

# Getting to know your data

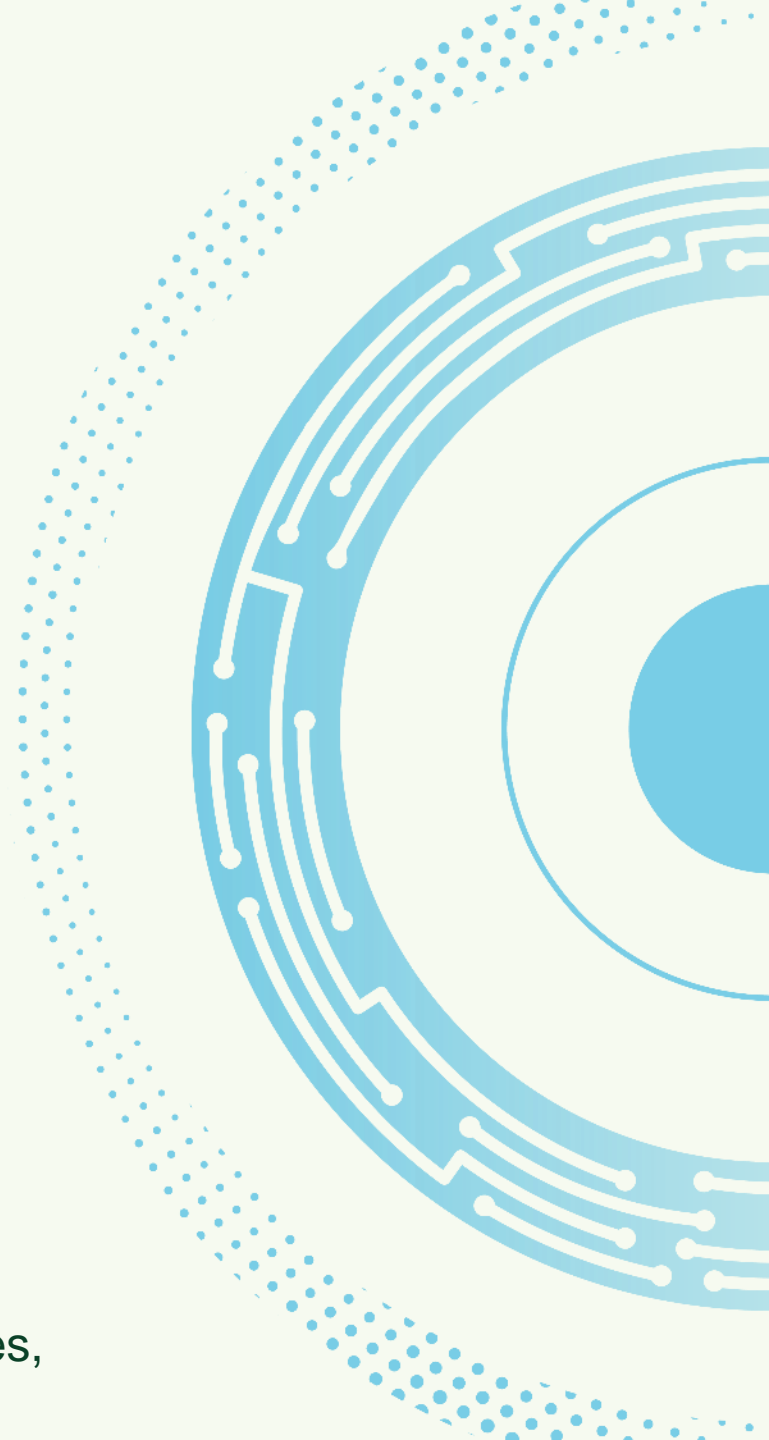
Data visualization

Measuring data similarity and dissimilarity



The Alexander Kofkin  
**Faculty of Engineering**  
Bar-Ilan University

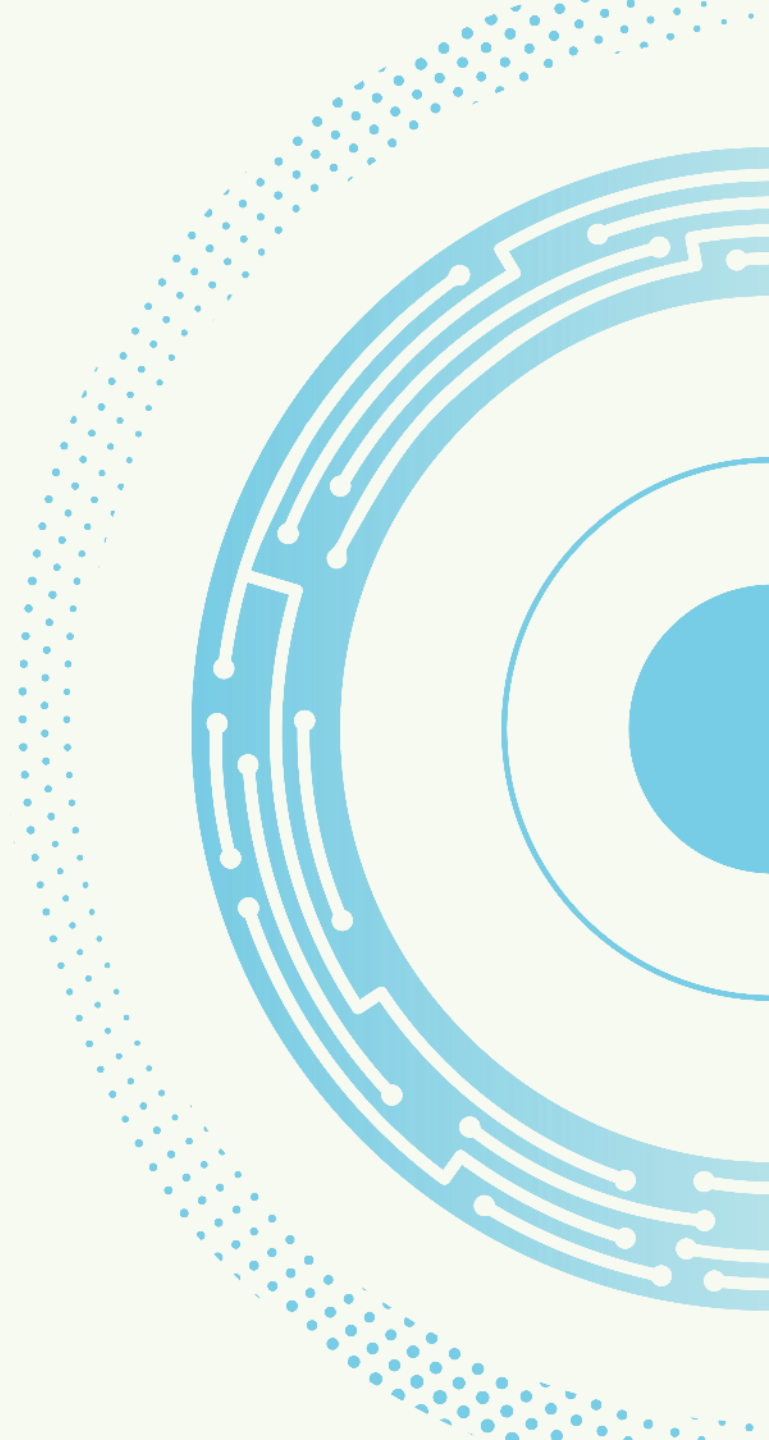
Based on Data Mining – Concept and Techniques,  
Jiawei Han, Micheline Kamber & Jian Pei



# Data visualization



The Alexander Kofkin  
**Faculty of Engineering**  
Bar-Ilan University



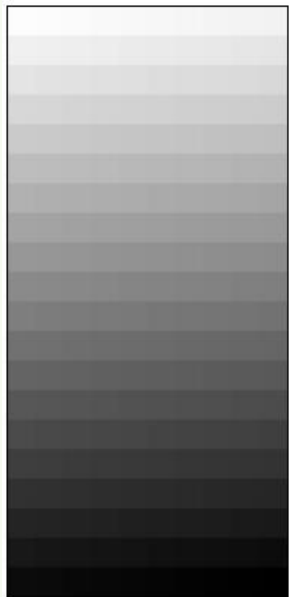
# What is data?

- **Why data visualization?**
  - Gain insight into an information space
  - Provide qualitative overview of large data sets
  - Search for patterns, trends, irregularities, relationships among data
