

CSE5243 INTRO. TO DATA MINING

Data & Data Preprocessing

Huan Sun, CSE@The Ohio State University

Data & Data Preprocessing

- What is Data: Data Objects and Attribute Types 
- Basic Statistical Descriptions of Data
- Measuring Data Similarity and Dissimilarity
- Data Preprocessing: An Overview
- Summary

What is Data?

- Collection of data objects and their attributes

What is Data?

- Collection of data objects and their attributes

Attributes

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Objects

Data Objects

- Data sets are made up of data objects
- A **data object** represents an entity
- Examples:
 - sales database: customers, store items, sales
 - medical database: patients, treatments
 - university database: students, professors, courses
- Also called *samples , examples, instances, data points, objects, tuples*

Data Objects

- Data sets are made up of data objects
- A **data object** represents an entity
- Examples:
 - sales database: customers, store items, sales
 - medical database: patients, treatments
 - university database: students, professors, courses
- Also called *samples , examples, instances, data points, objects, tuples*
- **Data objects are described by attributes**
- **Database rows → data objects; columns → attributes**

Attributes

- **Attribute (or dimensions, features, variables)**
 - A data field, representing a characteristic or feature of a data object.
 - E.g., *customer_ID*, *name*, *address*
- Types:
 - **Nominal** (e.g., red, blue)
 - **Binary** (e.g., {true, false})
 - **Ordinal** (e.g., {freshman, sophomore, junior, senior})
 - **Numeric**: quantitative

Attribute Types

- **Nominal:** categories, states, or “names of things”
 - *Hair_color* = {*auburn*, *black*, *blond*, *brown*, *grey*, *red*, *white*}
 - marital status, occupation, zip codes

Attribute Types

- **Nominal:** categories, states, or “names of things”
 - *Hair_color* = {*auburn*, *black*, *blond*, *brown*, *grey*, *red*, *white*}
 - marital status, occupation, zip codes
- **Binary**
 - **Nominal attribute with only 2 states (0 and 1)**
 - **Symmetric binary:** both outcomes equally important
 - e.g., gender
 - **Asymmetric binary:** outcomes not equally important.
 - e.g., medical test (positive vs. negative)
 - Convention: assign 1 to most important outcome (e.g., HIV positive)

Attribute Types

- **Nominal:** categories, states, or “names of things”
 - *Hair_color* = {*auburn, black, blond, brown, grey, red, white*}
 - marital status, occupation, ID numbers, zip codes
- **Binary**
 - Nominal attribute with only 2 states (0 and 1)
 - Symmetric binary: both outcomes equally important, e.g., gender
 - Asymmetric binary: outcomes not equally important.
 - e.g., medical test (positive vs. negative)
 - Convention: assign 1 to most important outcome (e.g., HIV positive)
- **Ordinal**
 - Values have a meaningful order (ranking) but magnitude between successive values is not known
 - Size = {*small, medium, large*}, grades, army rankings

Numeric Attribute Types

- Quantity (integer or real-valued)
- **Interval-scaled**
 - Measured on a scale of **equal-sized units**
 - Values have order
 - E.g., *temperature in C° or F°, calendar dates*
 - **No true zero-point**

Numeric Attribute Types

- Quantity (integer or real-valued)
- **Interval-scaled**
 - Measured on a scale of **equal-sized units**
 - Values have order
 - E.g., *temperature in C° or F°, calendar dates*
 - No true zero-point
- **Ratio-scaled**
 - Inherent zero-point
 - 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°).
 - e.g., *temperature in Kelvin, length, counts, monetary quantities*

- Q1: Is student ID a nominal, ordinal, or numerical attribute?
- Q2: What about eye color? Or color in the color spectrum of physics?



Discrete vs. Continuous Attributes

□ Discrete Attribute

- Has only a finite or countably infinite set of values
 - E.g., zip codes, profession, or the set of words in a collection of documents
- Sometimes, represented as integer variables
- **Note: Binary attributes are a special case of discrete attributes**



Discrete vs. Continuous Attributes

□ Discrete Attribute

- Has only a finite or countably infinite set of values
- Sometimes, represented as integer variables
- Note: Binary attributes are a special case of discrete attributes

□ Continuous Attribute

- Has real numbers as attribute values
 - E.g., temperature, height, or weight
- Practically, real values can only be measured and represented using a finite number of digits
- Continuous attributes are typically represented as floating-point variables

Types of Data Sets: (1) Record Data

- Relational records
 - Relational tables, highly structured

Person:

Pers_ID	Surname	First_Name	City
0	Miller	Paul	London
1	Ortega	Alvaro	Valencia
2	Huber	Urs	Zurich
3	Blanc	Gaston	Paris
4	Bertolini	Fabrizio	Rom

Car:

Car_ID	Model	Year	Value	Pers_ID
101	Bentley	1973	100000	0
102	Rolls Royce	1965	330000	0
103	Peugeot	1993	500	3
104	Ferrari	2005	150000	4
105	Renault	1998	2000	3
106	Renault	2001	7000	3
107	Smart	1999	2000	2

— no relation

Types of Data Sets: (1) Record Data

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

	China	England	France	Japan	USA	Total
Active Outdoors Crochet Glove		12.00	4.00	1.00	240.00	257.00
Active Outdoors Lycra Glove		10.00	6.00		323.00	339.00
InFlux Crochet Glove	3.00	6.00	8.00		132.00	149.00
InFlux Lycra Glove		2.00			143.00	145.00
Triumph Pro Helmet	3.00	1.00	7.00		333.00	344.00
Triumph Vertigo Helmet		3.00	22.00		474.00	499.00
Xtreme Adult Helmet	8.00	8.00	7.00	2.00	251.00	276.00
Xtreme Youth Helmet		1.00			76.00	77.00
Total	14.00	43.00	54.00	3.00	1,972.00	2,086.00

Types of Data Sets: (1) Record Data

- Transaction data

<i>TID</i>	<i>Items</i>
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

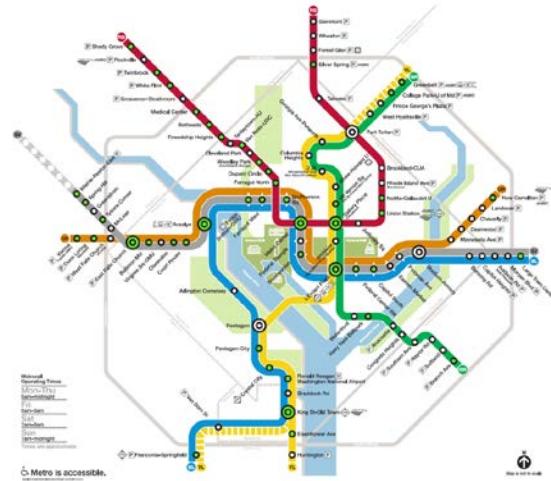
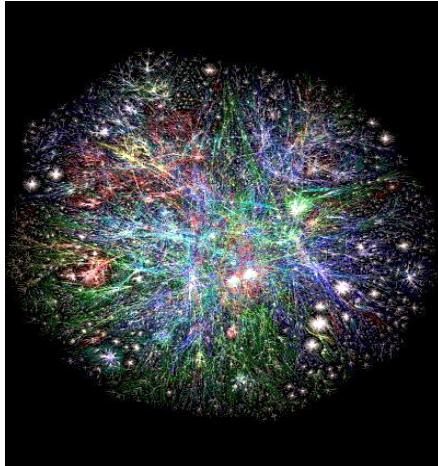
Types of Data Sets: (1) Record Data

- Document data: Term-frequency vector (matrix) of text documents

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

Types of Data Sets: (2) Graphs and Networks

- Transportation network
- World Wide Web

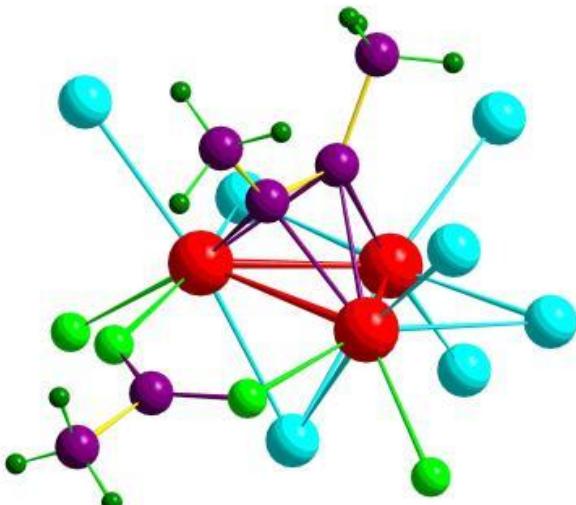
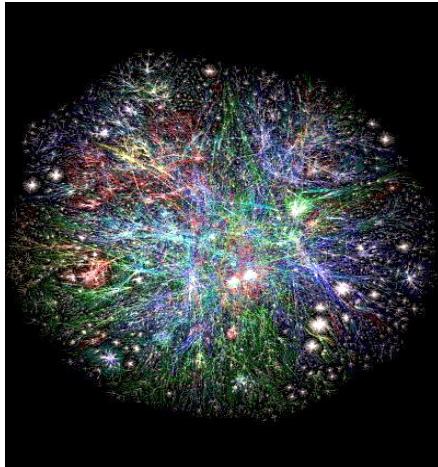


Types of Data Sets: (2) Graphs and Networks

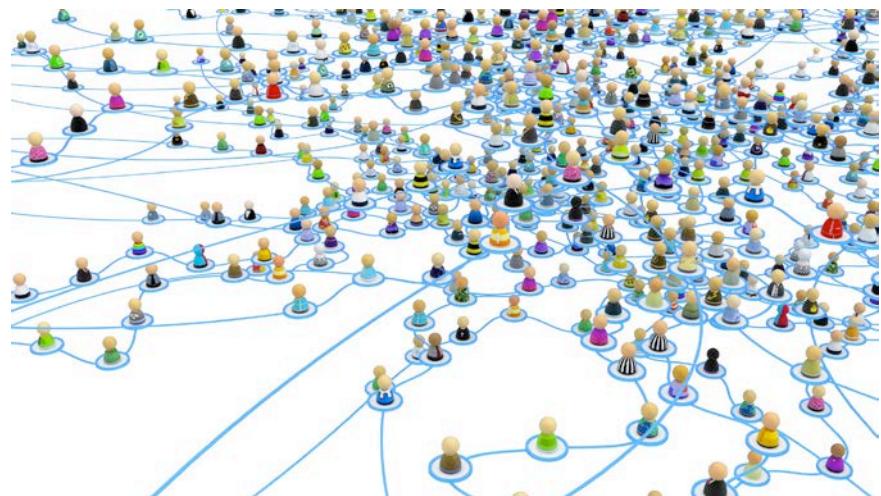
- Transportation network



- World Wide Web



- Molecular Structures



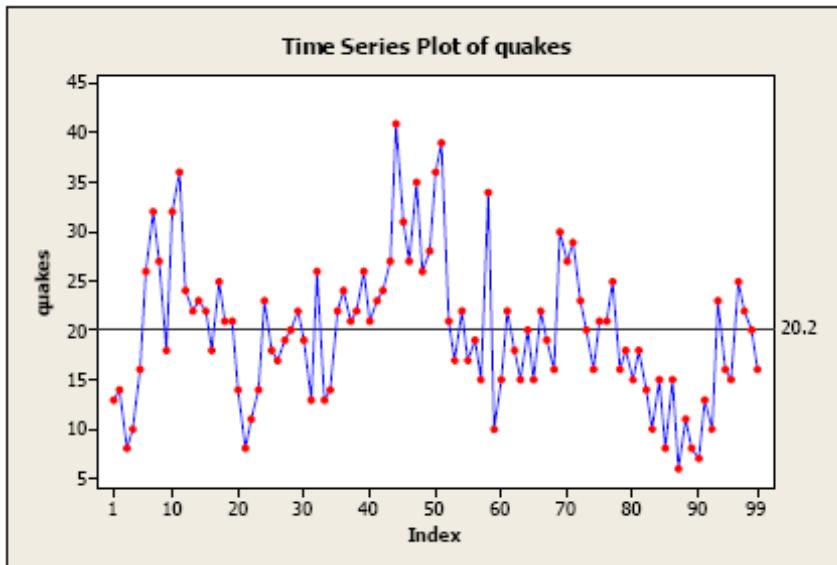
- Social or information networks

Types of Data Sets: (3) Ordered Data

- Video data: sequence of images



- Temporal data: time-series



- Sequential Data: transaction sequences

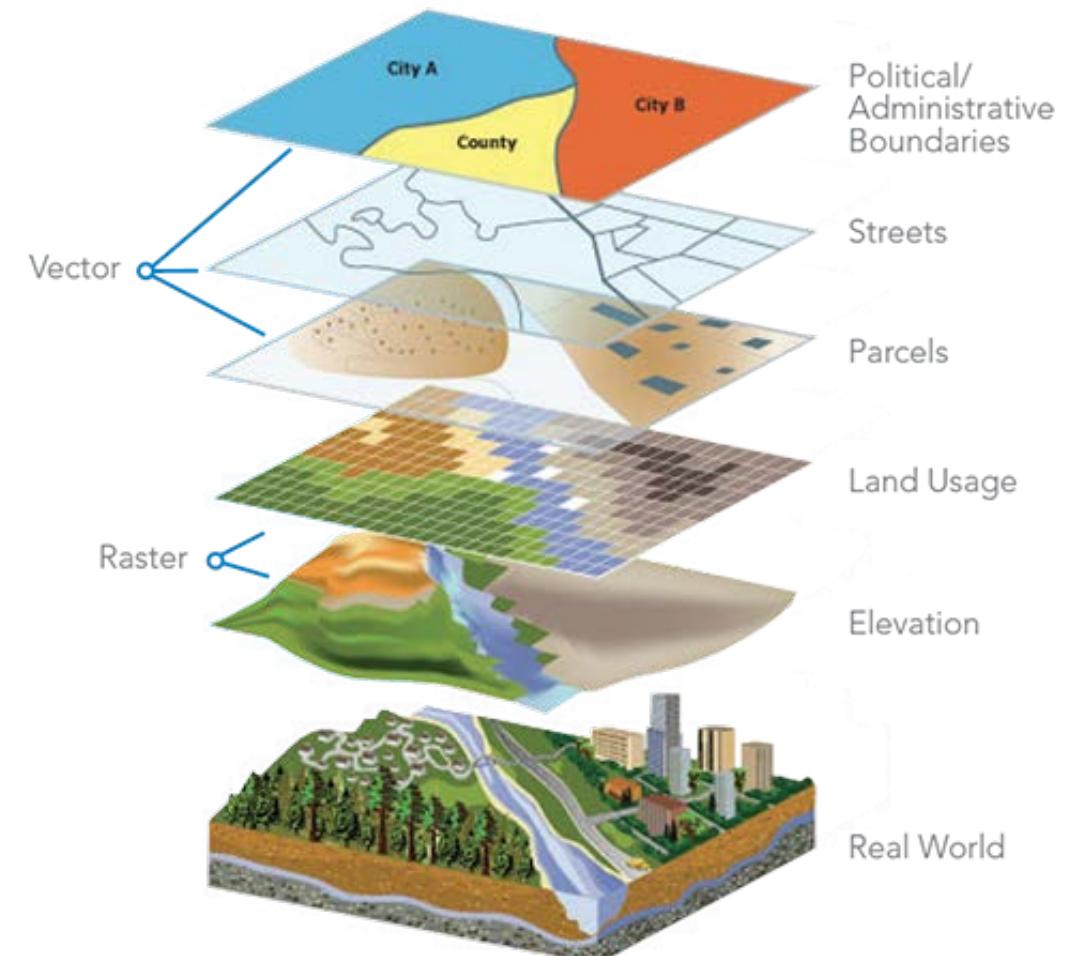
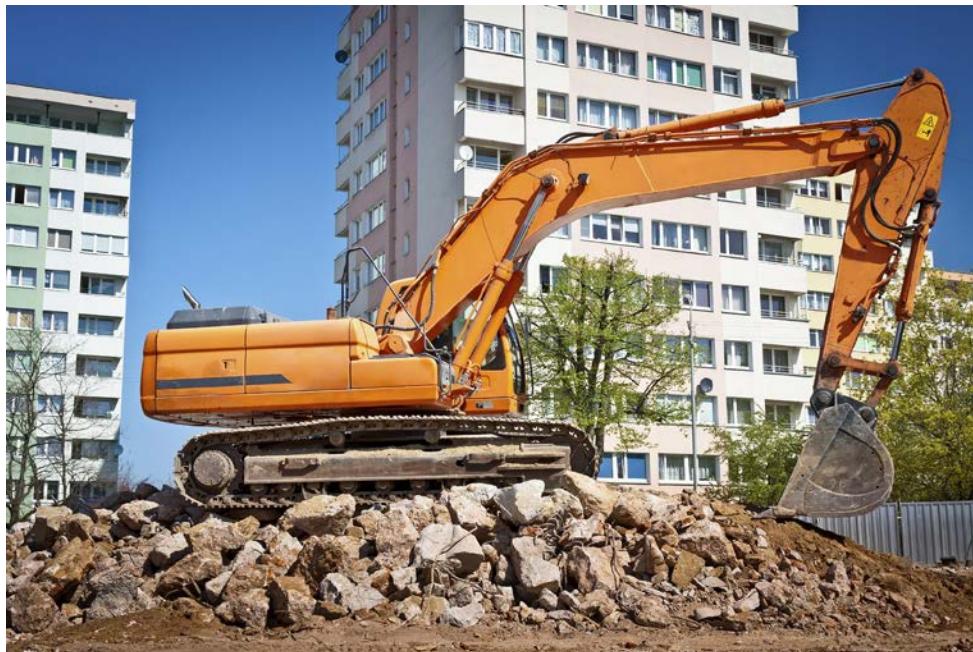
- Genetic sequence data

Start

Human	Chimpanzee	Macaque
GTTTGAGG - - -	ATGTTCAACAAATGCTCCTTCATTCTCTATTTACAGACCTGCCGCA	
GTTTGAGG - - -	ATGTTCAATAAAATGCTGCTTCACTCCCTATTTACAGACCTGCCGCA	
GTTTGAGG - - -	ATGCTCAATAAAATGCTCCTTCATTCTCCATTTACAAACTGCCGCA	
Human	Chimpanzee	Macaque
GACAATTCTGCTAGCAGCCTTGTGCTATTATCTGTTTCTAAACTTAGTAATTGAGTGT		
GACAATTCTGCTAGCAGCCTTGTGCTATTATCTGTTTCTAAACTTAGTAATTGAGTGT		
GACAATTCTGCTAGCAGCCTTGTGCTATTATCTGTTTCTAAACTTAGTAATTGAGTGT		
Human	Chimpanzee	Macaque
GATCTGGAGACTAA - CTCTGAAATAAAAGCTGATTATTATTTATTTCTCAAAACAA		
GATCTGGAGACTAAACTCTGAAATAAAAGCTGATTATTATTTATTTCTCAAAACAA		
TATCTGGAGACTAAACTCTGAAATAAAAGCTGATTATTATTTATTTCTCAAAACAA		
Human	Chimpanzee	Macaque
CAGAATACGATTAGCAAATTACTCTTAAGATATTATTTACATTCTATATTCTCCTA		
CAGAATACGATTAGCAAATTACTCTTAAGATACTATTTCACATTCTATATTCTCCTA		
CAGAATATGATTAGCAAATTACCTCTTAAGATATTATTTGCACCTCTATATTCTCCTA		
Human	Chimpanzee	Macaque
CCCTGAGTTGATGTGAGCAATATGTCACCTTCATAAAGGCCAGGTATACA - - - TTATG		
CCCTGAGTTGATGTGAGCGTATGTCACCTTCATAAAGGCCAGGTATACA - - - TTATG		
CCCTGAGTTGATGTGAGCAATATGTCACCTTCACACAAAGGCCAGGTATATAACATTACG		
Human	Chimpanzee	Macaque
GACAGGTAAGTAAAAAACATATTATTTATTCTACGTTTGTCCAAGAATTAAATTTC		
GACAGGTAAGTAAAAAACATATTATTTATTCTACGTTTGTCCAAGAATTAAATTTC		
GACAGGTAAGTAAAAA-CATATTATTTATTCTAGTTTGTCCAAGAGATTAAATTTC		
Human	Chimpanzee	Macaque
AAC TGTGCGCGTGTGGTAA - - - TGTAACAAACTCAAGTACA		
AAC TGTGCGCGTGTGGTAA - - - TGTAACAAACTCAAGTACA		
AAC TGTGCGCGTGTGGTAA - - - CGTAAACAAATTCAAGTACA		

Types of Data Sets: (4) Spatial, image and multimedia data

- Spatial data: maps



- Image data:
- Video data:

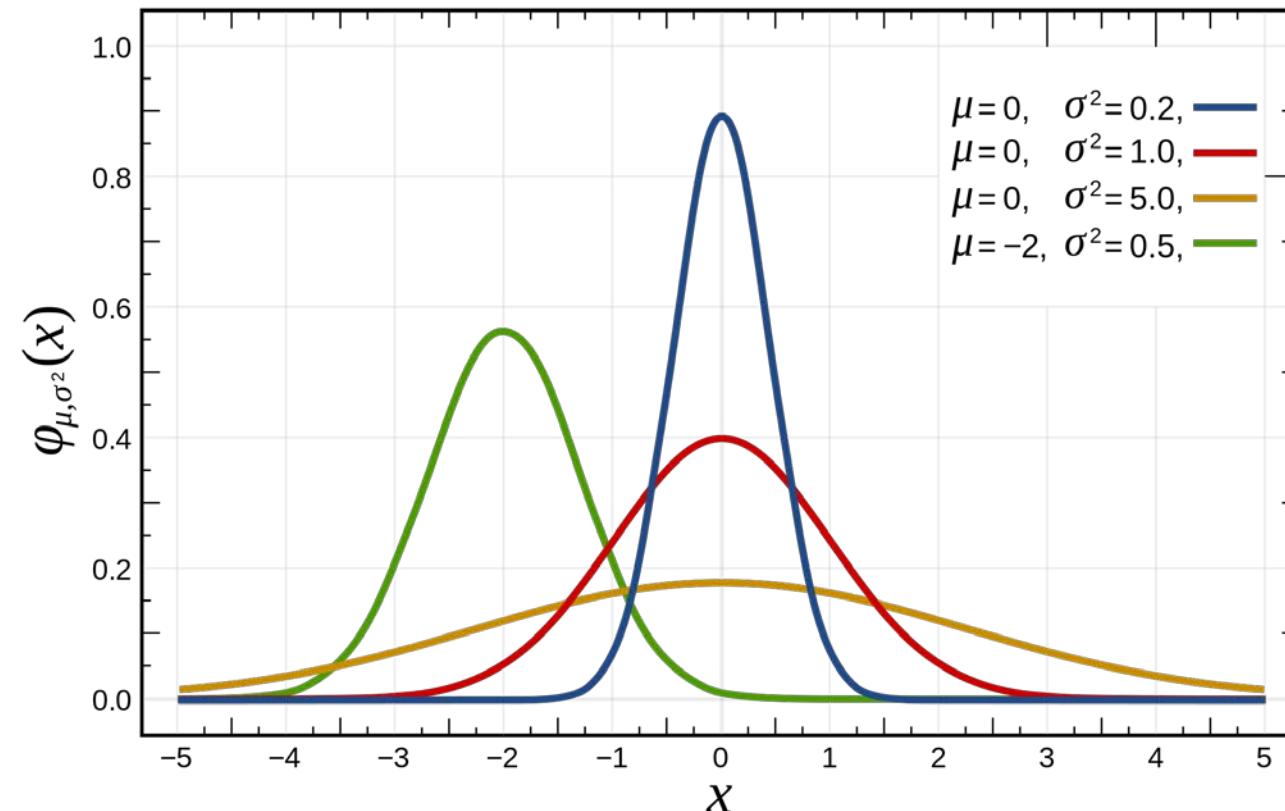
Chapter 2. Getting to Know Your Data

- Data Objects and Attribute Types
- Basic Statistical Descriptions of Data 
- Measuring Data Similarity and Dissimilarity
- Data Preprocessing: An Overview
- Summary

Basic Statistical Descriptions of Data

□ Motivation

- To better understand the data: central tendency, variation and spread



Measuring the Central Tendency: (1) Mean

- Mean (algebraic measure) (sample vs. population):

Note: n is sample size and N is population size.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \mu = \frac{\sum x}{N}$$


Measuring the Central Tendency: (1) Mean

□ Mean (algebraic measure) (sample vs. population):

Note: n is sample size and N is population size.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \mu = \frac{\sum x}{N}$$

□ **Weighted** arithmetic mean:

$$\bar{x} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

□ **Trimmed** mean:

- Chopping extreme values (e.g., Olympics gymnastics score computation)



Measuring the Central Tendency: (2) Median

Median:

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

Measuring the Central Tendency: (2) Median

□ Median:

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

□ Estimated by interpolation (for *grouped data*):

Approximate median

$$\text{median} = L_1 + \left(\frac{n/2 - (\sum \text{freq})_l}{\text{freq}_{\text{median}}} \right) \text{width}$$

Low interval limit

Sum before the median interval

Interval width ($L_2 - L_1$)

age	frequency
1–5	200
6–15	450
16–20	300
21–50	1500
51–80	700
81–110	44



Measuring the Central Tendency: (3) Mode

- Mode: Value that occurs most frequently in the data

- Unimodal

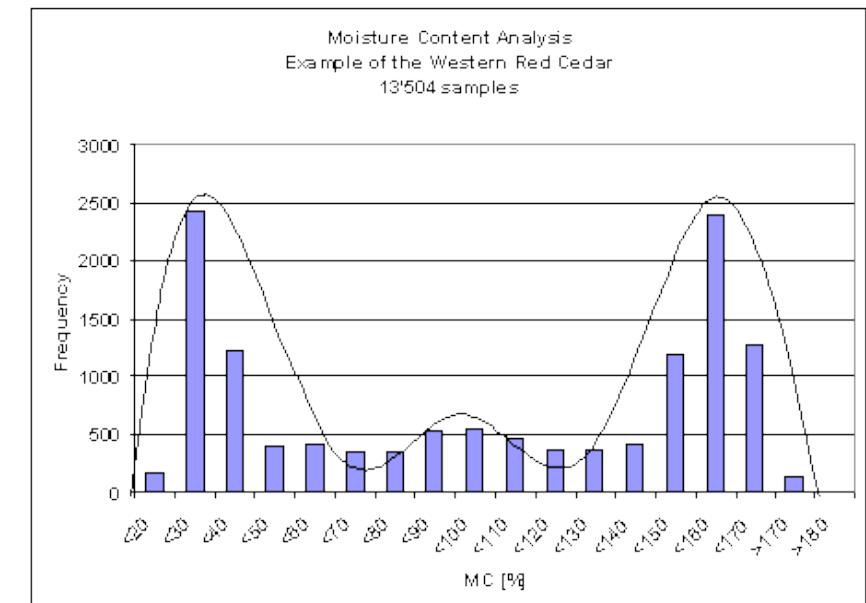
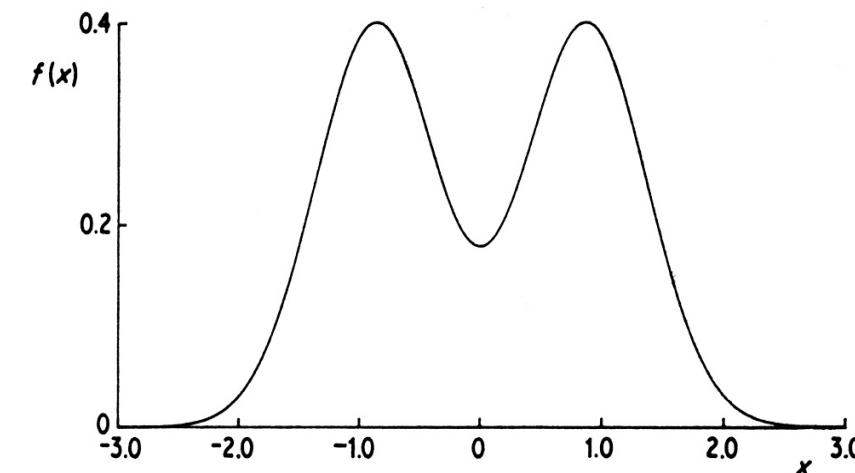
- Empirical formula:

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

- Multi-modal

- Bimodal

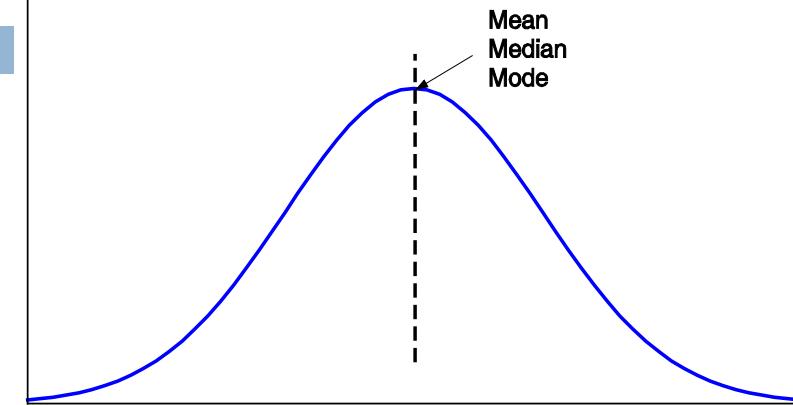
- Trimodal



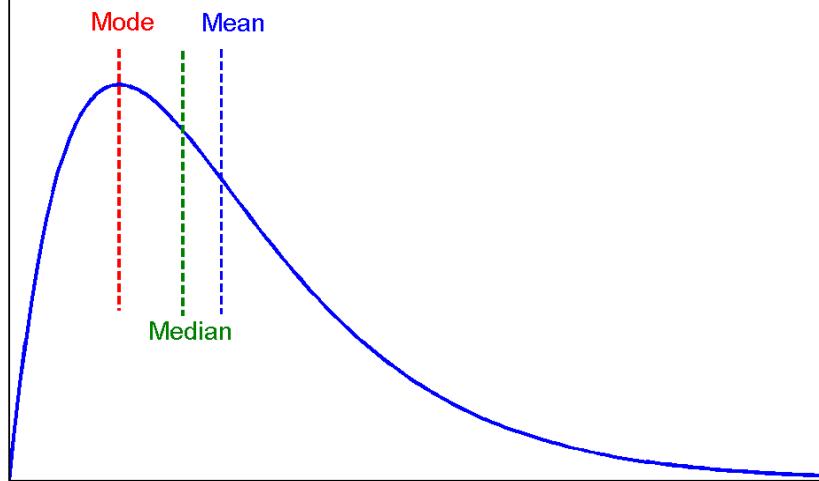
Symmetric vs. Skewed Data

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

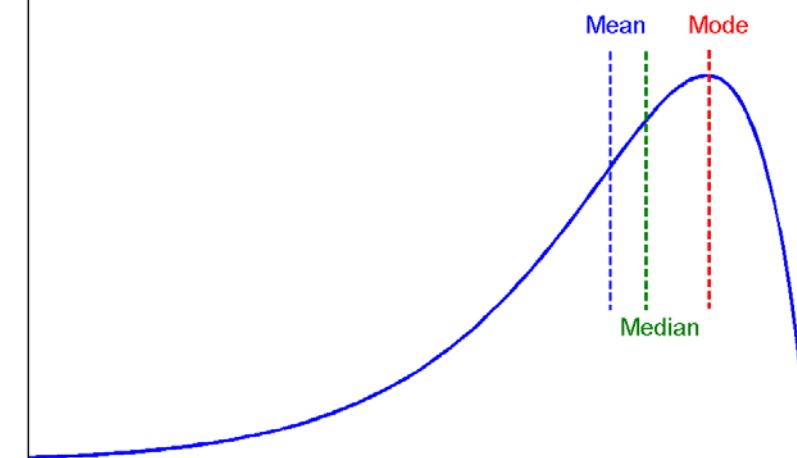
symmetric



positively skewed

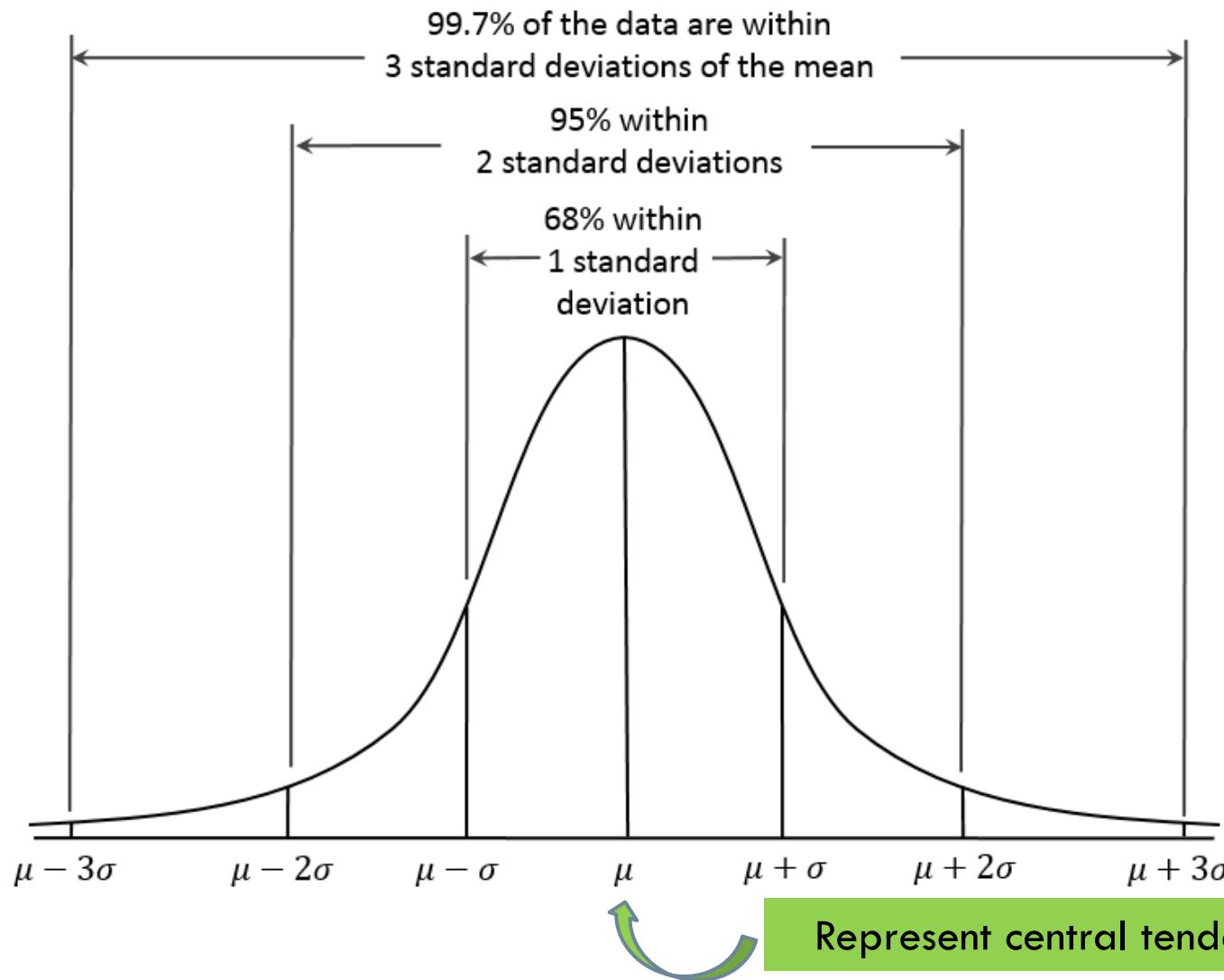


negatively skewed



Properties of Normal Distribution Curve

← ----- Represent data dispersion, spread ----- →



Measures Data Distribution: Variance and Standard Deviation

- Variance and standard deviation (*sample: s, population: σ*)
 - **Variance:** (algebraic, scalable computation)

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right]$$

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

- **Standard deviation s (or σ)** is the square root of variance s^2 (or σ^2)

Measures Data Distribution: Variance and Standard Deviation

- Variance and standard deviation (*sample: s, population: σ*)

- **Variance:** (algebraic, scalable computation)

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2 \right]$$

?

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

- **Standard deviation s (or σ)** is the square root of variance s^2 (or σ^2)

Standardizing Numeric Data

- Z-score:
$$z = \frac{x - \mu}{\sigma}$$
 - X: raw score to be standardized, μ : mean of the population, σ : standard deviation
 - **the distance between the raw score and the population mean in units of the standard deviation**
 - negative when the raw score is below the mean, positive when above

Standardizing Numeric Data

- Z-score:
$$z = \frac{x - \mu}{\sigma}$$
 - X: raw score to be standardized, μ : mean of the population, σ : standard deviation
 - the distance between the raw score and the population mean in units of the standard deviation
 - negative when the raw score is below the mean, positive when above

□ Mean absolute deviation:

$$s_f = \frac{1}{n}(|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|)$$

where

$$m_f = \frac{1}{n}(x_{1f} + x_{2f} + \dots + x_{nf}).$$

□ standardized measure (z-score):

$$z_{if} = \frac{x_{if} - m_f}{s_f}$$

□ Using mean absolute deviation is more robust than using standard deviation

Data & Data Preprocessing

- Data Objects and Attribute Types
- Basic Statistical Descriptions of Data
- Measuring Data Similarity and Dissimilarity
- Data Preprocessing: An Overview
- Summary



Similarity, Dissimilarity, and Proximity

- **Similarity measure or similarity function**
 - A **real-valued function** that quantifies the similarity between two objects
 - **Measure how two data objects are alike:** The higher value, the more alike
 - **Often falls in the range [0,1]:** 0: no similarity; 1: completely similar

Similarity, Dissimilarity, and Proximity

- **Similarity measure or similarity function**
 - A real-valued function that quantifies the similarity between two objects
 - Measure how two data objects are alike: The higher value, the more alike
 - Often falls in the range $[0,1]$: 0: no similarity; 1: completely similar
- **Dissimilarity (or distance) measure**
 - Numerical measure of how different two data objects are
 - In some sense, **the inverse of similarity**: The lower, the more alike
 - Minimum dissimilarity is often 0 (i.e., completely similar)
 - Range $[0, 1]$ or $[0, \infty)$, depending on the definition
- **Proximity** usually refers to either similarity or dissimilarity

Data Matrix and Dissimilarity Matrix

- Data matrix

- A data matrix of n data points with l dimensions

$$D = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1l} \\ x_{21} & x_{22} & \dots & x_{2l} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nl} \end{pmatrix}$$

Data Matrix and Dissimilarity Matrix

- Data matrix

- A data matrix of n data points with l dimensions

- **Dissimilarity (distance) matrix**

- n data points, but registers only the distance $d(i, j)$ (typically metric)

- Usually symmetric

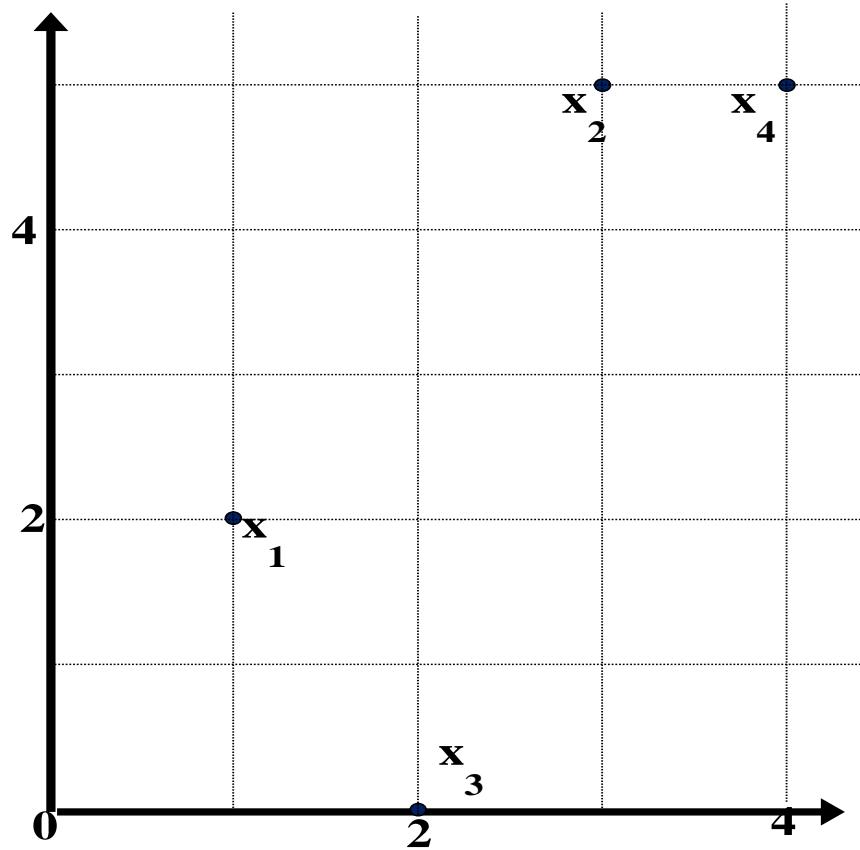
- **Distance functions** are usually different for real, boolean, categorical, ordinal, ratio, and vector variables

- Weights can be associated with different variables based on applications and data semantics

$$D = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1l} \\ x_{21} & x_{22} & \dots & x_{2l} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nl} \end{pmatrix}$$

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

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

Dissimilarity Matrix (by Euclidean Distance)

	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

Distance on Numeric Data: Minkowski Distance

- **Minkowski distance:** A popular distance measure

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

where $i = (x_{i1}, x_{i2}, \dots, x_{il})$ and $j = (x_{j1}, x_{j2}, \dots, x_{jl})$ are two l -dimensional data objects, and p is the order (the distance so defined is also called **L-p norm**)

Distance on Numeric Data: Minkowski Distance

- **Minkowski distance:** A popular distance measure

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

where $i = (x_{i1}, x_{i2}, \dots, x_{il})$ and $j = (x_{j1}, x_{j2}, \dots, x_{jl})$ are two l -dimensional data objects, and p is the order (the distance so defined is also called L- p norm)

- Properties
 - $d(i, j) > 0$ if $i \neq j$, and $d(i, i) = 0$ (Positivity)
 - $d(i, j) = d(j, i)$ (Symmetry)
 - $d(i, j) \leq d(i, k) + d(k, j)$ (Triangle Inequality)
- A distance that satisfies these properties is a **metric**
- Note: There are nonmetric dissimilarities, e.g., set differences

