

Chapter 2. Data, Measurements, and Data Preprocessing

- ❑ **Data Types**
- ❑ **Statics of Data**
- ❑ **Similarity and Distance Measures**
- ❑ **Data Quality, Data Cleaning and Data Integration**
- ❑ **Data Transformation**
- ❑ **Dimensionality Reduction**
- ❑ **Summary**

Types of Data Sets: (1) Record Data

- Relational records
 - Relational tables, highly structured
- Data matrix, e.g., numerical matrix, crosstabs

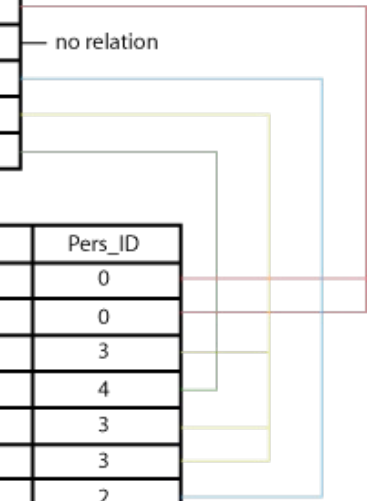
	China	England	France	Japan	USA	Total
Active Outdoors Crochet Glove		12.00	4.00	1.00	240.00	257.00
Active Outdoors Lycra Glove		10.00	6.00		323.00	339.00
InFlux Crochet Glove	3.00	6.00	8.00		132.00	149.00
InFlux Lycra Glove		2.00			143.00	145.00
Triumph Pro Helmet	3.00	1.00	7.00		333.00	344.00
Triumph 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

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



- Transaction data

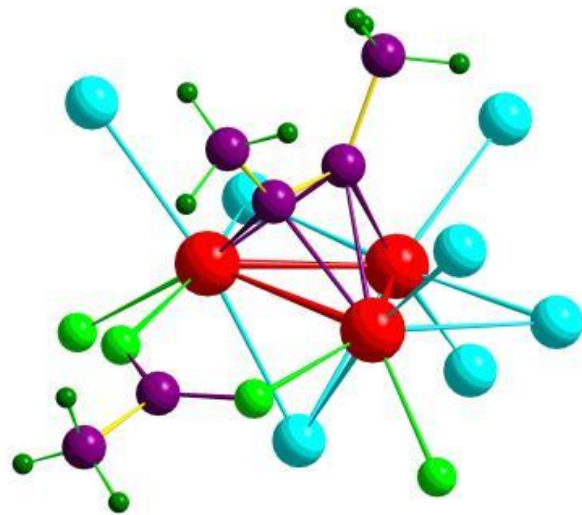
TID	Items
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

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

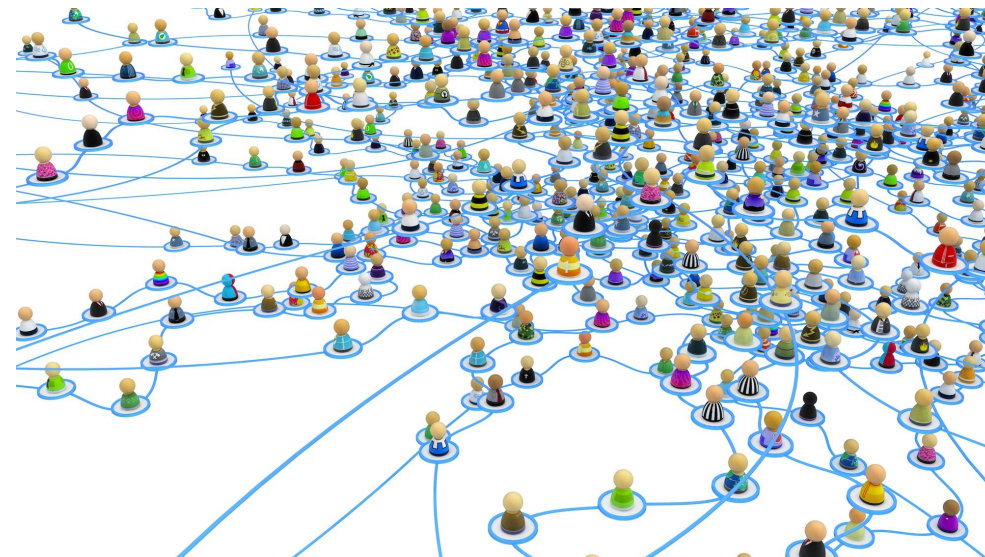
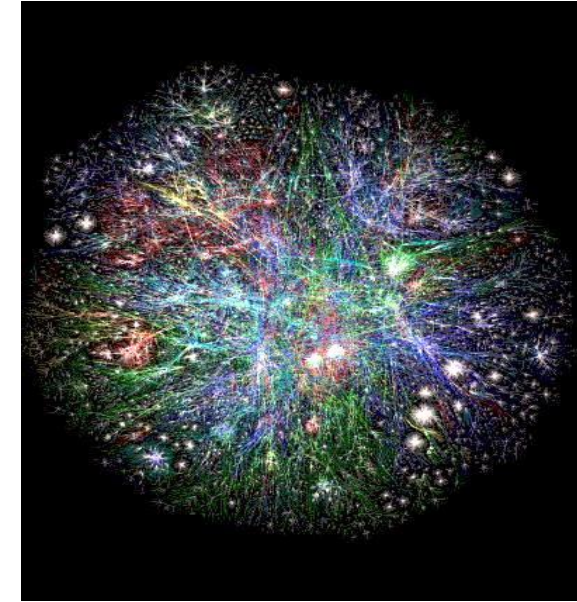
- Document data: Term-frequency vector (matrix) of text documents

Types of Data Sets: (2) Graphs and Networks

- ❑ Transportation network
- ❑ World Wide Web



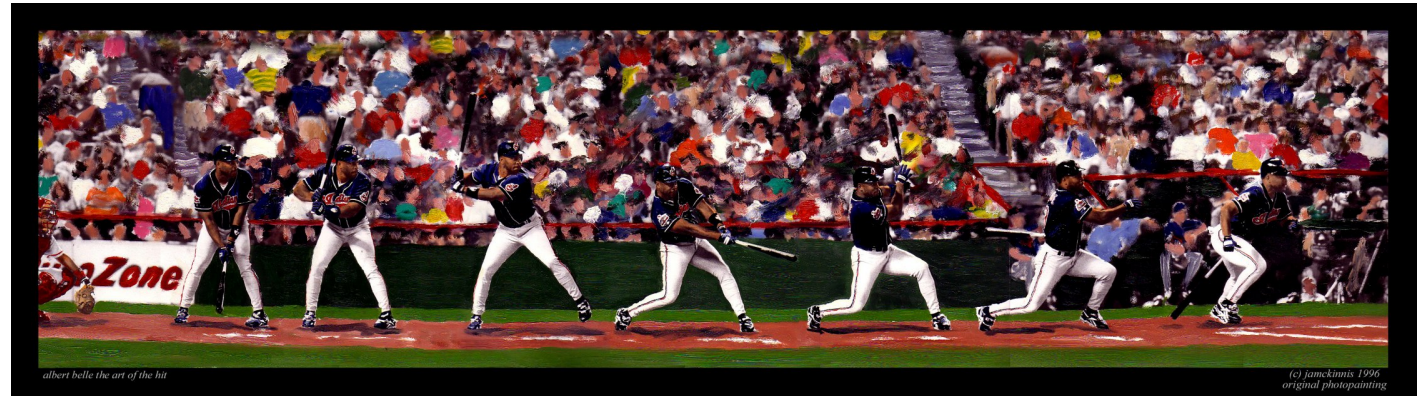
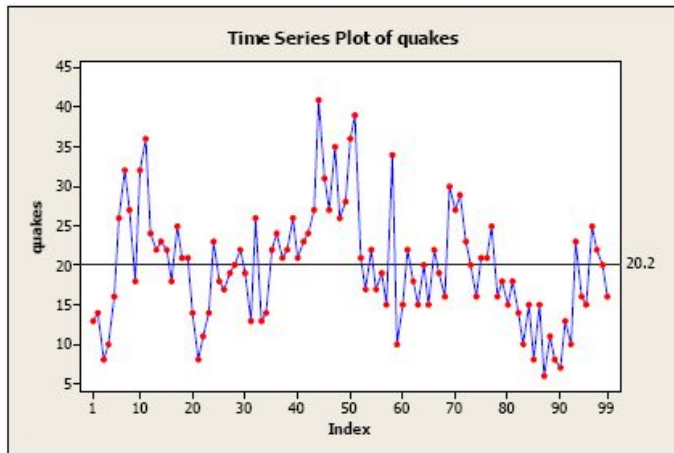
- ❑ Molecular Structures
- ❑ Social or information networks



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	GTTTGGAGG --- ATGTC AAC AAATGCTCCTTTCATTCCCTATTTACAGACC TGCCGCA
Chimpanzee	GTTTGGAGG --- ATGTC AAT AAATGCTGCTTTCACCTCCCTATTTACAGACC TGCCGCA
Macaque	GTTTGGAGG --- ATGC TCAAT AAATGCTCCTTTCATTCCCTATTTACA AACT TGCCGCA
Human	GACAATTCTGCTAGCAGCC TTTGTGCTATTATCTGTTTTCTAAACTTAGTAATTGAGTGT
Chimpanzee	GACAATTCTGCTAGCAGCC TTTGTGCTATTATCTGTTTTCTAAACTTAGTAATTGAGTGT
Macaque	GACAATTCTGCTAGCAGCC TTTGTGCTATTATCTGTTTTCTAAACTTAGTAATTGAGTGT
Human	GATCTGGAGACTAA - CTC TGAATAAATAAGCTGATTATTTATTTATTTCTCAAAACAA
Chimpanzee	GATCTGGAGACTAAACTCTG TGAATAAATAAGCTGATTATTTATTTATTTCTCAAAACAA
Macaque	TATCTGGAGACTAAACTCTG TGAATAAATAAGCTGATTATTTATTTATTTCTCAAAACAA
Human	CAGAATACGATTTAGCAAATTACTTCTTAAGATATATTTTACATTTCTATATTTCTCCTA
Chimpanzee	CAGAATACGATTTAGCAAATTACTTCTTAAGATATATTTTACATTTCTATATTTCTCCTA
Macaque	CAGAATA TGATTTAGCAAATTACTCTTAAAGATATATTTTGCAC TTCTATATTTCTCCTA
Human	CCCTGAGTTGATGTGTGAGCAATATGTCACCTTTCATAAAGCCAGGTATACA --- TTATG
Chimpanzee	CCCTGAGTTGATGTGTGAGCCG TATGTCACCTTTCATAAAGCCAGGTATACA --- TTATG
Macaque	CCCTGAGTTGATGTGTGAGCAATATGTCACCTTCCACA AAGCCAGGTATATATACATTACG
Human	GACAGGTAAGTAAAAAAC ATATTATTTATTCTACGTTTTGTCCAAAAATTTTAAATTTCT
Chimpanzee	GACAGGTAAGTAAAAAAC ATATTATTTATTCTACGTTTTGTCCAAAAATTTTAAATTTCT
Macaque	GACAGGTAAGTAAAAAAC - CATATTATTTATTCTAGBTTTTGTCCAAAGAG TTTTAAATTTCT
Human	AACTGTTGCGCGTGTGTTGGTAA --- TGTAAAACAAAC TCAGTACA
Chimpanzee	AACTGTTGCGCGTGTGTTGGTAA --- TGTAAAACAAAC TCAGTACA
Macaque	AACTGTTGTGCATGTGTTGGTAA --- CBTAAAACAAATTCAGTACG

