

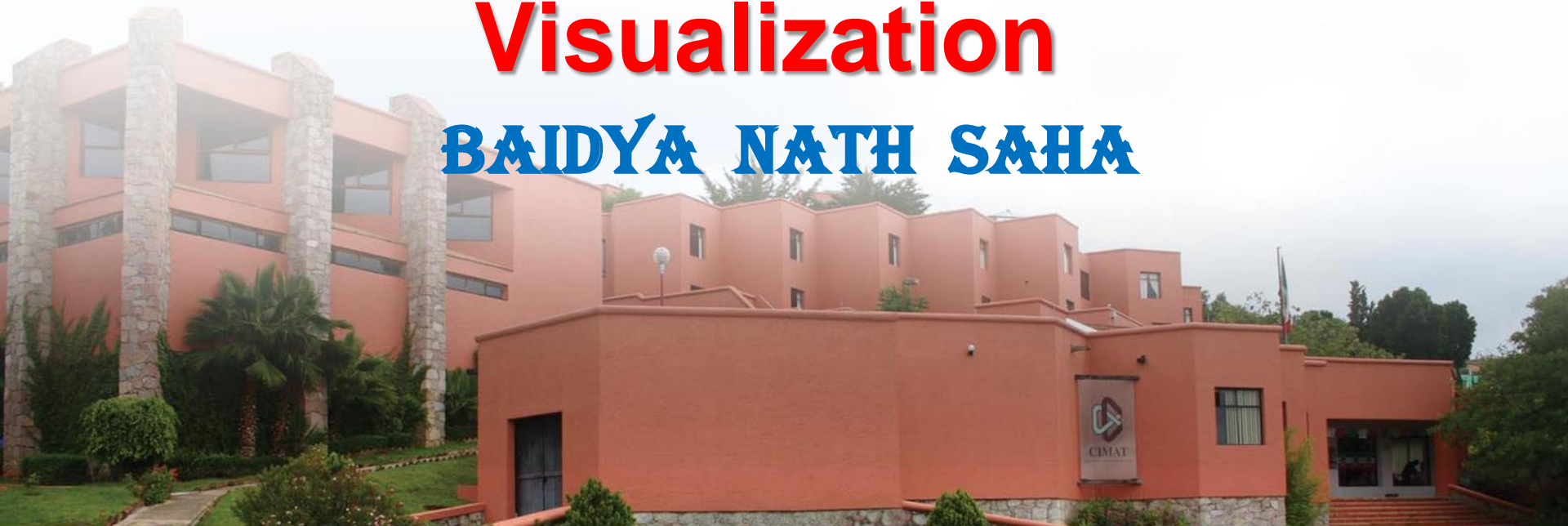


CIMAT Centro de Investigación en Matemáticas, A.C.



Computational Aspects of High Dimensional Large Scale Data Visualization

BAIDYA NATH SAHA



Outline

- Introduction
- History
- Modern Visualization Techniques and Computational Problems
 - Parallel Co-ordinates, heat map, SOM etc.
 - Space, data ordering, processing
- Conclusion

Outline

➤ Introduction

➤ History

➤ Modern Visualization Techniques and Computational Problems

- Parallel Co-ordinates, heat map, SOM etc.
- Space, data ordering, processing

➤ Conclusion

What is visualization?

- “any technique for creating images, diagrams, or animations to communicate a message”
(Wikipedia)
- “a systematic, rule-based, external, permanent, and graphic representation that depicts information in a way that is conducive to acquiring insights, developing an elaborate understanding, or communicating experiences”
(Lengler & Eppler)

Human Visual Perception

- ❑ Humans receive input from all five of their senses (sight, touch, hearing, smell, taste), however, Human visual cortex dominates our perception
- ❑ Accelerates the identification of hidden patterns in data
“A picture is worth a thousand words” – Chinese Proverb
- ❑ 50% of the human brain is dedicated to visual functions, and images are processed faster than text.
- ❑ The brain processes pictures all at once, but processes text in a linear fashion, meaning it takes much longer to obtain information from text.
- ❑ Furthermore, it is estimated that 65% of the population are visual learners.

Outline

- Introduction
- **History**
- Modern Visualization and Computational Problems
 - Parallel Co-ordinates, heat map, SOM etc.
 - Space, data ordering, processing
- Conclusion

Examples of Information Visualizations

- Pie chart
- Timeline
- Gantt Chart
- Metaphoric, e.g. iceberg
- Cartoon
- Org chart

A PERIODIC TABLE OF VISUALIZATION METHODS

C catalyst																G graphic utilization	
Tb table	Ca certain coordinates																Et cartoon
Pi pie chart	L line chart																Ri rich picture
B bar chart	Ac area chart	R radar chart cobweb	Pa parallel coordinates	Hy hyperbolic tree	Cy cycle diagram	T tree	Ve venn diagram	Mi mindmap	Sq square to opposites	Ce circular charts	Ar argument table	Sw swim lane diagram	Gc gantt chart	Pm perspective diagram	D diamond diagram	Pr process model	Hn knowledge map
Hi hierarchy	Sc strategy	Sa safety diagram	In information lines	E entity relationship diagram	Pt point set	Fl flow chart	Cl cleaning	Lc logic tree	Py pyramid model integrated	Ce cause-effect chain	Tl tunnel map	Dt decision tree	Gp goal setting path method	Cf concept map	Co concept map	Ic iceberg	Lm learning map
Tk task tree plot	Sp spring map	Da data map	Tp tree map	Cn cave tree	Sy system diagram simulation	Df data flow diagram	Se semantic network	So soft space mapping	Sn strategy map	Fo force field diagram	Ib ice berg diagram	Pr process flow chart	Pe pet chart	Ev evolution knowledge map	V ice berg diagram	Hh honey comb tree chart	I iceberg

- Cy** Process Visualization
- Hy** Structure Visualization
- ☉** Overview
- ⊙** Detail
- ☉⊙** Detail AND Overview
- < >** Divergent thinking
- > <** Convergent thinking

Note: Depending on your location and connection speed it can take some time to load a pop-up picture.
© Ralph Lengler & Martin J. Eppler. www.visual-literacy.org version 1.5

Su supply demand curve	Pe performance barometer	St strategy map	Oc organizational chart	Ho house of quality	Fd feedback diagram	Ft fishbone tree	Mq mind quadrant	Ld life-cycle diagram	Po poor man's flow chart	S siphon	Sm swimlane map	Is iceberg diagram	Yc youth strategy
Ed edge map tree	Pf particle diagram	Sg strategic game board	Mz matrix's organization	Z zoo's psychological tree	Ad advertising diagram	Be beaker discovery diagram	Bm big man's map	Stc strategy chart	Vc value chain	Hy hyper- cube	Sr swimlane strategy map	Ta top	Sd spring diagram

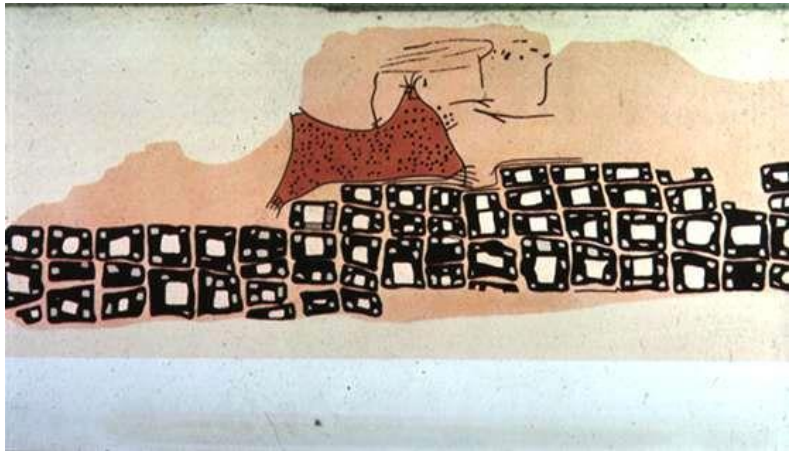
http://www.visual-literacy.org/periodic_table/periodic_table.html

A little history...

- ❑ Development of motion - about 100 years**
- ❑ Mobile sculpture – 20 years**
- ❑ Video art – 20 years**
- ❑ Animated maps – begin in the 1930's, but were not developed until 1959**
- ❑ Computer animated maps: 1990's**

Maps from the past

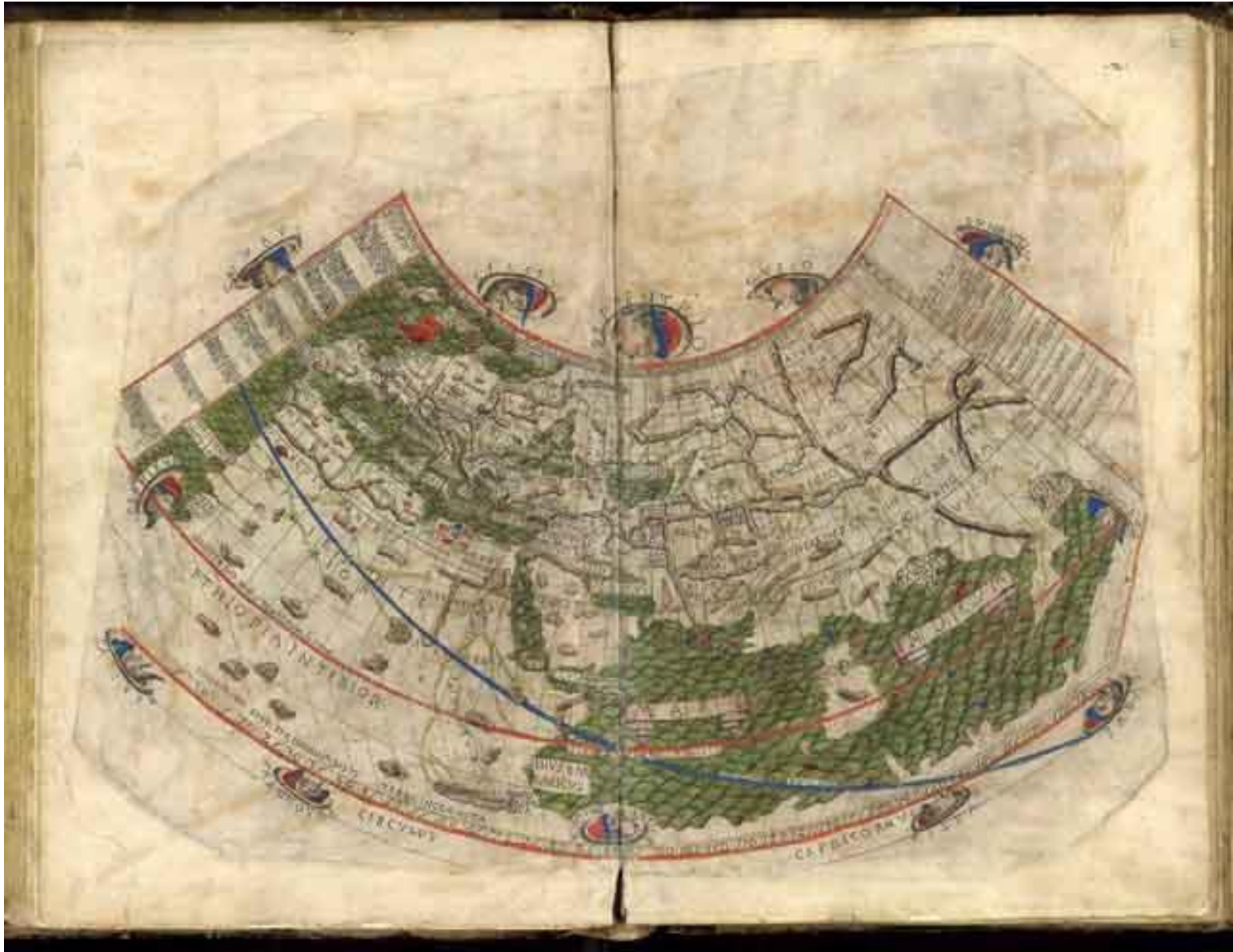
The oldest known map: Konya town, 6200 BC



Anaximander's Map of the World

Anaximander's Map of the World

The Ptolemy world map, written circa 150 BCE.



Early Representation



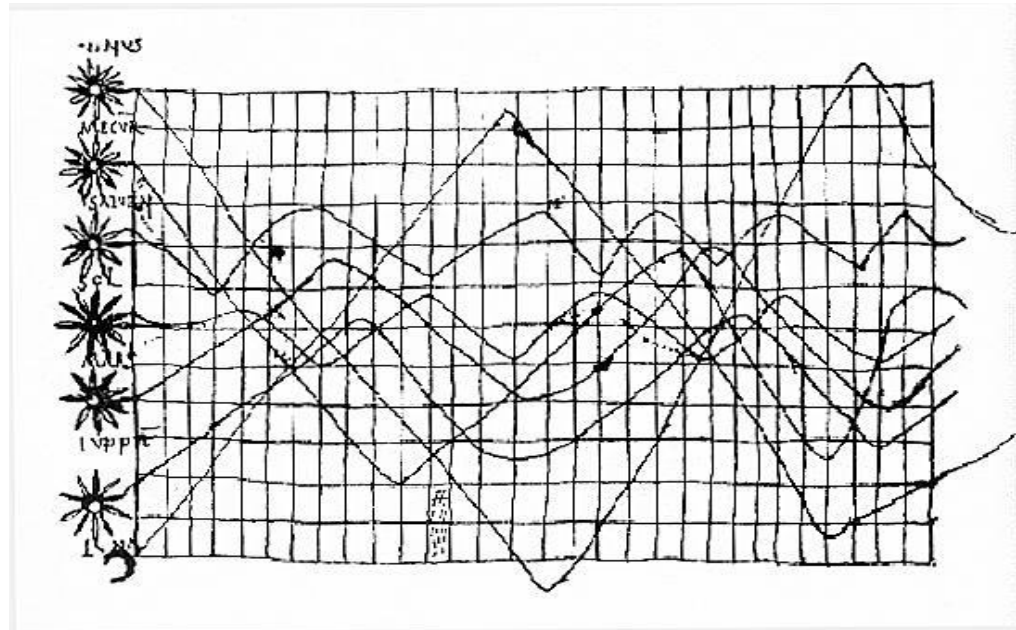
The Cave of Lascaux, France

~15,000 years old

- Tells a story

- ❑ Visualization existed before the invention of computers
- ❑ Representation of information allowing us to perceive such information visually

Planetary Orbits



- ❑ Tenth century
- ❑ Inclinations of the planetary orbits as a function of time.
- ❑ Oldest known attempt to show changing values graphically.