Special Cases of Minkowski Distance

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

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{il} - x_{jl}|$$

Special Cases of Minkowski Distance

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

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{il} - x_{jl}|$$

- $p = 2$: (L_2 norm) **Euclidean distance**

$$d(i, j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{il} - x_{jl}|^2}$$

Special Cases of Minkowski Distance

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

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{il} - x_{jl}|$$

- $p = 2$: (L_2 norm) **Euclidean distance**

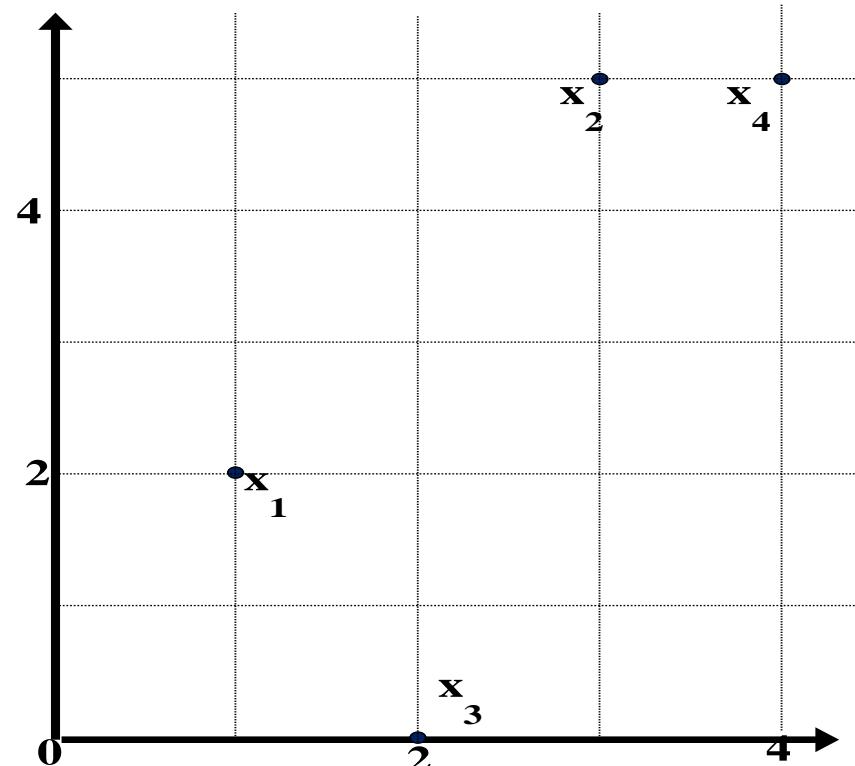
$$d(i, j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{il} - x_{jl}|^2}$$

- $p \rightarrow \infty$: (L_{\max} norm, L_∞ norm) **“supremum” distance**
 - ▣ The maximum difference between any component (attribute) of the vectors

$$d(i, j) = \lim_{p \rightarrow \infty} \sqrt[|p|]{|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{il} - x_{jl}|^p} = \max_{f=1}^l |x_{if} - x_{jf}|$$

Example: Minkowski Distance at Special Cases

point	attribute 1	attribute 2
x1	1	2
x2	3	5
x3	2	0
x4	4	5



L	x1	x2	x3	x4
x1	0			
x2	5	0		
x3	3	6	0	
x4	6	1	7	0

**Manhattan
(L_1)**

L2	x1	x2	x3	x4
x1	0			
x2	3.61	0		
x3	2.24	5.1	0	
x4	4.24	1	5.39	0

**Euclidean
(L_2)**

L_∞	x1	x2	x3	x4
x1	0			
x2	3	0		
x3	2	5	0	
x4	3	1	5	0

**Supremum
(L_∞)**

Proximity Measure for Binary Attributes

- A contingency table for binary data

		Object <i>j</i>		
		1	0	sum
Object <i>i</i>	1	<i>q</i>	<i>r</i>	<i>q+r</i>
	0	<i>s</i>	<i>t</i>	<i>s+t</i>
sum		<i>q+s</i>	<i>r+t</i>	<i>p</i>

- Distance measure for symmetric binary variables:

- Distance measure for **asymmetric** binary variables:

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

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

Proximity Measure for Binary Attributes

- A contingency table for binary data

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

- Distance measure for symmetric binary variables:

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

- Distance measure for asymmetric binary variables:

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

- Jaccard coefficient (*similarity* measure for *asymmetric*

binary variables): $sim_{Jaccard}(i, j) = \frac{q}{q + r + s}$

- Note: Jaccard coefficient is the same as “coherence”: (a concept discussed in Pattern Discovery)

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

Example: Dissimilarity between Asymmetric Binary Variables

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 (not counted in)
- The remaining attributes are asymmetric binary
- Let the values Y and P be 1, and the value N be 0

Example: Dissimilarity between Asymmetric Binary Variables

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4			Mary		
										1	0	Σ_{row}
Jack	M	Y	N	P	N	N	N			1	2	0
Mary	F	Y	N	P	N	P	N			0	1	3
Jim	M	Y	P	N	N	N	N			\sum_{co}	3	3
										1		6

- Gender is a symmetric attribute (not counted in)
- The remaining attributes are asymmetric binary
- Let the values Y and P be 1, and the value N be 0

		Jim		
		1	0	Σ_{row}
Jack		1	1	2
		0	1	3
		\sum_{co}	2	4
		1		6

		Mary		
		1	0	Σ_{row}
Jim		1	1	2
		0	2	4
		\sum_{col}	3	6

Contingency table

Example: Dissimilarity between Asymmetric Binary Variables

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4			Mary				Σ_{row}
Jack	M	Y	N	P	N	N	N			1	2	0	2	
Mary	F	Y	N	P	N	P	N			0	1	3	4	
Jim	M	Y	P	N	N	N	N			\sum_{co}	3	3	6	

- Gender is a symmetric attribute (not counted in)
- The remaining attributes are asymmetric binary
- Let the values Y and P be 1, and the value N be 0
- Distance: $d(i, j) = \frac{r + s}{q + r + s}$

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

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

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

			Jim	Mary										
				1	0	Σ_{row}								
				1	1	2								
				0	1	3								

Contingency table

Proximity Measure for Categorical Attributes

- Categorical data, also called nominal attributes
 - Example: Color (red, yellow, blue, green), profession, etc.
- 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 (e.g., freshman, sophomore, junior, senior)
- Can be treated like interval-scaled
 - Replace an ordinal variable value by its 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
$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$
- Example: freshman: 0; sophomore: 1/3; junior: 2/3; senior 1
 - Then distance: $d(\text{freshman}, \text{senior}) = 1$, $d(\text{junior}, \text{senior}) = 1/3$
- Compute the dissimilarity using methods for interval-scaled variables

Attributes of Mixed Type

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

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

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

Cosine Similarity of Two Vectors

- A **document** can be represented by a bag of terms or a long vector, with each attribute recording the *frequency* of a particular term (such as word, keyword, or phrase) in the document

Document	team	coach	hockey	baseball	soccer	penalty	score	win	loss	season
Document1	5	0	3	0	2	0	0	2	0	0
Document2	3	0	2	0	1	1	0	1	0	1
Document3	0	7	0	2	1	0	0	3	0	0
Document4	0	1	0	0	1	2	2	0	3	0

- Other vector objects: Gene features in micro-arrays
- Applications: Information retrieval, biologic taxonomy, gene feature mapping, etc.
- Cosine measure: If d_1 and d_2 are two vectors (e.g., term-frequency vectors), then

$$\cos(d_1, d_2) = \frac{d_1 \bullet d_2}{\|d_1\| \times \|d_2\|}$$

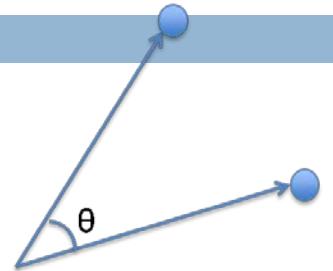
where \bullet indicates vector dot product, $\|d\|$: the length of vector d

Example: Calculating Cosine Similarity

- Calculating Cosine Similarity:

$$\cos(d_1, d_2) = \frac{d_1 \bullet d_2}{\|d_1\| \times \|d_2\|} \quad sim(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$

where \bullet indicates vector dot product, $\|d\|$: the length of vector d



Example: Calculating Cosine Similarity

- Calculating Cosine Similarity:

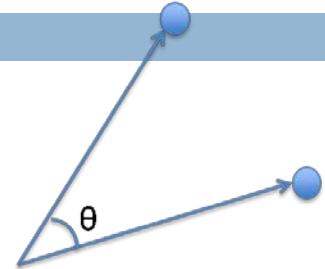
$$\cos(d_1, d_2) = \frac{d_1 \bullet d_2}{\|d_1\| \times \|d_2\|} \quad sim(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$

where \bullet indicates vector dot product, $\|d\|$: the length of vector d

- Ex: Find the **similarity** between documents 1 and 2.

$$d_1 = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)$$

$$d_2 = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1)$$



- First, calculate vector dot product

$$\begin{aligned} d_1 \bullet d_2 &= 5 \times 3 + 0 \times 0 + 3 \times 2 + 0 \times 0 + 2 \times 1 + 0 \times 1 + 0 \times 1 + 2 \times 1 + 0 \\ &\quad \times 0 + 0 \times 1 = 25 \end{aligned}$$

- Then, calculate $\|d_1\|$ and $\|d_2\|$

$$\|d_1\| = \sqrt{5 \times 5 + 0 \times 0 + 3 \times 3 + 0 \times 0 + 2 \times 2 + 0 \times 0 + 0 \times 0 + 2 \times 2 + 0 \times 0 + 0 \times 0} = 6.481$$

$$\|d_2\| = \sqrt{3 \times 3 + 0 \times 0 + 2 \times 2 + 0 \times 0 + 1 \times 1 + 1 \times 1 + 0 \times 0 + 1 \times 1 + 0 \times 0 + 1 \times 1} = 4.12$$

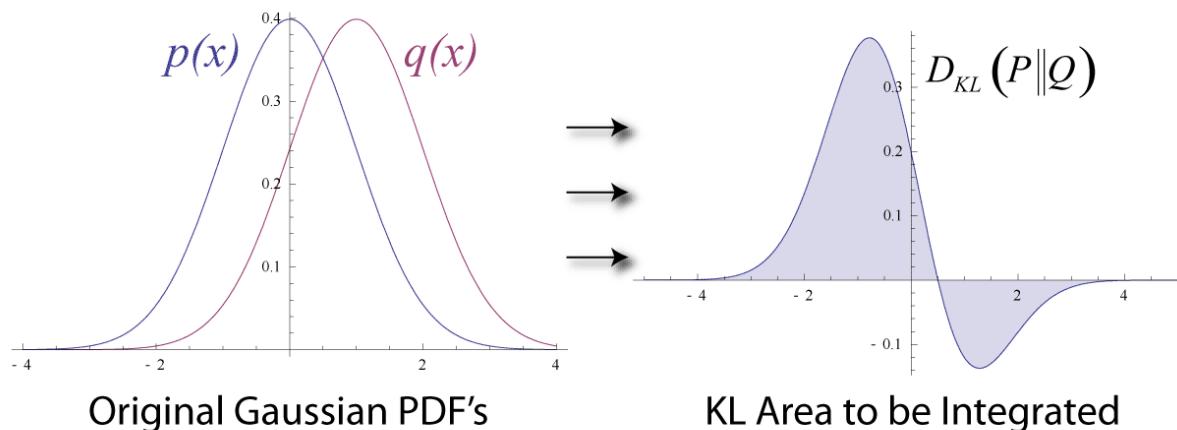
- Calculate cosine similarity: $\cos(d_1, d_2) = 25 / (6.481 \times 4.12) = 0.94$

KL Divergence: Comparing Two Probability Distributions

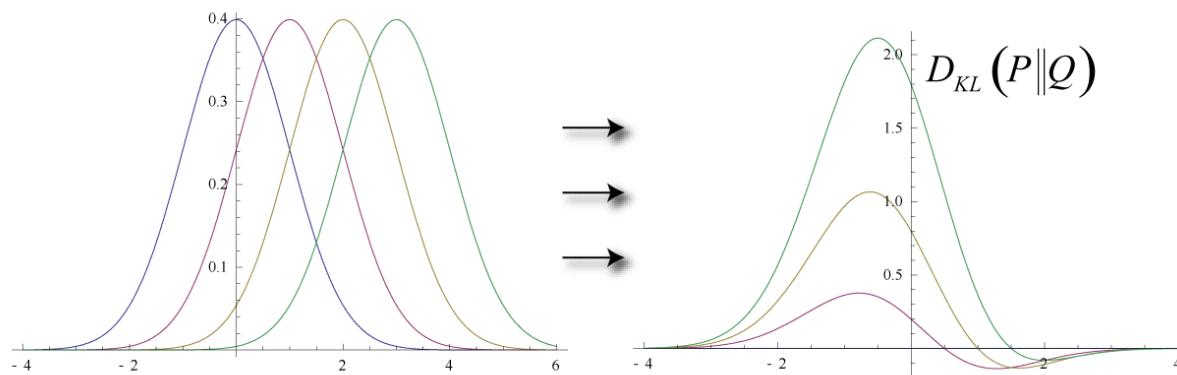
- The Kullback-Leibler (KL) divergence:
Measure the difference between two probability distributions over the same variable x
- From information theory, closely related to *relative entropy*, *information divergence*, and *information for discrimination*
- $D_{KL}(p(x) \parallel q(x))$: divergence of $q(x)$ from $p(x)$, measuring the information lost when $q(x)$ is used to approximate $p(x)$

Discrete form

$$D_{KL}(p(x) \parallel q(x)) = \sum_{x \in X} p(x) \ln \frac{p(x)}{q(x)}$$
$$D_{KL}(p(x) \parallel q(x)) = \int_{-\infty}^{\infty} p(x) \ln \frac{p(x)}{q(x)} dx$$



Original Gaussian PDF's



KL Area to be Integrated

Ack.: Wikipedia entry: The Kullback-Leibler (KL) divergence

Continuous form

More on KL Divergence

$$D_{KL}(p(x)||q(x)) = \sum_{x \in X} p(x) \ln \frac{p(x)}{q(x)}$$

- The KL divergence measures the expected number of extra bits required to code samples from $p(x)$ ("true" distribution) when using a code based on $q(x)$, which represents a theory, model, description, or approximation of $p(x)$
- **The KL divergence is not a distance measure, not a metric: asymmetric ($D_{KL}(P||Q)$ does not equal $D_{KL}(Q||P)$)**
- In applications, P typically represents the "true" distribution of data, observations, or a precisely calculated theoretical distribution, while Q typically represents a theory, model, description, or approximation of P .
- The Kullback–Leibler divergence from Q to P , denoted $D_{KL}(P||Q)$, is a measure of the information gained when one revises one's beliefs from the prior probability distribution Q to the posterior probability distribution P . In other words, it is the amount of information lost when Q is used to approximate P .
- The KL divergence is sometimes also called the information gain achieved if P is used instead of Q . It is also called the relative entropy of P with respect to Q .

Subtlety at Computing the KL Divergence

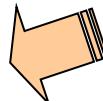
- Base on the formula, $D_{KL}(P||Q) \geq 0$ and $D_{KL}(P || Q) = 0$ if and only if $P = Q$
- How about when $p = 0$ or $q = 0$?
 - $\lim_{p \rightarrow 0} p \log p = 0$
 - when $p \neq 0$ but $q = 0$, $D_{KL}(p || q)$ is defined as ∞ , i.e., if one event e is possible (i.e., $p(e) > 0$), and the other predicts it is absolutely impossible (i.e., $q(e) = 0$), then the two distributions are absolutely different
- However, in practice, P and Q are derived from frequency distributions, not counting the possibility of unseen events. Thus smoothing is needed.

$$D_{KL}(p(x)||q(x)) = \sum_{x \in X} p(x) \ln \frac{p(x)}{q(x)}$$

Subtlety at Computing the KL Divergence

- Base on the formula, $D_{KL}(P, Q) \geq 0$ and $D_{KL}(P || Q) = 0$ if and only if $P = Q$
- How about when $p = 0$ or $q = 0$?
 - $\lim_{p \rightarrow 0} p \log p = 0$
 - when $p \neq 0$ but $q = 0$, $D_{KL}(p || q)$ is defined as ∞ , i.e., if one event e is possible (i.e., $p(e) > 0$), and the other predicts it is absolutely impossible (i.e., $q(e) = 0$), then the two distributions are absolutely different
- However, in practice, P and Q are derived from frequency distributions, not counting the possibility of unseen events. Thus *smoothing* is needed
- Example: $P : (a : 3/5, b : 1/5, c : 1/5)$. $Q : (a : 5/9, b : 3/9, d : 1/9)$
 - need to introduce a small constant ϵ , e.g., $\epsilon = 10^{-3}$
 - The sample set observed in P, $SP = \{a, b, c\}$, $SQ = \{a, b, d\}$, $SU = \{a, b, c, d\}$
 - Smoothing, add missing symbols to each distribution, with probability ϵ
 - $P' : (a : 3/5 - \epsilon/3, b : 1/5 - \epsilon/3, c : 1/5 - \epsilon/3, d : \epsilon)$
 - $Q' : (a : 5/9 - \epsilon/3, b : 3/9 - \epsilon/3, c : \epsilon, d : 1/9 - \epsilon/3)$
 - $D_{KL}(P' || Q')$ can then be computed easily

Data & Data Preprocessing

- Data Objects and Attribute Types
- Basic Statistical Descriptions of Data
- Measuring Data Similarity and Dissimilarity
- Data Preprocessing: An Overview 
- Summary

Why Preprocess the Data? -Data Quality Issues

- Measures for data quality: A multidimensional view
 - Accuracy: correct or wrong, accurate or not
 - Completeness: not recorded, unavailable, ...
 - Consistency: some modified but some not, dangling, ...
 - Timeliness: timely update?
 - Believability: how trustable the data are correct?
 - Interpretability: how easily the data can be understood?

Data Quality Issues - Examples

- Data in the Real World Is Dirty: Lots of potentially incorrect data, e.g., instrument faulty, human or computer error, and transmission error
 - Incomplete: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
 - e.g., *Occupation* = “ ” (missing data)
 - Noisy: containing noise, errors, or outliers
 - e.g., *Salary* = “−10” (an error)
 - Inconsistent: containing discrepancies in codes or names, e.g.,
 - *Age* = “42”, *Birthday* = “03/07/2010”
 - Was rating “1, 2, 3”, now rating “A, B, C”
 - discrepancy between duplicate records
 - Intentional (e.g., disguised missing data)
 - Jan. 1 as everyone’s birthday?

Missing (Incomplete) Values

- Reasons for missing values
 - Information is not collected
(e.g., people decline to give their age and weight)
 - Attributes may not be applicable to all cases
(e.g., annual income is not applicable to children)

Missing (Incomplete) Values

- Reasons for missing values
 - Information is not collected
(e.g., people decline to give their age and weight)
 - Attributes may not be applicable to all cases
(e.g., annual income is not applicable to children)

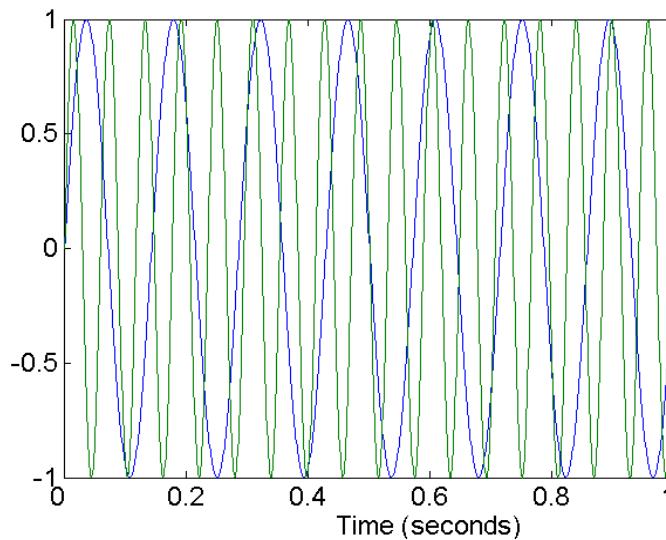
How to handle them?

How to Handle Missing Data?

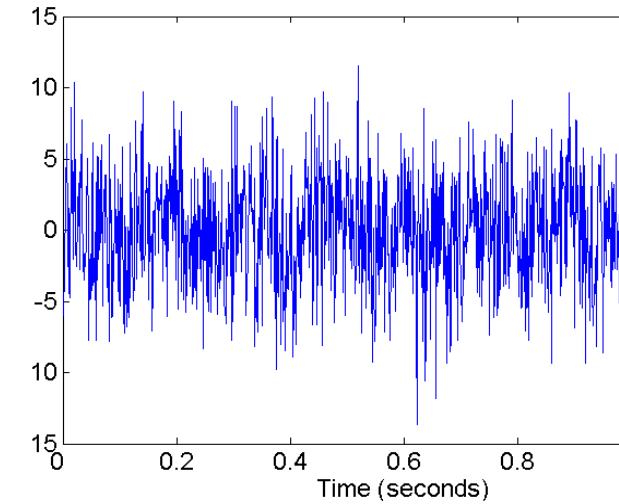
- **Ignore the tuple**: usually done when class label is missing (when doing classification)—not effective when the % of missing values per attribute varies considerably
- Fill in the missing value **manually**: tedious + infeasible?
- Fill in it automatically with
 - a global constant : e.g., “unknown”, a new class?!
 - **the attribute mean**
 - the attribute mean for all samples belonging to the same class: smarter
 - **the most probable value**: inference-based such as Bayesian formula or decision tree

Noise

- Noise refers to modification of original values
 - Examples: distortion of a person's voice when talking on a poor phone and "snow" on television screen



Two Sine Waves



Two Sine Waves + Noise



How to Handle Noisy Data?