  - Help find interesting regions
- **Categorization of visualization methods:**
  - Pixel-oriented visualization techniques
  - Geometric projection visualization techniques
  - Icon-based visualization techniques
  - Hierarchical visualization techniques
  - Visualizing complex data and relations

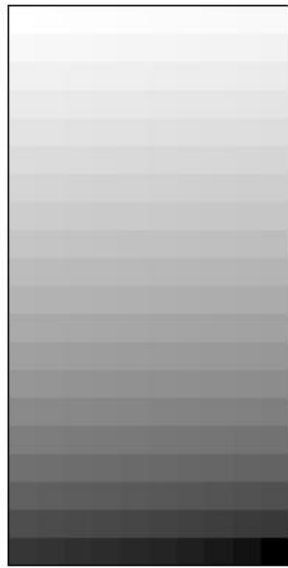


# Pixel-oriented visualization techniques

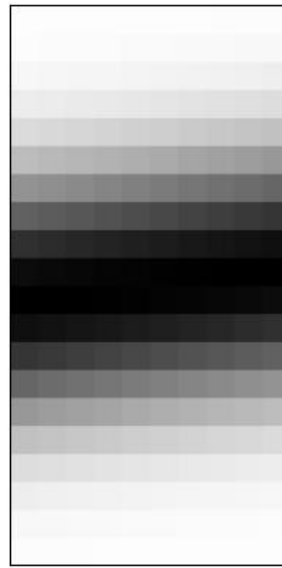
- For a data set of  $m$  dimensions, create  $m$  windows on the screen, one for each dimension
- The  $m$  dimension values of a record are mapped to  $m$  pixels at the corresponding positions in the windows
- The colors of the pixels reflect the corresponding values