Message Communication - Example

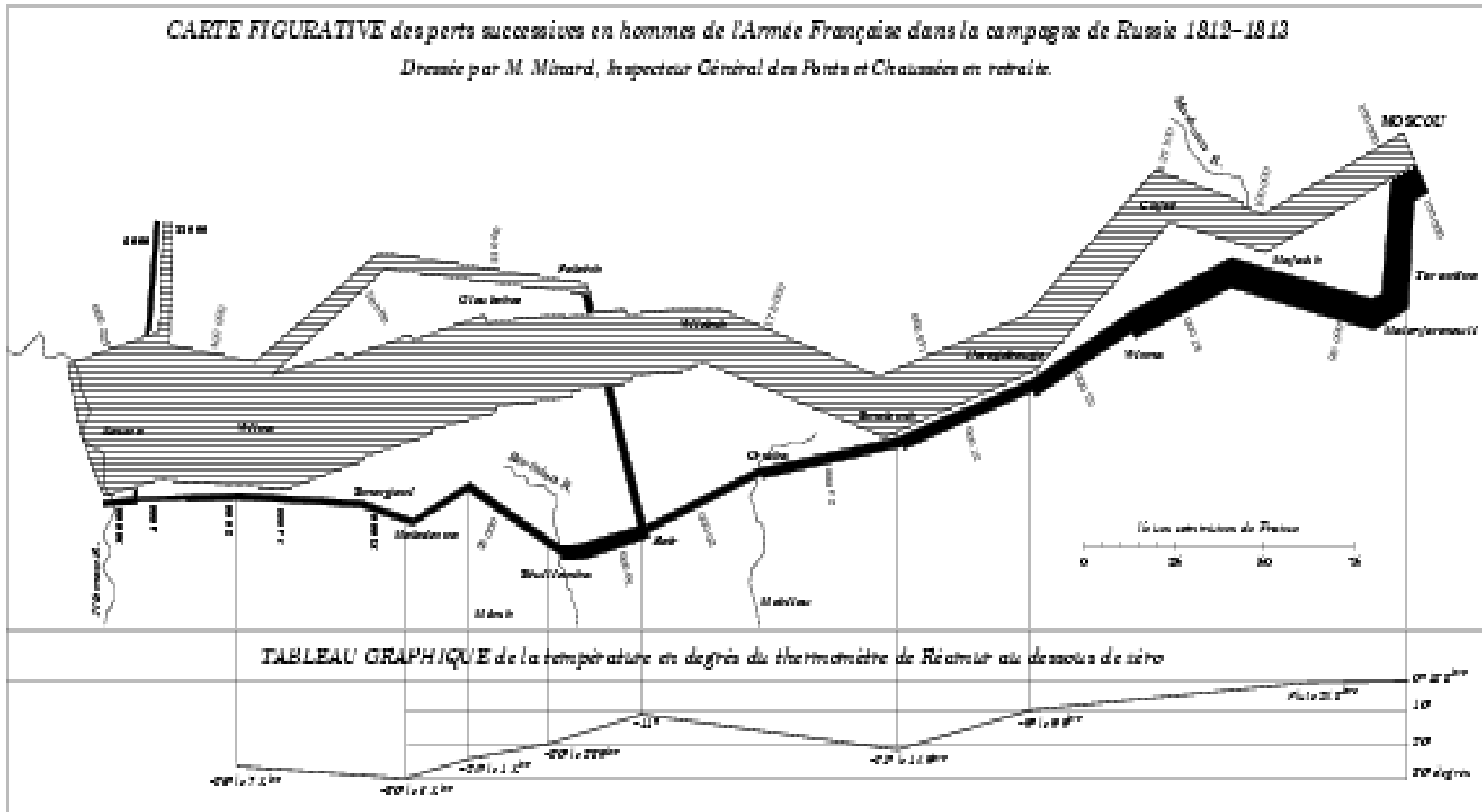
□ Napoleon's 1812 campaign on Russia

□ Input data

- Size of army
 - at the start of the campaign = 442,000
 - at the end of the campaign = 10,000
- Location of the army (2 dimensions)
- Direction of the army's movement
- Temperature and
- Time



Minard's Drawing



Created in 1861 by French engineer Charles Joseph Minard

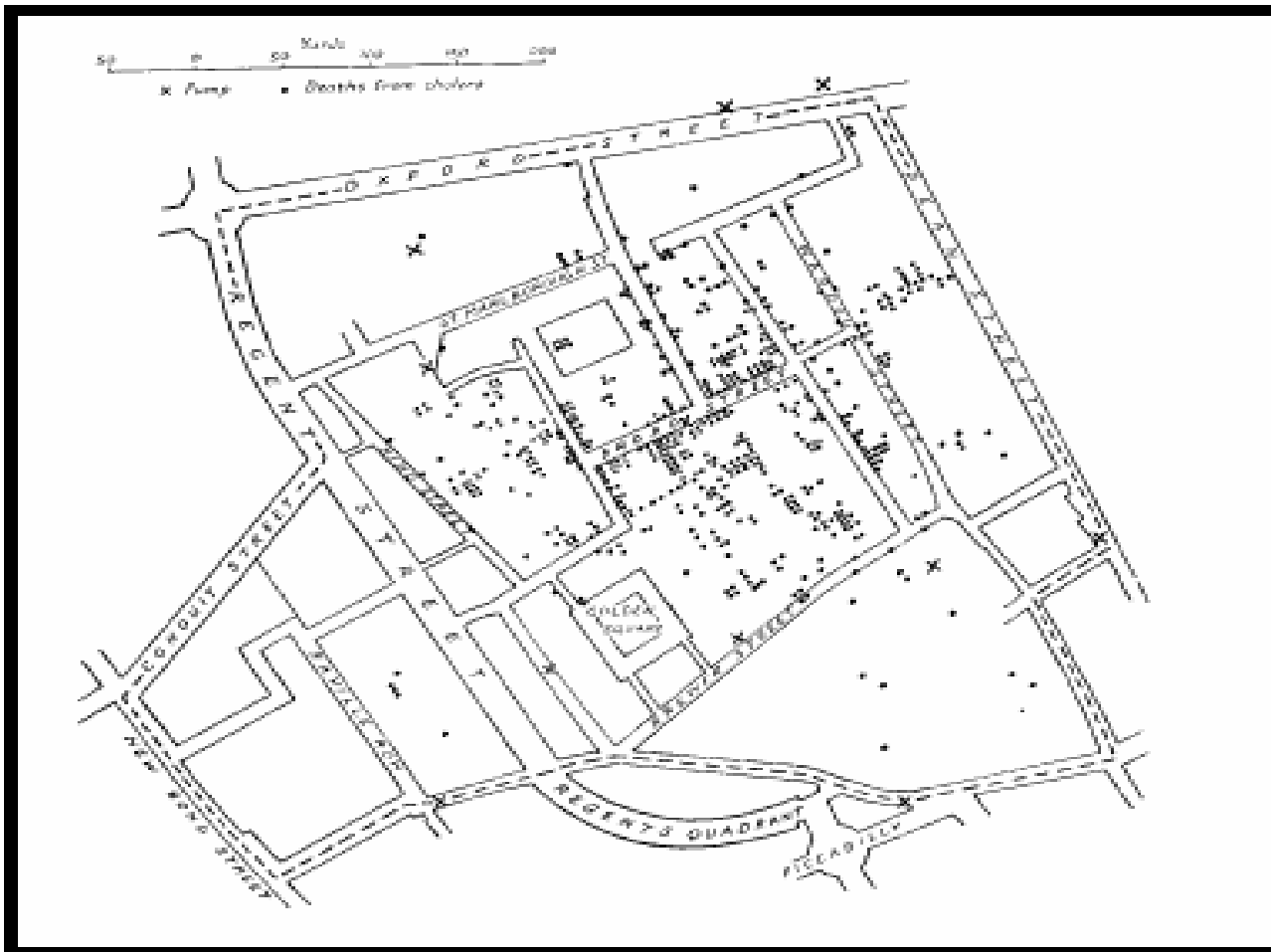
Minard's Drawing (2)

- ❑ Considered the best graphic ever produced
 - Inspiration for modern researchers
- ❑ Plots all the data corresponding to all the six input variables
- ❑ Clearly shows the message underlying the input data
 - Gradual reduction in the size of the army
 - Linked to the gradual fall in temperatures
- ❑ Input data is complex
- ❑ Yet, most important information abstracted out and presented in a simple graphic

Problem Solving - Example

- London cholera epidemic of 1854
- At that time, two hypotheses of causes of cholera:
 - Cholera is related to miasmas concentrated in the swampy areas of the city
 - Cholera is related to ingestion of contaminated water
- Input Data
 - Locations of deaths due to cholera
 - Locations of water pumps

Dr Snow's Cholera Map



Dots locate deaths due to cholera

Crosses locate water pumps

Dr Snow's Cholera Map (2)

- Plotting the input data on the map helped Dr Snow
 - to detect the epicentre of the epidemic
 - Close to a pump on Broad Street
- Considered a classic case of visualization helping reasoning with data

Outline

- Introduction
- History
- Modern Visualization and Computational Problems
 - Parallel Co-ordinates, heat map, SOM etc.
 - Space, data ordering, processing
- Conclusion

Key points in Visualization

- ❑ Visualization can enhance understanding
- ❑ Different visualization methods are appropriate to different types of data
- ❑ Free online tools such as ManyEyes, Swivel & Google Docs:
 - Make it (fairly) easy to create a visualization
 - Encourage visual literacy
 - Build dialogue and community

Information Visualization (IV)

- ❑ Visual presentation of abstractions or relationships underlying input data

- ❑ IV has two goals
 - Communication
 - to communicate a rich **message**
 - Problem solving/ reasoning/ analysis
 - to display large amount of information to facilitate reasoning to uncover **new** facts or relationships

- ❑ Limited screen sizes pose a serious challenge for using IV on very large data sets

- ❑ Therefore the main task is to pack large information into a simple graphic
 - Highlighting all the required (important) information

- ❑ **Creative art?**

Gospel like guidelines for IV

As you create graphics keep the following in mind.

- ❑ Avoid **distortion** of the true story, Don't tell lie.
- ❑ Induce the viewer to think about the **substance**, not the graph.

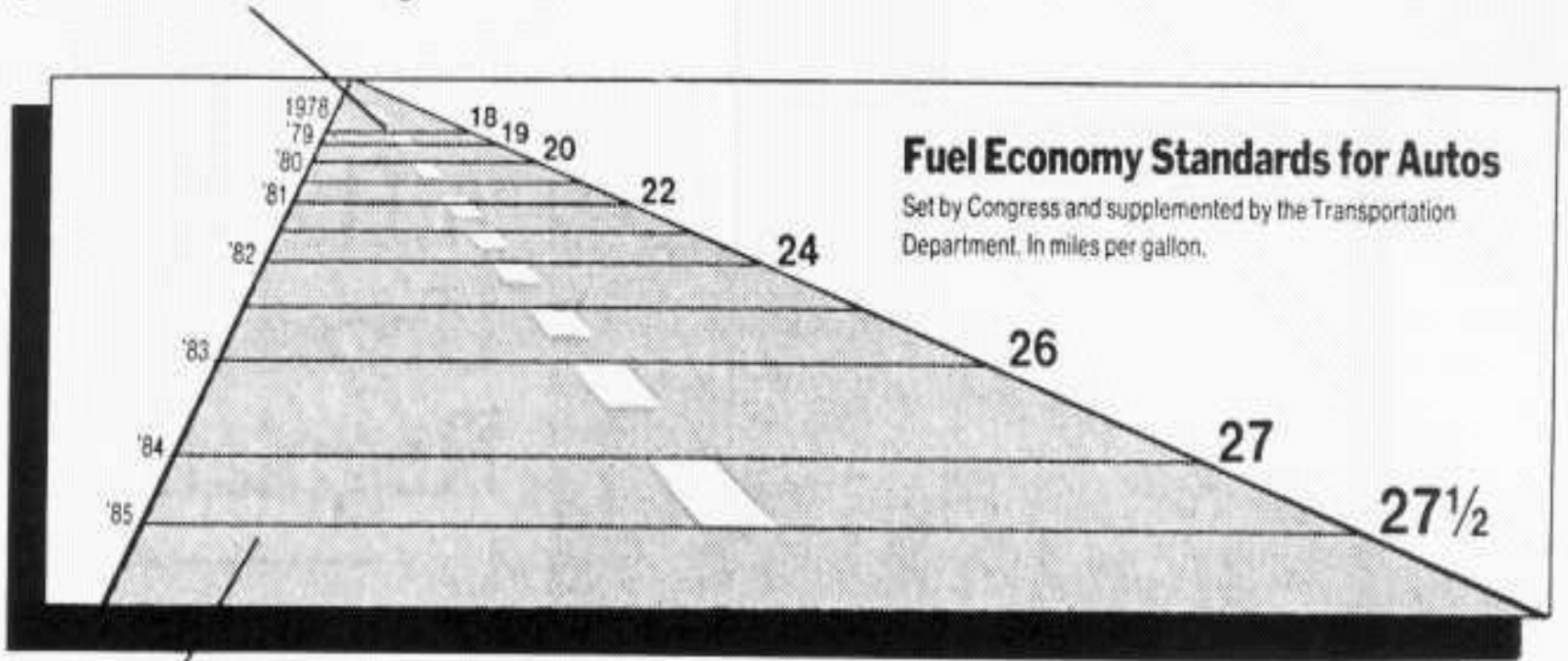
Richard Hamming: "The purpose of computing is insight, not numbers."