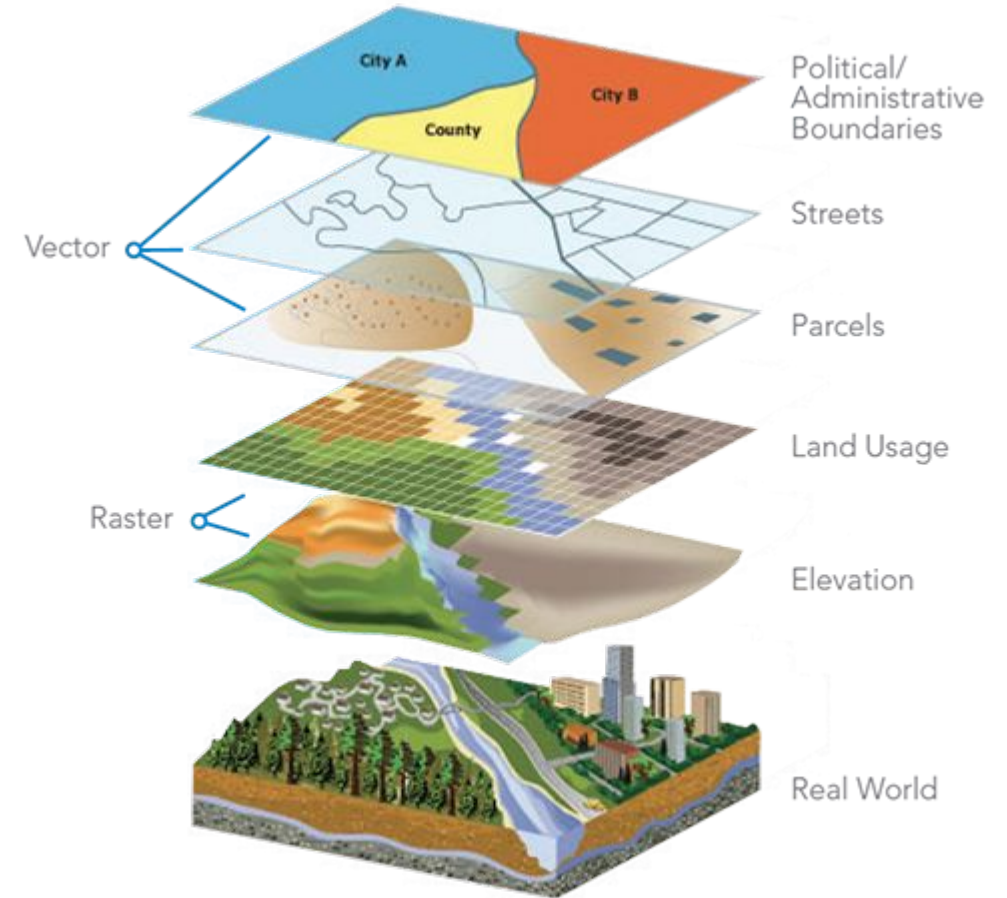
Types of Data Sets: (4) Spatial, image and multimedia Data

- Spatial data: maps



- Image data:

- Video data:



Important Characteristics of Structured Data

- Dimensionality
 - Curse of dimensionality
- Sparsity
 - Only presence counts
- Resolution
 - Patterns depend on the scale
- Distribution
 - Centrality and dispersion

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
 - ❑ Interval-scaled: 100°C is interval scales
 - ❑ Ratio-scaled: 100°K is ratio scaled since it is twice as high as 50°K
 - ❑ Discrete vs. Continuous Attributes

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**

- 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**

- 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*

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

□ 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

Statics of Data

- ❑ Measuring the Central Tendency
- ❑ Measuring the Dispersion of Data
- ❑ Covariance and Correlation Analysis
- ❑ Graphic Displays of Basic Statics of Data

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:

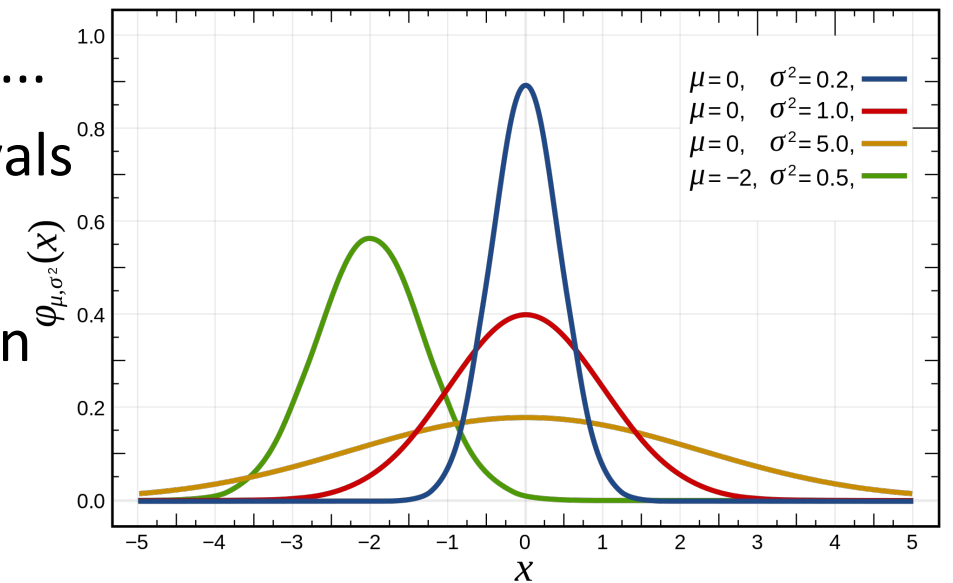
- Analyzed with multiple granularities of precision

- 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




Measuring the Central Tendency: (1) Mean

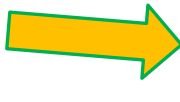
- Mean (algebraic measure) (sample vs. population):

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


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


$$\mu = \frac{\sum x}{N}$$

- 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)

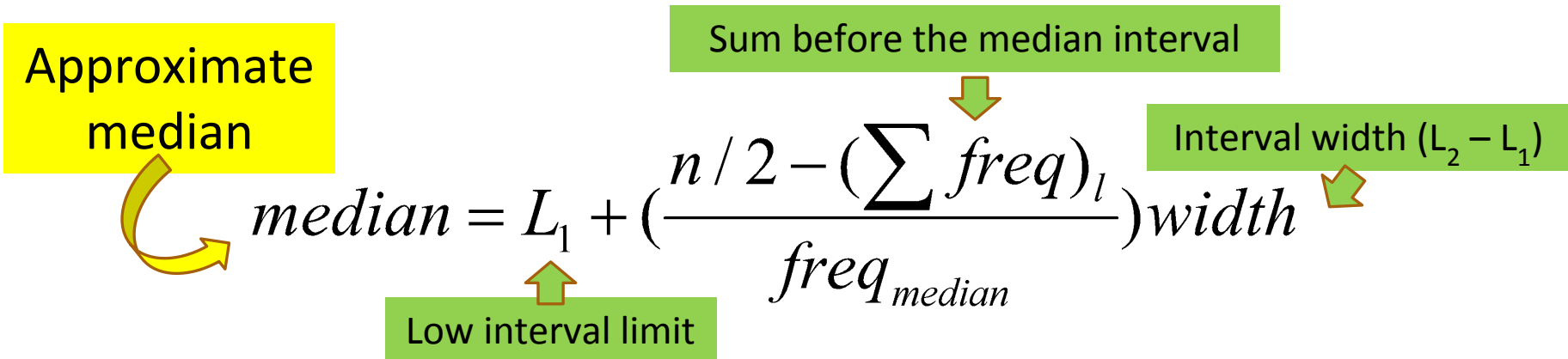
Measuring the Central Tendency: (2) Median

□ Median:

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

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

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



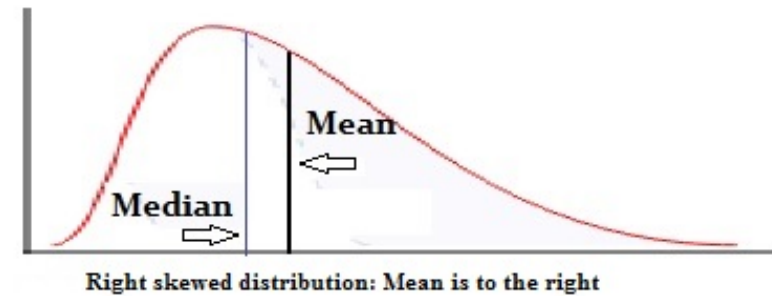
Measuring the Central Tendency: (3) Mode

□ Mode: Value that occurs most frequently in the data

□ Unimodal

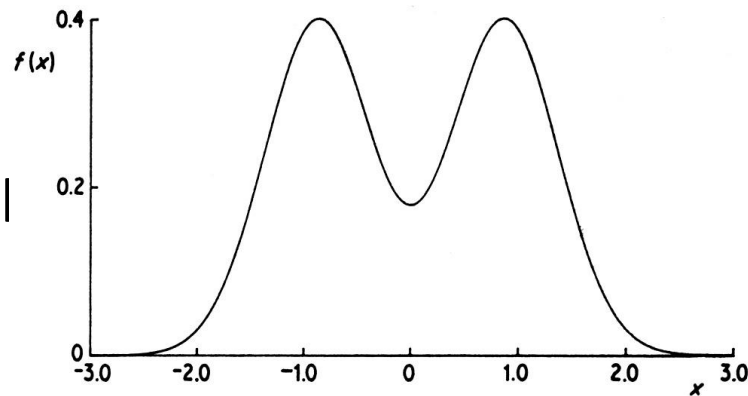
□ Empirical formula:

$$mean - mode = 3 \times (mean - median)$$

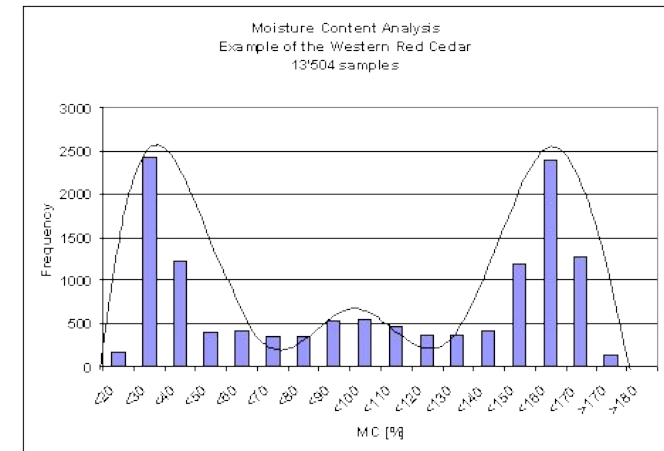
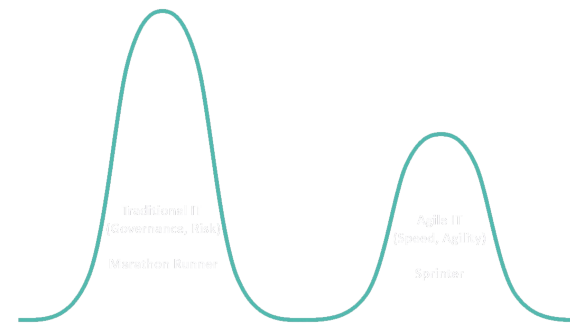


