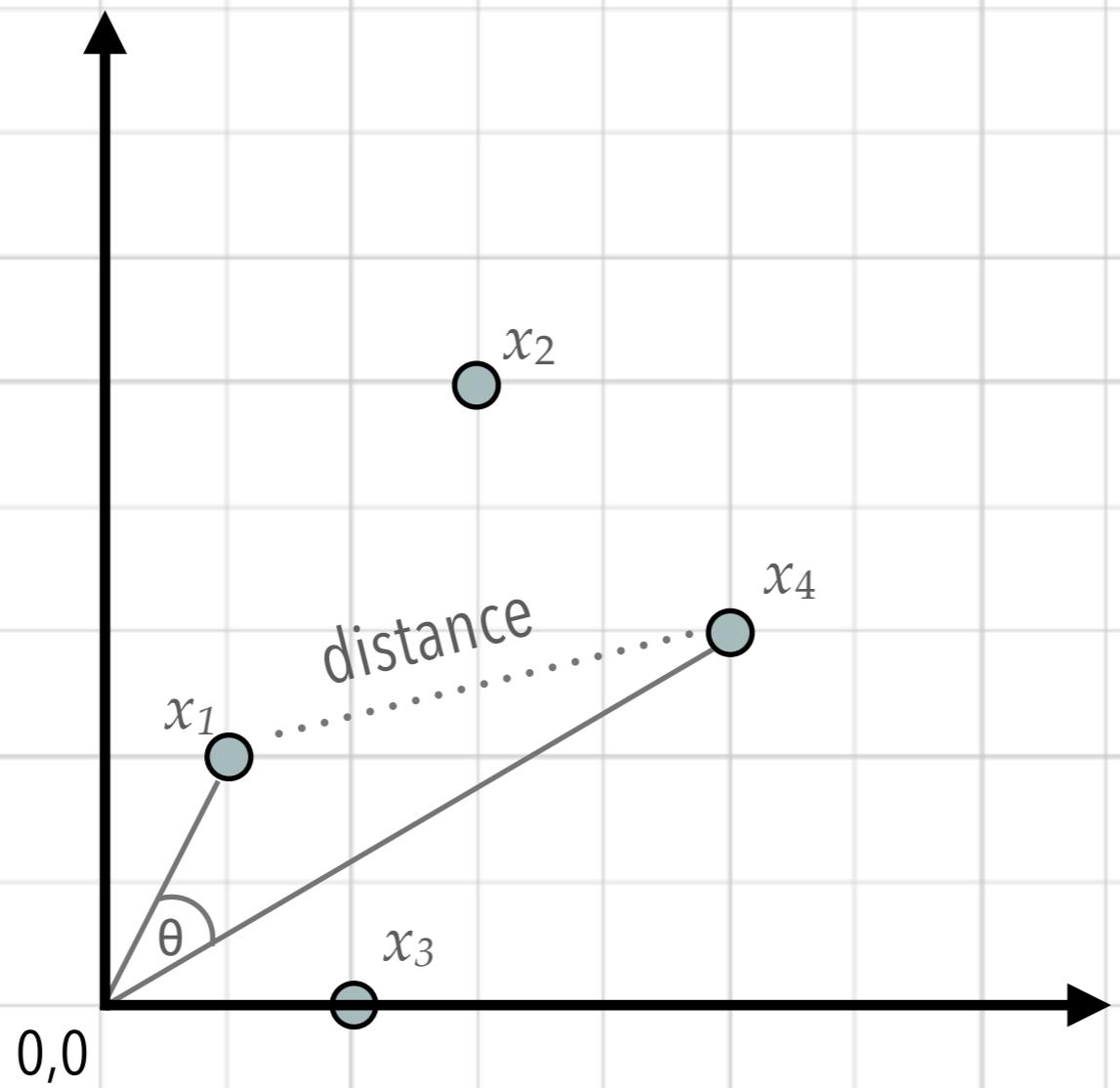
Micro-array may have tens of thousands of dimensions

<i>Document</i>	<i>teamcoach</i>	<i>hockey</i>	<i>baseball</i>	<i>soccer</i>	<i>penalty</i>	<i>score</i>	<i>win</i>	<i>loss</i>	<i>season</i>
Document1	5	0	3	0	2	0	2	0	0
Document2	3	0	2	0	1	1	1	0	1
Document3	0	7	0	2	1	0	3	0	0
Document4	0	1	0	0	1	2	0	3	0

A document can be represented by thousands of attributes, each recording the frequency of a particular word (such as keywords) or phrase in the document.

cosine similarity is used for high dimensional data

instead of measuring length, we measure the angle between two objects



cosine similarity $s(x, y) = \cos(\theta) = \frac{\langle x, y \rangle}{\|x\| \|y\|}$

dot product

cosine distance $d(x, y) = 1 - s(x, y)$

L₂ norm

$$x = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)$$

$$y = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1)$$

compute the
cosine similarity

$$\begin{aligned} \langle x, y \rangle &= 5 \times 3 + 0 \times 0 + 3 \times 2 + 0 \times 0 + 2 \times 1 + 0 \times 1 + 2 \times 1 + 0 \times 0 + 0 \times 1 \\ &= 25 \end{aligned}$$