- Binning
 - First sort data and partition into (equal-frequency) bins
 - Then one can **smooth by bin means, smooth by bin median, smooth by bin boundaries**, etc.
- Regression
 - Smooth by fitting the data into regression functions
- Clustering
 - **Detect and remove outliers**
- Semi-supervised: Combined computer and human inspection
 - Detect suspicious values and check by human (e.g., deal with possible outliers)

Data Cleaning as a Process

□ Data discrepancy detection

- Use metadata (e.g., domain, range, dependency, distribution)
- Check field overloading
- Check uniqueness rule, consecutive rule and null rule
- Use commercial tools
 - Data scrubbing: use simple domain knowledge (e.g., postal code, spell-check) to **detect errors and make corrections**
 - Data auditing: by analyzing data to discover rules and relationship to **detect violators** (e.g., correlation and clustering to find outliers)

□ Data migration and integration

- Data migration tools: allow transformations to be specified
- ETL (Extraction/Transformation>Loading) tools: allow users to specify transformations through a graphical user interface

□ Integration of the two processes

- Iterative and interactive

Data & Data Preprocessing

- Data Objects and Attribute Types
- Basic Statistical Descriptions of Data
- Measuring Data Similarity and Dissimilarity
- Data Preprocessing: An Overview
- Summary 



Summary

- Data attribute types: nominal, binary, ordinal, interval-scaled, ratio-scaled
- Many types of data sets, e.g., numerical, text, graph, Web, image.
- Gain insight into the data by:
 - Basic statistical data description: central tendency, dispersion
 - Measure data similarity
- Data quality issues and preprocessing
- Many methods have been developed but still an active area of research

References

- W. Cleveland, *Visualizing Data*, Hobart Press, 1993
- T. Dasu and T. Johnson. *Exploratory Data Mining and Data Cleaning*. John Wiley, 2003
- U. Fayyad, G. Grinstein, and A. Wierse. *Information Visualization in Data Mining and Knowledge Discovery*, Morgan Kaufmann, 2001
- L. Kaufman and P. J. Rousseeuw. *Finding Groups in Data: an Introduction to Cluster Analysis*. John Wiley & Sons, 1990.
- H. V. Jagadish et al., Special Issue on Data Reduction Techniques. *Bulletin of the Tech. Committee on Data Eng.*, 20(4), Dec. 1997
- D. A. Keim. Information visualization and visual data mining, *IEEE trans. on Visualization and Computer Graphics*, 8(1), 2002
- D. Pyle. *Data Preparation for Data Mining*. Morgan Kaufmann, 1999
- S. Santini and R. Jain," Similarity measures", *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 21(9), 1999
- E. R. Tufte. *The Visual Display of Quantitative Information*, 2nd ed., Graphics Press, 2001
- C. Yu, et al., Visual data mining of multimedia data for social and behavioral studies, *Information Visualization*, 8(1), 2009

References

- W. Cleveland, *Visualizing Data*, Hobart Press, 1993
- T. Dasu and T. Johnson. *Exploratory Data Mining and Data Cleaning*. John Wiley, 2003
- U. Fayyad, G. Grinstein, and A. Wierse. *Information Visualization in Data Mining and Knowledge Discovery*, Morgan Kaufmann, 2001
- L. Kaufman and P. J. Rousseeuw. *Finding Groups in Data: an Introduction to Cluster Analysis*. John Wiley & Sons, 1990.
- H. V. Jagadish et al., Special Issue on Data Reduction Techniques. *Bulletin of the Tech. Committee on Data Eng.*, 20(4), Dec. 1997
- D. A. Keim. Information visualization and visual data mining, *IEEE trans. on Visualization and Computer Graphics*, 8(1), 2002
- D. Pyle. *Data Preparation for Data Mining*. Morgan Kaufmann, 1999
- S. Santini and R. Jain," Similarity measures", *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 21(9), 1999
- E. R. Tufte. *The Visual Display of Quantitative Information*, 2nd ed., Graphics Press, 2001
- C. Yu, et al., Visual data mining of multimedia data for social and behavioral studies, *Information Visualization*, 8(1), 2009

Next Class: Data Preprocessing

Homework #1 is Out!

Due Date: 09/12/2017 (Two weeks from now)

Online Submissions ONLY!

Backup slides

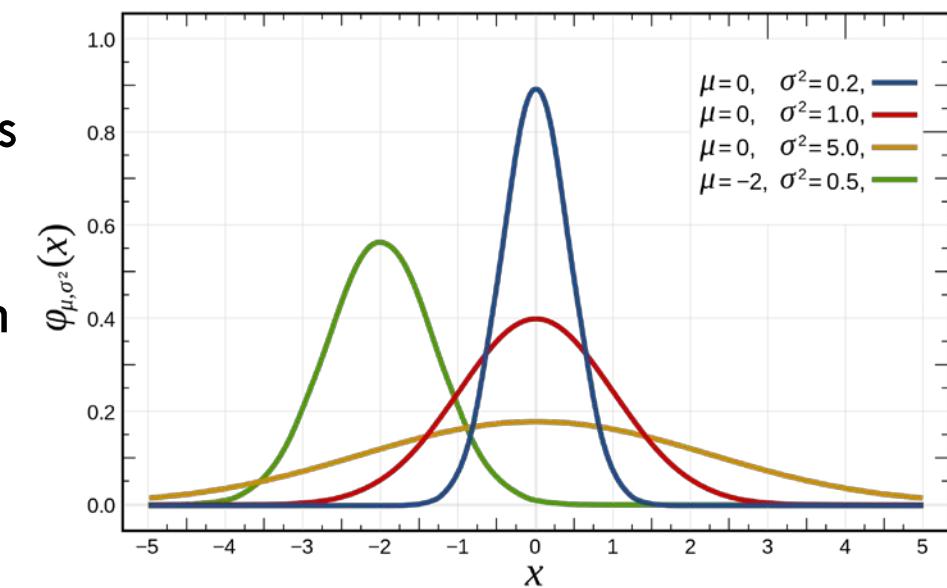
Basic Statistical Descriptions of Data

Motivation

- To better understand the data: central tendency, variation and spread

Data dispersion characteristics

- Median, max, min, quantiles, outliers, variance, ...
- Numerical dimensions correspond to sorted intervals
- Data dispersion:
 - Analyzed with multiple granularities of precision
 - Boxplot or quantile analysis on sorted intervals
- Dispersion analysis on computed measures
 - Folding measures into numerical dimensions
 - Boxplot or quantile analysis on the transformed cube



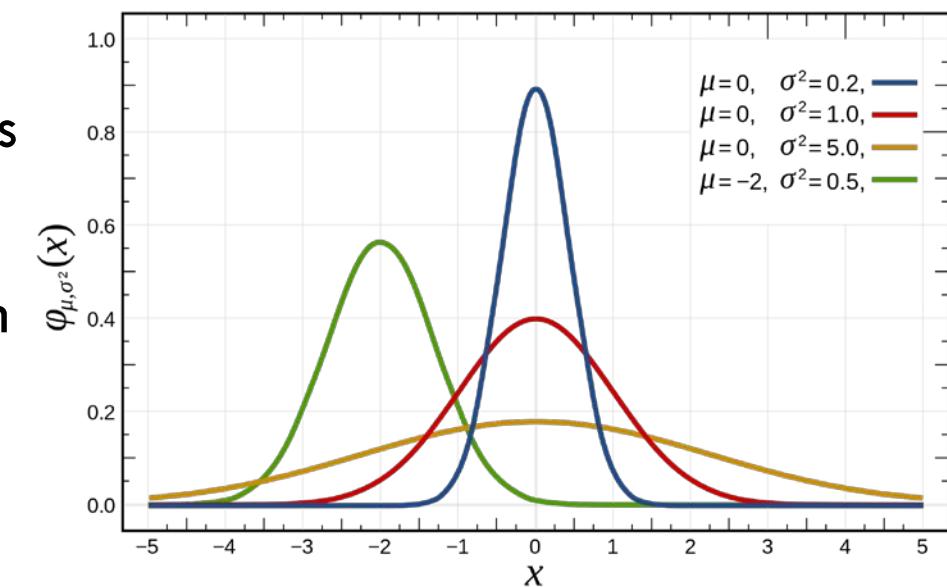
Basic Statistical Descriptions of Data

Motivation

- To better understand the data: central tendency, variation and spread

Data dispersion characteristics

- Median, max, min, quantiles, outliers, variance, ...
- Numerical dimensions correspond to sorted intervals
- Data dispersion:
 - Analyzed with multiple granularities of precision
 - Boxplot or quantile analysis on sorted intervals
- Dispersion analysis on computed measures
 - Folding measures into numerical dimensions
 - Boxplot or quantile analysis on the transformed cube

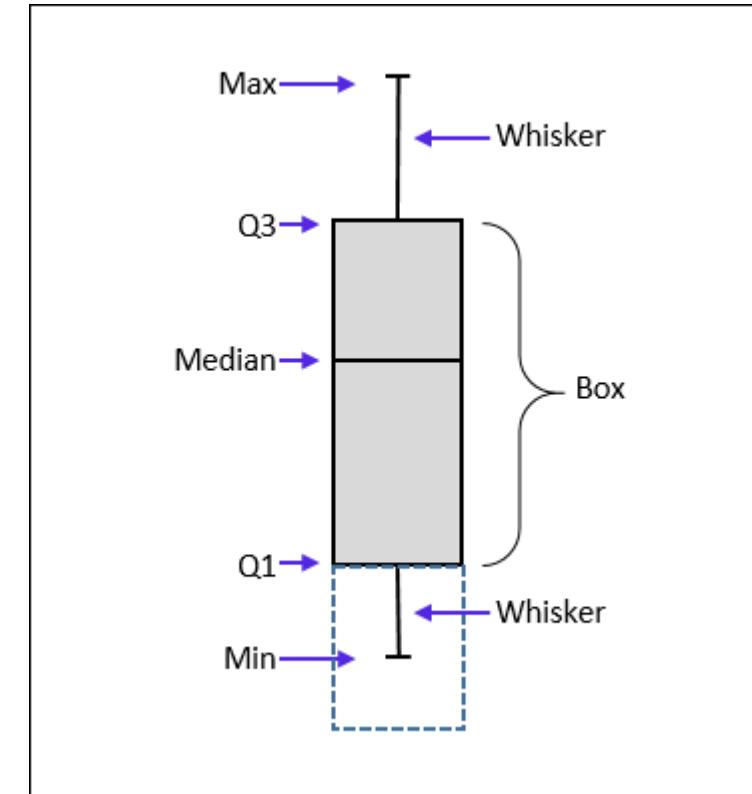


Graphic Displays of Basic Statistical Descriptions

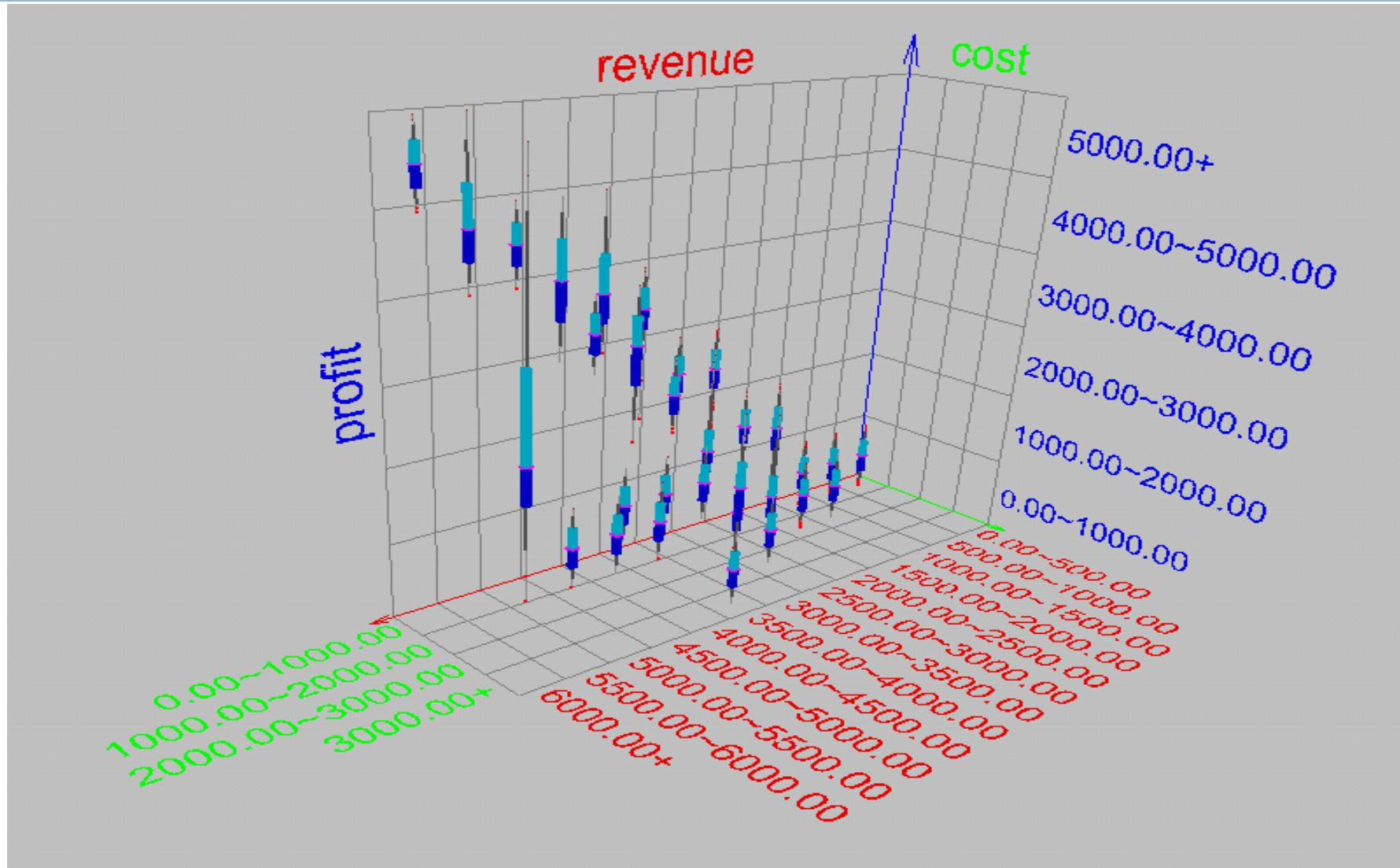
- **Boxplot:** graphic display of five-number summary
- **Histogram:** x-axis are values, y-axis represents frequencies
- **Quantile plot:** each value x_i is paired with f_i , indicating that approximately $100\% * f_i$ of data are $\leq x_i$
- **Quantile-quantile (q-q) plot:** graphs the quantiles of one univariant distribution against the corresponding quantiles of another
- **Scatter plot:** each pair of values is a pair of coordinates and plotted as points in the plane

Measuring the Dispersion of Data: Quartiles & Boxplots

- **Quartiles:** Q_1 (25th percentile), Q_3 (75th percentile)
- **Inter-quartile range:** $IQR = Q_3 - Q_1$
- **Five number summary:** min, Q_1 , median, Q_3 , max
- **Boxplot:** Data is represented with a box
 - Q_1 , Q_3 , IQR: The ends of the box are at the first and third quartiles, i.e., the height of the box is IQR
 - Median (Q_2) is marked by a line within the box
 - Whiskers: two lines outside the box extended to Minimum and Maximum
- Outliers: points beyond a specified outlier threshold, plotted individually
 - **Outlier:** usually, a value higher/lower than $1.5 \times IQR$

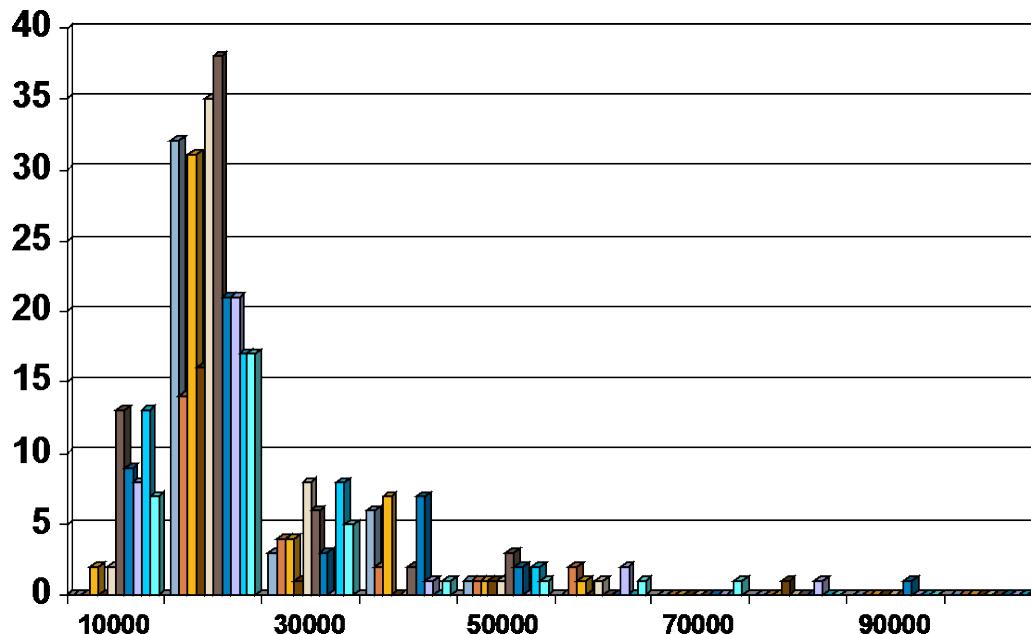


Visualization of Data Dispersion: 3-D Boxplots

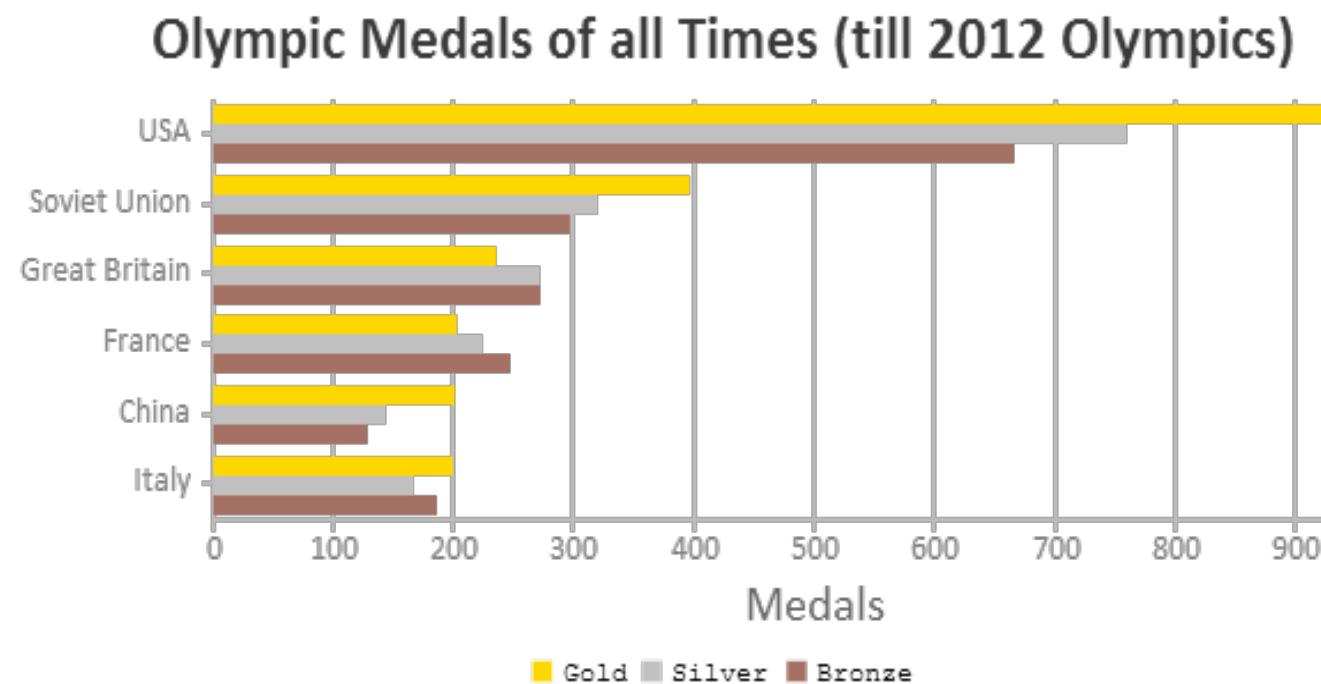


Histogram Analysis

- Histogram: Graph display of tabulated frequencies, shown as bars
- Differences between histograms and bar chart



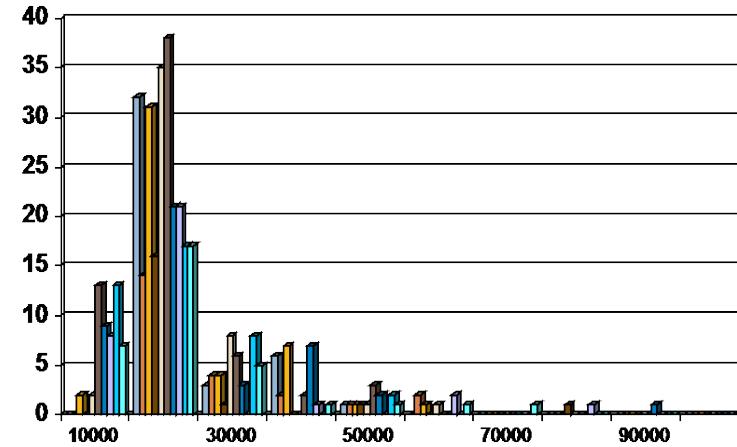
Histogram



Bar chart

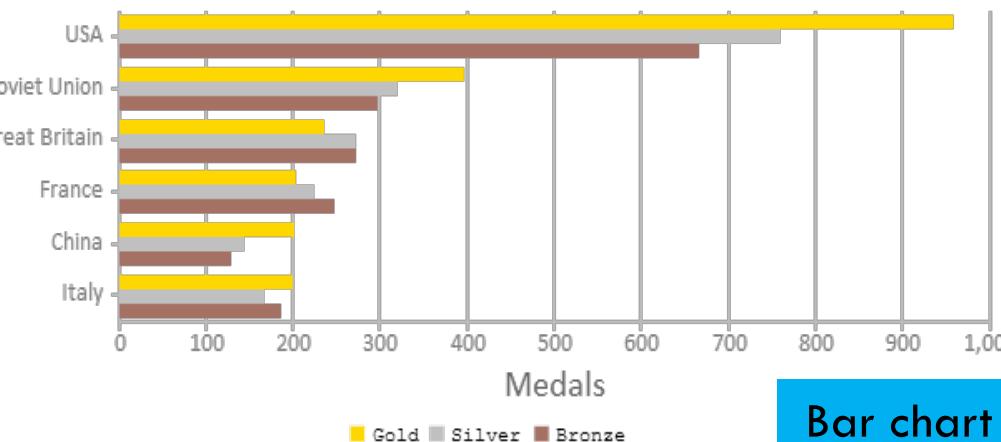
Histogram Analysis

- 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



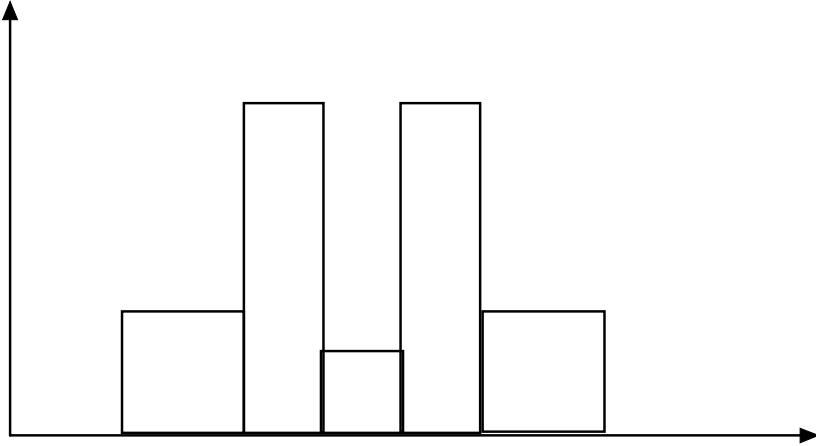
Histogram

Olympic Medals of all Times (till 2012 Olympics)