- ❑ Reveal the data at several **layers of detail**.
- ❑ Encourage the eye to **compare** different pieces.
- ❑ Support the statistical and verbal **descriptions** of the data.
- ❑ Maximize the data-ink ratio (Edward Tufte, www.edwardtufte.com)

Data-ink ratio= data-ink/total ink used on the graphic

Tufte's Principles of Graphical Excellence

This line, representing 18 miles per gallon in 1978, is 0.6 inches long.



Fuel Economy Standards for Autos

Set by Congress and supplemented by the Transportation Department. In miles per gallon.

This line, representing 27.5 miles per gallon in 1985, is 5.3 inches long.

Lie Factor=14.8

New York Times, August 9, 1978, p. D-2.

(E.R. Tufte, "The Visual Display of Quantitative Information", 2nd edition)

Lie Factor

$$\text{Lie Factor} = \frac{\text{size of effect shown in graphic}}{\text{size of effect in data}} =$$

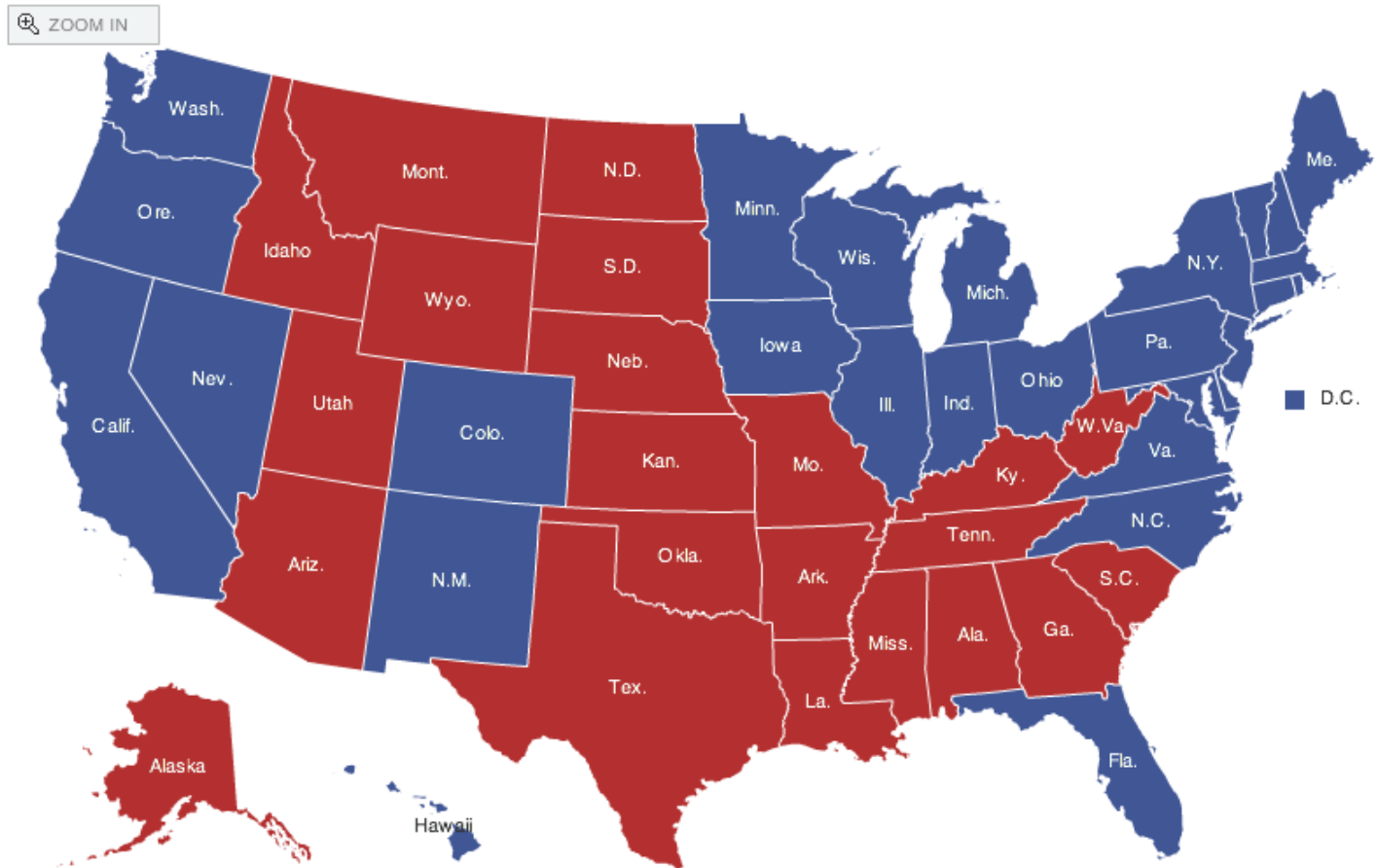
$$= \frac{\frac{(5.3 - 0.6)}{0.6}}{\frac{(27.5 - 18.0)}{18}} = \frac{7.833}{0.528} = 14.8$$

Tufte requirement: $0.95 < \text{Lie Factor} < 1.05$

Scientific Visualization

- ❑ *Visualization for scientific computing*, shortened to *scientific visualization*, was coined in 1987 and refers to the science or methodology of quickly and effectively displaying scientific data.
 - 1987 NSF report: “Visualization in Scientific Computing”
- ❑ Formal name given to the field in computer science that includes user interface, data representation and processing algorithms, visual representations and other sensory presentation such as sound or touch.
- ❑ Visual modelling of scientific data using computer graphics
- ❑ Examples
 - visualizations of protein docking
 - Molecular structural data is hard to understand without visualization
 - Visualization of brain models
- ❑ **Focus** is
 - on modelling (visually) the input data as close to reality as possible
 - Not on presenting abstractions or relationships from the input data

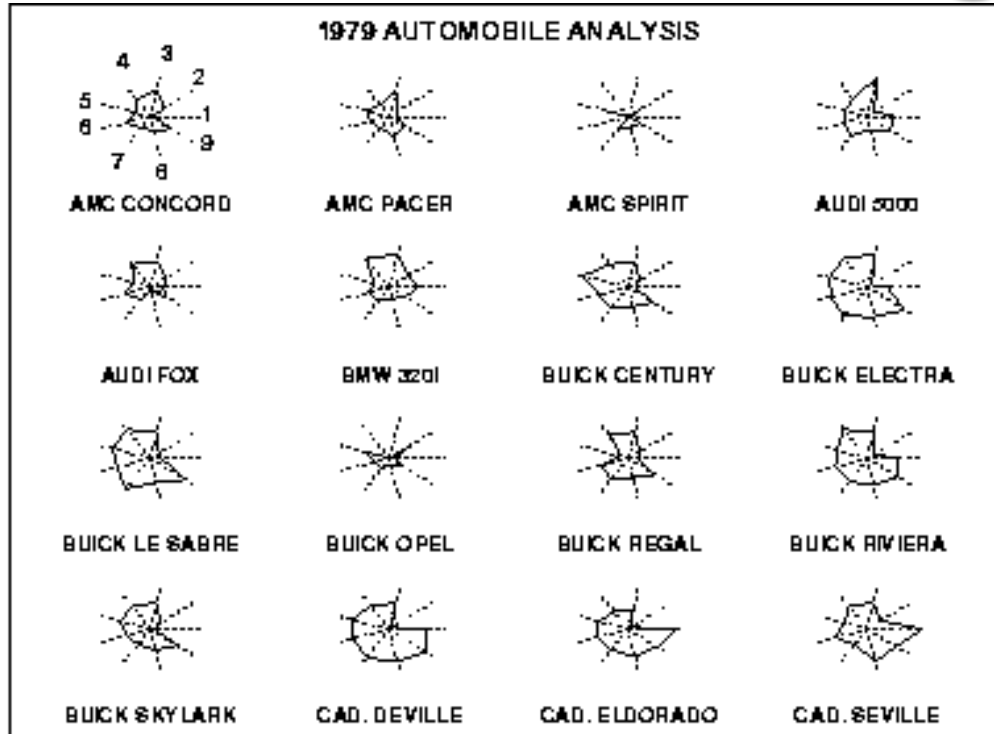
Spatial data: choropleth Maps



- ❑ Maps using color shadings to represent numerical values are called choropleth maps

<http://elections.nytimes.com/2008/results/president/map.html>

Using Icons to Encode Information, e.g., Star Plots



1 Price

2 Mileage (MPG)

3 1978 Repair Record (1 = Worst, 5 = Best)

4 1977 Repair Record (1 = Worst, 5 = Best)

5 Headroom

6 Rear Seat Room

7 Trunk Space

8 Weight

9 Length

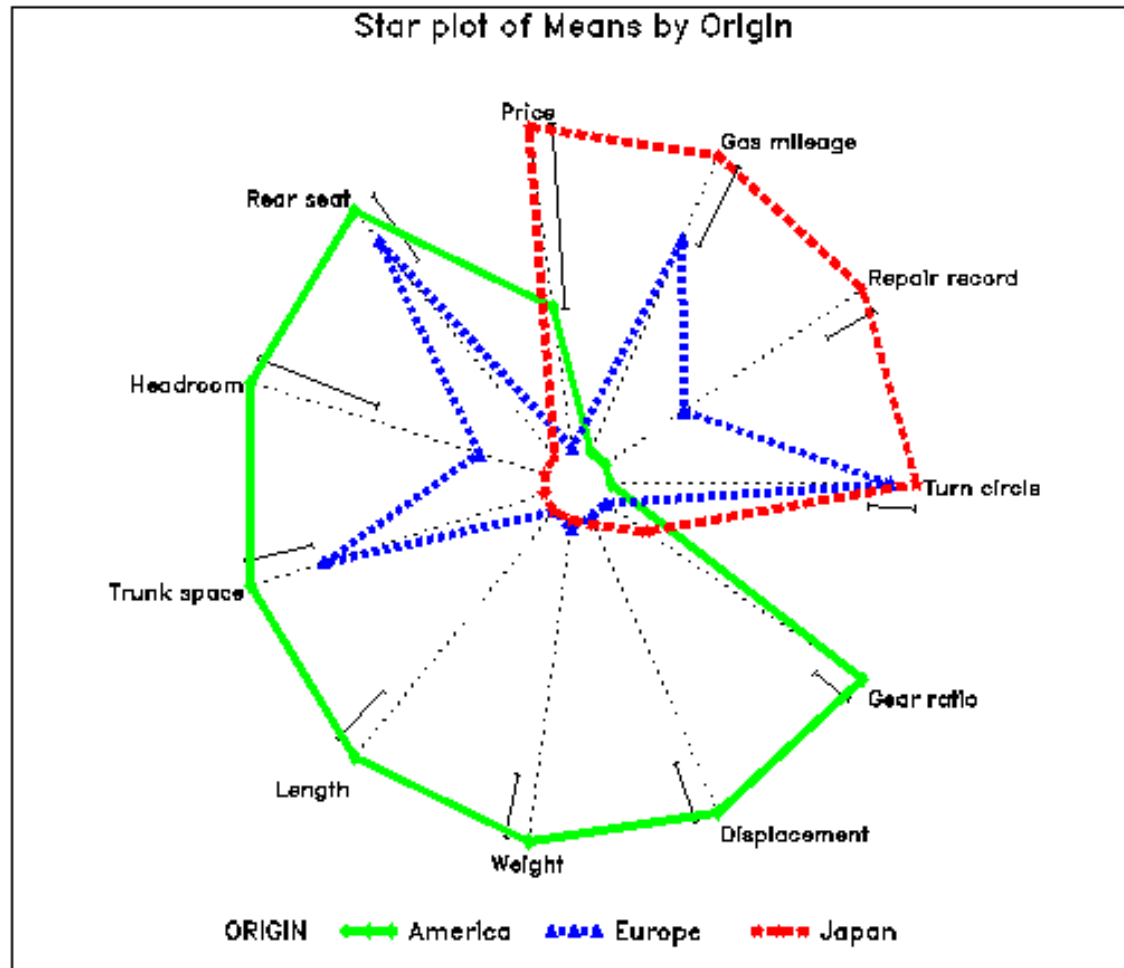
- Each star represents a single observation. Star plots are used to examine the relative values for a single data point
- The star plot consists of a sequence of equi-angular spokes, called radii, with each spoke representing one of the variables.

Useful for small data sets with up to 10 or so variables

□ Limitations?