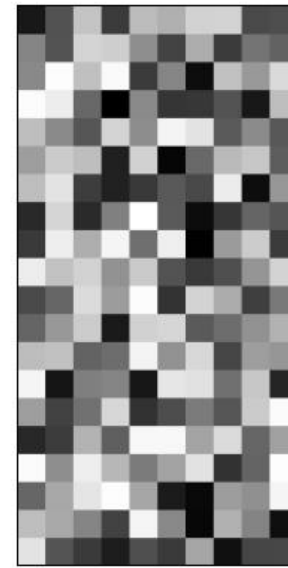
(a) Income



(b) Credit Limit



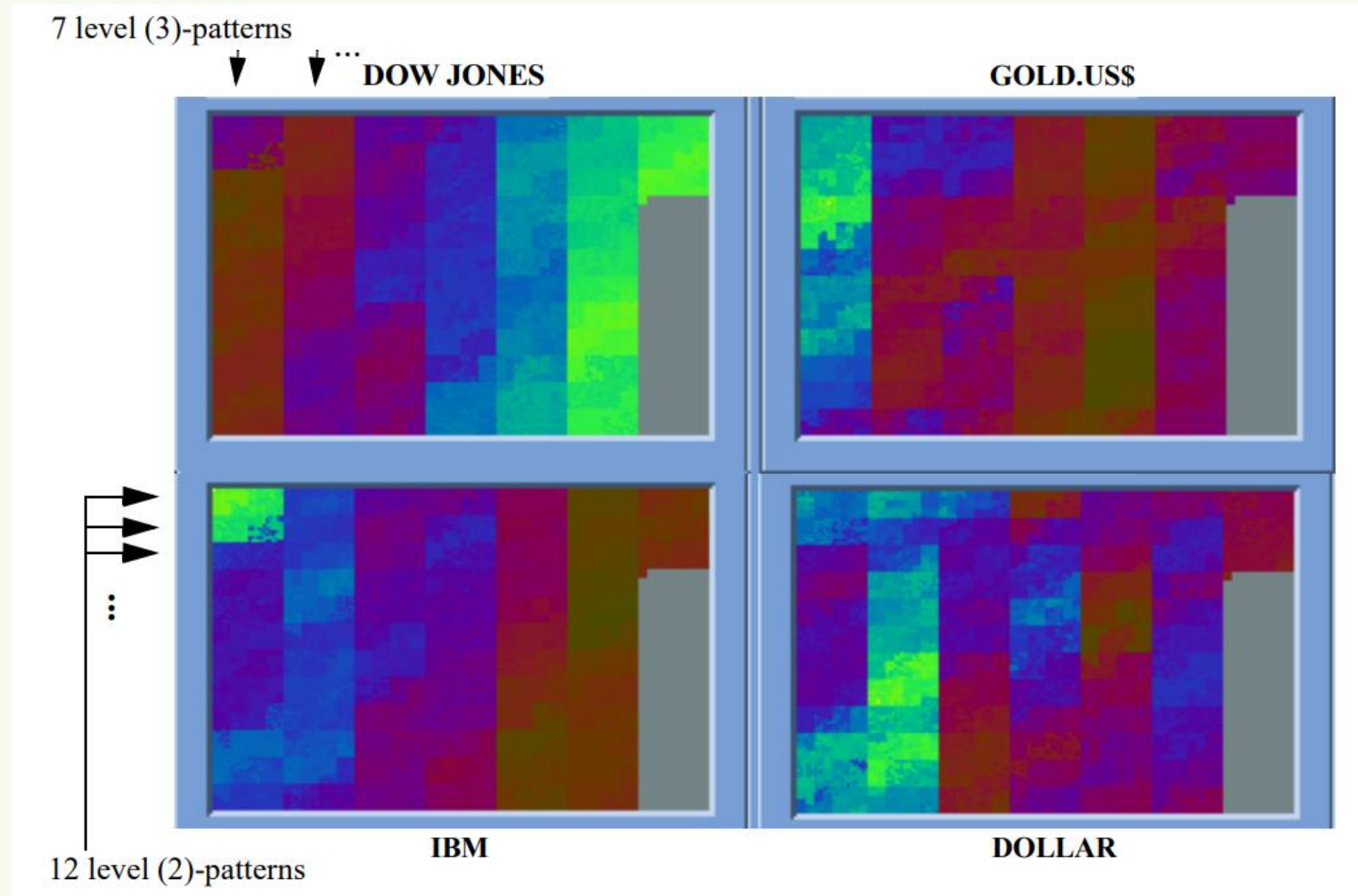
(c) transaction volume



(d) age



# Pixel-oriented visualization techniques (Recursive pattern technique)

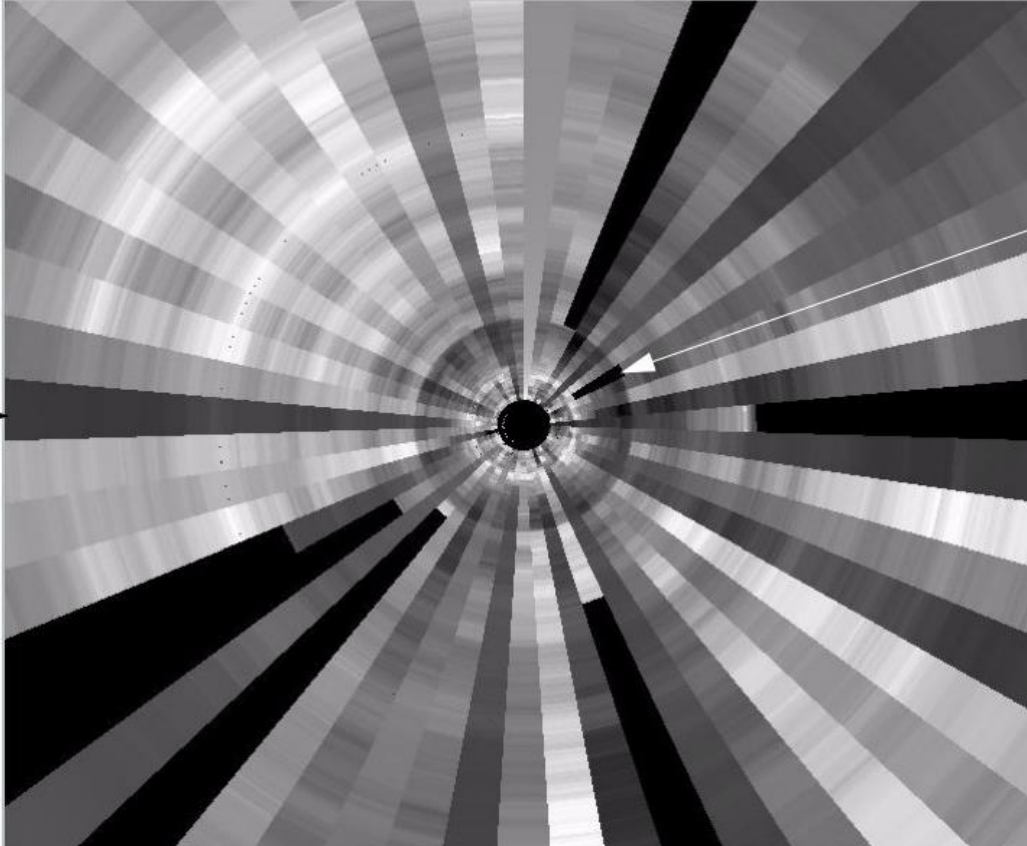


From: Visual Data Mining with Pixel-oriented Visualization Techniques, Mihael Ankerst

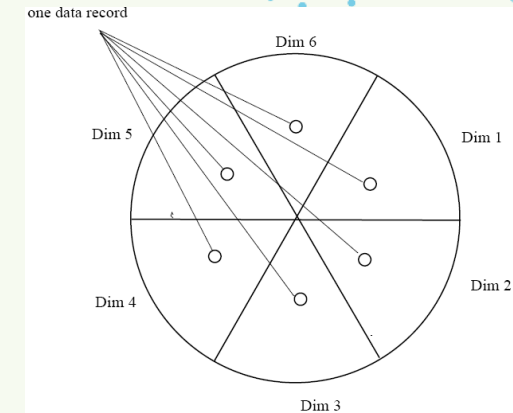