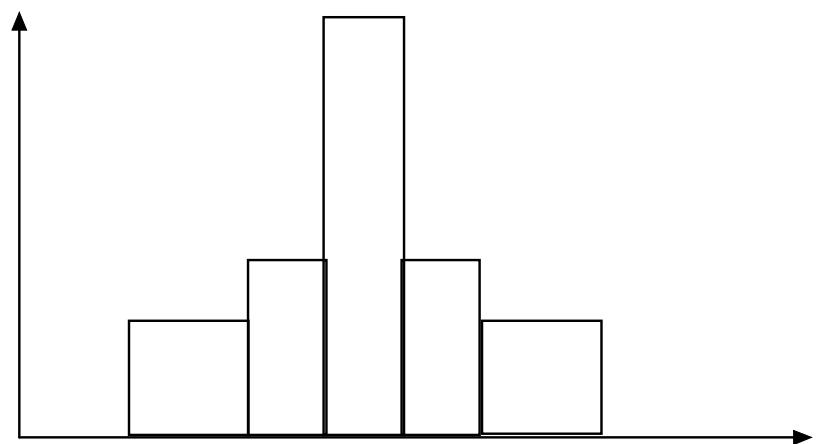


Bar chart

Histograms Often Tell More than Boxplots



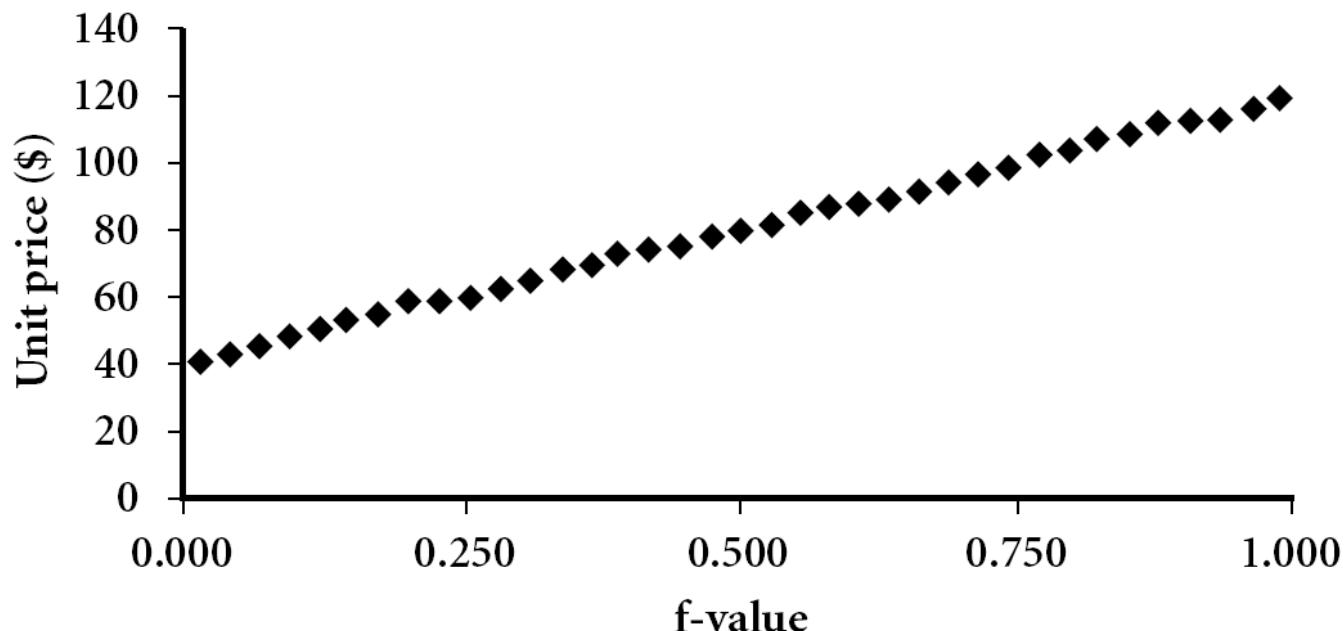
- ❑ The two histograms shown in the left may have the same boxplot representation
- ❑ The same values for: min, Q1, median, Q3, max



- ❑ But they have rather different data distributions

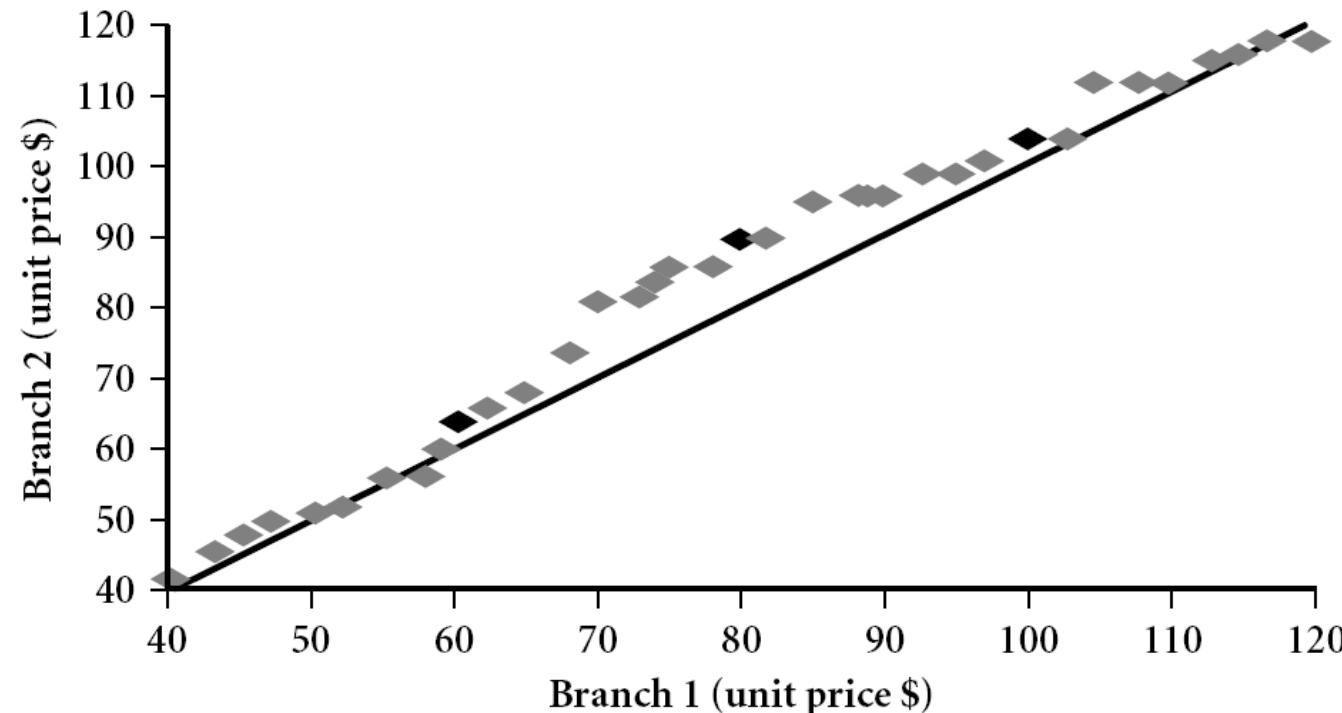
Quantile Plot

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



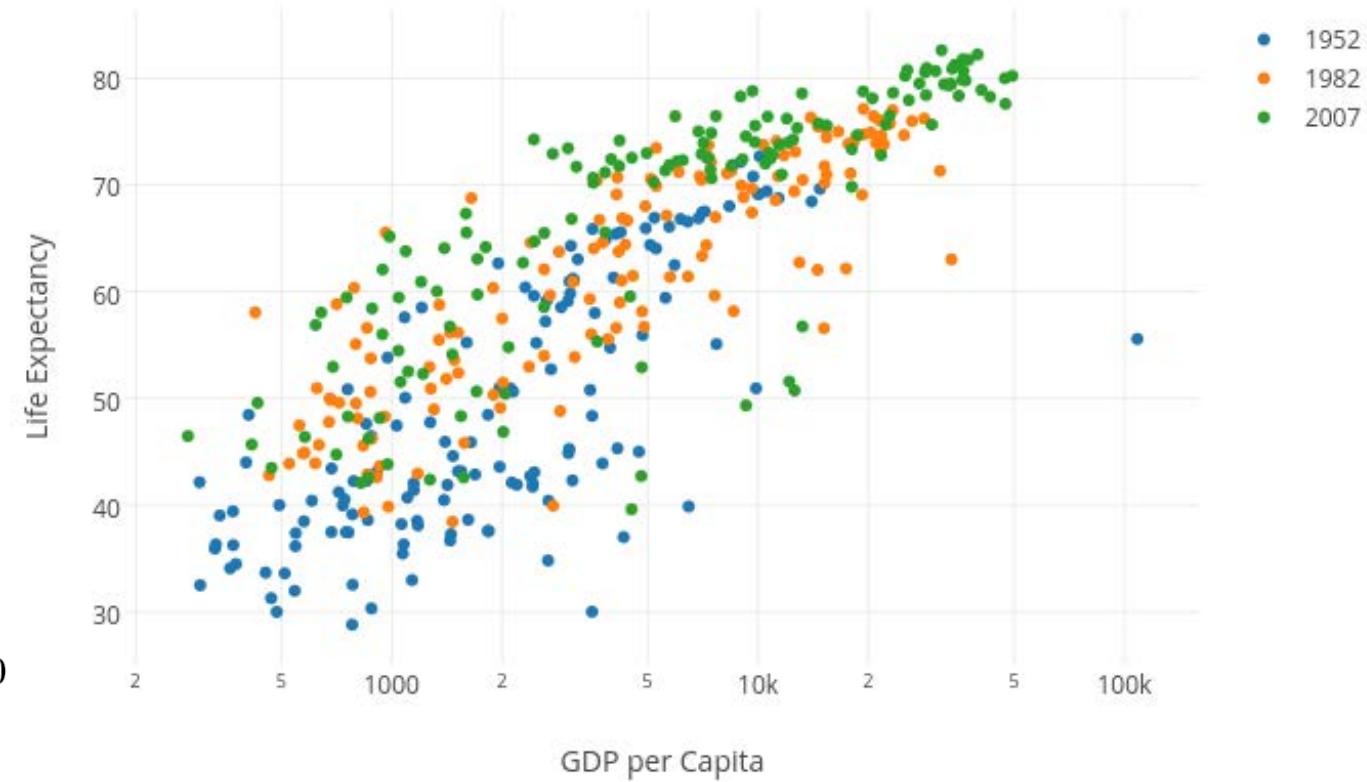
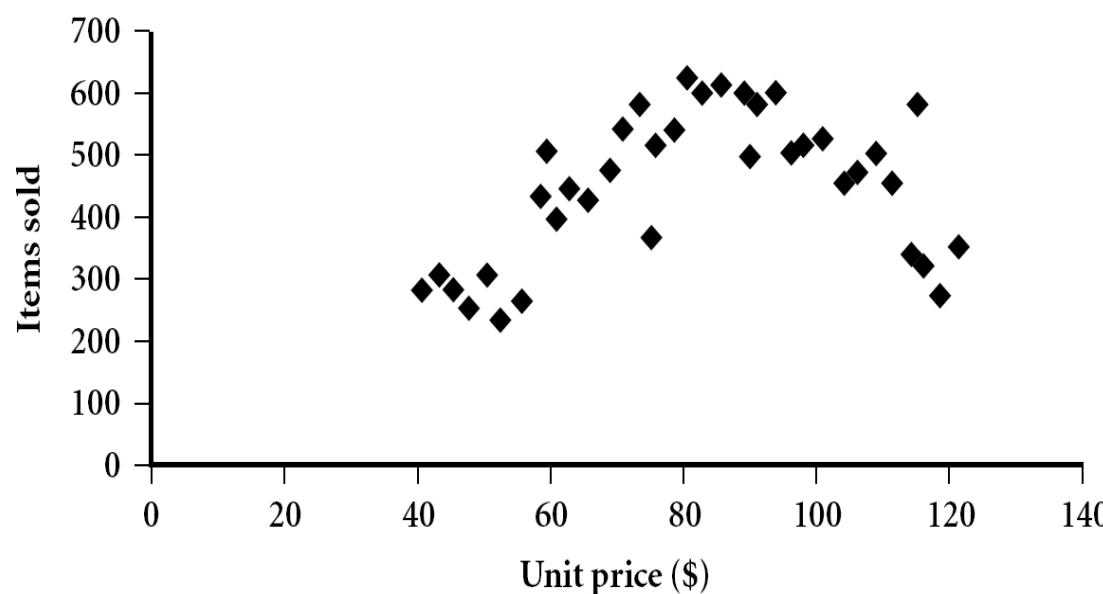
Quantile-Quantile (Q-Q) Plot

- Graphs the quantiles of one univariate distribution against the corresponding quantiles of another
- View: Is there 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

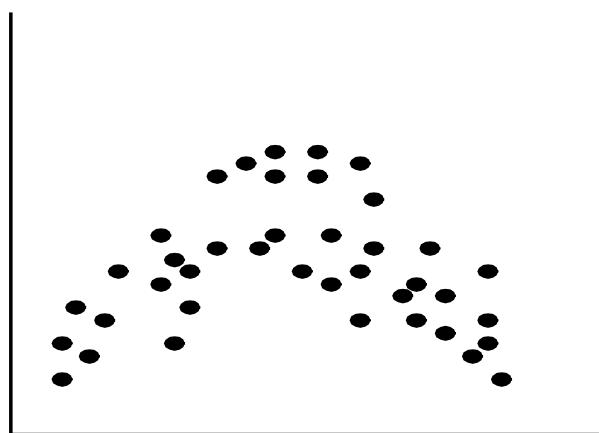
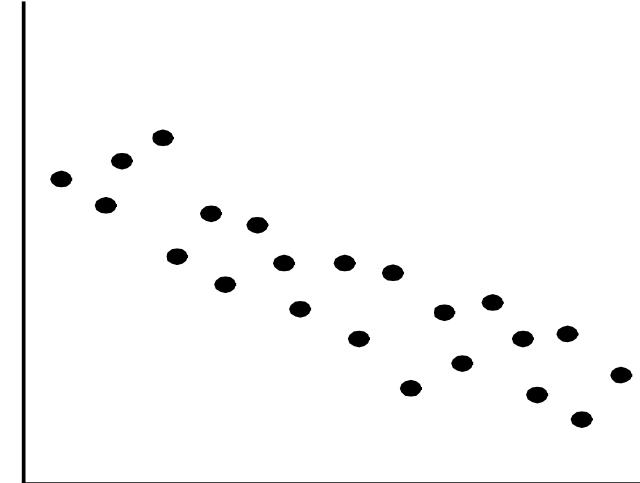
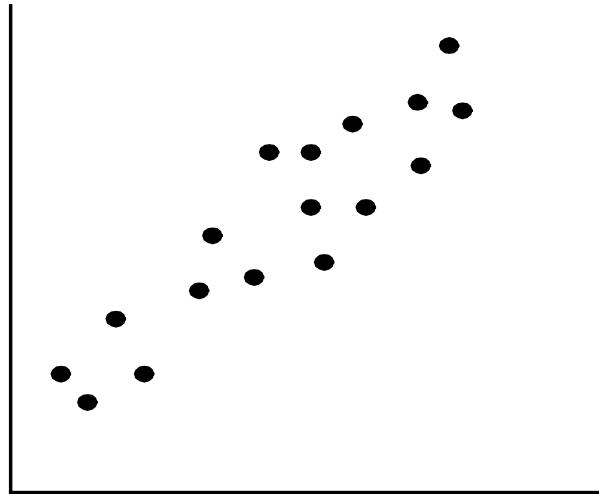


Scatter plot

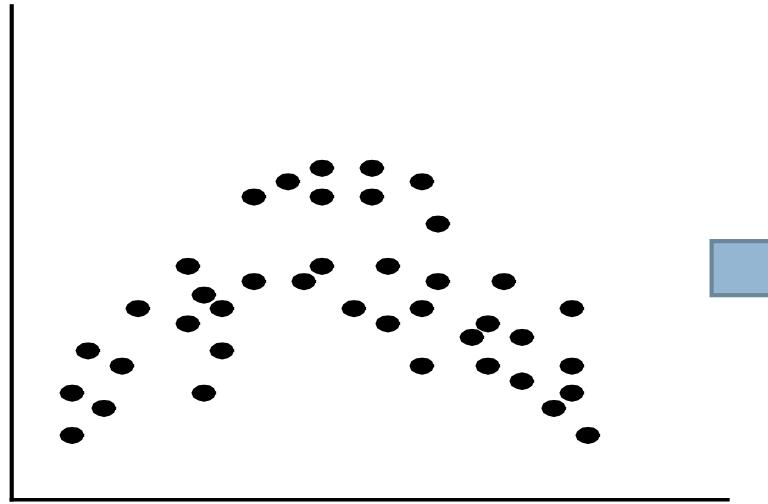
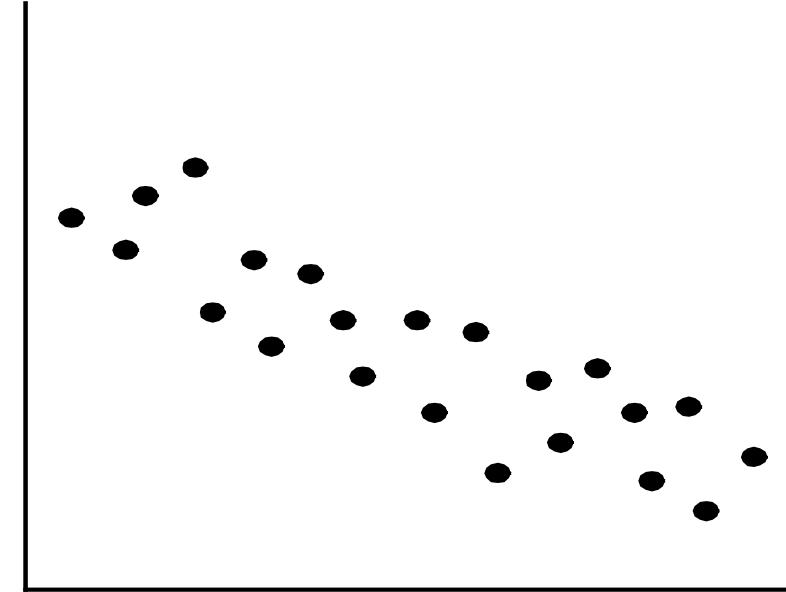
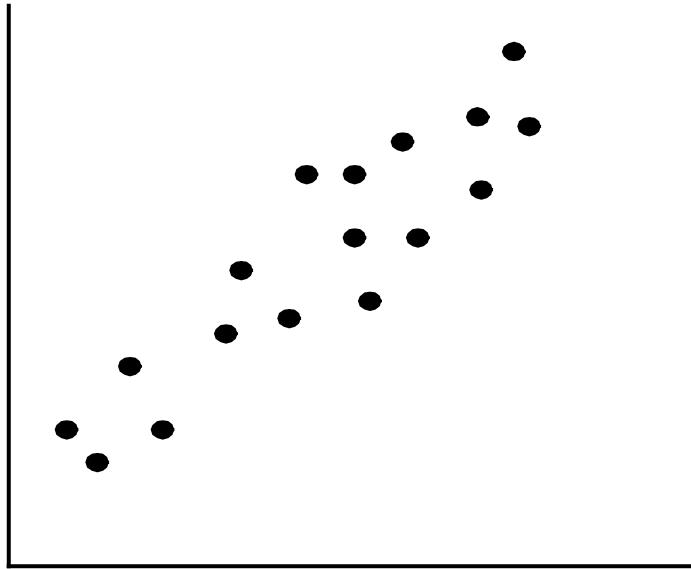
- Provides a first look at bivariate data to see clusters of points, outliers, etc.
- Each pair of values is treated as a pair of coordinates and plotted as points in the plane



Positively and Negatively Correlated Data ?