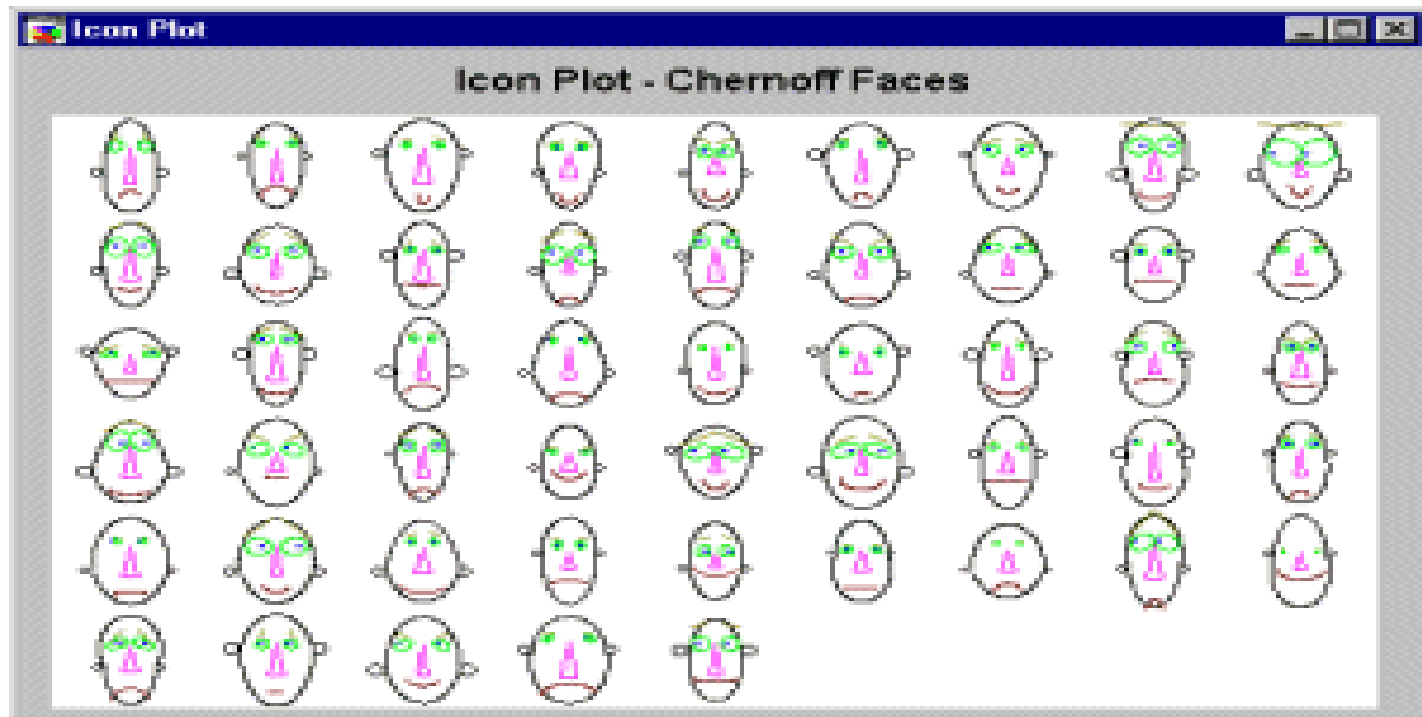
- Small data sets, small dimensions
- Ordering of variables may affect perception

Starplots



Chernoff's Faces

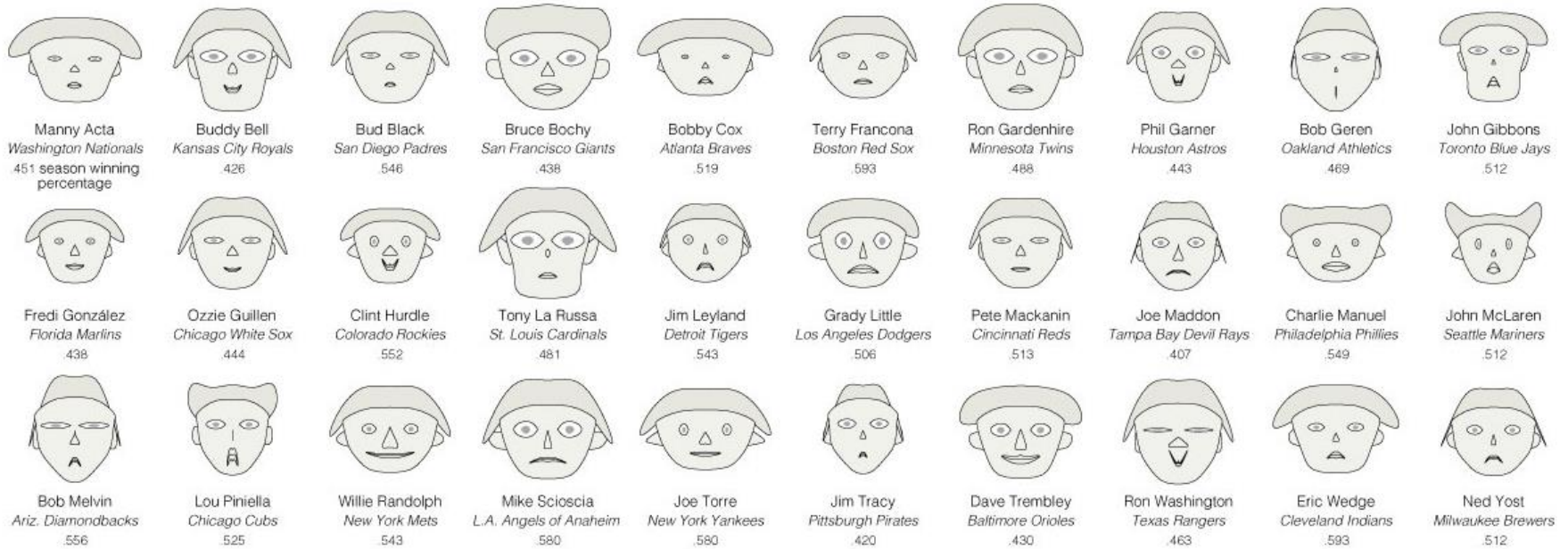
- described by ten facial characteristic parameters: head eccentricity, eye eccentricity, pupil size, eyebrow slant, nose size, mouth shape, eye spacing, eye size, mouth length and degree of mouth opening
- Much derided in statistical circles



Chernoff faces

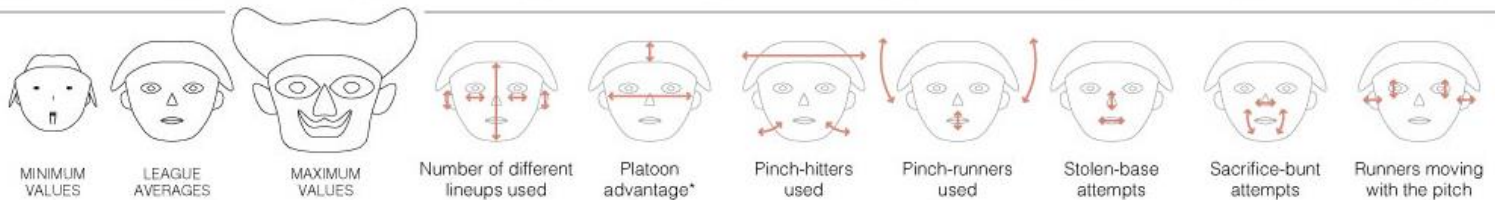
The New York Times

April 1, 2008



SMILE IF YOU BUNT

Steve C. Wang, an associate professor of statistics at Swarthmore College, charted baseball managers from the 2007 season as Chernoff faces, a method of using the heights, widths and angles of facial features to represent different sets of numbers.



*Percentage of players who had the advantage of batting against an opposite-handed pitcher at the start of the game.

Note: Because different rules cause National League managers to use more pinch-hitters, for example, each manager's rates are compared with his league's average.

JONATHAN CORUM/
THE NEW YORK TIMES

Chernoff faces

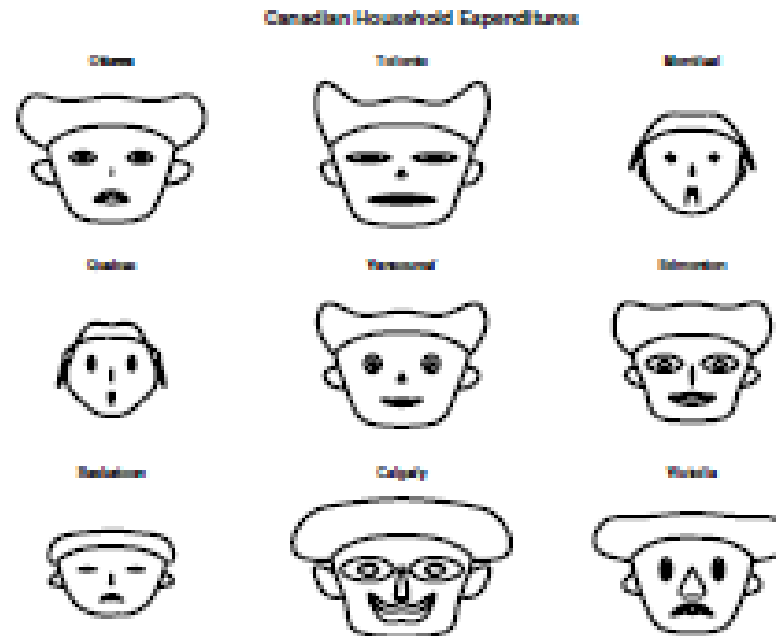
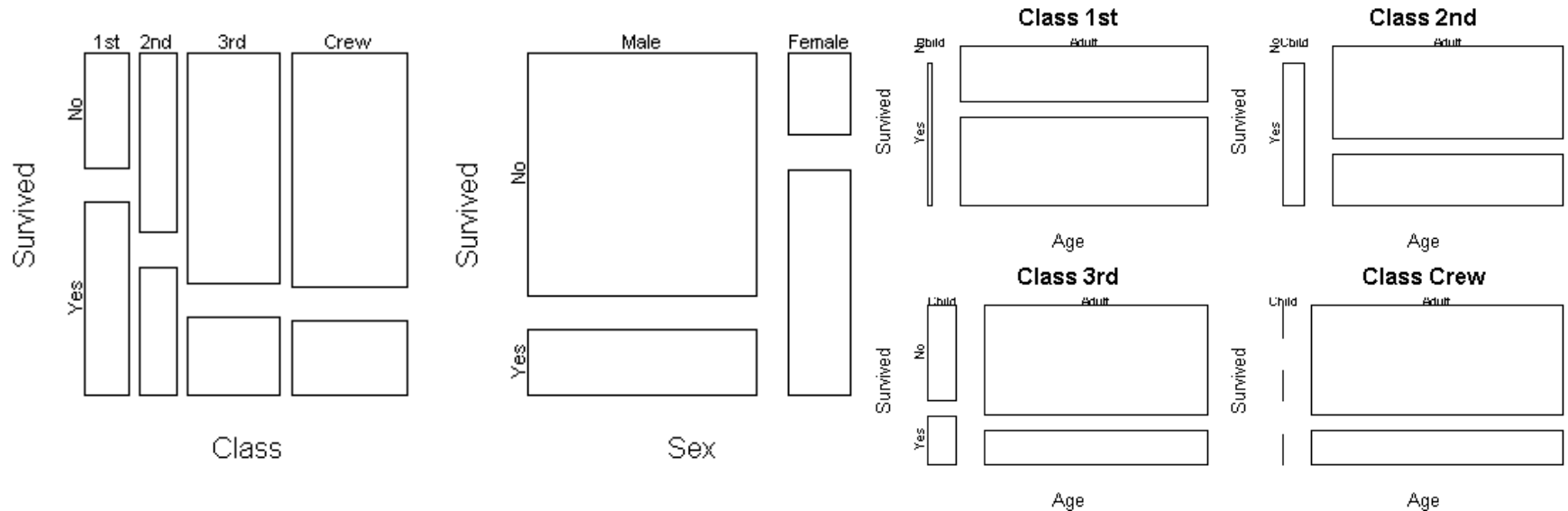


Fig. 9.9 Chernoff faces for household expenditures for nine Canadian cities. Note that food and shelter primarily determine the size of the face. It is great fun to try to match the expression on the face to the character of the city. This shows how exquisitely sensitive we are to facial expressions, underscoring the strength of the method.