□ Multi-modal

□ Bimodal

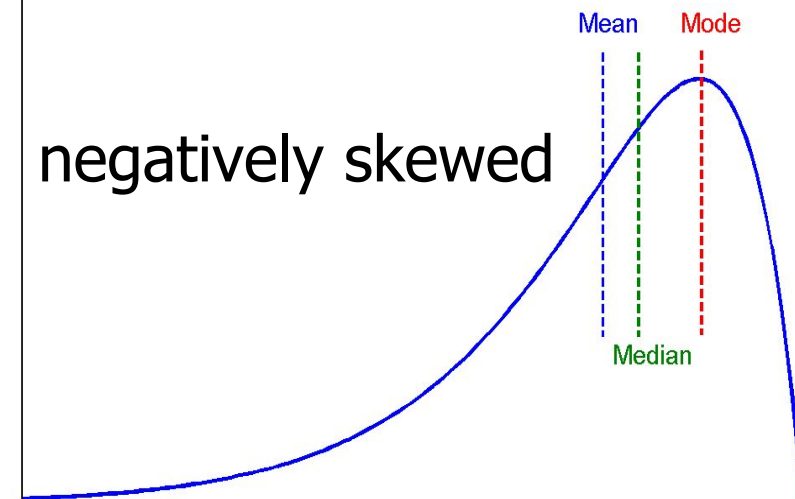
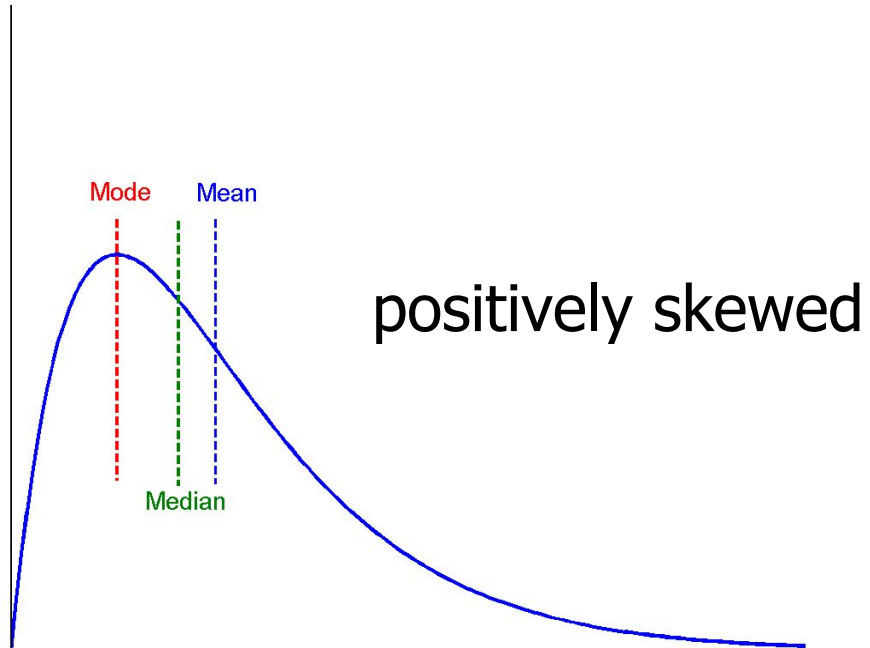
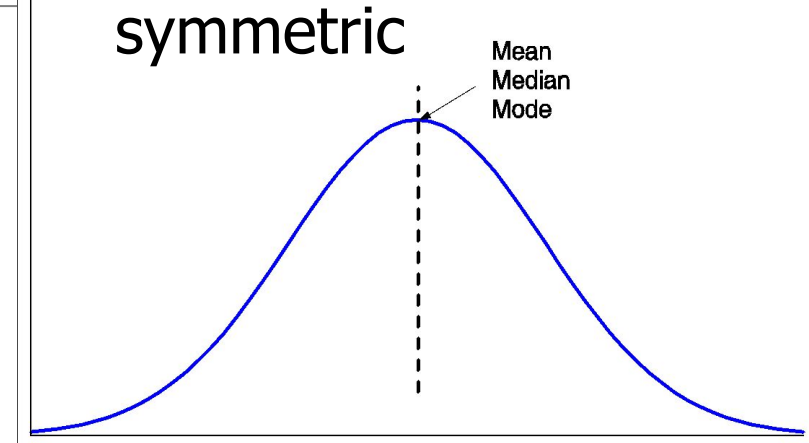


□ Trimodal



Symmetric vs. Skewed Data

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





Measures Data Distribution: Variance and Standard Deviation

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

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

- Q: Can you compute it incrementally and efficiently?


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


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

Note: The subtle difference of formulae for sample vs. population

- n : the size of the sample
- N : the size of the population

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

Correlation Analysis (for Categorical Data)

□ χ^2 (chi-square) test:

$$\chi^2 = \sum_i^n \frac{\overset{\text{observed}}{\downarrow} (O_i - E_i)^2}{\underset{\text{expected}}{E_i}}$$

- Null hypothesis: The two distributions are independent
- The cells that contribute the most to the χ^2 value are those whose actual count is very different from the expected count
 - The larger the χ^2 value, the more likely the variables are related
- Note: Correlation does not imply causality
 - # of hospitals and # of car-theft in a city are correlated
 - Both are causally linked to the third variable: population

Chi-Square Calculation: An Example

	Play chess	Not play chess	Sum (row)
Like science fiction	250 (X1)	200 (X2)	450
Not like science fiction	50 (X3)	1000 (X4)	1050
Sum(col.)	300	1200	1500

- ❑ Null hypothesis: The two distributions are independent
 - ❑ What does that mean?
 - ❑ The ratio between people who play chess vs not play chess is the same for both groups of like science fiction and not like science fiction
 - ❑ $X1:X2=X3:X4=300:1200$
 - ❑ $X1:X3=X2:X4=450:1050$
 - ❑ $X1+X2=450$ $X3+X4=1050$
 - ❑ $X1+X3=300$ $X2+X4=1200$

Chi-Square Calculation: An Example

	Play chess	Not play chess	Sum (row)
Like science fiction	250 (90)	200 (360)	450
Not like science fiction	50 (210)	1000 (840)	1050
Sum(col.)	300	1200	1500

How to derive 90?
 $450/1500 * 300 = 90$

- χ^2 (chi-square) calculation (numbers in parenthesis are expected counts calculated based on the data distribution in the two categories)

$$\chi^2 = \frac{(250 - 90)^2}{90} + \frac{(50 - 210)^2}{210} + \frac{(200 - 360)^2}{360} + \frac{(1000 - 840)^2}{840} = 507.93$$

We can reject the null hypothesis of independence at a confidence level of 0.001

- It shows that like_science_fiction and play_chess are correlated in the group

Correlation between Two Numerical Variables

- **Correlation** between two variables X_1 and X_2 is the standard covariance, obtained by normalizing the covariance with the standard deviation of each variable

$$\rho_{12} = \frac{\sigma_{12}}{\sigma_1\sigma_2} = \frac{\sigma_{12}}{\sqrt{\sigma_1^2\sigma_2^2}}$$

- **Sample correlation** for two attributes X_1 and X_2 :
$$\hat{\rho}_{12} = \frac{\hat{\sigma}_{12}}{\hat{\sigma}_1\hat{\sigma}_2} = \frac{\sum_{i=1}^n (x_{i1} - \hat{\mu}_1)(x_{i2} - \hat{\mu}_2)}{\sqrt{\sum_{i=1}^n (x_{i1} - \hat{\mu}_1)^2 \sum_{i=1}^n (x_{i2} - \hat{\mu}_2)^2}}$$

where n is the number of tuples, μ_1 and μ_2 are the respective means of X_1 and X_2 ,
 σ_1 and σ_2 are the respective standard deviation of X_1 and X_2

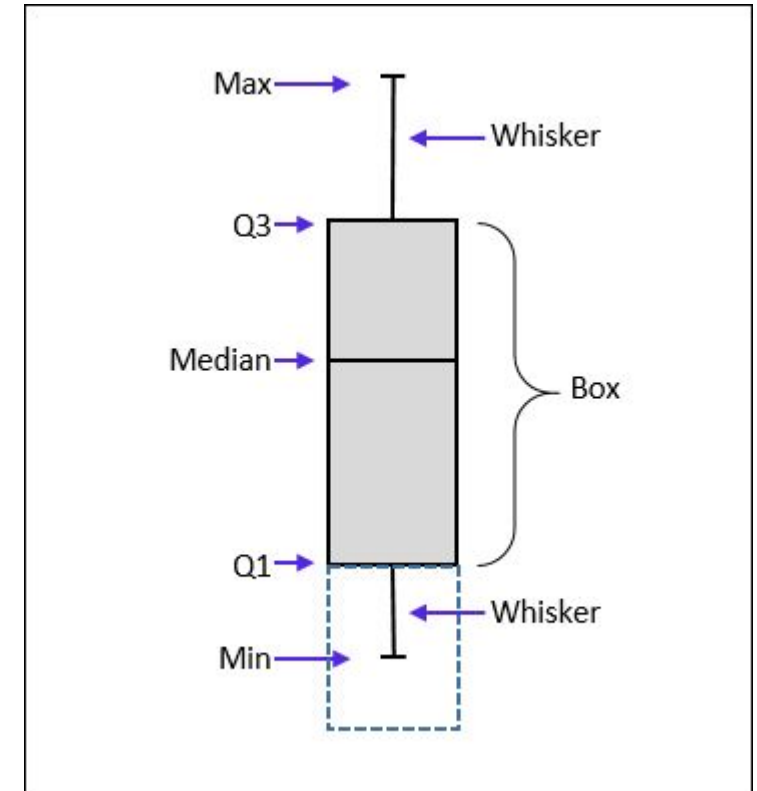
- If $\rho_{12} > 0$: A and B are positively correlated (X_1 's values increase as X_2 's)
 - The higher, the stronger correlation
- If $\rho_{12} = 0$: independent (under the same assumption as discussed in co-variance)
- If $\rho_{12} < 0$: negatively correlated

Graphic Displays of Basic Statistical Descriptions

- ❑ **Boxplot:** graphic display of five-number summary
- ❑ **Histogram:** x-axis are values, y-axis repres. 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