Positively and Negatively Correlated Data



- The left half fragment is positively correlated
- The right half is negative correlated

Chapter 2. Getting to Know Your Data

- Data Objects and Attribute Types
- Basic Statistical Descriptions of Data
- Data Visualization 
- Measuring Data Similarity and Dissimilarity
- Summary

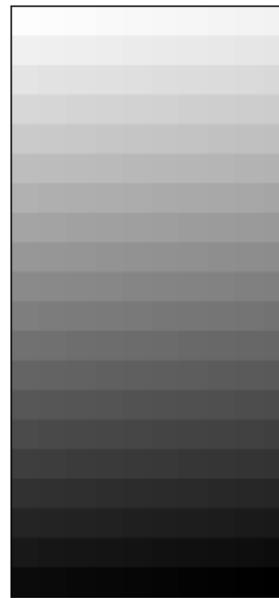


Data Visualization

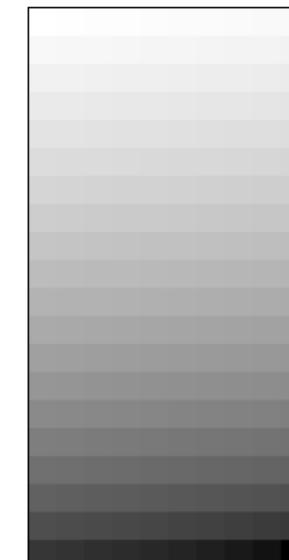
- Why data visualization?
 - Gain insight into an information space by mapping data onto graphical primitives
 - Provide qualitative overview of large data sets
 - Search for patterns, trends, structure, irregularities, relationships among data
 - Help find interesting regions and suitable parameters for further quantitative analysis
 - Provide a visual proof of computer representations derived
- 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

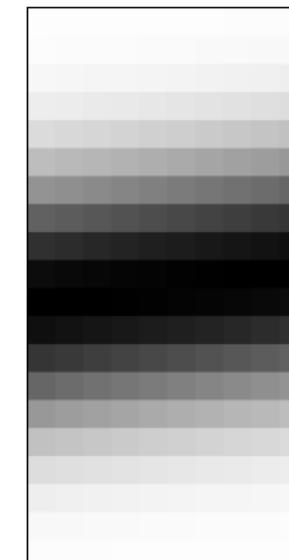
- For a data set of m dimensions, create m windows on the screen, one for each dimension
- The m dimension values of a record are mapped to m pixels at the corresponding positions in the windows
- The colors of the pixels reflect the corresponding values



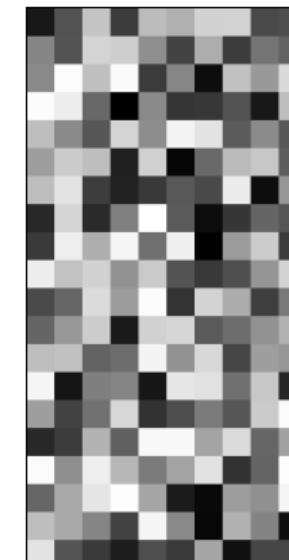
(a) Income



(b) Credit Limit



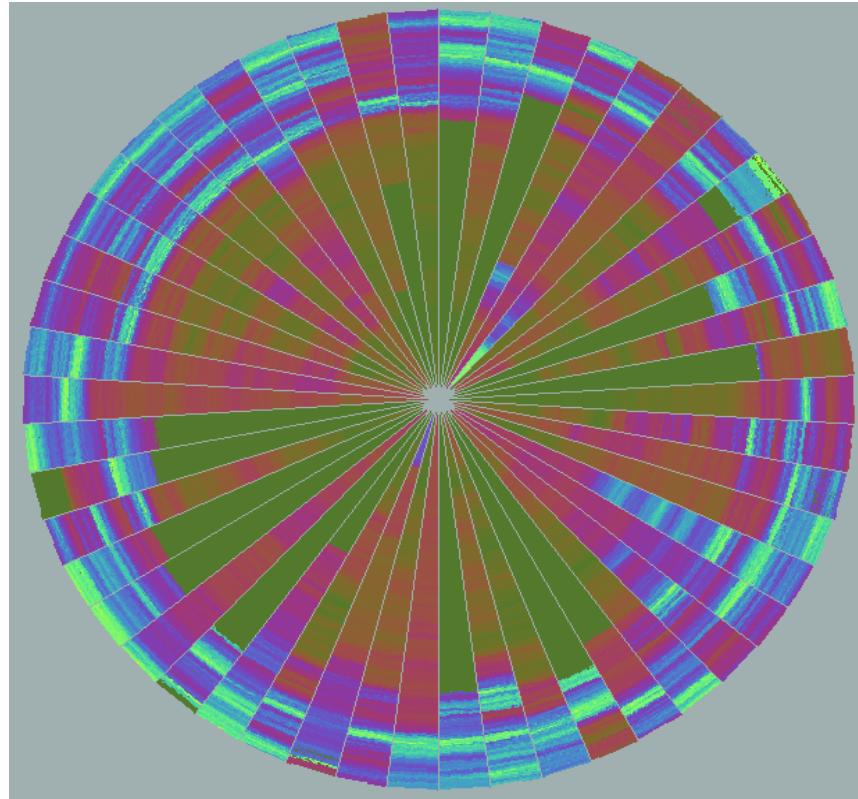
(c) transaction volume



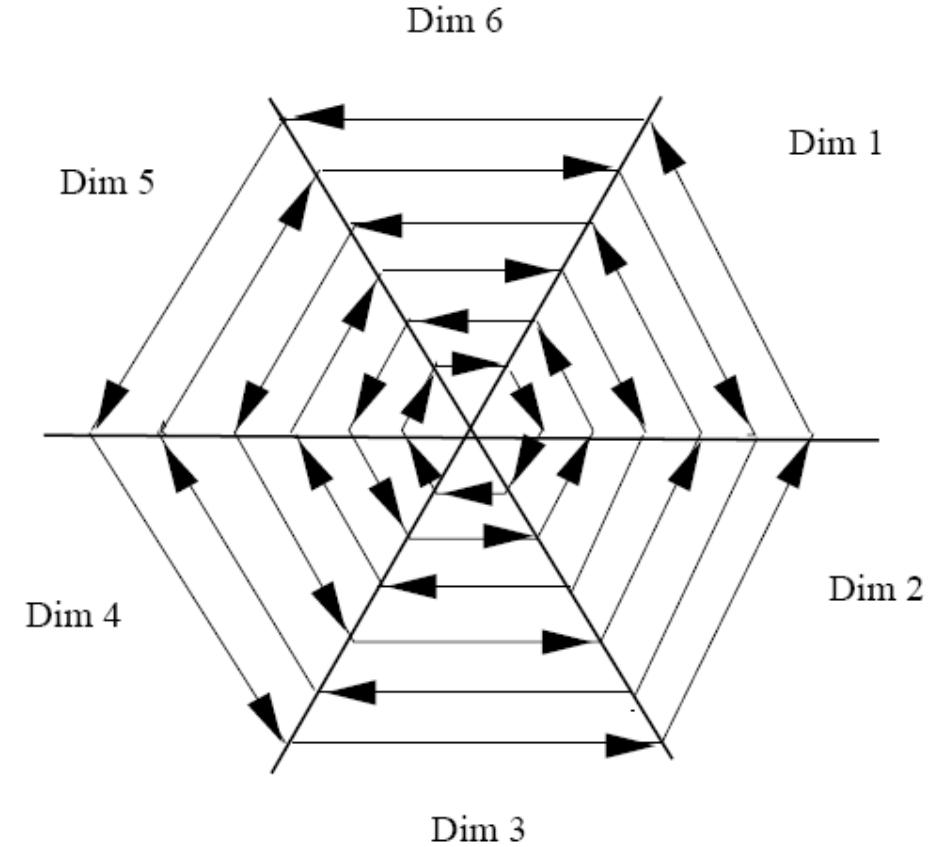
(d) age

Laying Out Pixels in Circle Segments

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



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



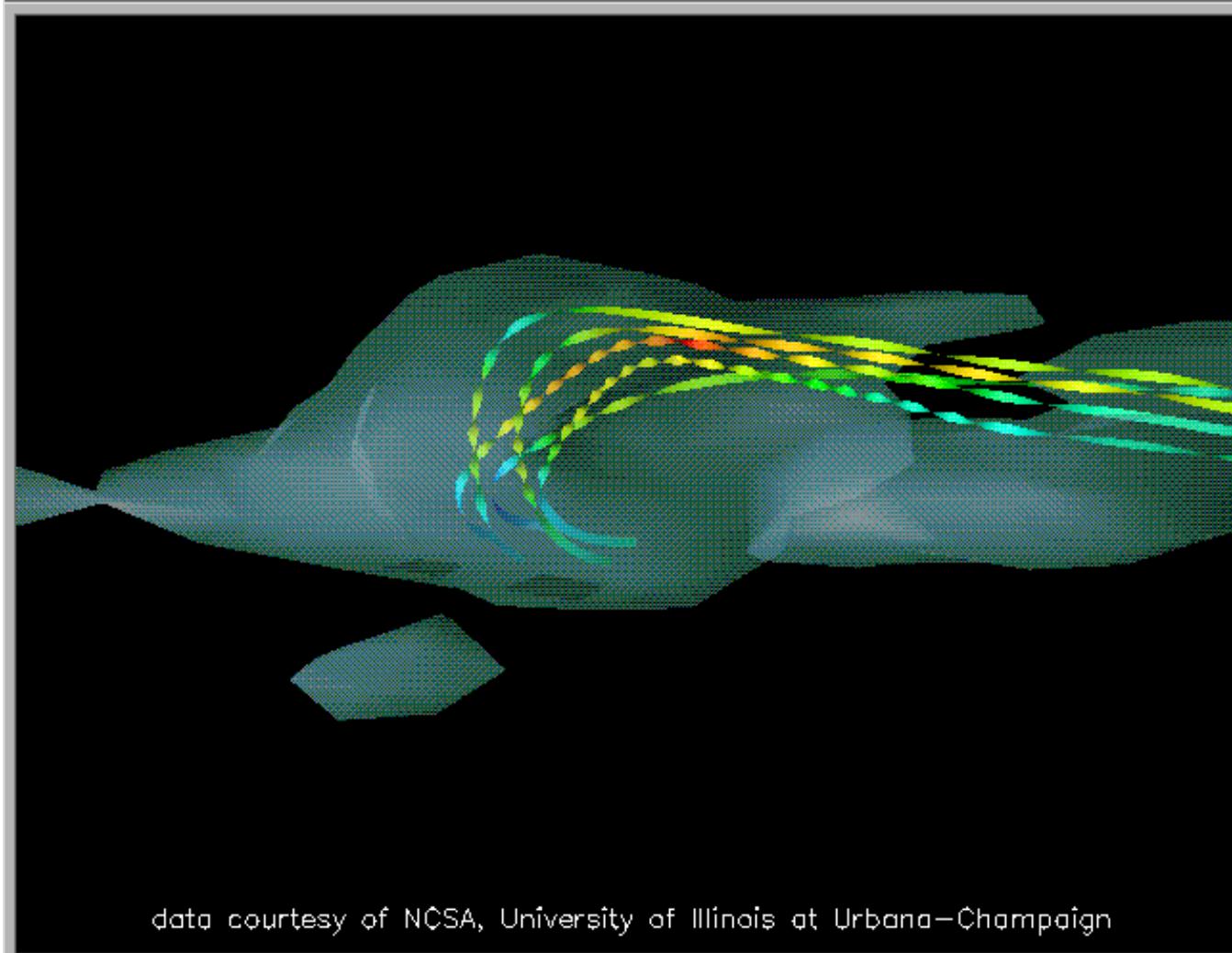


Geometric Projection Visualization Techniques

- Visualization of geometric transformations and projections of the data
- Methods
 - Direct visualization
 - Scatterplot and scatterplot matrices
 - Landscapes
 - Projection pursuit technique: Help users find meaningful projections of multidimensional data
 - Prosection views
 - Hyperslice
 - Parallel coordinates

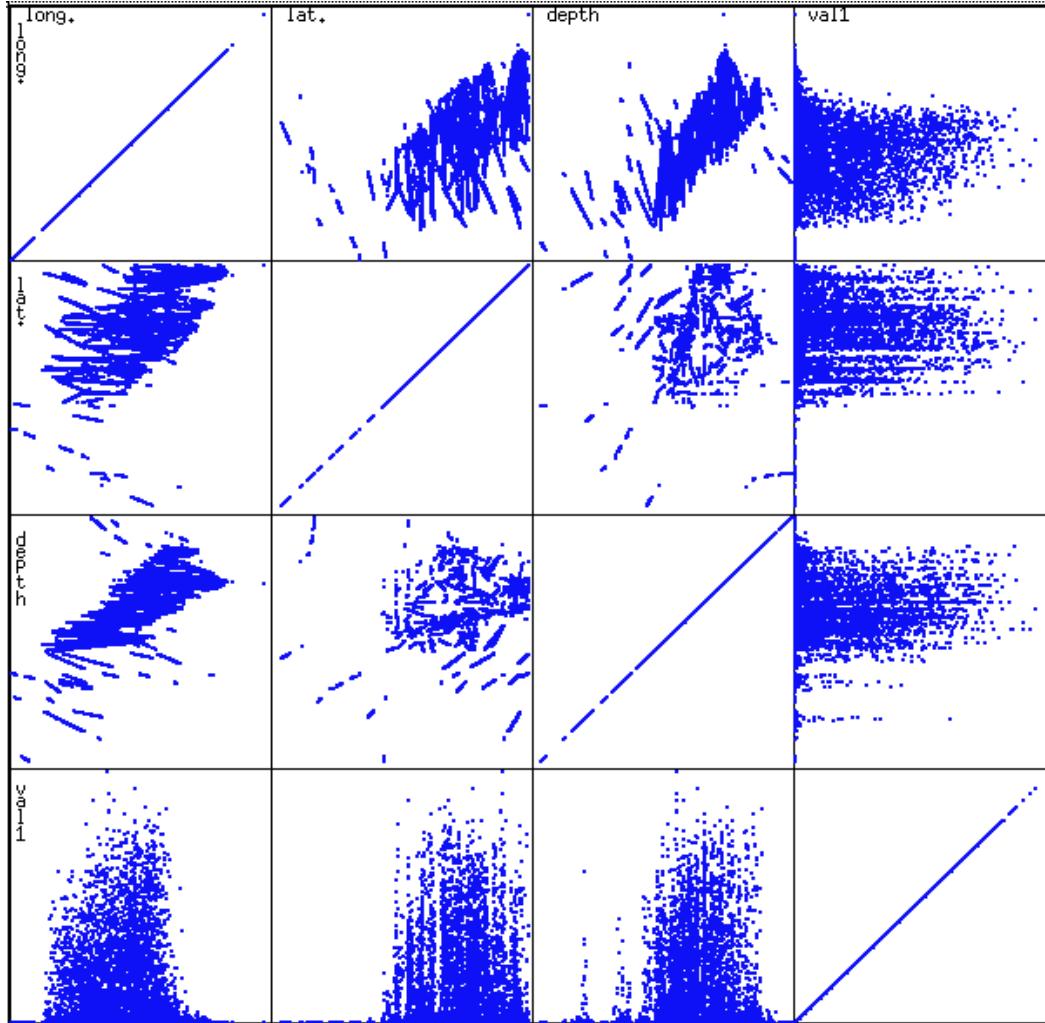
Direct Data Visualization

Ribbons with Twists Based on Vorticity



Scatterplot Matrices

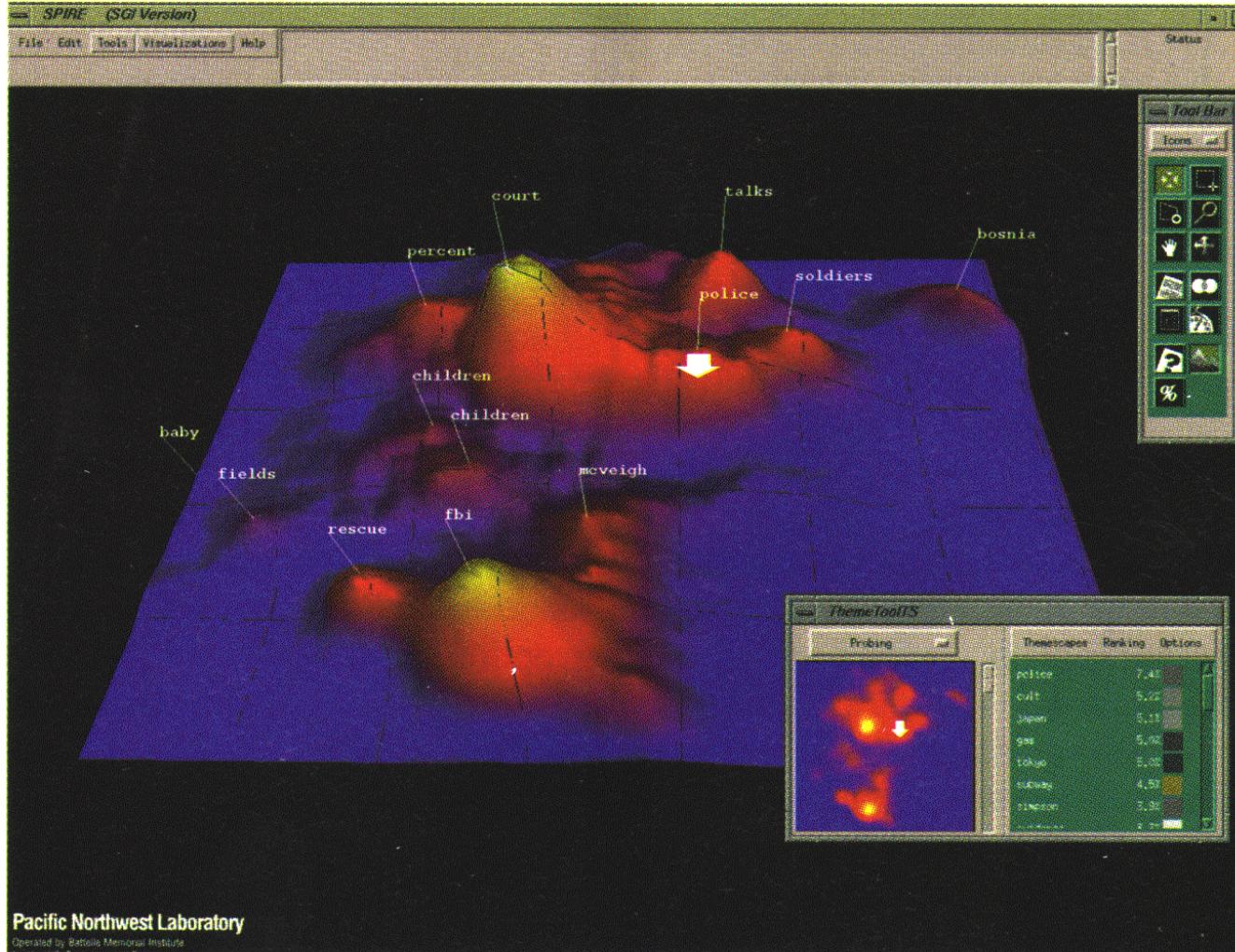
Used by permission of M. Ward, Worcester Polytechnic Institute



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

Landscapes

Used by permission of B. Wright, Visible Decisions Inc.