Source: Principles and Theory for Data Mining and Machine Learning, Bertrand Clarke et al., Springer, 2012

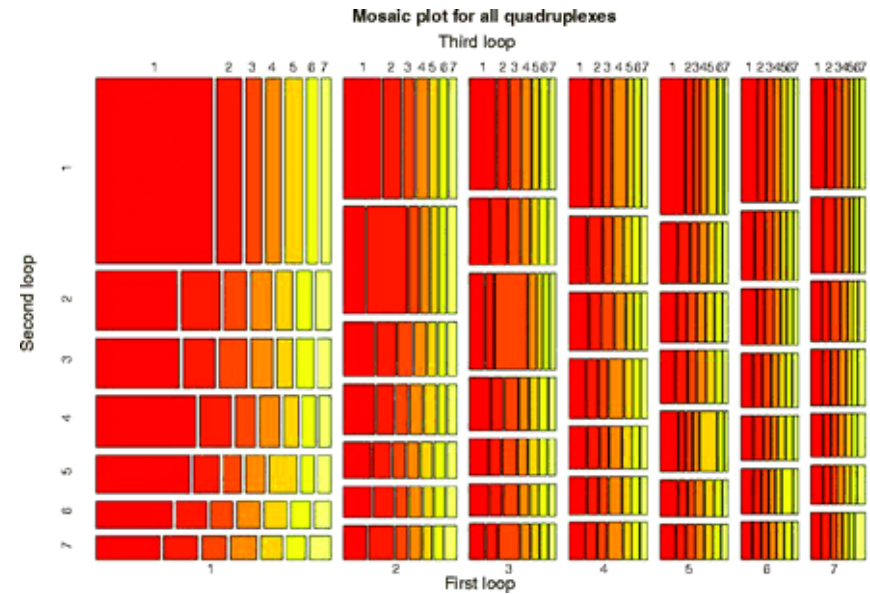
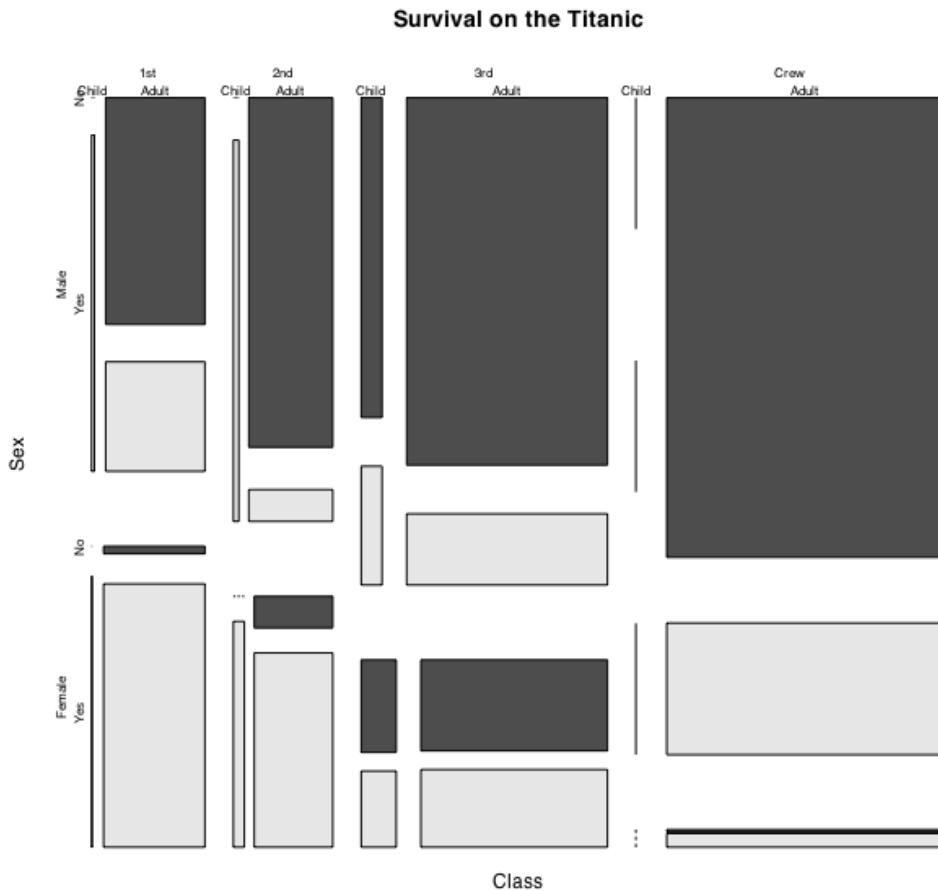
Mosaic Plots

- ❑ Very useful for many categorical variables
- ❑ sensitive to the order which they are applied
- ❑ Titanic Data:



Mosaic plots

Can be effective, but can get out of hand:



Heatmap

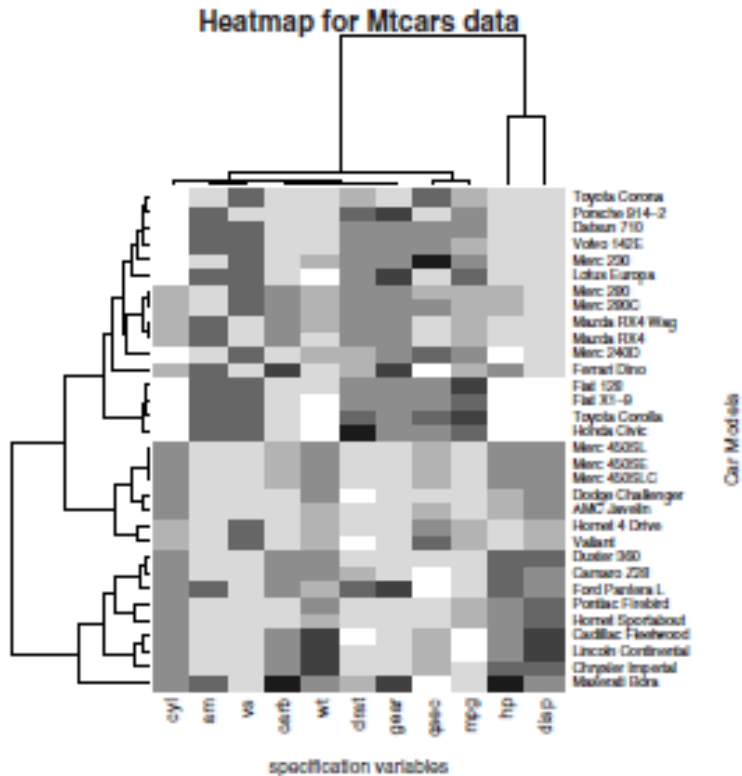


Fig. 9.5 Heatmap for the MTCars data, with clustering on the models and variables done separately. Darker regions correspond to higher values, lighter regions to lower values.

Inventor: Cormac Kinney, 1991,

Source: Principles and Theory for Data Mining and Machine Learning, Bertrand Clarke et al., Springer, 2012

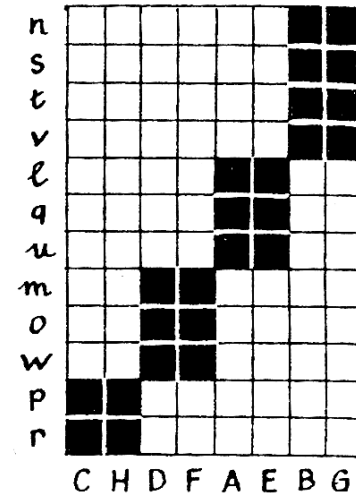
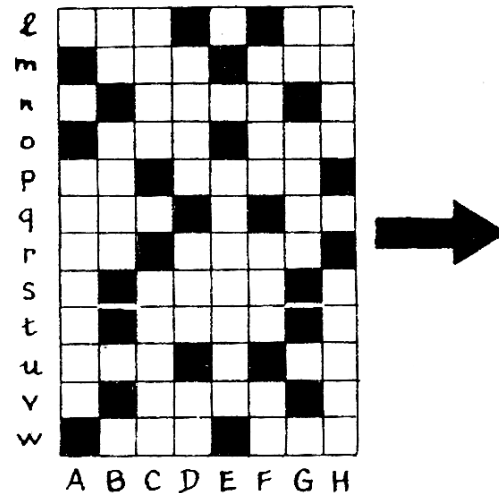
- A heatmap is a matrix of values that have been color coded, usually so that higher values are brighter and lower values are darker, in analogy with temperature.

- The columns and rows of the matrix typically have an interpretation.

- In genomics, The columns and rows represent gene and brighter color represent higher gene expression.

Heatmap

- To rearrange the rows and columns so that two dark blocks place in the main diagonal and two lighter blocks in off the main diagonal is the ideal way to represent a Heat Map



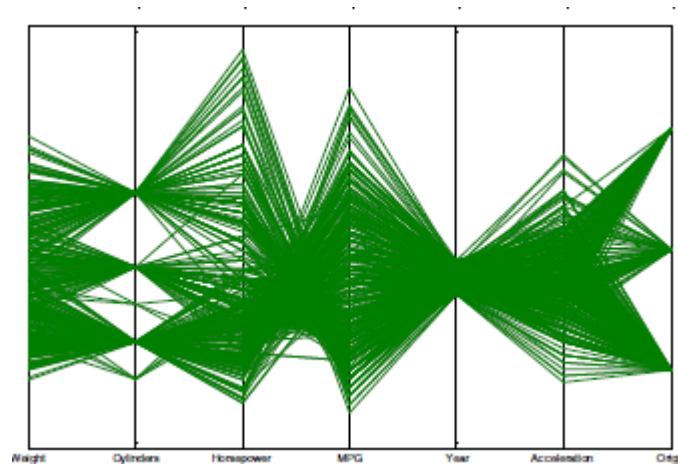
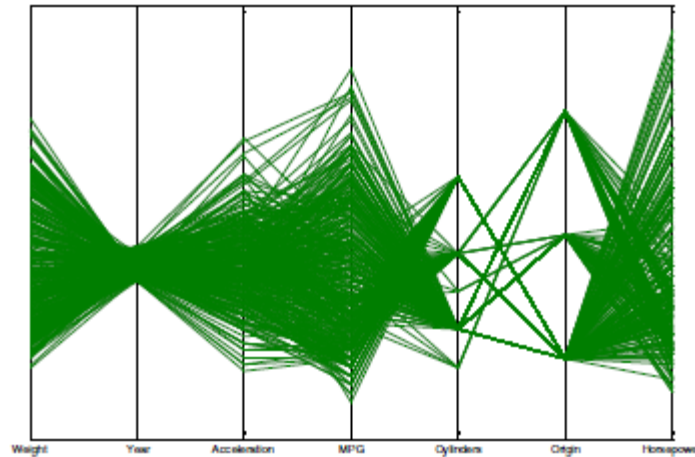
Multi-dimensional data

- Two solutions
 - Plot all possible pairs of variables as 2D scatter plots
 - Simple but not helpful to visualize the data as a whole
 - Parallel coordinates
 - A novel way of plotting multi-dimensional data proposed by **Alfred Inselberg**

Parallel Coordinates

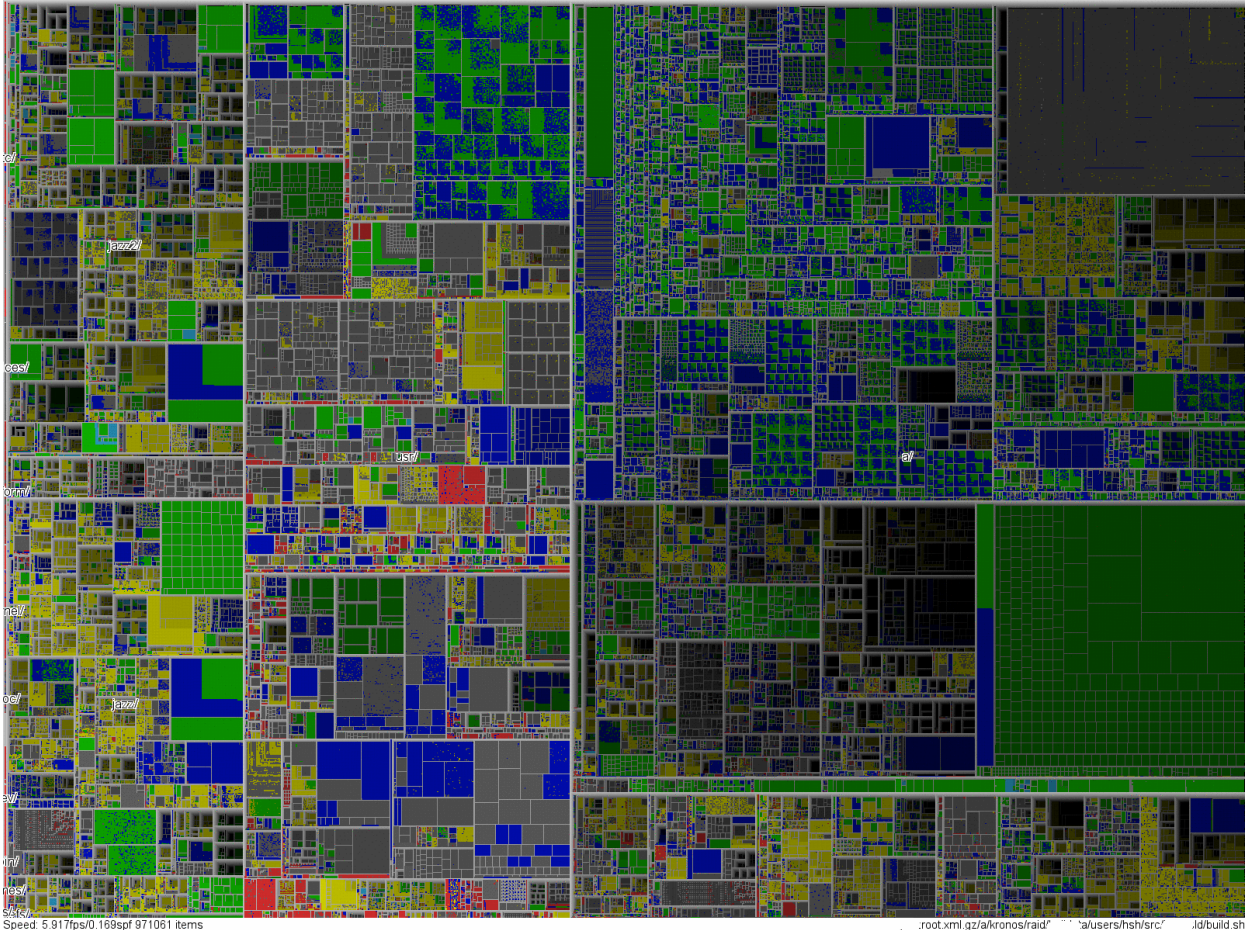
- ❑ One vertical bar per dimension drawn in parallel
- ❑ Each point is represented by a set of lines
- ❑ Revolutionary representation for multi-dimensional data
- ❑ But users may need long time to learn how to understand the graphs