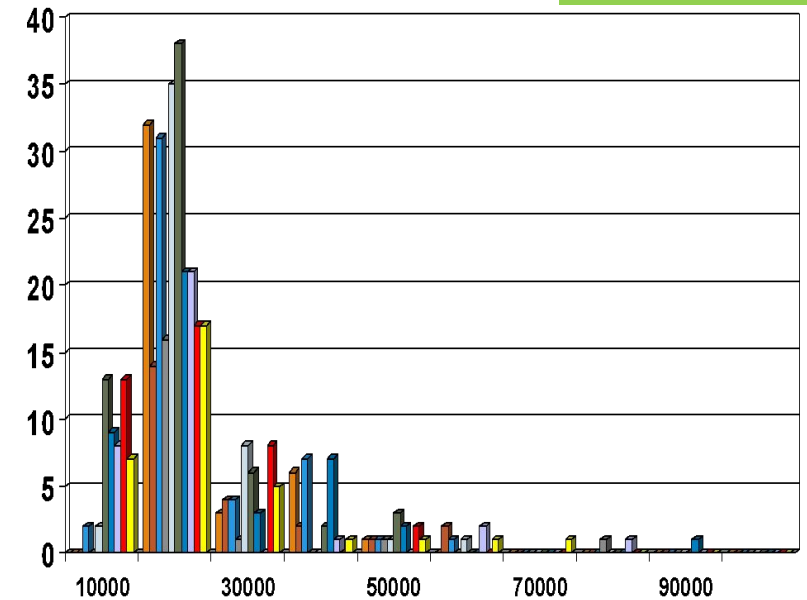
- ❑ **Quartiles:** Q_1 (25th percentile), Q_3 (75th percentile)
- ❑ **Inter-quartile range:** $IQR = Q_3 - Q_1$
- ❑ **Five number summary:** min, Q_1 , median, Q_3 , max
- ❑ **Boxplot:** Data is represented with a box
 - ❑ Q_1 , Q_3 , IQR: The ends of the box are at the first and third quartiles, i.e., 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 $1.5 \times IQR$



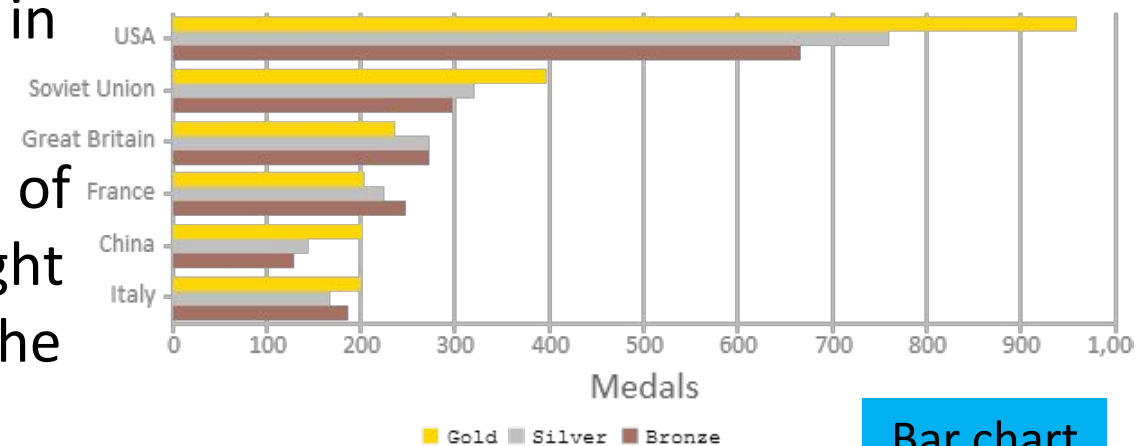
Histogram Analysis

- ❑ 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

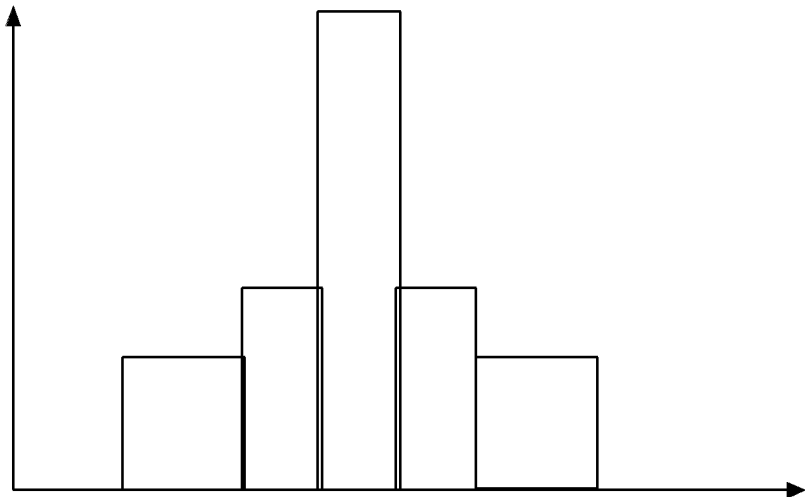
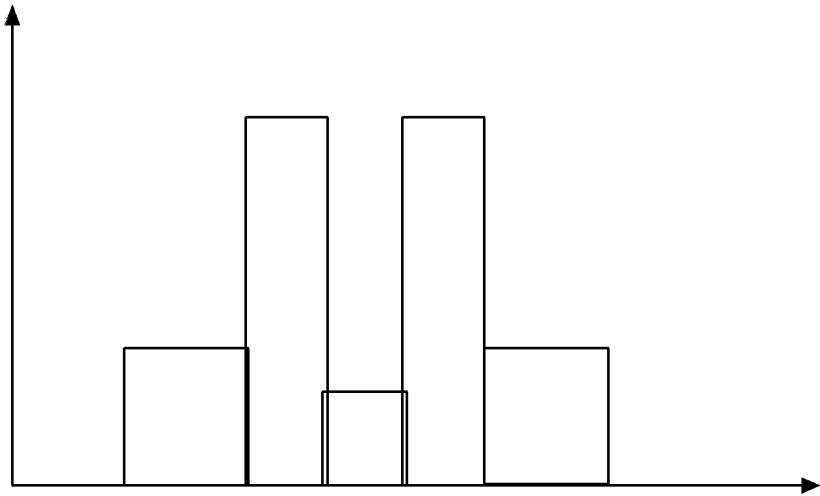


Olympic Medals of all Times (till 2012 Olympics)



Bar chart

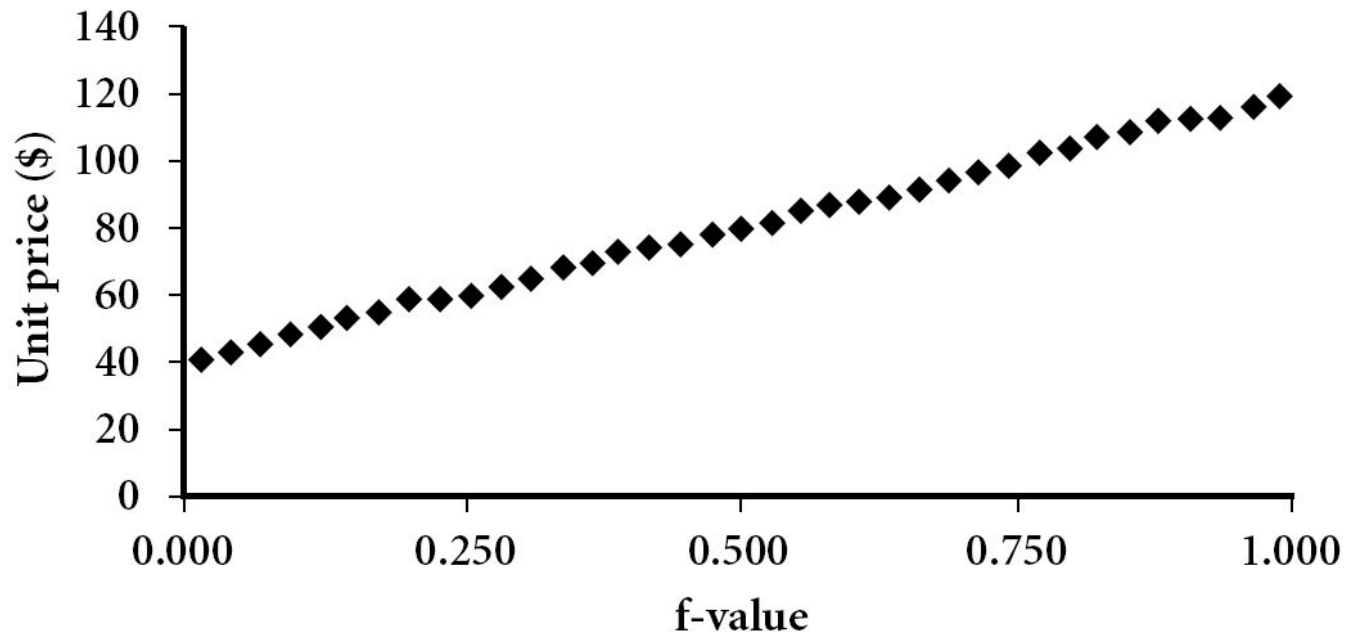
Histograms Often Tell More than Boxplots