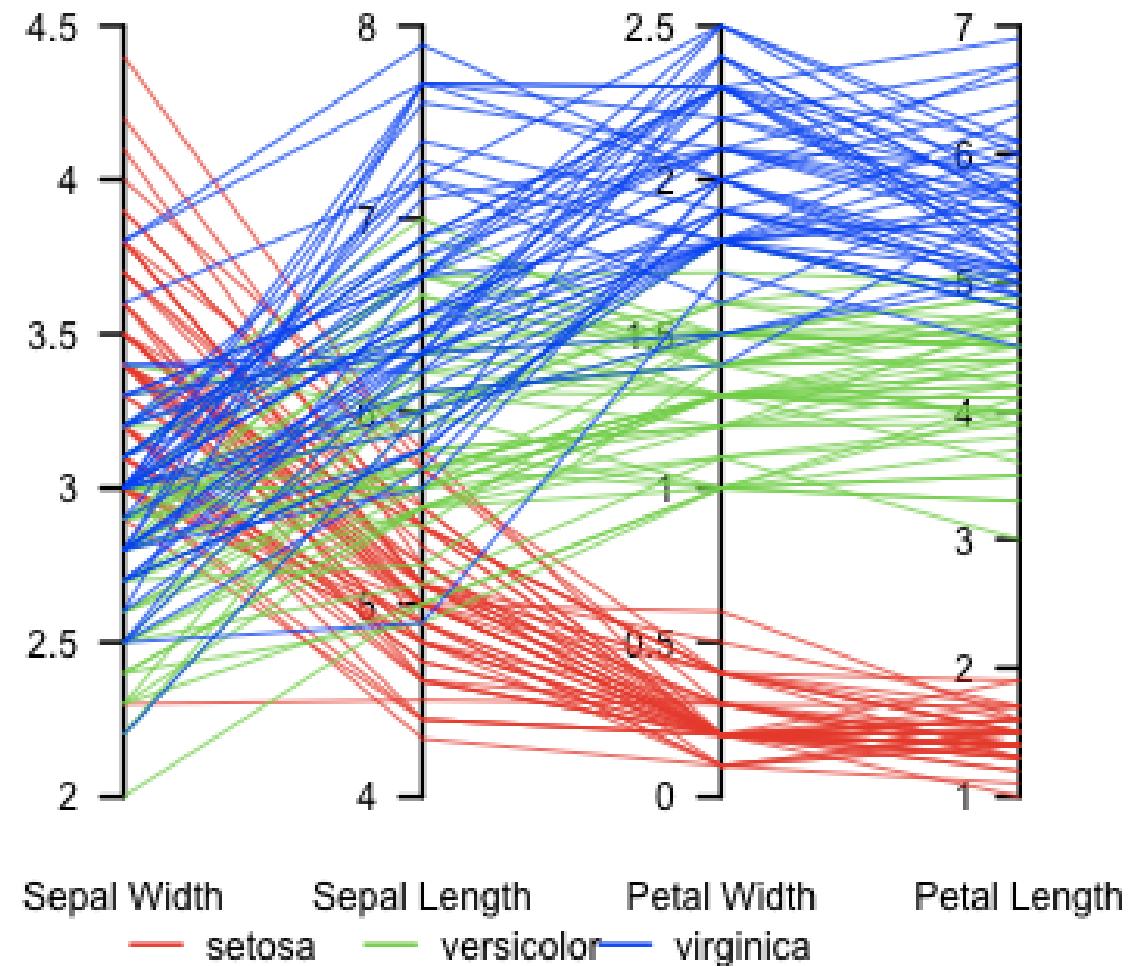
news articles visualized as a landscape

- 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

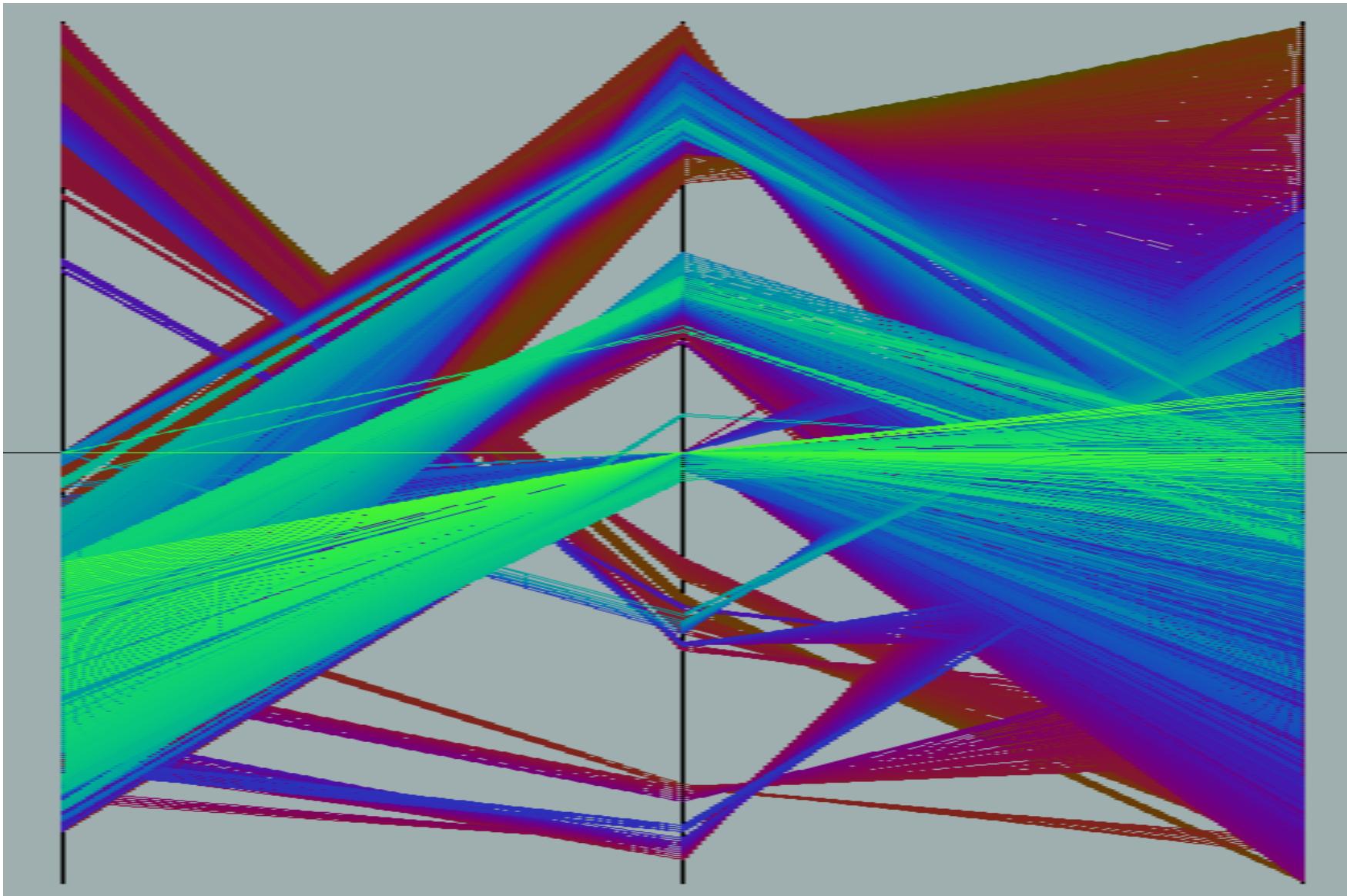
Parallel Coordinates

- 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

Parallel coordinate plot, Fisher's Iris data



Parallel Coordinates of a Data Set

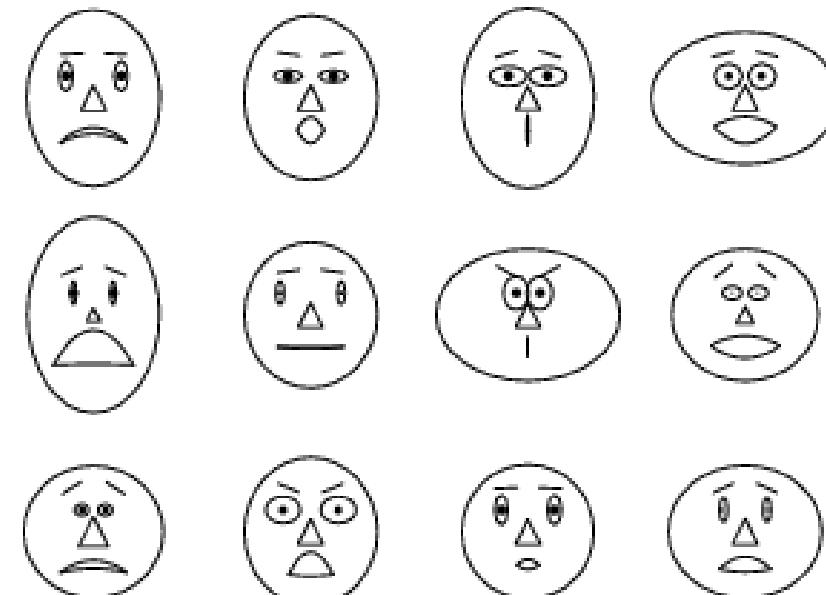


Icon-Based Visualization Techniques

- 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

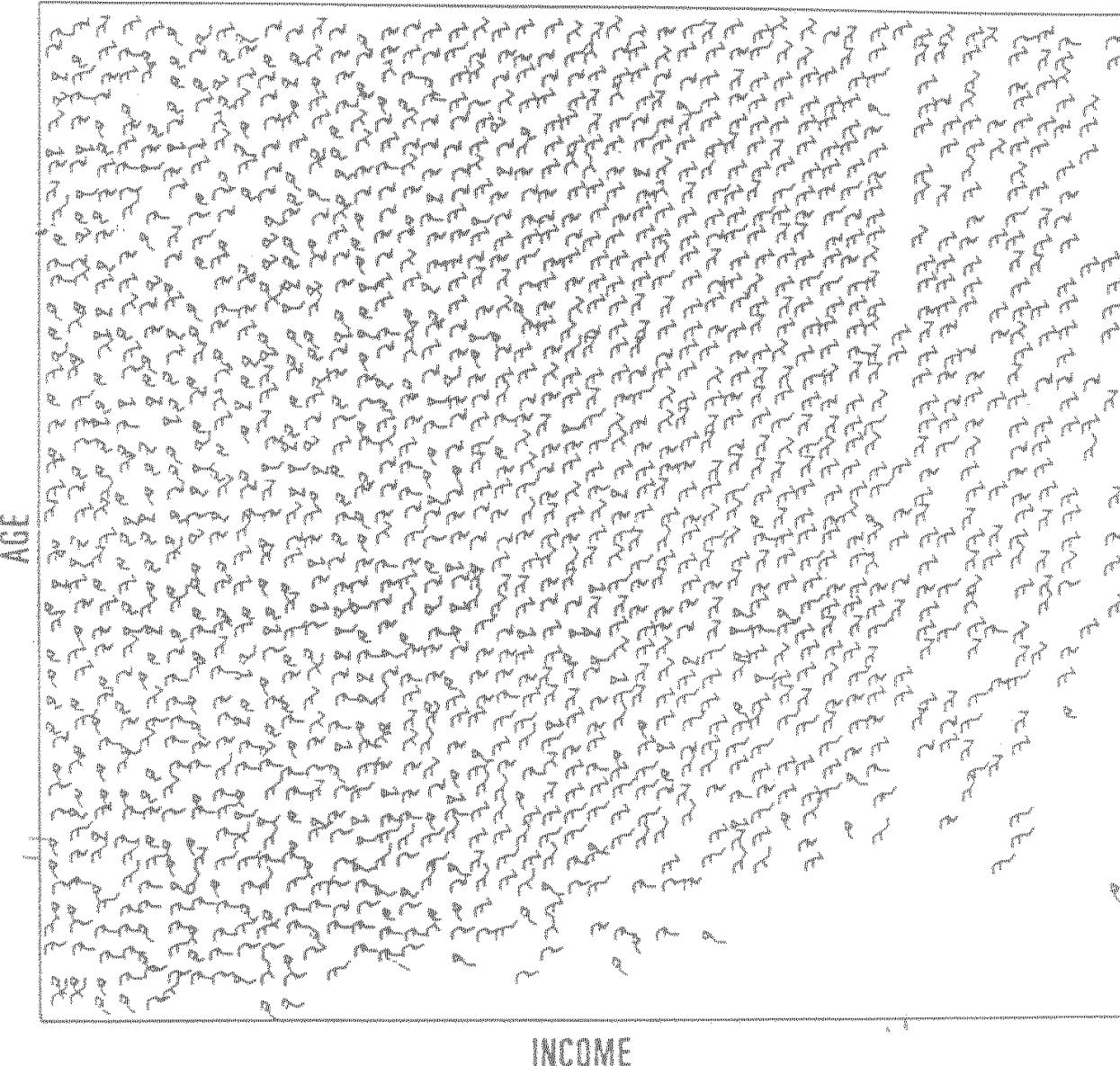
- A way to display variables on a two-dimensional surface, e.g., let x be eyebrow slant, y be eye size, z be nose length, etc.
- The figure shows faces produced using 10 characteristics--head eccentricity, eye size, eye spacing, eye eccentricity, pupil size, eyebrow slant, nose size, mouth shape, mouth size, and mouth opening): Each assigned one of 10 possible values, generated using *Mathematica* (S. Dickson)
- 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





Stick Figure

used by permission of G. Grinstein, University of Massachusetts at Lowell

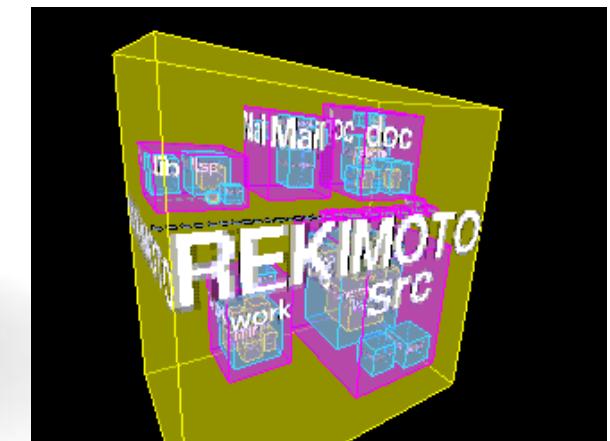
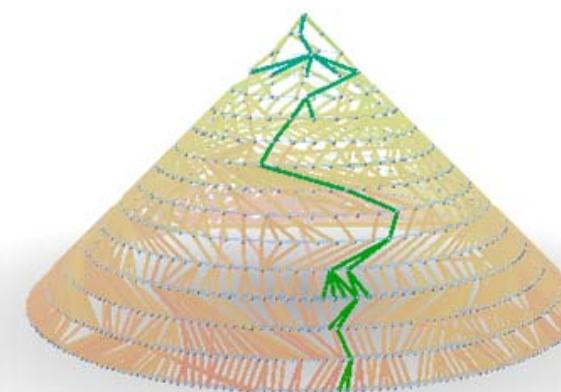
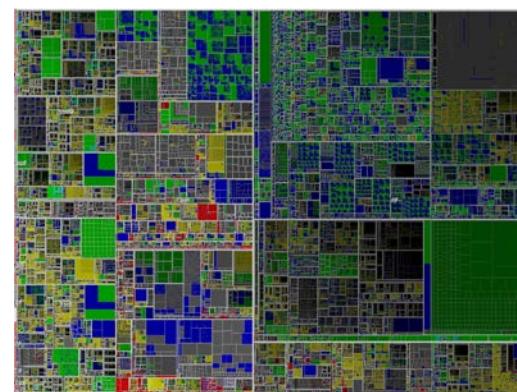
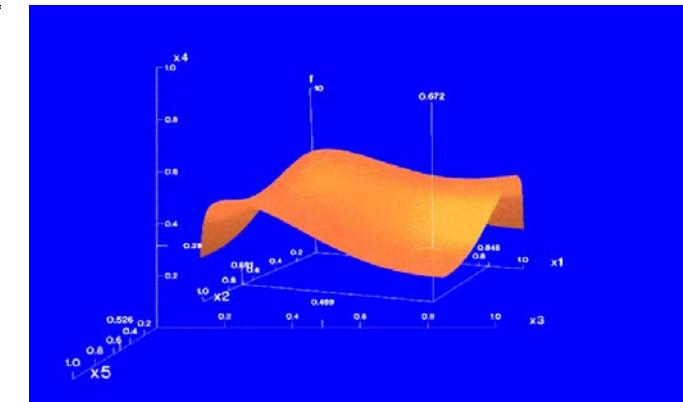
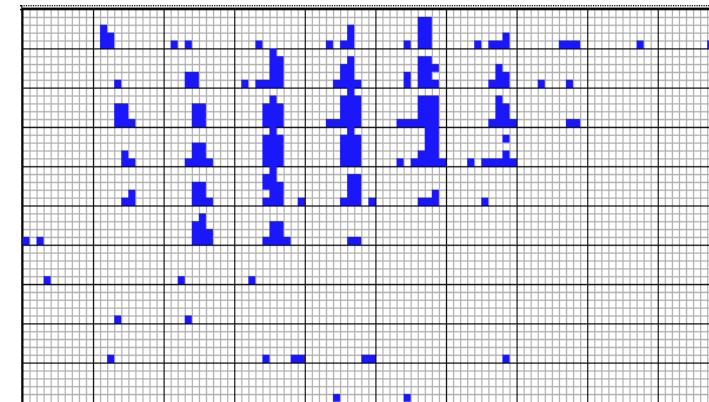


- A census data figure showing age, income, gender, education, etc.

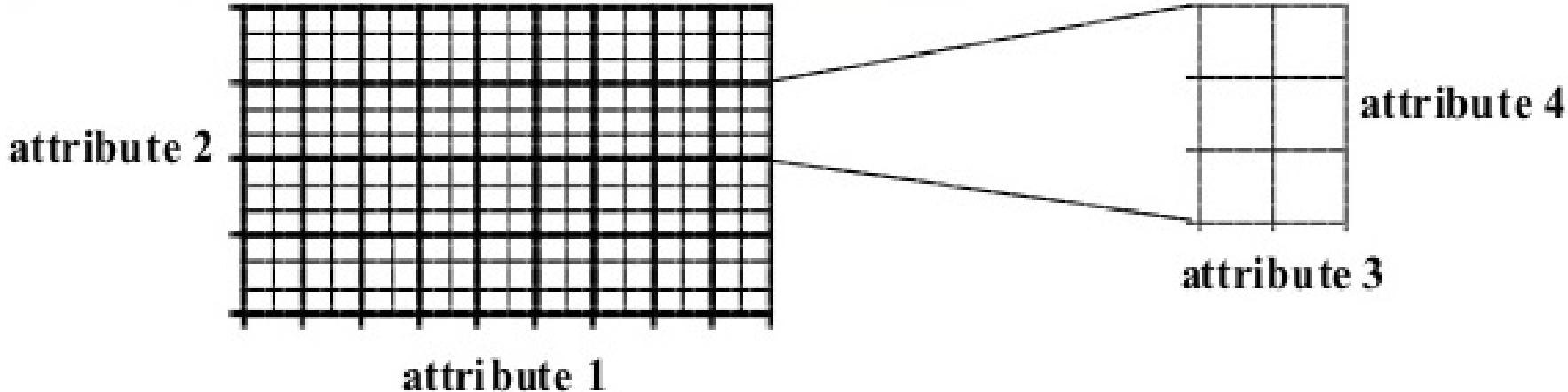
- A 5-piece stick figure (1 body and 4 limbs w. different angle/length)

Hierarchical Visualization Techniques

- Visualization of the data using a hierarchical partitioning into subspaces
- Methods
 - Dimensional Stacking
 - Worlds-within-Worlds
 - Tree-Map
 - Cone Trees
 - InfoCube



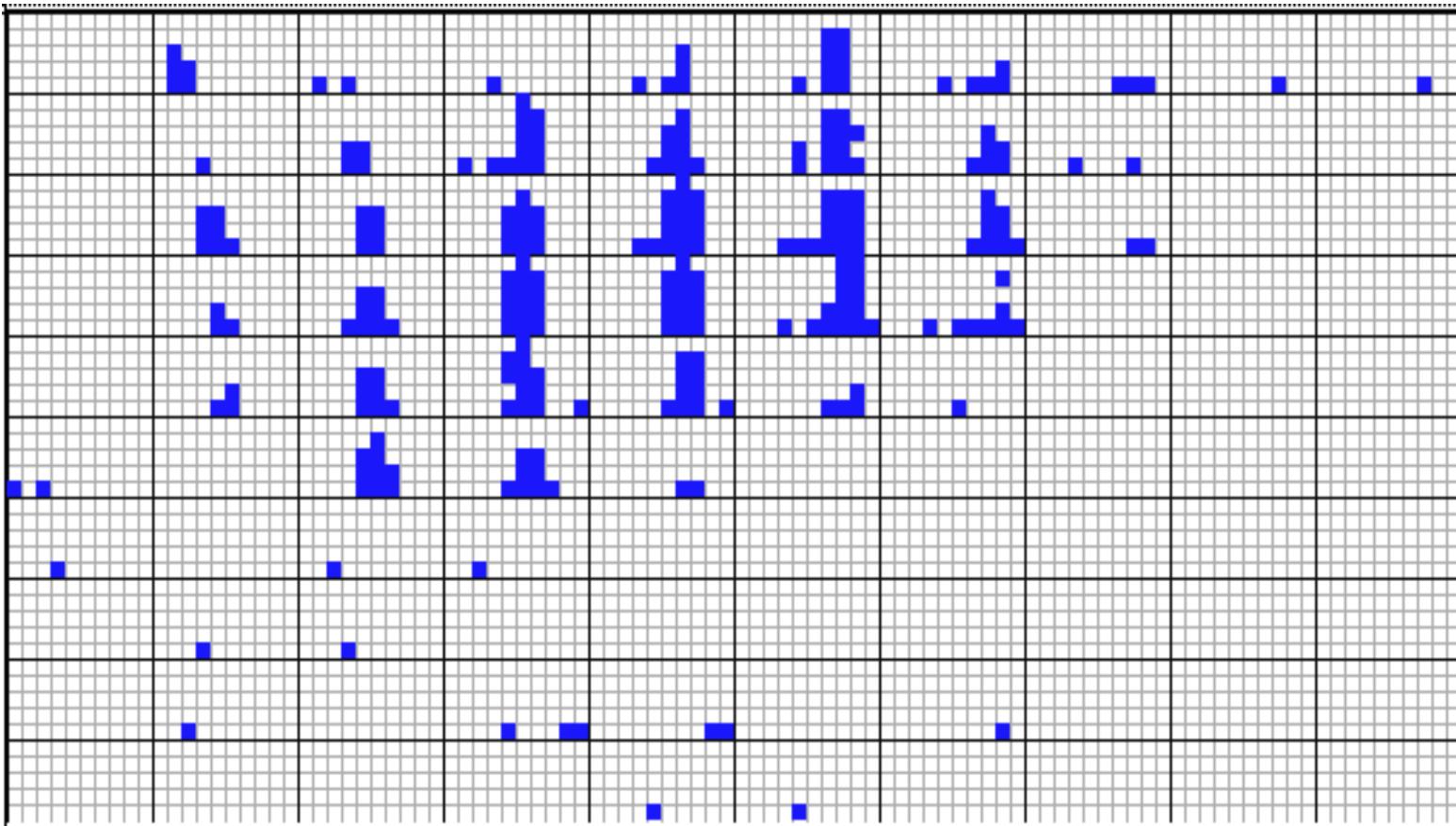
Dimensional Stacking



- 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

Dimensional Stacking

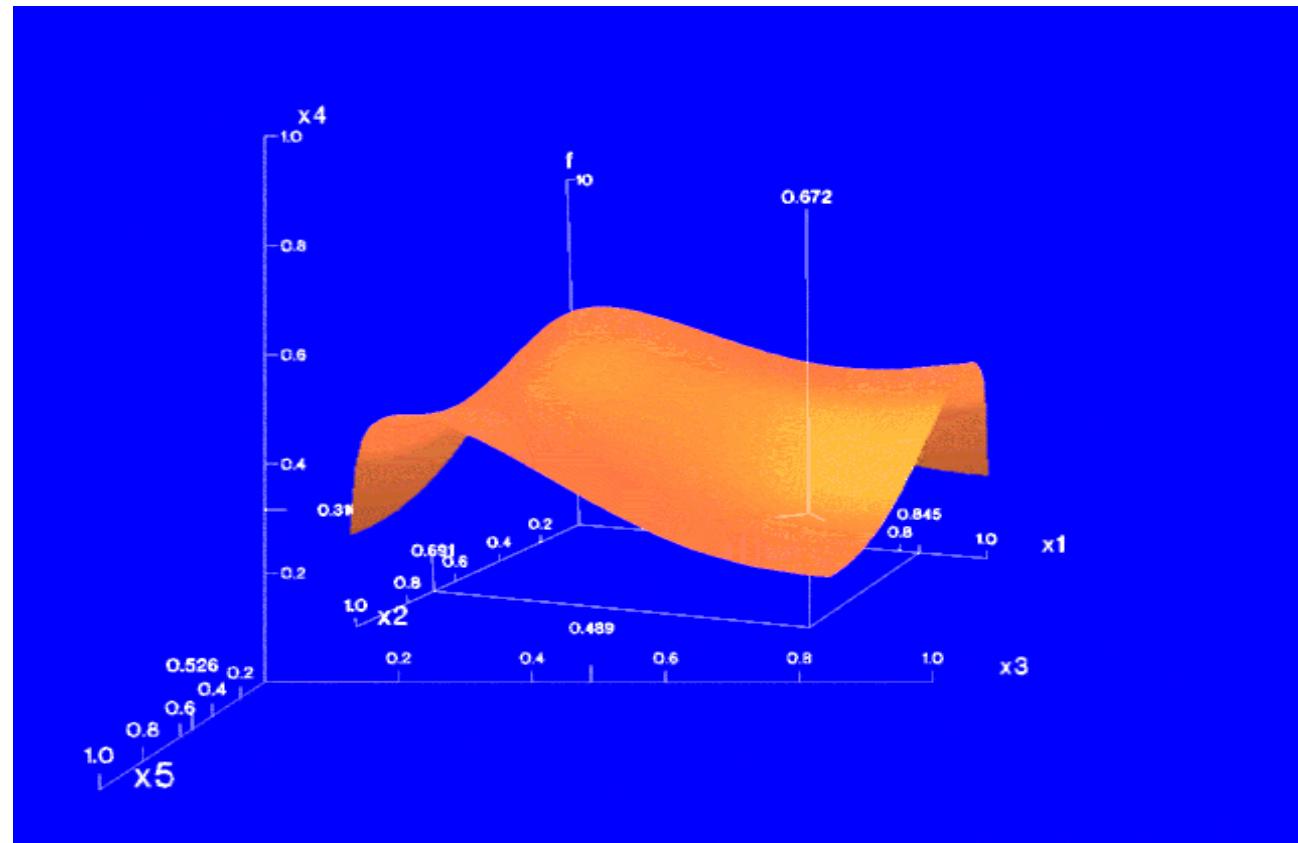
Used by permission of M. Ward, Worcester Polytechnic Institute