Parallel Coordinates



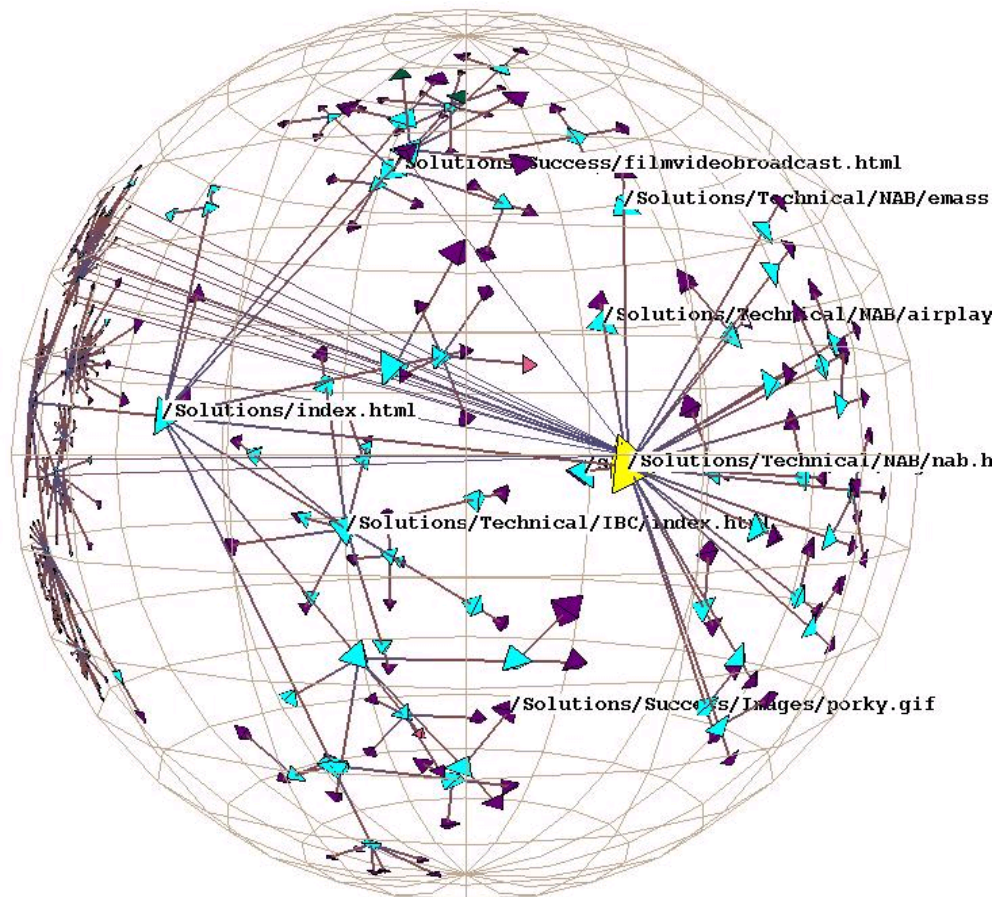
- Reordering of the axes so that crossing angle is minimized increases visualization efficiency

TreeMap of One million items



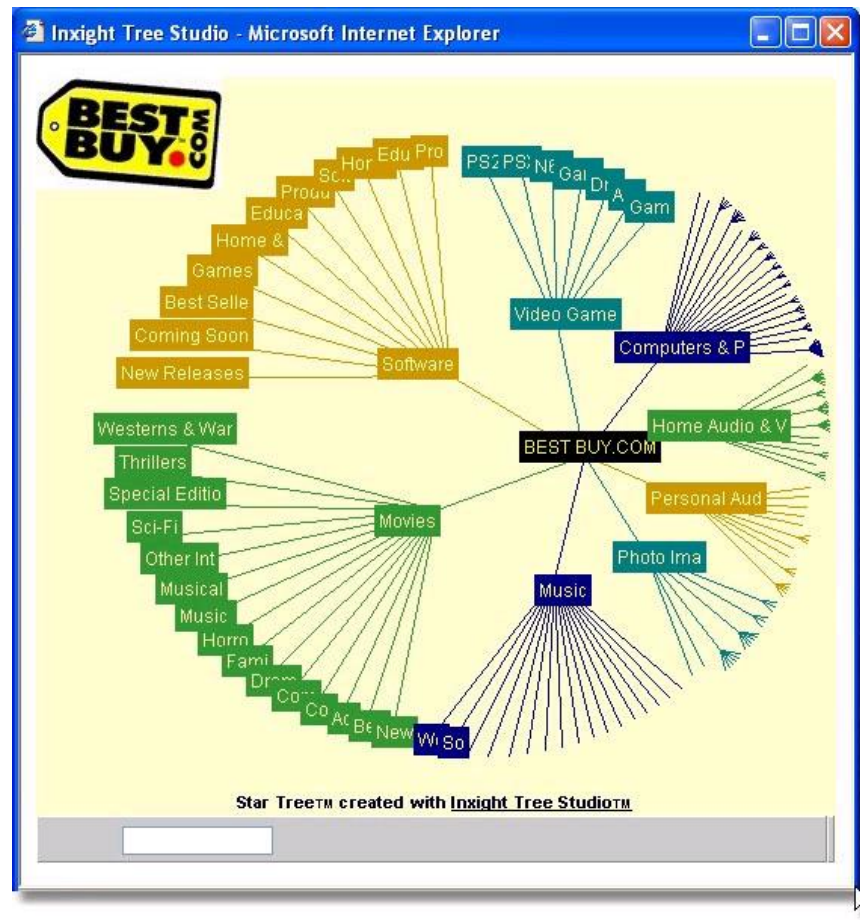
Tree

3-D hyperbolic space (Munzner, 2000)



Visualization of a single Website

StarTree (by InXight)



The Curse of Dimensionality

- The term *curse of dimensionality* was coined by **Richard E. Bellman**

$$\mathbf{X} \sim \text{MVN}_p(\mathbf{0}, \mathbf{I})$$

- Gaussian kernel density estimation
- Bandwidth chosen to minimize MSE at the mean
- Suppose want:

$$\frac{E[(\hat{p}(x) - p(x))^2]}{p(x)^2} < 0.1 \Big|_{x=0}$$

- In order to obtain a statistically sound and reliable result, the amount of data needed to support the result often grows exponentially with the dimensionality.

Dimension **# data points**

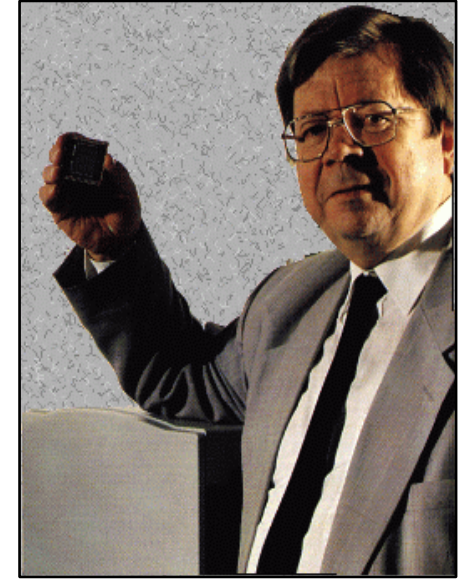
1	4
2	19
3	67
6	2,790
10	842,000

- The volume of the spaces increases very rapidly
- The Availability of the data becomes sparse
- *This sparsity is problematic* for any method that requires statistical significance.
- Organizing and searching data often relies on detecting areas where objects form groups with similar properties
- In high dimensional data however all objects appear to be sparse and dissimilar in many ways.

Self-Organizing Maps

Self-Organizing Maps

- Ideas first introduced by C. von der Malsburg (1973), developed and refined by T. Kohonen (1982)
- Neural network algorithm using unsupervised competitive learning
- Primarily used for organization and visualization of complex data
- Biological basis: 'brain maps'

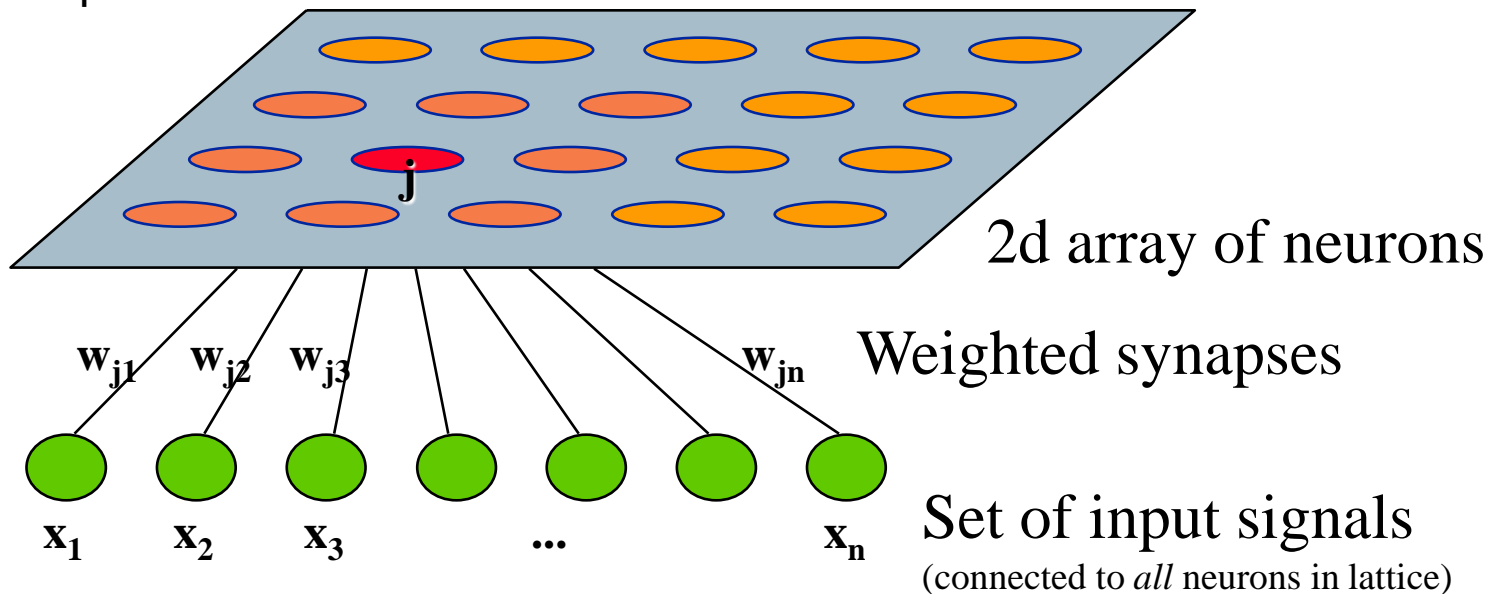


Teuvo Kohonen

Self-Organizing Maps

SOM - Architecture

- Lattice of neurons ('nodes') accepts and responds to set of input signals
- Responses compared; 'winning' neuron selected from lattice
- Selected neuron activated together with 'neighbourhood' neurons
- Adaptive process changes weights to more closely resemble inputs



Self-Organizing Maps

- ❑ The brain cells are self organizing themselves in groups, according to incoming information.
- ❑ This incoming information is not only received by a single neural cell, but also influences other cells in its neighbourhood. This organisation results in some kind of map, where Neural cells with similar functions are arranged close together.
- ❑ SOM mechanism is also based on this principle

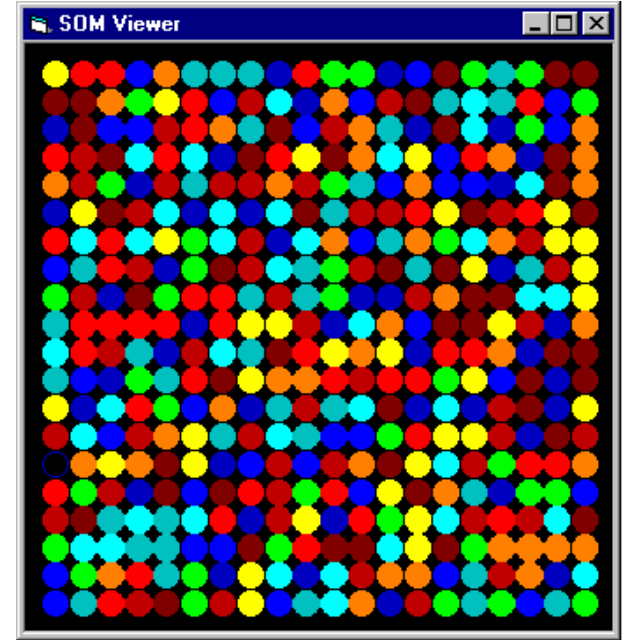
Self-Organizing Maps

SOM – Algorithm Overview

1. Randomly initialise all weights
2. Select input vector $\mathbf{x} = [x_1, x_2, x_3, \dots, x_n]$
3. Compare \mathbf{x} with weights \mathbf{w}_j for each neuron j to determine winner
4. Update winner so that it becomes more like \mathbf{x} , together with the winner's *neighbours*
5. Adjust parameters: learning rate & 'neighbourhood function'
6. Repeat from (2) until the map has converged (i.e. no noticeable changes in the weights) or pre-defined no. of training cycles have passed