- 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

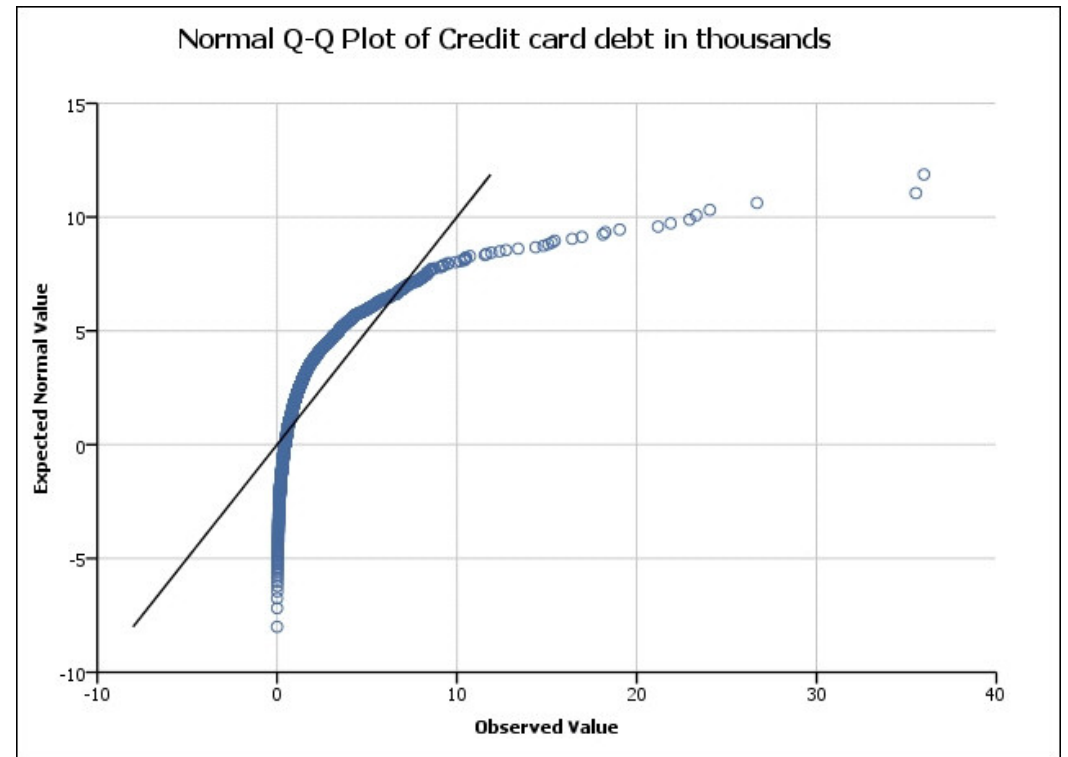
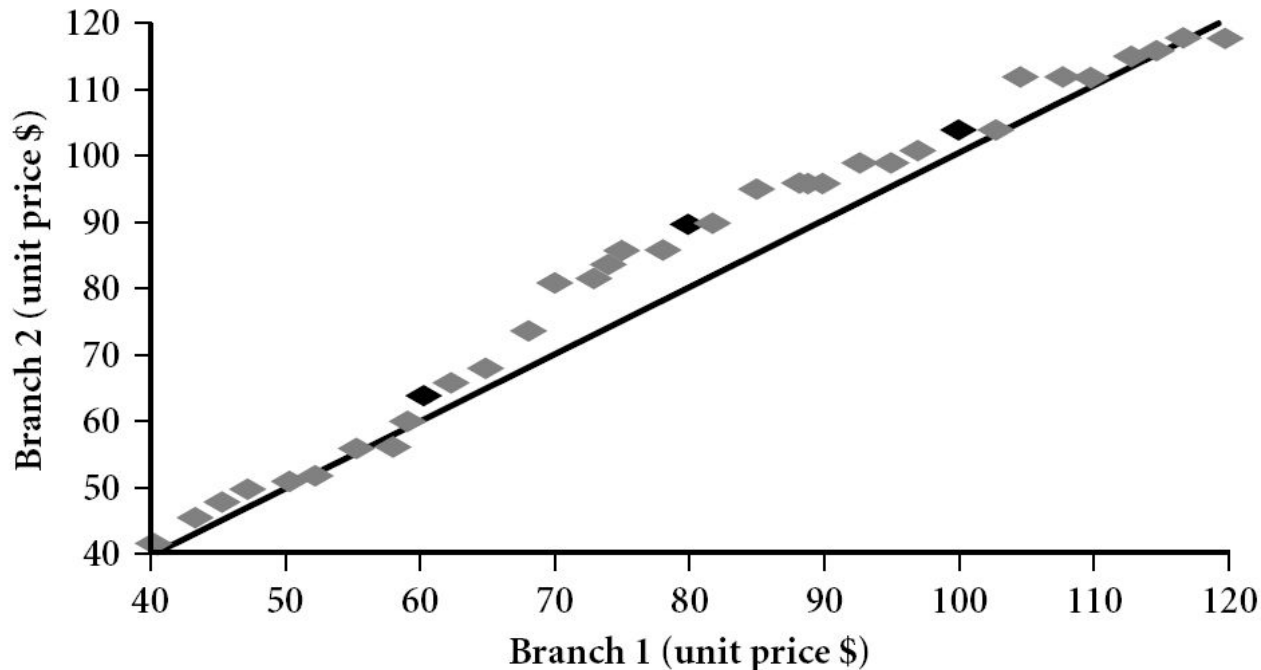
Quantile Plot

- 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 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



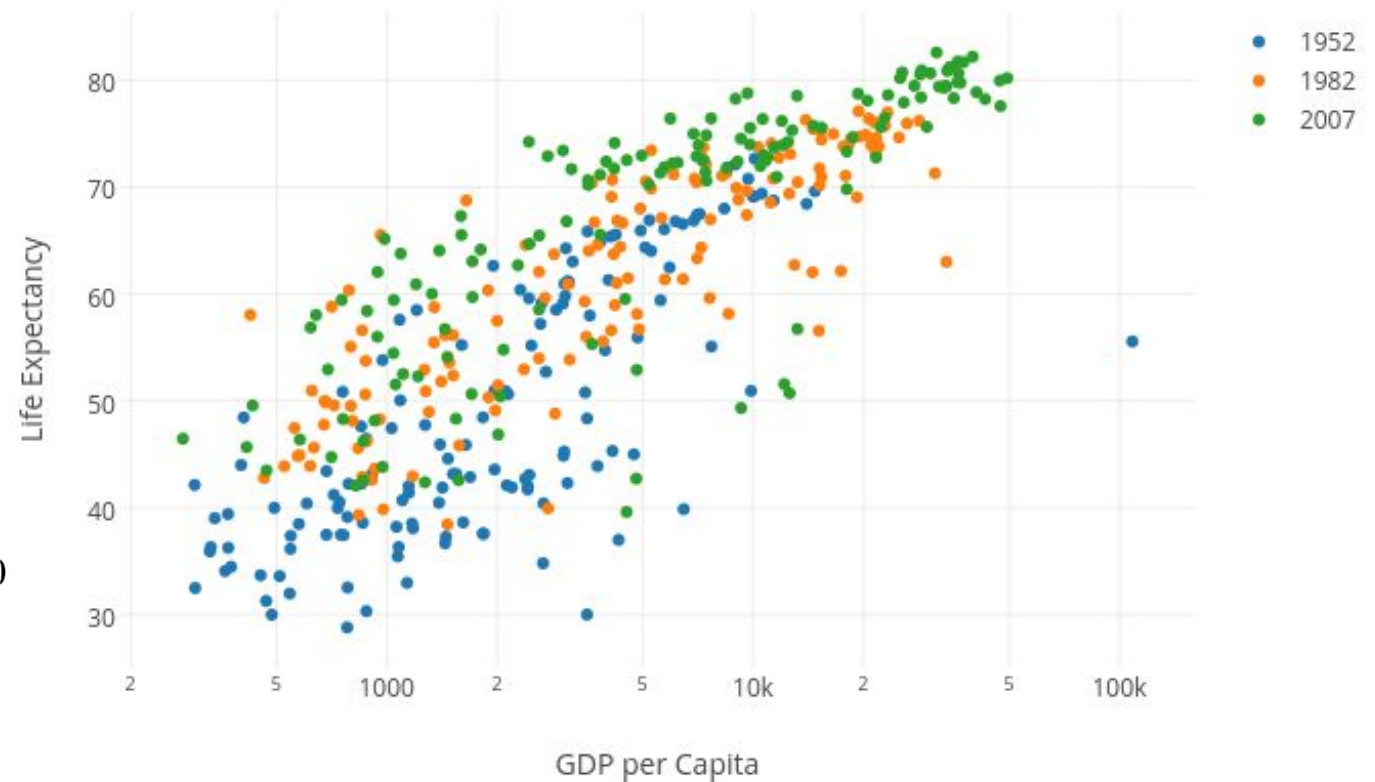
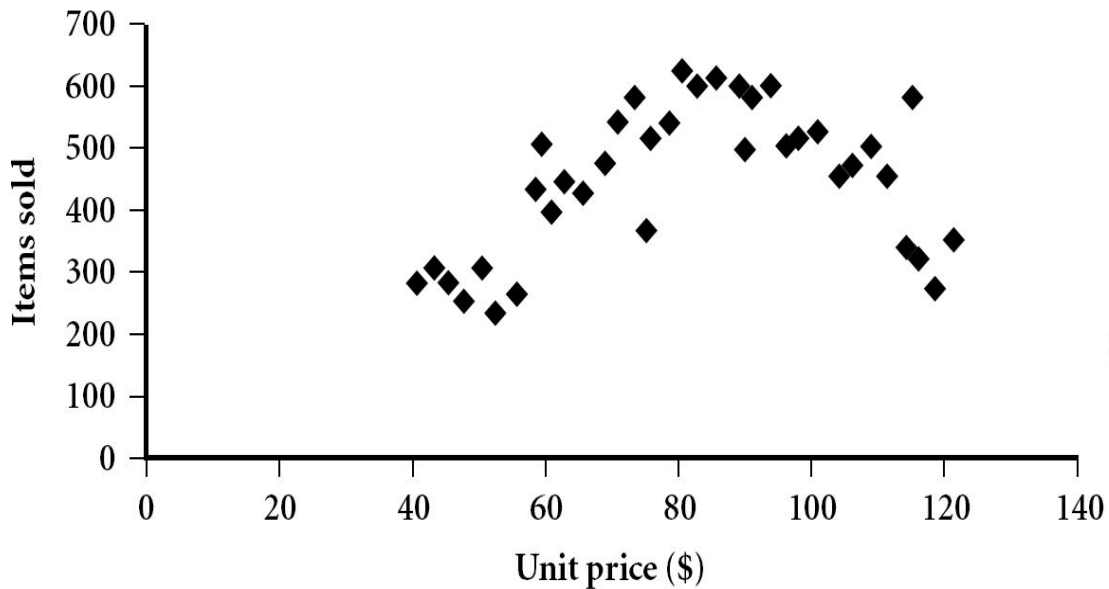
Quantile-Quantile (Q-Q) Plot

- Graphs the quantiles of one univariate distribution against the corresponding quantiles of another
- View: Is there is a shift in going from one distribution to another?
- Example shows unit price of items sold at Branch 1 vs. Branch 2 for each quantile. Unit prices of items sold at Branch 1 tend to be lower than those at Branch 2