# Pixel-oriented visualization techniques (Laying Out Pixels in Circle Segments)

Exceptional  
stock



Crash



Representing a data  
record in circle segment

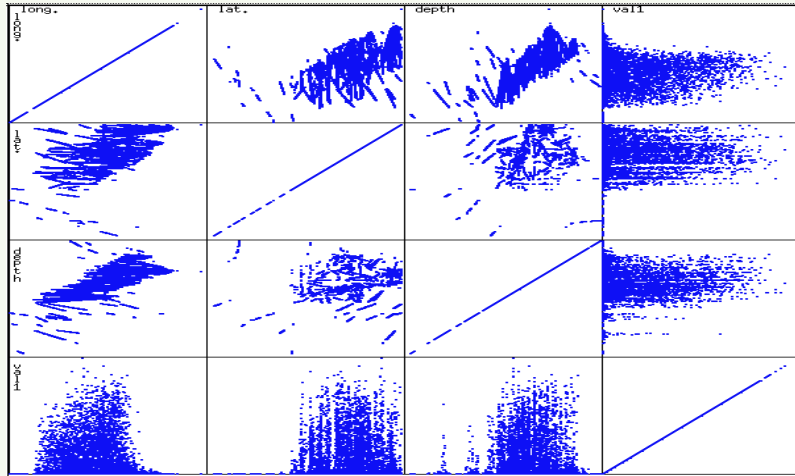




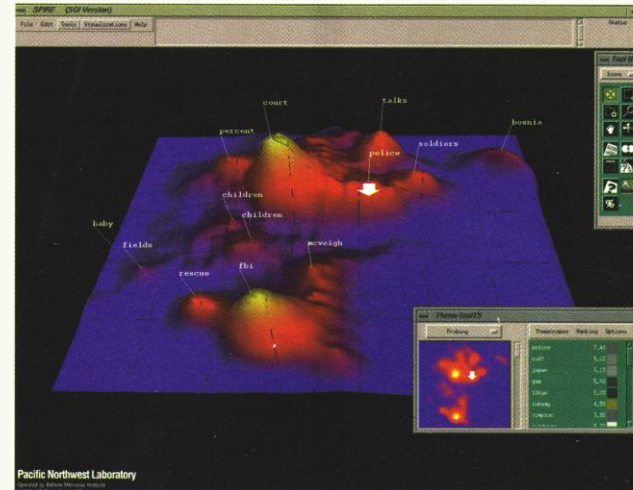
# Geometric projection visualization techniques

Matrix of scatterplots (x-y-diagrams) of the k-dim.

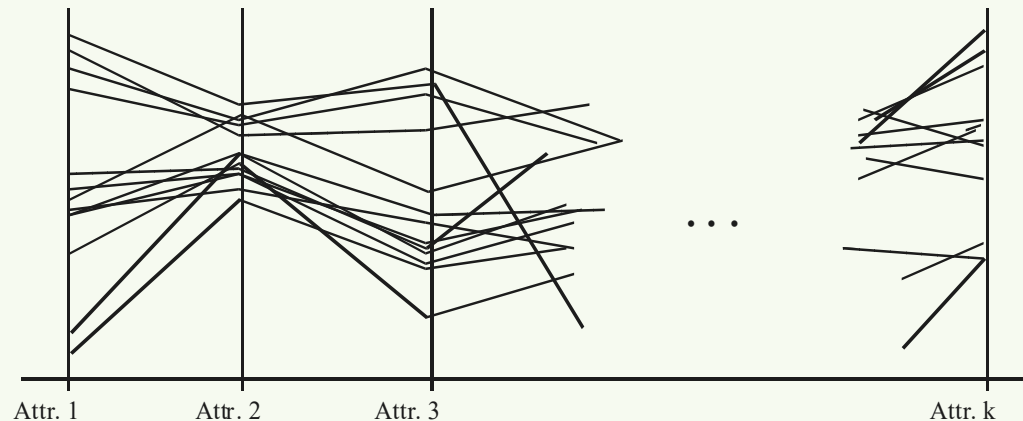
data [total of  $\frac{k^2 - k}{2}$  scatterplots]