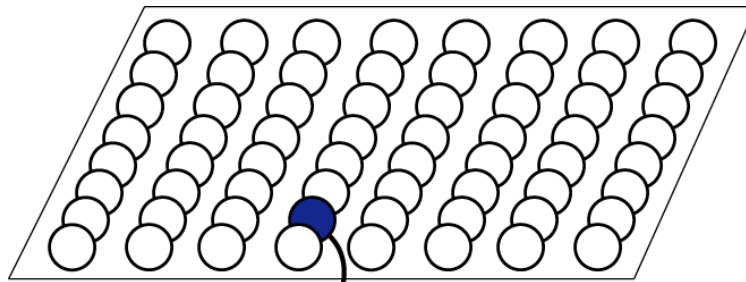
Initialization

(i) Randomly initialize the weight vectors \mathbf{w}_j for all nodes j



Input vector

(ii) Choose an input vector \mathbf{x} from the training set



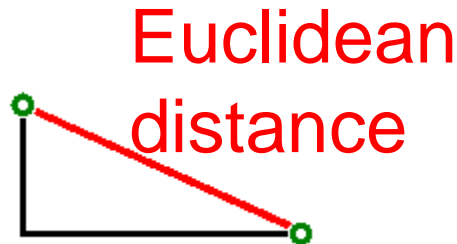
Dokumentenkarten	1
Tiere	0
Texte	1
Michael	1
Schuhmacher	1
Dokumentenvektoren	1

Finding a Winner

(iii) Find the best-matching neuron $w(\mathbf{x})$, usually the neuron whose weight vector has **smallest Euclidean distance** from the input vector \mathbf{x}

The winning node is that which is in some sense 'closest' to the input vector
'Euclidean distance' is the straight line distance between the data points, if they were plotted on a (multi-dimensional) graph

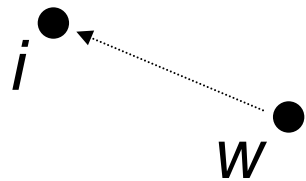
Euclidean distance between two vectors \mathbf{a} and \mathbf{b} , $\mathbf{a} = (a_1, a_2, \dots, a_n)$, $\mathbf{b} = (b_1, b_2, \dots, b_n)$, is calculated as:



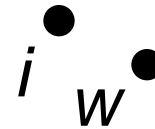
$$d_{\mathbf{a}, \mathbf{b}} = \sqrt{\sum_i (a_i - b_i)^2}$$

Learning the SOM

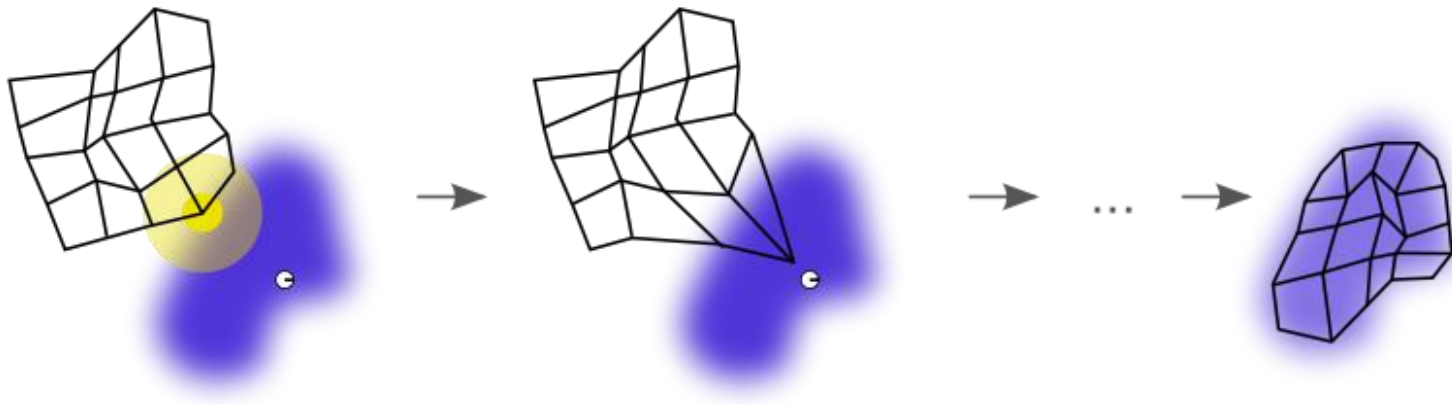
- Determine the winner (the neuron of which the weight vector has the smallest distance to the input vector)
- Move the weight vector w of the winning neuron towards the input i



Before learning



After learning



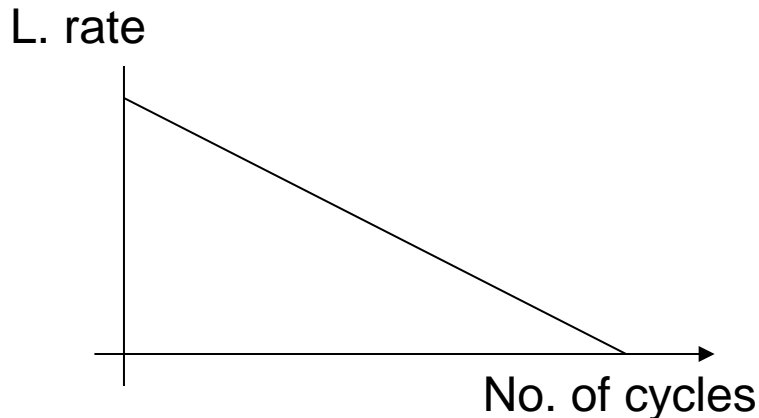
Weight Update

SOM Weight Update Equation

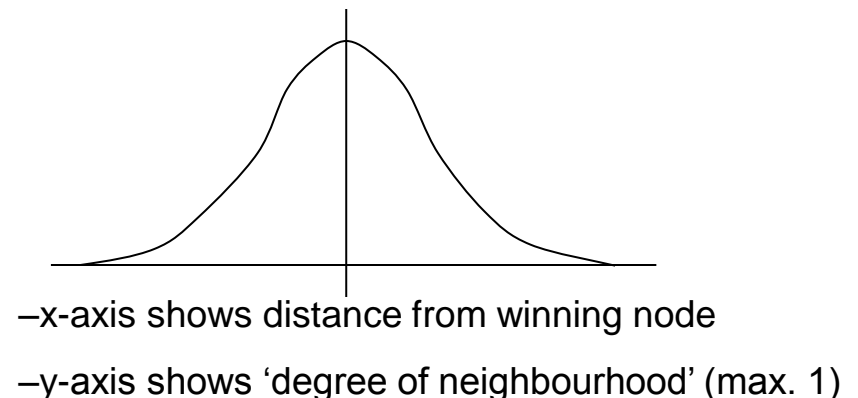
$$\mathbf{w}_j(t+1) = \mathbf{w}_j(t) + \mu(t) \lambda_{\omega(x)}(j,t) [\mathbf{x} - \mathbf{w}_j(t)]$$

“The weights of every node are updated at each cycle by adding
Current learning rate \times Degree of neighbourhood with respect to winner \times
Difference between current weights and input vector
to the current weights”

Example of $\mu(t)$

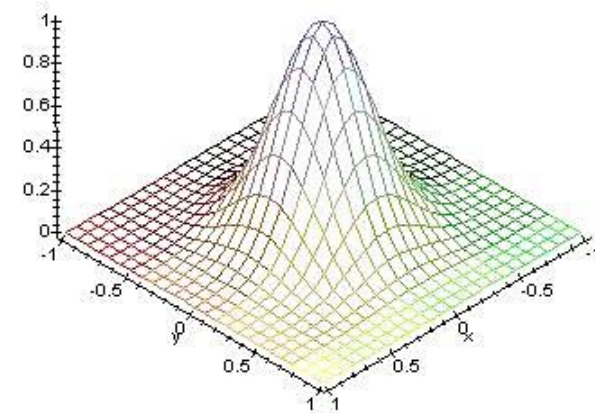
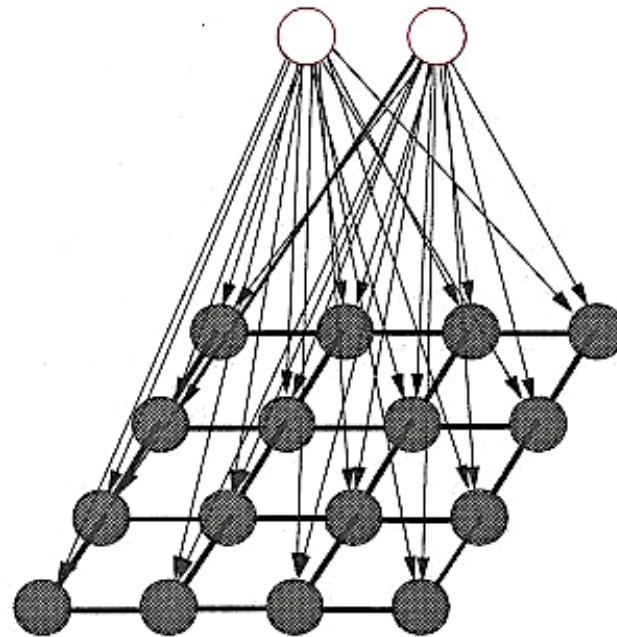


Example of $\lambda_{\omega(x)}(j,t)$



Self Organizing Map

- ❑ Neighborhood function to preserve topological properties of the input space
- ❑ Neighbors share the prize (postcode lottery principle)



Self Organizing Map

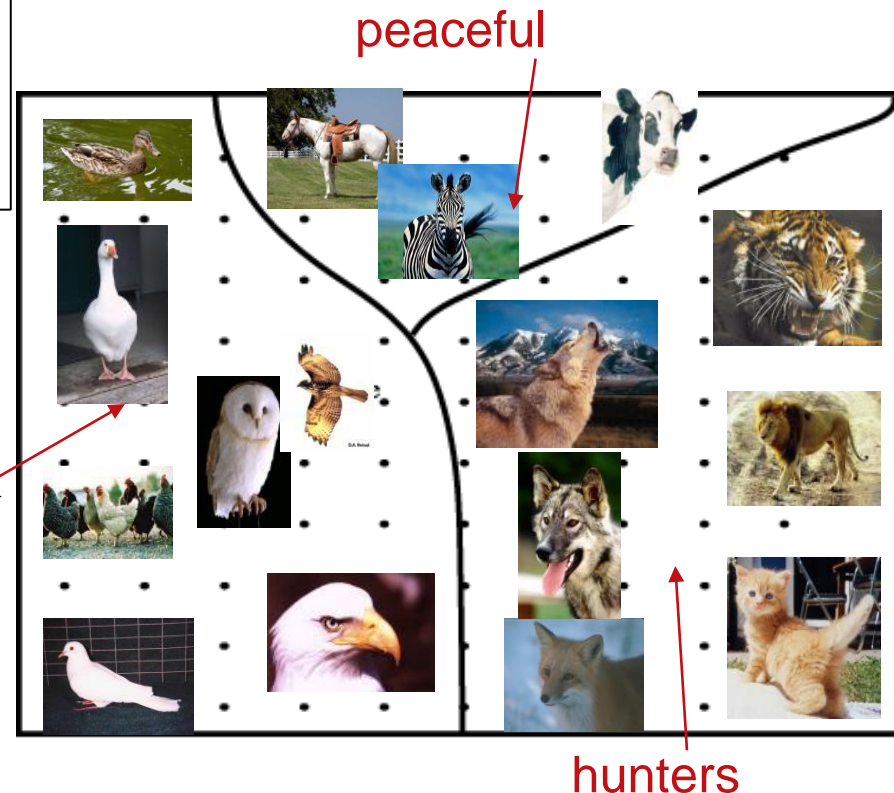
- ❑ Input: high-dimension input space
- ❑ Output: low dimensional (typically 2 or 3)
 - ❑ network topology
- ❑ Training
 - ❑ Starting with a large learning rate and neighborhood size, both are gradually decreased to facilitate convergence
- ❑ After learning, neurons with similar weights tend to cluster on the map

Example: Self-Organizing Maps

Animal names and their attributes

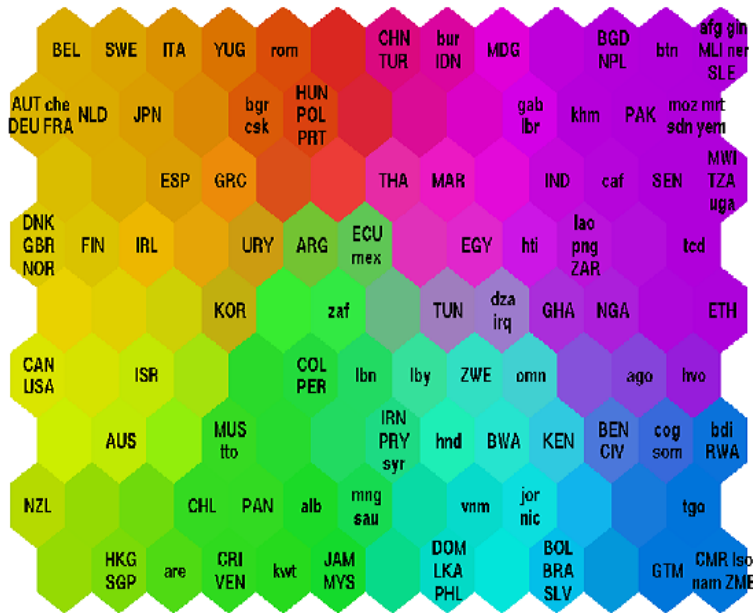
	Dove	Hen	Duck	Goose	Owl	Hawk	Eagle	Fox	Dog	Wolf	Cat	Tiger	Lion	Horse	Zebra	Cow
is	Small	1	1	1	1	1	1	0	0	0	1	0	0	0	0	0
	Medium	0	0	0	0	0	1	1	1	1	0	0	0	0	0	0
	Big	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1
has	2 legs	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
	4 legs	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
	Hair	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
	Hooves	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1
	Mane	0	0	0	0	0	0	0	0	1	0	0	1	1	1	0
likes	feathers	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
to	Hunt	0	0	0	0	1	1	1	1	0	1	1	1	0	0	0
	Run	0	0	0	0	0	0	0	1	1	0	1	1	1	1	0
	Fly	1	0	0	1	1	1	1	0	0	0	0	0	0	0	0
	Swim	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0