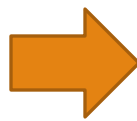
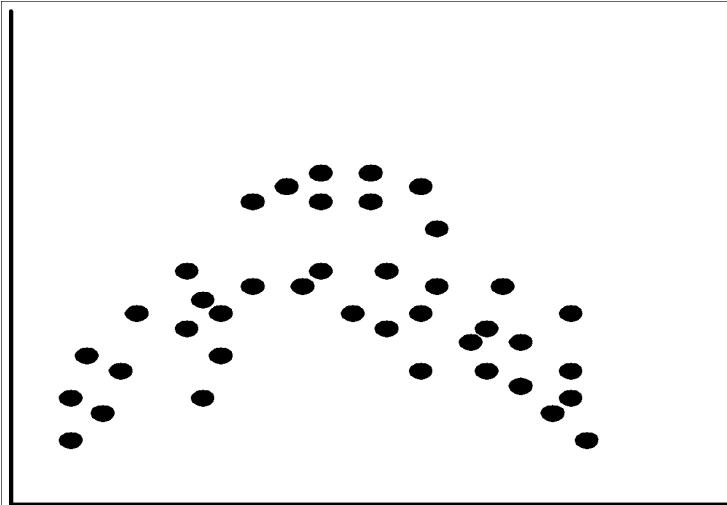
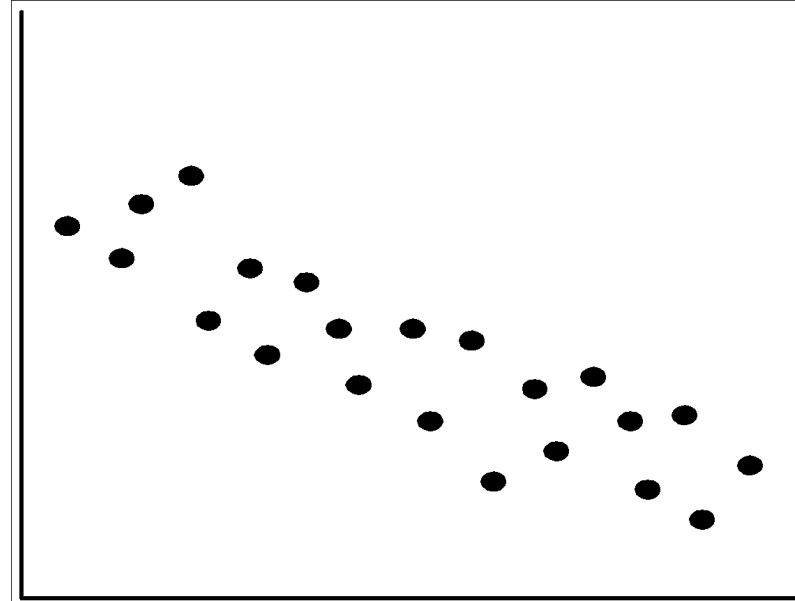
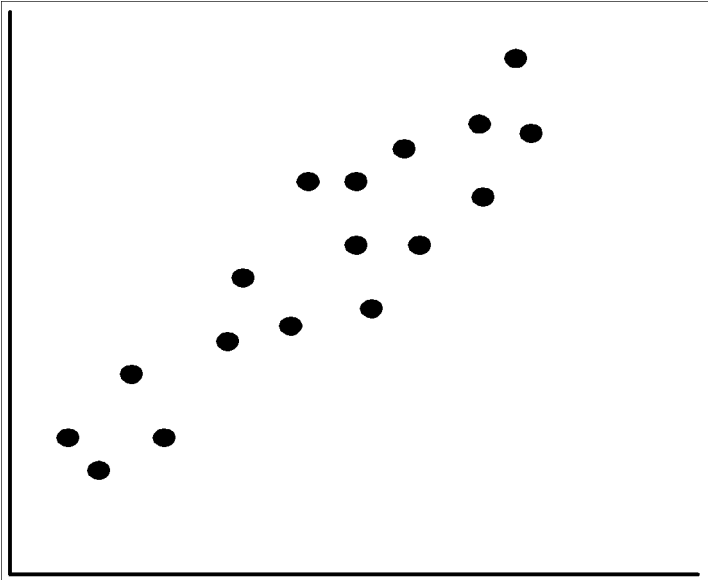


Scatter plot

- ❑ 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

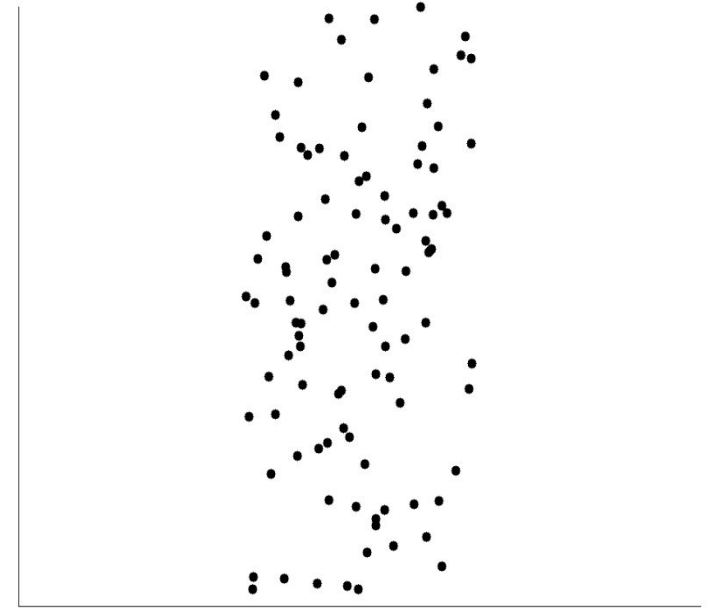
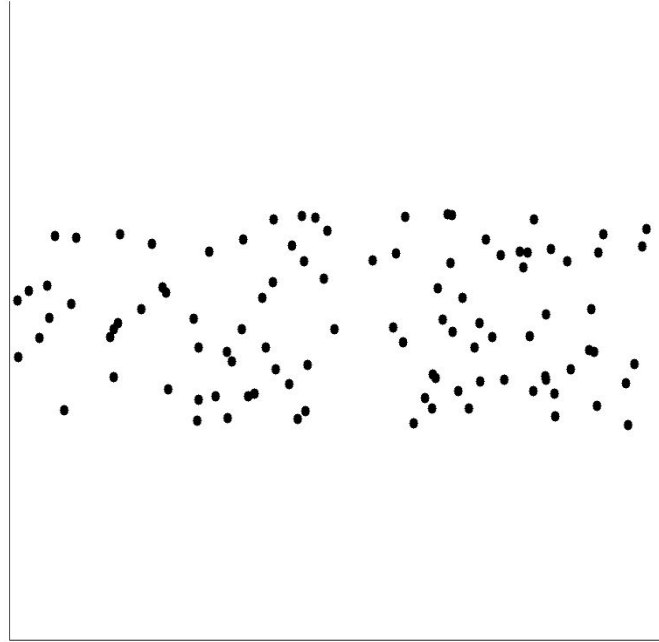
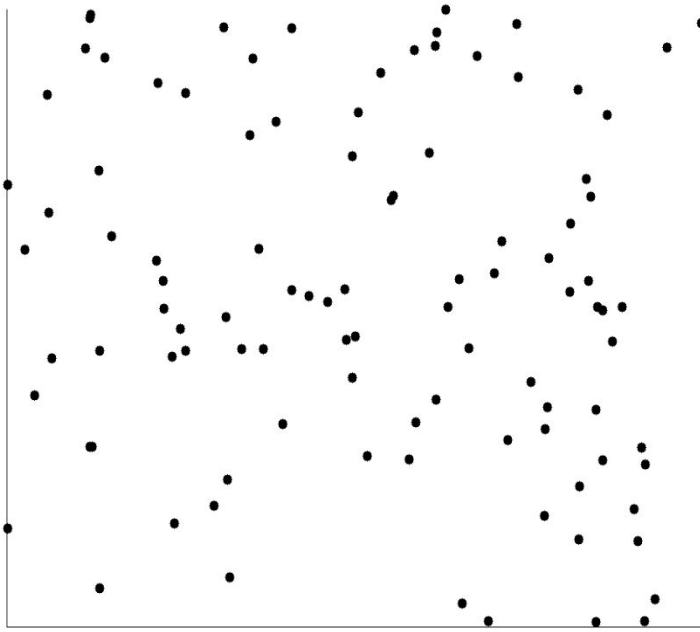


Positively and Negatively Correlated Data



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

Uncorrelated Data



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, “+” 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

Data Quality, Data Cleaning and Data Integration

- ❑ Data Quality Measures
- ❑ Data Cleaning
- ❑ Data Integration

What is Data Preprocessing? — Major Tasks

□ **Data cleaning**

- Handle missing data, smooth noisy data, identify or remove outliers, and resolve inconsistencies

□ **Data integration**

- Integration of multiple databases, data cubes, or files

□ **Data reduction**

- Dimensionality reduction
- Numerosity reduction
- Data compression

□ **Data transformation and data discretization**

- Normalization
- Concept hierarchy generation

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 Cleaning

- ❑ 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?

Incomplete (Missing) Data

- ❑ Data is not always available
 - ❑ E.g., many tuples have no recorded value for several attributes, such as customer income in sales data
- ❑ Missing data may be due to
 - ❑ Equipment malfunction
 - ❑ Inconsistent with other recorded data and thus deleted
 - ❑ Data were not entered due to misunderstanding
 - ❑ Certain data may not be considered important at the time of entry
 - ❑ Did not register history or changes of the data
- ❑ Missing data may need to be inferred

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**

Noisy Data

- ❑ **Noise:** random error or variance in a measured variable
- ❑ **Incorrect attribute values** may be due to
 - ❑ Faulty data collection instruments
 - ❑ Data entry problems
 - ❑ Data transmission problems
 - ❑ Technology limitation
 - ❑ Inconsistency in naming convention
- ❑ **Other data problems**
 - ❑ Duplicate records
 - ❑ Incomplete data
 - ❑ Inconsistent data

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 (e.g., Potter's Wheels)

Data Integration

- Data integration
 - Combining data from multiple sources into a coherent store
- Why data integration?
 - Help reduce/avoid noise
 - Get a more complete picture
 - Improve mining speed and quality
- **Schema integration:**
 - e.g., A.cust-id \equiv B.cust-#
 - Integrate metadata from different sources
- **Entity identification:**
 - Identify real world entities from multiple data sources, e.g., Bill Clinton = William Clinton

Handling Noise in Data Integration

- ❑ Detecting data value conflicts
 - ❑ For the same real world entity, attribute values from different sources are different
 - ❑ Possible reasons: no reason, different representations, different scales, e.g., metric vs. British units
- ❑ Resolving conflict information
 - ❑ Take the mean/median/mode/max/min
 - ❑ Take the most recent
 - ❑ Truth finding: consider the source quality
- ❑ Data cleaning + data integration

Handling Redundancy in Data Integration

- 📖 Redundant data occur often when integration of multiple databases
 - ❑ *Object identification*: The same attribute or object may have different names in different databases
 - ❑ *Derivable data*: One attribute may be a “derived” attribute in another table, e.g., annual revenue
- ❑ What’s the problem?
 - ❑ $Y = 2X \rightarrow Y = X_1 + X_2 \quad Y = 3X_1 - X_2 \quad Y = -1291X_1 + 1293X_2$
- ❑ Redundant attributes may be detected by correlation analysis and covariance analysis

Data Transformation

- Normalization
- Discretization
- Data Compression
- Sampling

Data Transformation

- ❑ A function that maps the entire set of values of a given attribute to a new set of replacement values s.t. each old value can be identified with one of the new values
- ❑ Methods
 - ❑ Smoothing: Remove noise from data
 - ❑ Attribute/feature construction
 - ❑ New attributes constructed from the given ones
 - ❑ Aggregation: Summarization, data cube construction
 - ❑ Normalization: Scaled to fall within a smaller, specified range
 - ❑ min-max normalization
 - ❑ z-score normalization
 - ❑ normalization by decimal scaling
 - ❑ Discretization: Concept hierarchy climbing

Normalization

- **Min-max normalization:** to $[\text{new_min}_A, \text{new_max}_A]$

$$v' = \frac{v - \text{min}_A}{\text{max}_A - \text{min}_A} (\text{new_max}_A - \text{new_min}_A) + \text{new_min}_A$$

- Ex. Let income range \$12,000 to \$98,000 normalized to [0.0, 1.0]

$$\frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$$

- Then \$73,000 is mapped to

- **Z-score normalization** (μ : mean, σ : s

$$v' = \frac{v - \mu_A}{\sigma_A}$$

Z-score: The distance between the raw score and the population mean in the unit of the standard deviation

$$\frac{73,600 - 54,000}{16,000} = 1.225$$

- Ex. Let $\mu = 54,000$, $\sigma = 16,000$. Then

- **Normalization by decimal scaling:** $v' = \frac{v}{10^j}$ where j is the smallest integer such that $\text{Max}(|v'|) < 1$