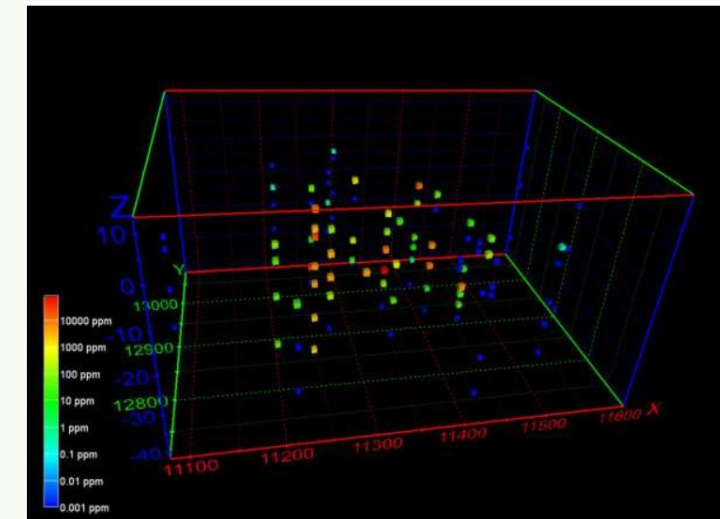
Landscape



Parallel Coordinates

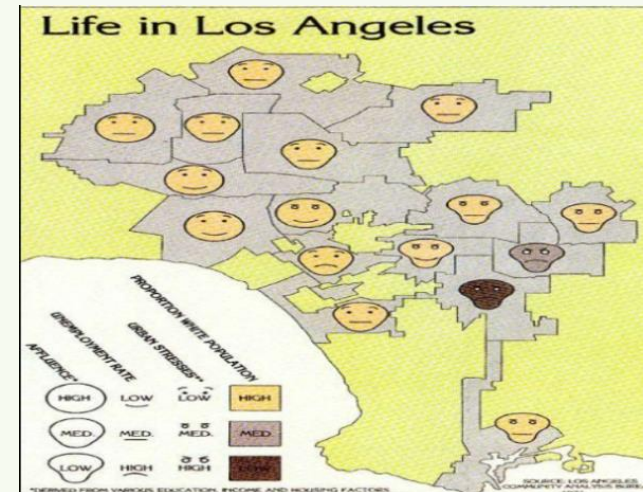
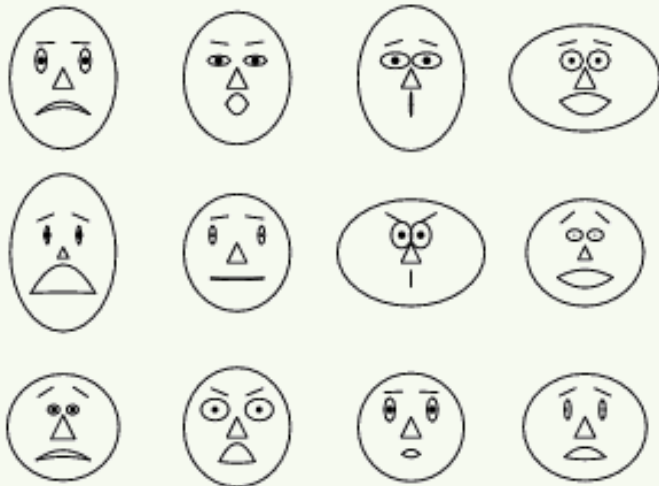