A grouping according to similarity has emerged



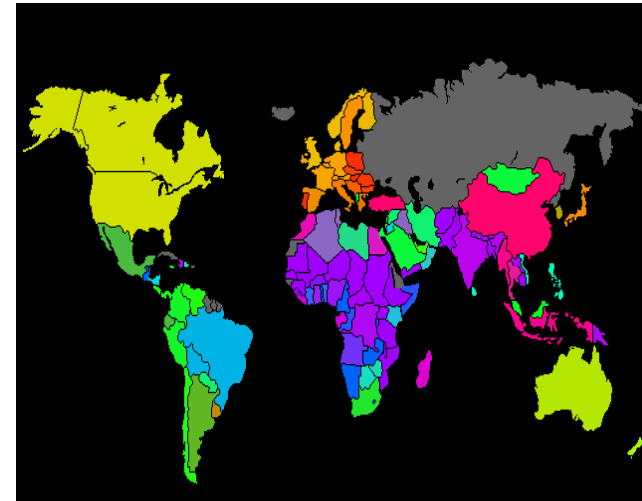
Self-Organizing Maps

SOM – Result Example Classifying World Poverty



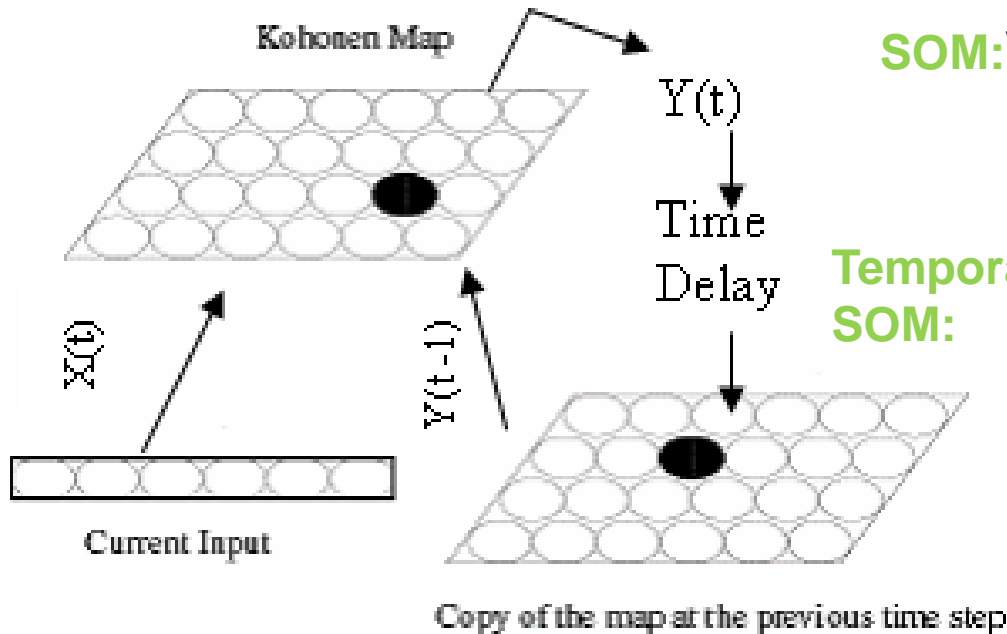
The Country Names

AFG	Afghanistan	GTM	Guatemala	NZL	New Zealand
AGO	Angola	HKG	Hong Kong	OMN	Oman
ALB	Albania	IDN	Indonesia	OMR	Oman
ARE	United Arab Emirates	IND	India	PAK	Pakistan
ARG	Argentina	IRN	Iran	PAN	Panama
ARM	Armenia	IRQ	Iraq	PAR	Paraguay
AUT	Austria	ISR	Israel	PER	Peru
AZE	Azerbaijan	ITA	Italy	PHL	Philippines
BAN	Banladesh	JAM	Jamaica	PRT	Portugal
BEL	Belgium	JOR	Jordan	ROM	Romania
BEN	Benin	KAZ	Kazakhstan	RUS	Russia
BGD	Bangladesh	KEN	Kenya	SAC	Swaziland
BGR	Bulgaria	KHM	Khmer Rep.	SAR	Saudi Arabia
BHS	Bahamas	KOR	South Korea	SEN	Senegal
BOL	Bolivia	KWT	Kuwait	SGP	Singapore
BRA	Brazil	LAO	Laos PDR	SLV	El Salvador
BRE	Brexit	LBN	Lebanon	SOM	Somalia
BUL	Bulgaria	LES	Lesotho	SWE	Sweden
BUR	Burkina Faso	LKA	Sri Lanka	SVK	Slovakia
BUR	Burundi	LUX	Luxembourg	THA	Thailand
BUR	Burundi	MDG	Madagascar	TIM	Timor-Leste
BUR	Burundi	MEX	Mexico	TUN	Tunisia
BUR	Burundi	MIL	Maldives	TUR	Turkey
BUR	Burundi	MNI	Mongolia	UGA	Uganda
BUR	Burundi	MOR	Morocco	USA	United States
BUR	Burundi	MRT	Mauritius	VEN	Venezuela
BUR	Burundi	MUS	Mauritius	VNM	Vietnam
BUR	Burundi	MWI	Malawi	YEM	Yemen
BUR	Burundi	MYS	Malaysia	YUG	Yugoslavia
BUR	Burundi	NAM	Namibia	ZAF	South Africa
BUR	Burundi	NER	Niger	ZAR	Zimbabwe
BUR	Burundi	NEA	Nepal	ZMB	Zambia
BUR	Burundi	NLD	Netherlands	ZWE	Zimbabwe
BUR	Burundi	NOR	Norway		
BUR	Burundi	NPL	Nepal		



‘Poverty map’ based on 39 indicators describing various quality-of-life factors, such as state of health, nutrition, educational services, etc.) from World Bank statistics (1992)

Temporal SOM



$$\text{SOM: } \mathbf{w}_j(t+1) = \mathbf{w}_j(t) + \mu(t) \lambda_{\omega(x)}(j,t) [\mathbf{x}(t) - \mathbf{w}_j(t)]$$

Temporal SOM:

$$\mathbf{y}_j(t) = (1-\alpha) \mathbf{y}_j(t-1) + \alpha [\mathbf{x}(t) - \mathbf{w}_j(t)]$$

$$\mathbf{w}_j(t+1) = \mathbf{w}_j(t) + \mu(t) \lambda_{\omega(x)}(j,t) \mathbf{y}_j(t)$$

$[\mathbf{x}(t) - \mathbf{w}_j(t)]$ is called quantization error

$\mathbf{y}_j(t)$ is called recursive of the neuron j at time t

α is the leaking coefficient (value is between 0 and 1)

Recursive SOM architecture. The original SOM algorithm is used recursively, on both the input vector $X(t)$ and the representation derived at the previous time step, $y(t-1)$

Self-Organizing Maps

– Advantages

- SOM is Algorithm that projects high-dimensional data onto a two-dimensional map.
- The projection preserves the topology of the data so that similar data items will be mapped to nearby locations on the map.
- SOM still have many practical applications in pattern recognition, speech analysis, industrial and medical diagnostics, data mining

– Disadvantages

- Large quantity of good quality representative training data required
- No generally accepted measure of 'quality' of a SOM
e.g. Average quantization error (how well the data is classified)

Outline

- Introduction
- History
- Modern Visualization and Computational Problems
 - Parallel Co-ordinates, heat map, SOM etc.
 - Space, data ordering, processing
- Conclusion

Conclusion

- ❑ The advent of computer capacity and power push the envelope of computational sciences and scientific visualization (SciVis).
- ❑ SciVis has revolutionized the way we do sciences.
- ❑ SciVis provides scientists a process to probe into enormously large data sets, perceive incredible details of the domain, and discover unexpected insights.
- ❑ Challenging issues in SciVis evolve, but we will continue to face them, solve the problems, and face future challenges.

Scientists make music from DNA

By Daniel Woolls
The Associated Press

MADRID, Spain - Imagine the human genome as music. Unravel DNA's double helix, picture its components lined up like piano keys and assign a note to each. Run your finger along the keys.

Spanish scientists did that just for fun and recorded what they call an audio version of the blueprint for life.

The team at Madrid's Ramon y Cajal Hospital was intrigued by music's lure - how it can make toddlers dance and adults cry - and looked for hints in the genetic material that makes us what we are. They also had some microbial genes wax melodic.

The end product is *Genoma Music*, a 10-tune CD due out in February. "It's a way to bring science and music closer together," said Dr. Aurora Sanchez Sousa, a piano-playing microbiologist who specializes in fungi.

DNA, or deoxyribonucleic acid, is composed of long strings of molecules called nucleotides, which

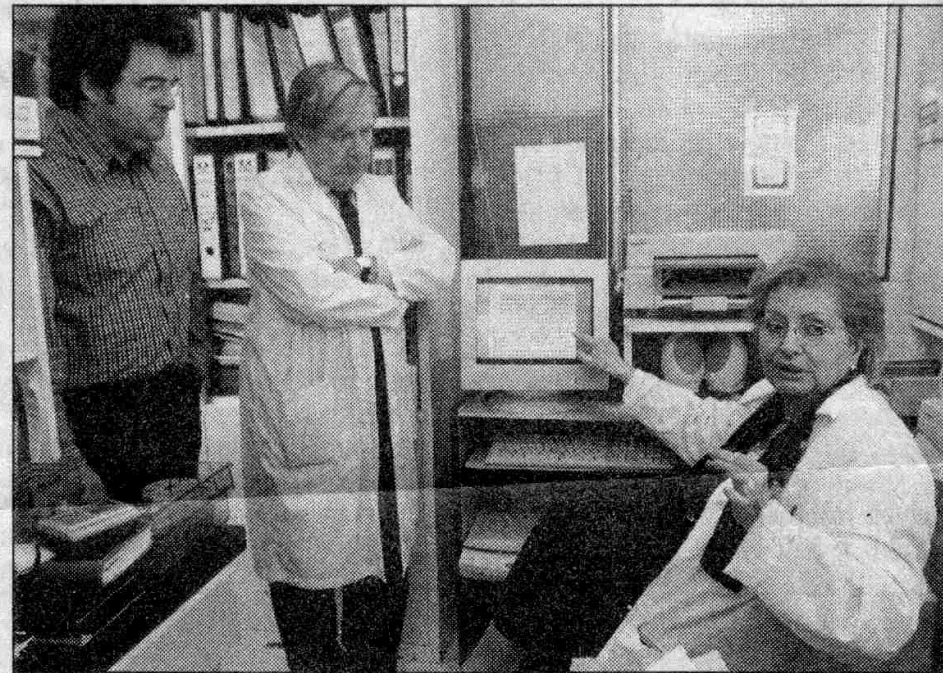
are distinguished by which of four nitrogen-containing bases they contain: adenine, guanine, thymine or cytosine, represented as A, G, T and C. These became the musical notes.

French-born composer Richard Krull turned DNA sequences - a snippet of a gene might look like AGCGTATACGAGT - into sheet music. He arbitrarily assigned tones of the eight-note, do-re-mi scale to each letter. Thymine became re, for instance. Guanine is so, adenine la and cytosine do.

Played solo on percussion, classical guitar or the other instruments used on the CD, the sequences would sound cute but rudimentary, the musical equivalent of PacMan in an era of Microsoft Xbox.

So the alphabet soup of bases served as just that, base lines to accompany melodies composed by Mr. Krull and his scientific colleague. They say the melodies were influenced, even dictated, by the mood and rhythm of the underlying genetic code.

In general, the genome music is



The Associated Press/PAUL WHITE

Dr. Aurora Sanchez Sousa (right) points to some sheet music that has been translated from the DNA code into easy-listening music, inside the Ramon y Cajal hospital in Madrid, Spain.

an easy-listening sound that is vaguely New Age. One of the prettiest songs is based on Connexin 26, a human gene that causes deafness when it mutates. The DNA skeleton is expressed with tinkling bells and a flute melody does the rest.

Another song draws on a yeast gene known as SLT2. Dr. Sanchez Sousa, the main author of the project, is fond of the sequence because it features a stretch in which one triplet of nitrogen bases ap-

pears several times in rapid succession - a repetitive phenomenon that has a musical equivalent called ostinato.

"This is a very sad part, but a beautiful one," Dr. Sanchez Sousa said, wearing a white lab coat and waving her arms like a musical conductor as she played the segment for a visitor.

Her team's plans for future music include having the hospital choir sing a vocal piece based on DNA from a bacteria.



CIMAT

Centro de Investigación en Matemáticas, A.C.

Unidad Monterrey

CIIDIT, UANL, Planta Baja, PIIT Autopista al Aeropuerto Km. 10,
C.P. 66600 Apodaca N.L., México
Tel. +52 (81) 1340 4000

CIDCS, Campus de la Salud, UANL Ave. Carlos Canseco s/n
Col Mitras Centro, Monterrey N.L., México, C.P.66460 México
Tel +52 (81) 8329 4000 Ext 1778
cimat@cimat.mx / www.cimat.mx