Discretization

- ❑ Three types of attributes
 - ❑ Nominal—values from an unordered set, e.g., color, profession
 - ❑ Ordinal—values from an ordered set, e.g., military or academic rank
 - ❑ Numeric—real numbers, e.g., integer or real numbers
- ❑ Discretization: Divide the range of a continuous attribute into intervals
 - ❑ Interval labels can then be used to replace actual data values
 - ❑ Reduce data size by discretization
 - ❑ Supervised vs. unsupervised
 - ❑ Split (top-down) vs. merge (bottom-up)
 - ❑ Discretization can be performed recursively on an attribute
 - ❑ Prepare for further analysis, e.g., classification

Data Discretization Methods

- Binning
 - Top-down split, unsupervised
- Histogram analysis
 - Top-down split, unsupervised
- Clustering analysis
 - Unsupervised, top-down split or bottom-up merge
- Decision-tree analysis
 - Supervised, top-down split
- Correlation (e.g., χ^2) analysis
 - Unsupervised, bottom-up merge
- Note: All the methods can be applied recursively

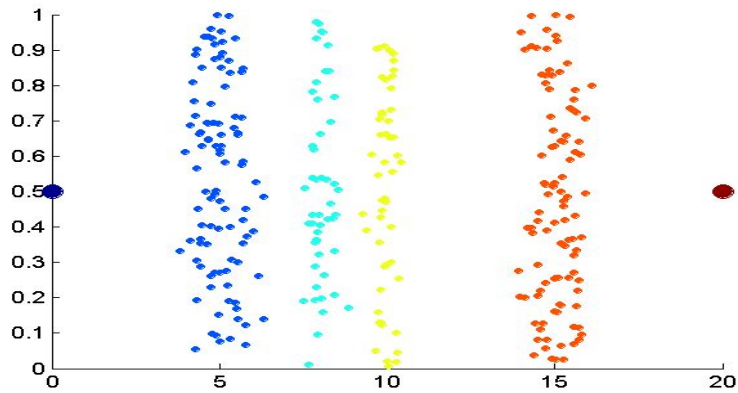
Simple Discretization: Binning

- ❑ **Equal-width** (distance) partitioning
 - ❑ Divides the range into N intervals of equal size: uniform grid
 - ❑ if A and B are the lowest and highest values of the attribute, the width of intervals will be: $W = (B - A)/N$.
 - ❑ The most straightforward, but outliers may dominate presentation
 - ❑ Skewed data is not handled well
- ❑ **Equal-depth** (frequency) partitioning
 - ❑ Divides the range into N intervals, each containing approximately same number of samples
 - ❑ Good data scaling
 - ❑ Managing categorical attributes can be tricky

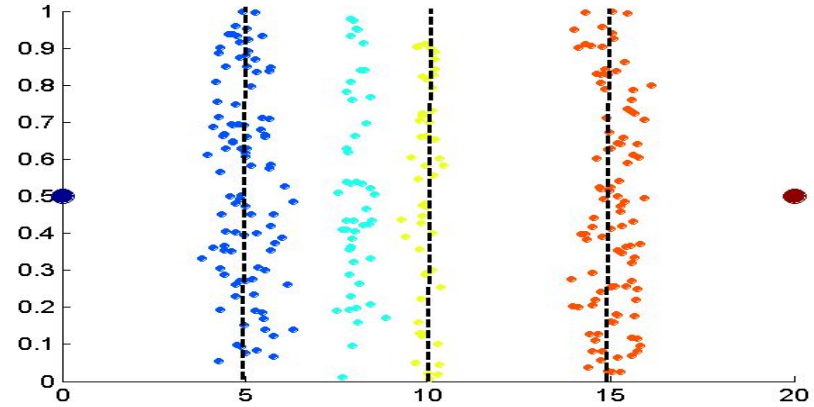
Example: Binning Methods for Data Smoothing

- Sorted data for price (in dollars): 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34
- * Partition into equal-frequency (**equal-depth**) bins:
 - Bin 1: 4, 8, 9, 15
 - Bin 2: 21, 21, 24, 25
 - Bin 3: 26, 28, 29, 34
- * Smoothing by **bin means**:
 - Bin 1: 9, 9, 9, 9
 - Bin 2: 23, 23, 23, 23
 - Bin 3: 29, 29, 29, 29
- * Smoothing by **bin boundaries**:
 - Bin 1: 4, 4, 4, 15
 - Bin 2: 21, 21, 25, 25
 - Bin 3: 26, 26, 26, 34

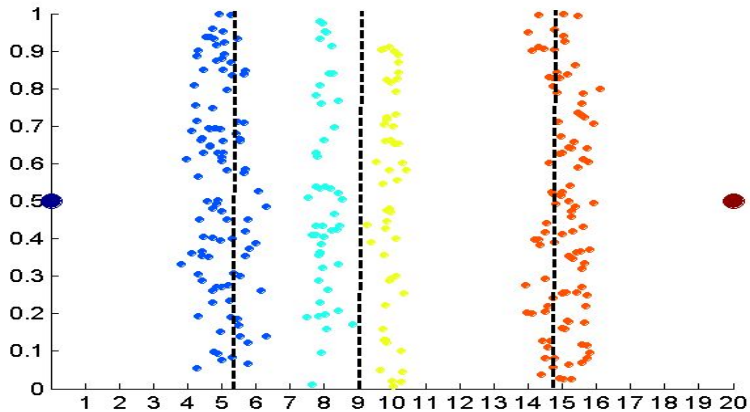
Discretization Without Supervision: Binning vs. Clustering



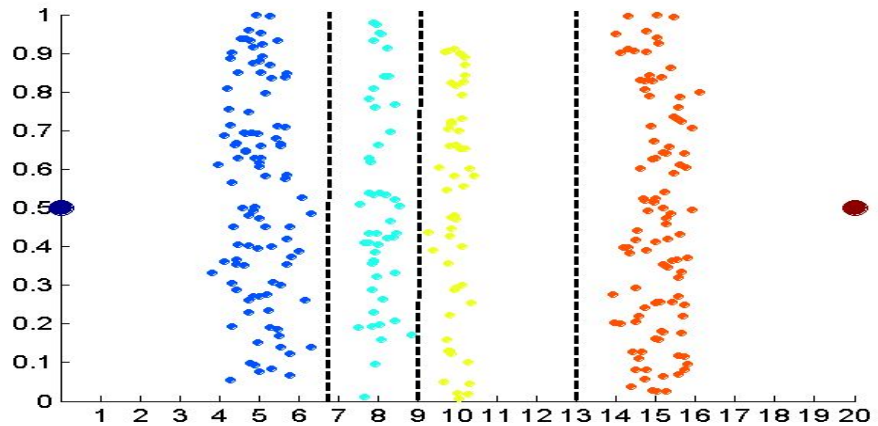
Data



Equal width (distance) binning



Equal depth (frequency) (binning)



K-means clustering leads to better results

Discretization by Classification & Correlation Analysis

- ❑ Classification (e.g., decision tree analysis)
 - ❑ Supervised: Given class labels, e.g., cancerous vs. benign
 - ❑ Using *entropy* to determine split point (discretization point)
 - ❑ Top-down, recursive split
 - ❑ Details to be covered in Chapter “Classification”
- ❑ Correlation analysis (e.g., Chi-merge: χ^2 -based discretization)
 - ❑ Supervised: use class information
 - ❑ Bottom-up merge: Find the best neighboring intervals (those having similar distributions of classes, i.e., low χ^2 values) to merge
 - ❑ Merge performed recursively, until a predefined stopping condition

Concept Hierarchy Generation

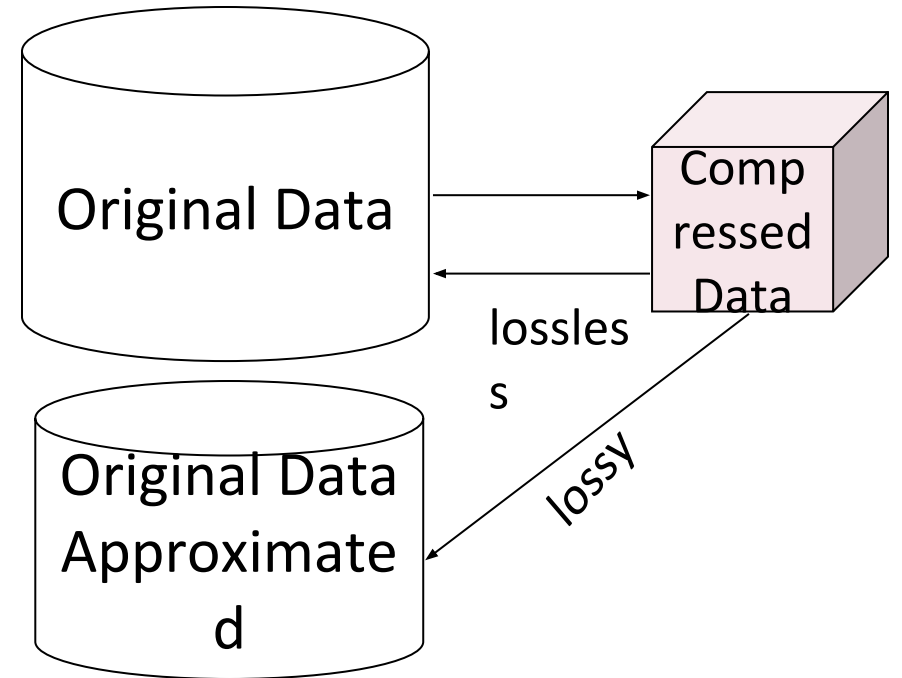
- ❑ **Concept hierarchy** organizes concepts (i.e., attribute values) hierarchically and is usually associated with each dimension in a data warehouse
- ❑ Concept hierarchies facilitate drilling and rolling in data warehouses to view data in multiple granularity
- ❑ Concept hierarchy formation: Recursively reduce the data by collecting and replacing low level concepts (such as numeric values for *age*) by higher level concepts (such as *youth, adult, or senior*)
- ❑ Concept hierarchies can be explicitly specified by domain experts and/or data warehouse designers
- ❑ Concept hierarchy can be automatically formed for both numeric and nominal data—For numeric data, use discretization methods shown

Concept Hierarchy Generation for Nominal Data

- Specification of a partial/total ordering of attributes explicitly at the schema level by users or experts
 - *street < city < state < country*
- Specification of a hierarchy for a set of values by explicit data grouping
 - {Urbana, Champaign, Chicago} < Illinois
- Specification of only a partial set of attributes
 - E.g., only *street < city*, not others
- Automatic generation of hierarchies (or attribute levels) by the analysis of the number of distinct values
 - E.g., for a set of attributes: {*street, city, state, country*}

Data Compression

- ❑ String compression
 - ❑ There are extensive theories and well-tuned algorithms
 - ❑ Typically lossless, but only limited manipulation is possible without expansion
- ❑ Audio/video compression
 - ❑ Typically lossy compression, with progressive refinement
 - ❑ Sometimes small fragments of signal can be reconstructed without reconstructing the whole
- ❑ Time sequence is not audio
 - ❑ Typically short and vary slowly with time
- ❑ Data reduction and dimensionality reduction may also be considered as forms of data compression