Direct visualization



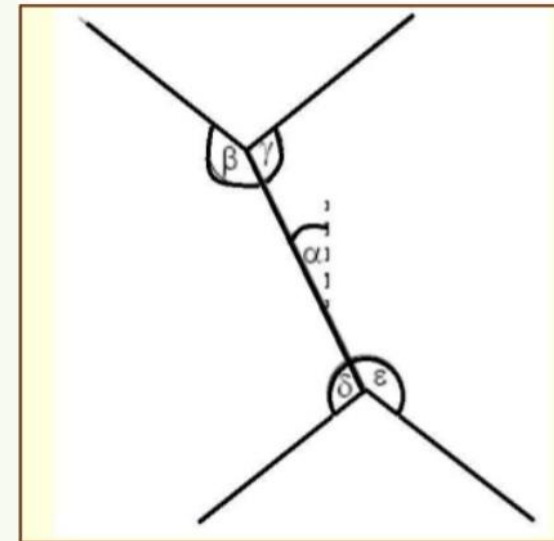
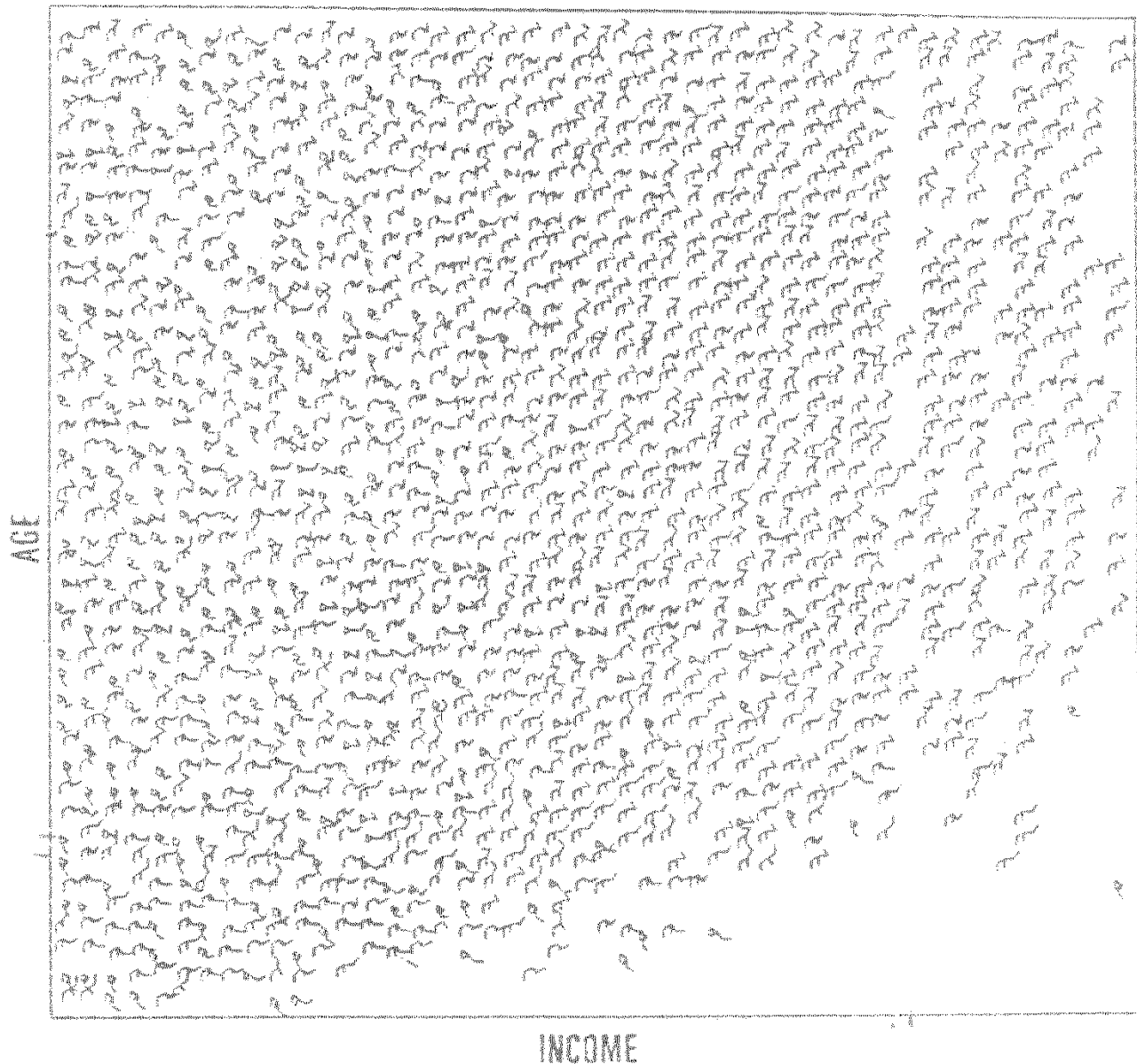
# Icon-Based Visualization Techniques (chernoff faces)

- 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. (10 characteristics)
- assigned one of 10 possible values, generated using
- REFERENCE: Gonick, L. and Smith, W. *The Cartoon Guide to Statistics*. New York: Harper Perennial, p. 212, 1993
- Weisstein, Eric W. "Chernoff Face." From *MathWorld*--A Wolfram Web Resource. [mathworld.wolfram.com/ChernoffFace.html](http://mathworld.wolfram.com/ChernoffFace.html)

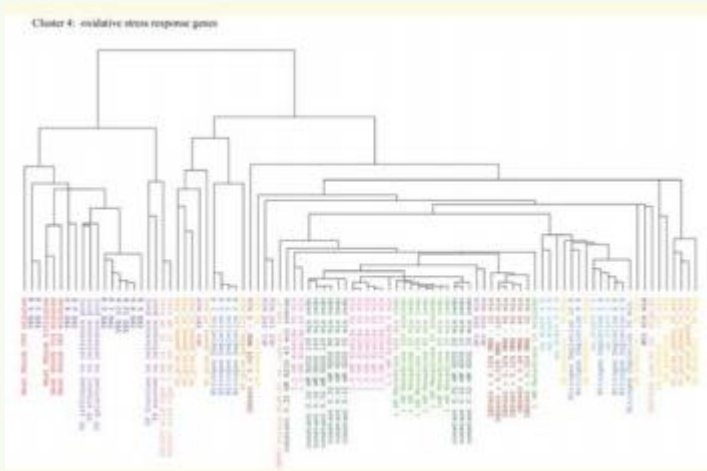




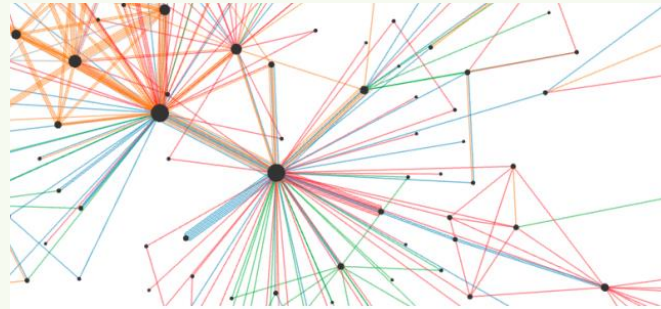
# Icon-Based Visualization Techniques (Stick figure)