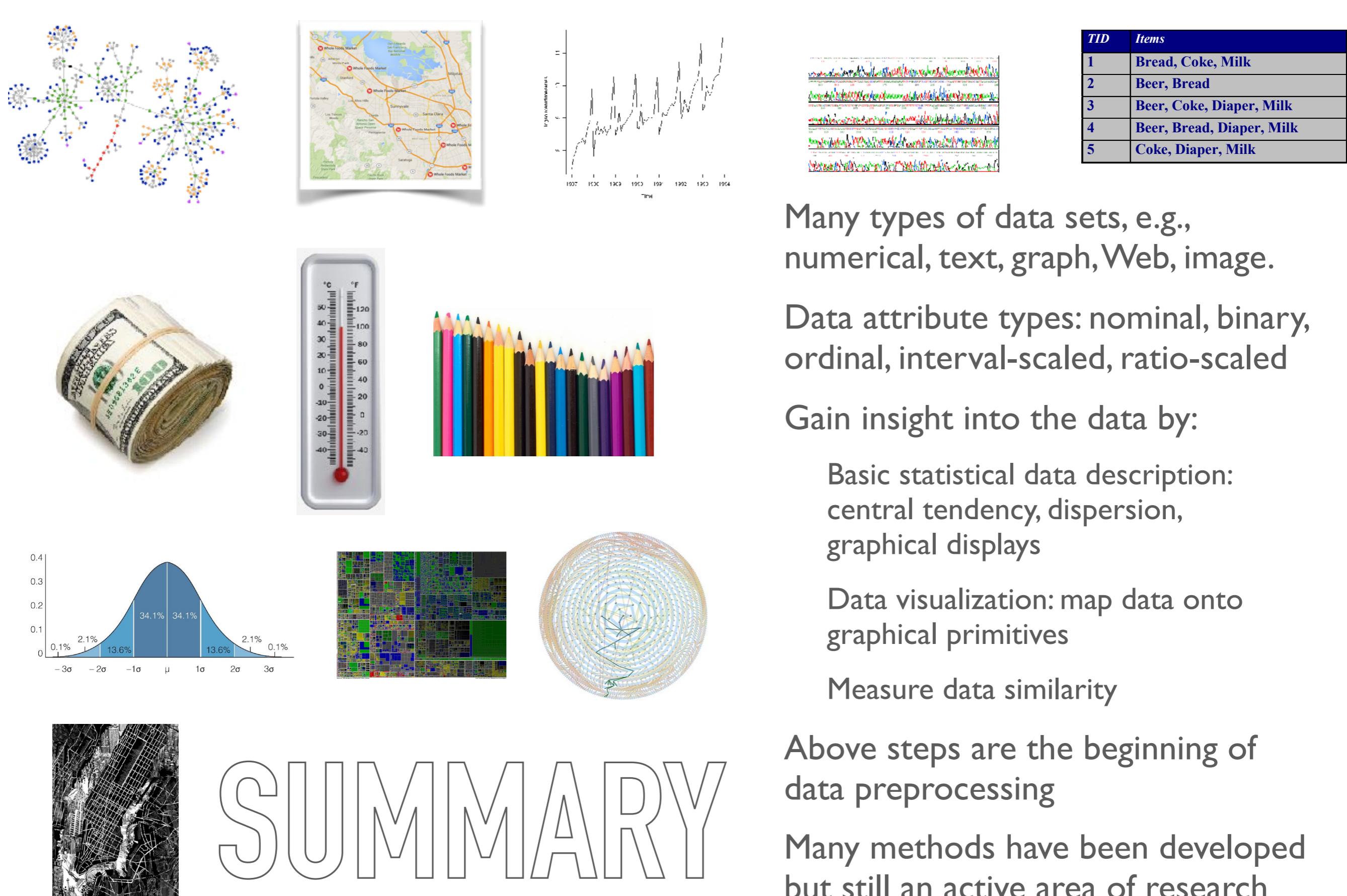
$$\|x\| = 6.481$$

$$\|y\| = 4.12$$

$$s(x, y) = \cos(\theta) = \frac{\langle x, y \rangle}{\|x\| \|y\|}$$

$$d(x, y) = 1 - s(x, y)$$

What distance
measure should we
use and when?



Many types of data sets, e.g., numerical, text, graph, Web, image.

Data attribute types: nominal, binary, ordinal, interval-scaled, ratio-scaled

Gain insight into the data by:

Basic statistical data description:
central tendency, dispersion,
graphical displays

Data visualization: map data onto graphical primitives

Measure data similarity

Above steps are the beginning of data preprocessing

Many methods have been developed
but still an active area of research

$$d(x, y) = \left(\sum_i^p |x_i - y_i|^h \right)^{\frac{1}{h}}$$

$$d(x, y) = \cos(\theta) = \frac{\langle x, y \rangle}{\|x\| \|y\|}$$