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

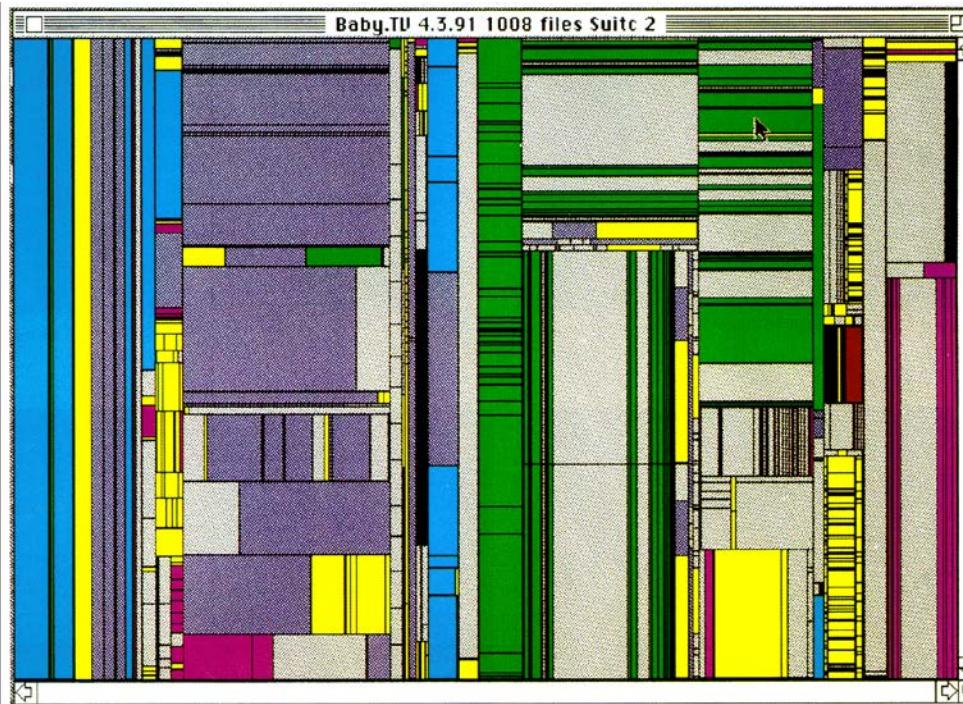
Worlds-within-Worlds

- 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)
- Software that uses this paradigm
- N-vision: Dynamic interaction through data glove and stereo displays, including rotation, scaling (inner) and translation (inner/outer)
- Auto Visual: Static interaction by means of queries

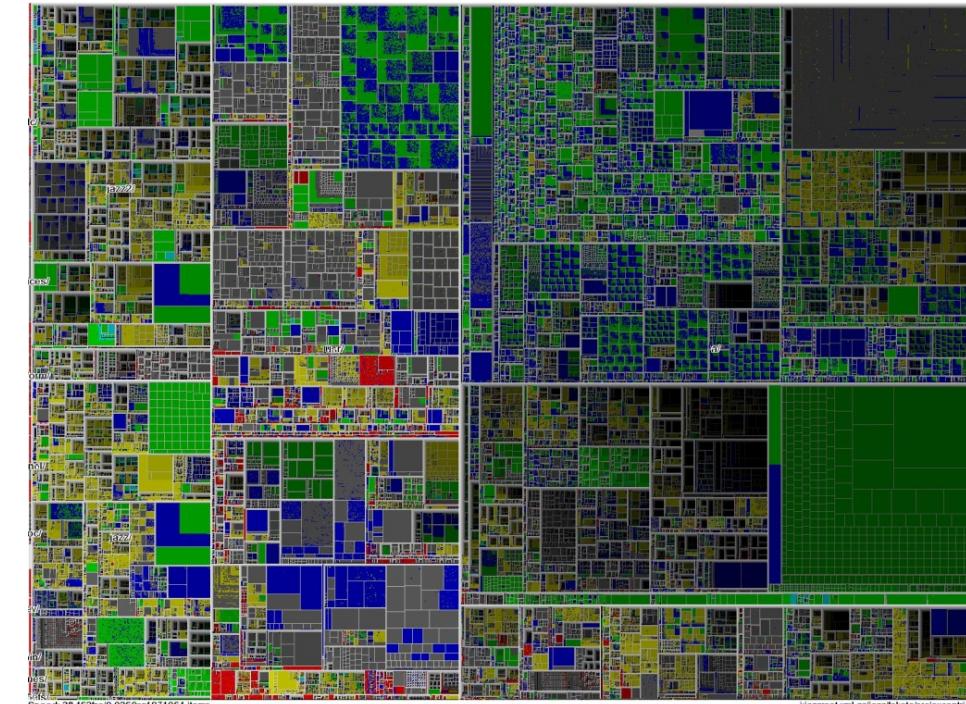


Tree-Map

- 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)



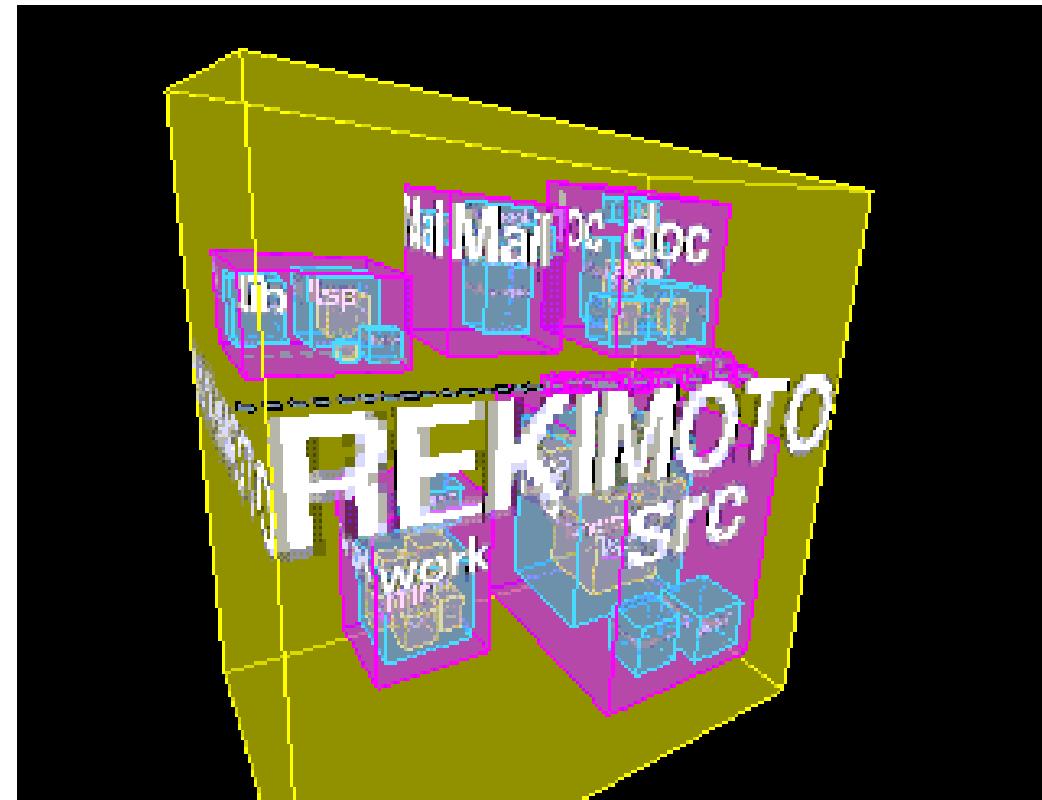
Schneiderman@UMD: Tree-Map of a File System



Schneiderman@UMD: Tree-Map to support large data sets of a million items

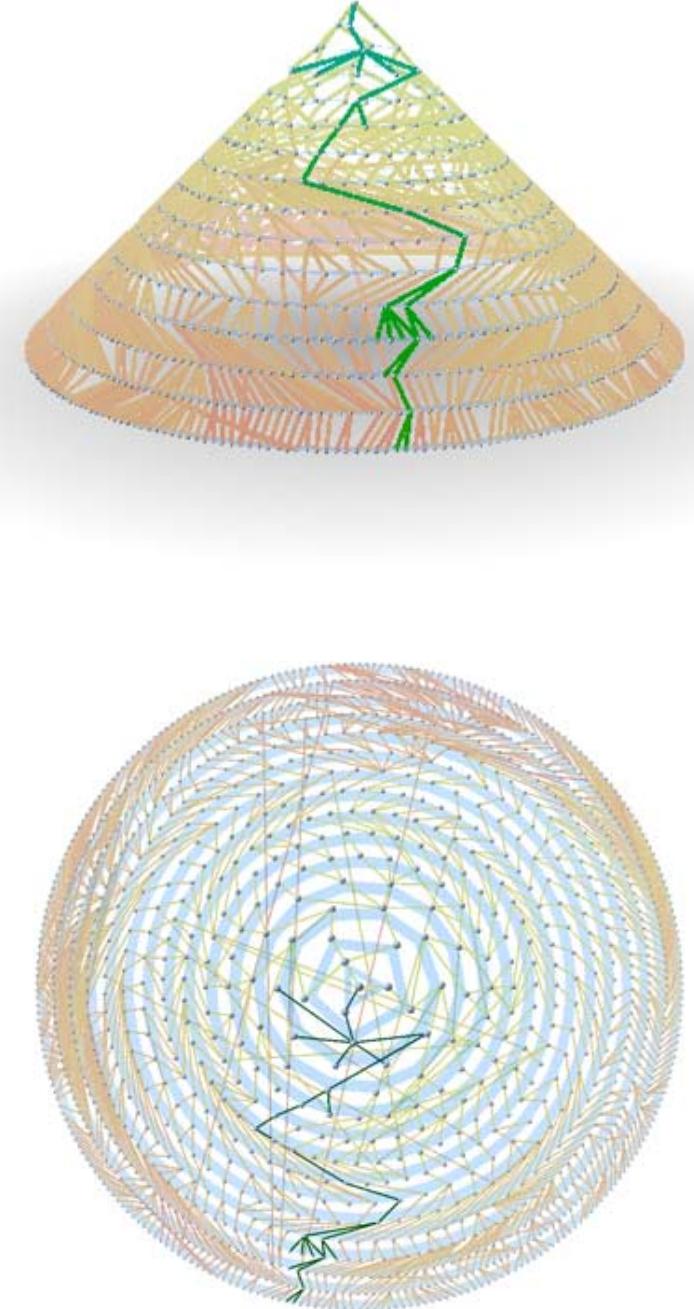
InfoCube

- 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 subnodes or the lower level data are represented as smaller cubes inside the outermost cubes, etc.



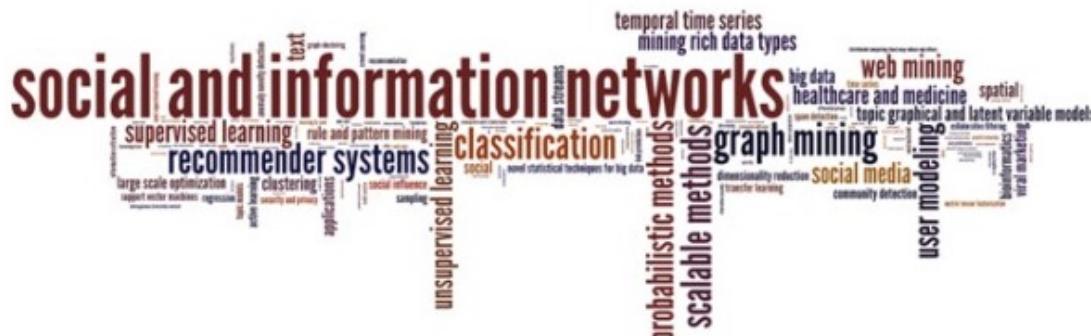
Three-D Cone Trees

- 3D cone tree visualization technique works well 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
- G. Robertson, J. Mackinlay, S. Card. “Cone Trees: Animated 3D Visualizations of Hierarchical Information”, ACM SIGCHI’91
- Graph from Nadeau Software Consulting website: Visualize a social network data set that models the way an infection spreads from one person to the next



Visualizing Complex Data and Relations: Tag Cloud

- **Tag cloud:** Visualizing user-generated tags
 - The importance of tag is represented by font size/color
 - Popularly used to visualize word/phrase distributions



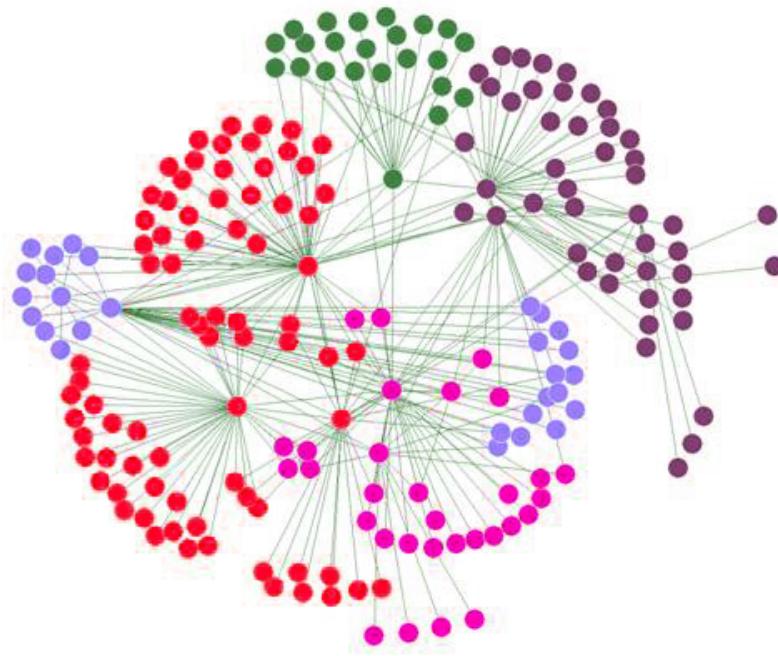
KDD 2013 Research Paper Title Tag Cloud



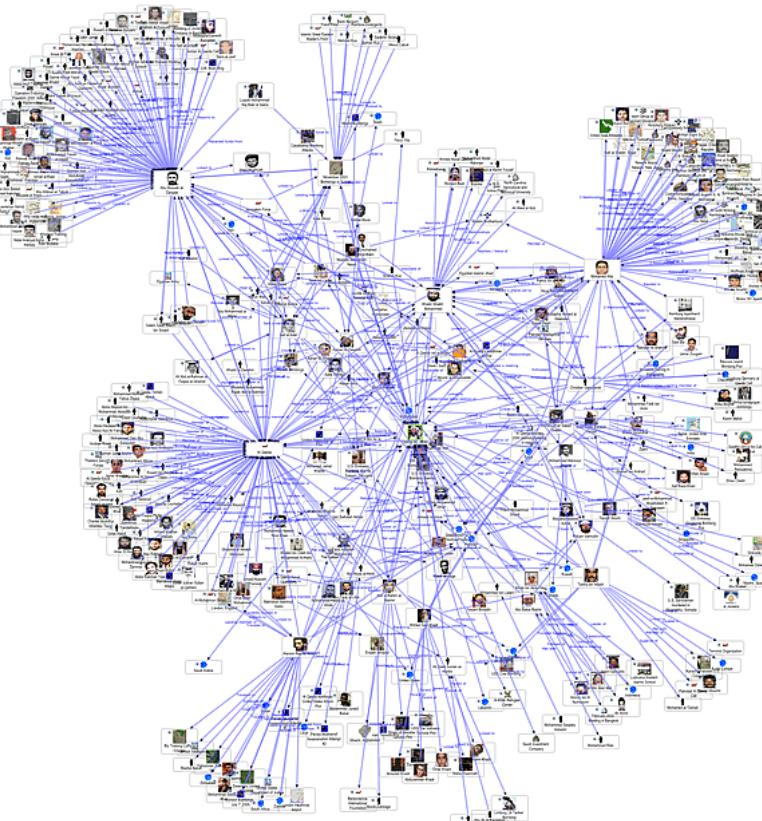
Newsmap: Google News Stories in 2005

Visualizing Complex Data and Relations: Social Networks

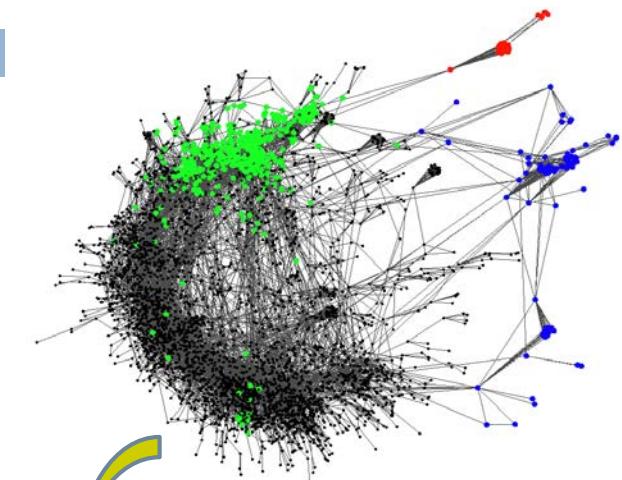
- Visualizing non-numerical data: **social and information networks**



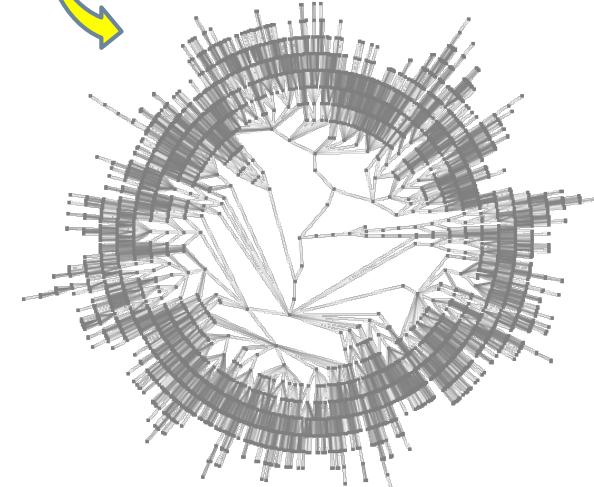
A typical network structure



A social network



organizing
information networks



Standardizing Numeric Data

- Z-score: $z = \frac{x - \mu}{\sigma}$
 - X: raw score to be standardized, μ : mean of the population, σ : standard deviation
 - the distance between the raw score and the population mean in units of the standard deviation
 - negative when the raw score is below the mean, positive when above
- An alternative way: Calculate the mean absolute deviation

$$s_f = \frac{1}{n}(|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|)$$

where

$$m_f = \frac{1}{n}(x_{1f} + x_{2f} + \dots + x_{nf}).$$

□ standardized measure (z-score):

$$z_{if} = \frac{x_{if} - m_f}{s_f}$$

- Using mean absolute deviation is more robust than using standard deviation