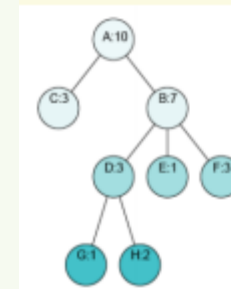
# Hierarchical visualization techniques



Dendrogram



Node-link diagram



Treemaps, voronoi treemaps



# Measuring Data Similarity and Dissimilarity



The Alexander Kofkin  
**Faculty of Engineering**  
Bar-Ilan University

# Similarity and dissimilarity

- **Similarity**
  - Numerical measure of how alike two data objects are
  - Value is higher when objects are more alike
  - Often falls in the range  $[0,1]$
- **Dissimilarity** (e.g., distance)
  - Numerical measure of how different two data objects are
  - Lower when objects are more alike
  - Minimum dissimilarity is often 0
  - Upper limit varies
- **Proximity** refers to a similarity or dissimilarity



# Data matrix and dissimilarity matrix

- **Data matrix**

- n data points with p dimensions

$$\begin{bmatrix} x_{11} & \dots & x_{1f} & \dots & x_{1p} \\ \dots & \dots & \dots & \dots & \dots \\ x_{i1} & \dots & x_{if} & \dots & x_{ip} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & \dots & x_{nf} & \dots & x_{np} \end{bmatrix}$$

- **Dissimilarity matrix**

- n data points, but registers only the distance
- A triangular matrix

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



# Proximity measures for binary attributes

- A contingency table for binary data
- Distance measure for symmetric binary variables:

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

- Distance measure for asymmetric binary variables

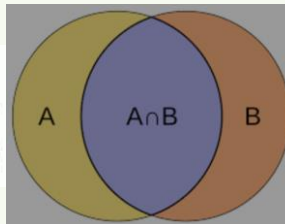
$$d(i, j) = \frac{r + s}{q + r + s + t}$$

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

- Jaccard coefficient (*similarity* measure for *asymmetric* binary variables):

$$\text{sim}_{\text{Jaccard}}(i, j) = \frac{q}{q + r + s}$$

- Note: Jaccard coefficient is the same as “coherence”



$$\text{coherence}(i, j) = \frac{\text{sup}(i, j)}{\text{sup}(i) + \text{sup}(j) - \text{sup}(i, j)} = \frac{q}{(q + r) + (q + s) - q}$$





# Dissimilarity between binary variables

- Example

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

- Gender is a symmetric attribute
- The remaining attributes are asymmetric binary
- Let the values Y and P be 1, and the value N 0

$$d(jack, mary) = \frac{0 + 1}{2 + 0 + 1} = 0.33$$

$$d(jack, jim) = \frac{1 + 1}{1 + 1 + 1} = 0.67$$

$$d(jim, mary) = \frac{1 + 2}{1 + 1 + 2} = 0.75$$



# Proximity measures for nominal attributes

- Can take 2 or more states, e.g., red, yellow, blue, green (generalization of a binary attribute)

- Method 1: Simple matching

$m$ : # of matches,  $p$ : total # of variables

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

- Method 2: Use a large number of binary attributes

creating a new binary attribute for each of the  $M$  nominal states



# Ordinal Variables

- An ordinal variable can be discrete or continuous
- Order is important, e.g., rank
- Can be treated like interval-scaled
  - replace  $x_{if}$  by their rank  $r_{if} \in \{1, \dots, M_f\}$
  - map the range of each variable onto  $[0, 1]$  by replacing  $i$ -th object in the  $f$ -th variable by

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

- compute the dissimilarity using methods for interval-scaled variables



# Distance on Numeric Data: Minkowski Distance

- *Minkowski distance*: A popular distance measure

$$d(i, j) = \sqrt[h]{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \dots + |x_{ip} - x_{jp}|^h}$$

- where  $i = (x_{i1}, x_{i2}, \dots, x_{ip})$  and  $j = (x_{j1}, x_{j2}, \dots, x_{jp})$  are two  $p$ -dimensional data objects, and  $h$  is the order (the distance so defined is also called  $L$ - $h$  norm)
- Properties
  - $d(i, j) > 0$  if  $i \neq j$ , and  $d(i, i) = 0$  (Positive definiteness)
  - $d(i, j) = d(j, i)$  (Symmetry)
  - $d(i, j) \leq d(i, k) + d(k, j)$  (Triangle Inequality)
- A distance that satisfies these properties is a **metric**



# Special Cases of Minkowski Distance

- $h = 1$ : Manhattan (city block,  $L_1$  norm) distance
  - E.g., the Hamming distance: the number of bits that are different between two binary vectors

$$d(i, j) = |x_{i_1} - x_{j_1}| + |x_{i_2} - x_{j_2}| + \dots + |x_{i_p} - x_{j_p}|$$

- $h = 2$ : ( $L_2$  norm) Euclidean distance

$$d(i, j) = \sqrt{(|x_{i_1} - x_{j_1}|^2 + |x_{i_2} - x_{j_2}|^2 + \dots + |x_{i_p} - x_{j_p}|^2)}$$

- $h \rightarrow \infty$ . “supremum” ( $L_{\max}$  norm,  $L_{\infty}$  norm) distance. Chebyshev distance
  - This is the maximum difference between any component (attribute) of the vectors

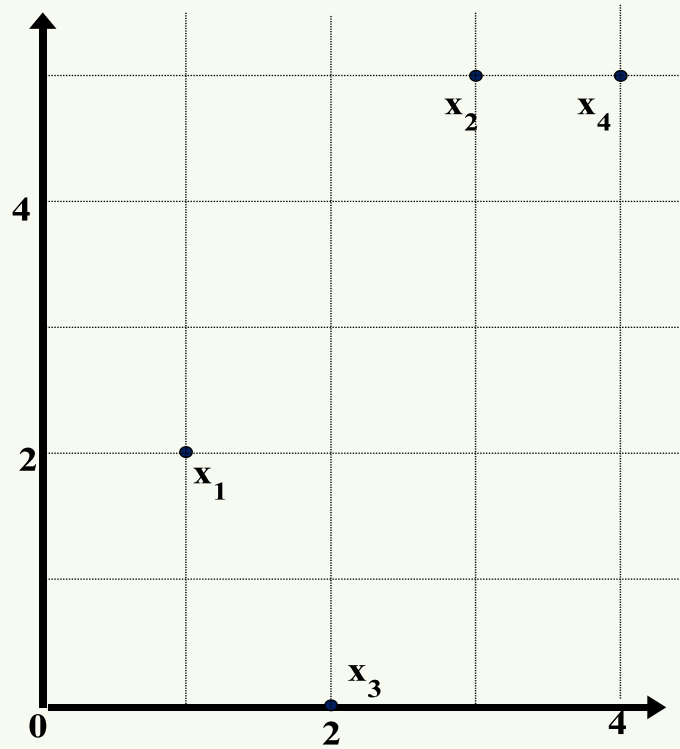
$$d(i, j) = \lim_{h \rightarrow \infty} \left( \sum_{f=1}^p |x_{if} - x_{jf}|^h \right)^{\frac{1}{h}} = \max_f |x_{if} - x_{jf}|$$



# Example: Data Matrix and Dissimilarity Matrix

Data Matrix

point	attribute1	attribute2
$x_1$	1	2
$x_2$	3	5
$x_3$	2	0
$x_4$	4	5



Manhattan ( $L_1$ )

L	$x_1$	$x_2$	$x_3$	$x_4$
$x_1$	0			
$x_2$	5	0		
$x_3$	3	6	0	
$x_4$	6	1	7	0

Euclidean ( $L_2$ )

$L_2$	$x_1$	$x_2$	$x_3$	$x_4$
$x_1$	0			
$x_2$	3.61	0		
$x_3$	2.24	5.1	0	
$x_4$	4.24	1	5.39	0

Supremum

$L_\infty$	$x_1$	$x_2$	$x_3$	$x_4$
$x_1$	0			
$x_2$	3	0		
$x_3$	2	5	0	
$x_4$	3	1	5	0





# Attributes of Mixed Type

- A database may contain all attribute types
  - Nominal, symmetric binary, asymmetric binary, numeric, ordinal
- One may use a weighted formula to combine their effects

$$d(i, j) = \frac{\sum_{f=1}^p \delta_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^p \delta_{ij}^{(f)}}$$

- $f$  is binary or nominal:
  - $d_{ij}^{(f)} = 0$  if  $x_{if} = x_{jf}$ , or  $d_{ij}^{(f)} = 1$  otherwise
- $f$  is numeric: use the normalized distance
- $f$  is ordinal
  - Compute ranks  $r_{if}$  and
  - Treat  $z_{if}$  as interval-scaled  $z_{if} = \frac{r_{if} - 1}{M_f - 1}$



# Cosine Similarity

- 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.

Document	teamcoach	hockey	baseball	soccer	penalty	score	win	loss	season
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

- Other vector objects: gene features in micro-arrays, ...
- Applications: information retrieval, biologic taxonomy, gene feature mapping, ...
- Cosine measure: If  $d_1$  and  $d_2$  are two vectors (e.g., term-frequency vectors), then
- $$\cos(d_1, d_2) = (d_1 \bullet d_2) / \|d_1\| \|d_2\| ,$$
  - where  $\bullet$  indicates vector dot product,  $\|d\|$ : the length of vector  $d$



# Example: Cosine Similarity

- $\cos(d_1, d_2) = (d_1 \bullet d_2) / \|d_1\| \|d_2\|$  ,
  - where  $\bullet$  indicates vector dot product,  $\|d\|$  : the length of vector  $d$
- Ex: Find the **similarity** between documents 1 and 2.
  - $d_1 = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)$
  - $d_2 = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1)$
  - $d_1 \bullet d_2 = 5*3+0*0+3*2+0*0+2*1+0*1+0*1+2*1+0*0+0*1 = 25$
  - $\|d_1\| = (5*5+0*0+3*3+0*0+2*2+0*0+0*0+2*2+0*0+0*0)^{0.5} = (42)^{0.5} = 6.481$
  - $\|d_2\| = (3*3+0*0+2*2+0*0+1*1+1*1+0*0+1*1+0*0+1*1)^{0.5} = (17)^{0.5} = 4.12$
  - $\cos(d_1, d_2) = 0.94$



# Exploratory Data Analysis (EDA) & Data Visualization

**Knowledge about your data is useful for data preprocessing**

- What are the types of attributes?
- What kind of values does each attribute have?
- What do the data look like?
  - Exploring measures of central tendency (symmetric or skewed)
  - Mean, Median, Mode, Variance, Weighted Mean, Trimmed Mean etc.
- How are the values distributed?
- Are there ways we can visualize the data to get a better sense of it all?
  - Identify relations, trends, biases etc.
  - Simple techniques such as scatter-plot matrices, histograms, Q-Q plots etc.
- Measure data similarity



# Thank you!



The Alexander Kofkin  
**Faculty of Engineering**  
Bar-Ilan University