Lossy vs. lossless compression

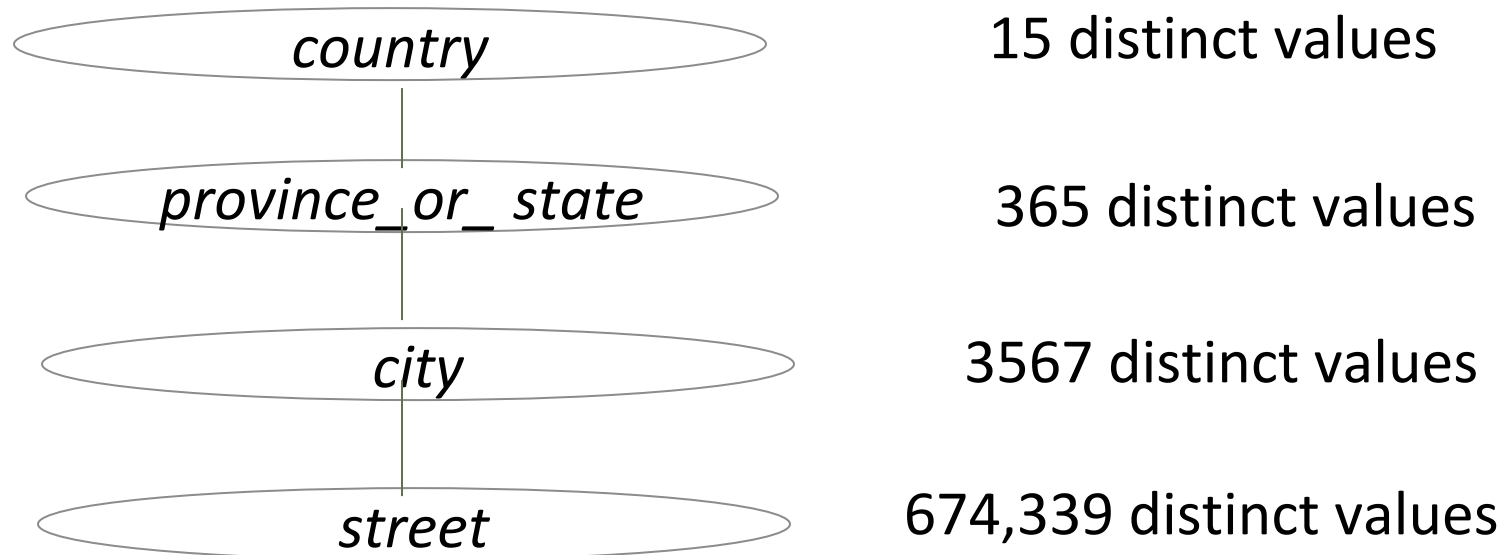
Data Cube Aggregation

- The lowest level of a data cube (base cuboid)
 - The aggregated data for an **individual entity of interest**
 - E.g., a customer in a phone calling data warehouse
- Multiple levels of aggregation in data cubes
 - Further reduce the size of data to deal with
- Reference appropriate levels
 - Use the smallest representation which is enough to solve the task
- Queries regarding aggregated information should be answered using data cube, when possible



Automatic Concept Hierarchy Generation

- ❑ Some hierarchies can be automatically generated based on the analysis of the number of distinct values per attribute in the data set
 - ❑ The attribute with the most distinct values is placed at the lowest level of the hierarchy
 - ❑ Exceptions, e.g., weekday, month, quarter, year

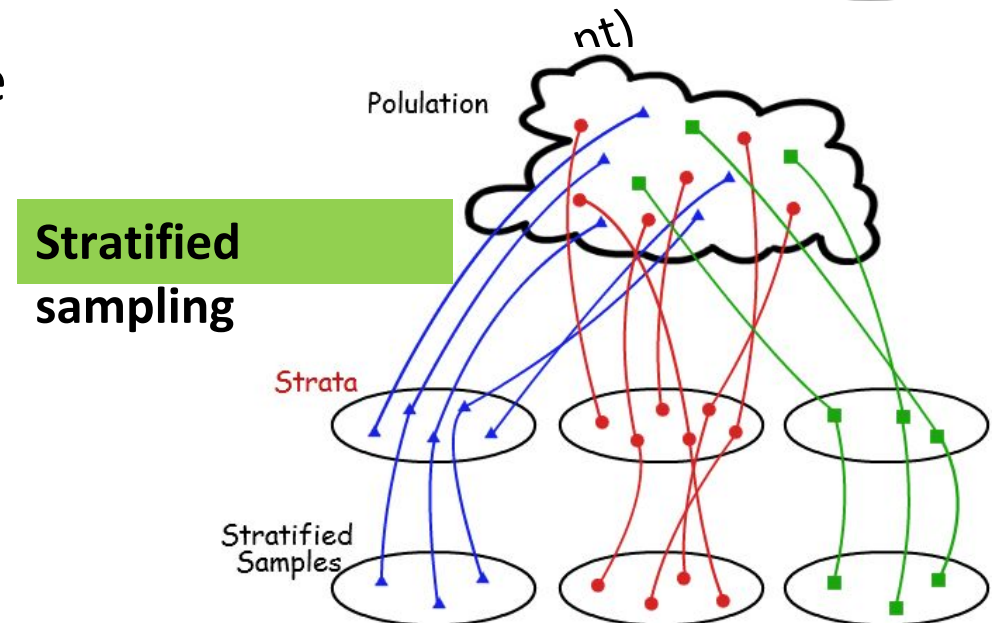
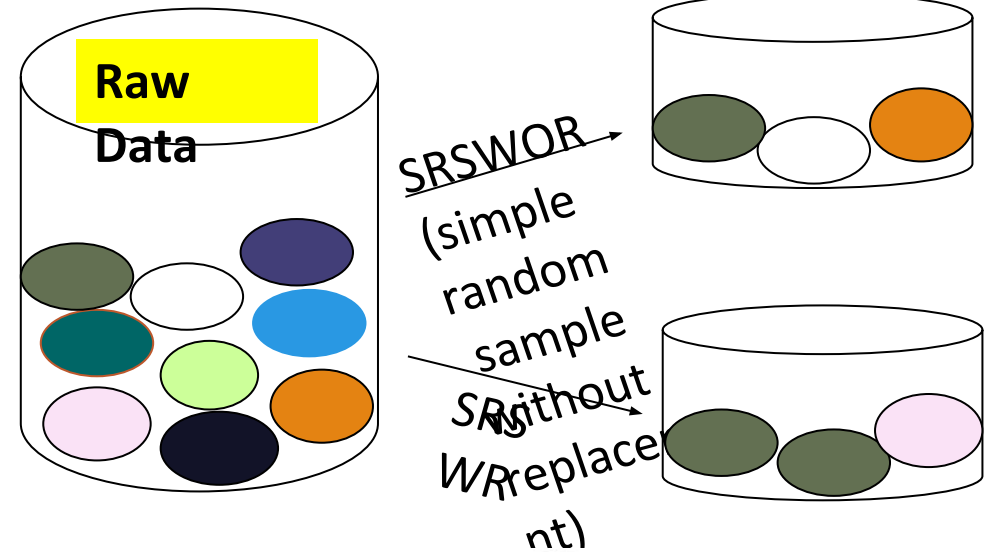


Sampling

- ❑ Sampling: obtaining a small sample s to represent the whole data set N
- ❑ Allow a mining algorithm to run in complexity that is potentially sub-linear to the size of the data
- ❑ Key principle: Choose a **representative** subset of the data
 - ❑ Simple random sampling may have very poor performance in the presence of skew
 - ❑ Develop adaptive sampling methods, e.g., stratified sampling:
- ❑ Note: Sampling may not reduce database I/Os (page at a time)

Types of Sampling

- ❑ **Simple random sampling:** equal probability of selecting any particular item
- ❑ **Sampling without replacement**
 - ❑ Once an object is selected, it is removed from the population
- ❑ **Sampling with replacement**
 - ❑ A selected object is not removed from the population
- ❑ **Stratified sampling**
 - ❑ Partition (or cluster) the data set, and draw samples from each partition (proportionally, i.e., approximately the same percentage of the data)

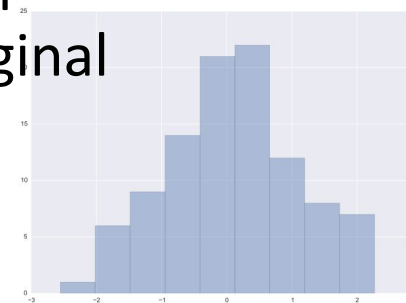
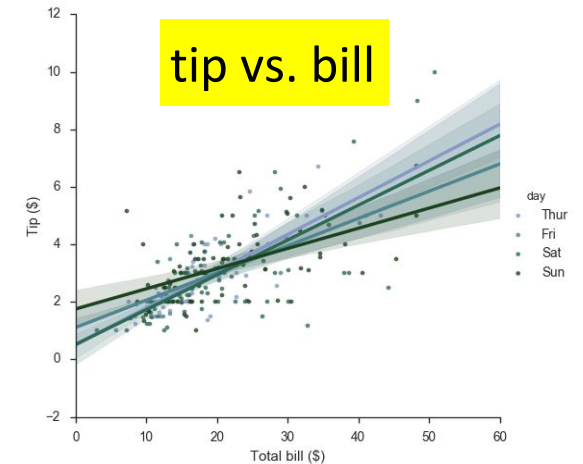


Data Reduction

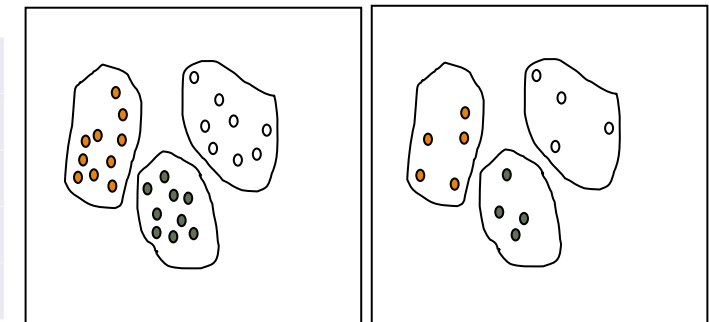
- **Data reduction:**
 - Obtain a reduced representation of the data set
 - much smaller in volume but yet produces *almost* the same analytical results
- Why data reduction?—A database/data warehouse may store terabytes of data
 - Complex analysis may take a very long time to run on the complete data set
- **Methods for data reduction** (also *data size reduction* or *numerosity reduction*)
 - Regression and Log-Linear Models
 - Histograms, clustering, sampling
 - Data cube aggregation
 - Data compression

Data Reduction: Parametric vs. Non-Parametric Methods

- ❑ Reduce data volume by choosing alternative, *smaller forms* of data representation
- ❑ **Parametric methods** (e.g., regression)
 - ❑ Assume the data fits some model, estimate model parameters, store only the parameters, and discard the data (except possible outliers)
 - ❑ Ex.: Log-linear models—obtain value at a point in m -D space as the product on appropriate marginal subspaces
- ❑ **Non-parametric methods**
 - ❑ Do not assume models
 - ❑ Major families: histograms, clustering, sampling, ...



Histogram

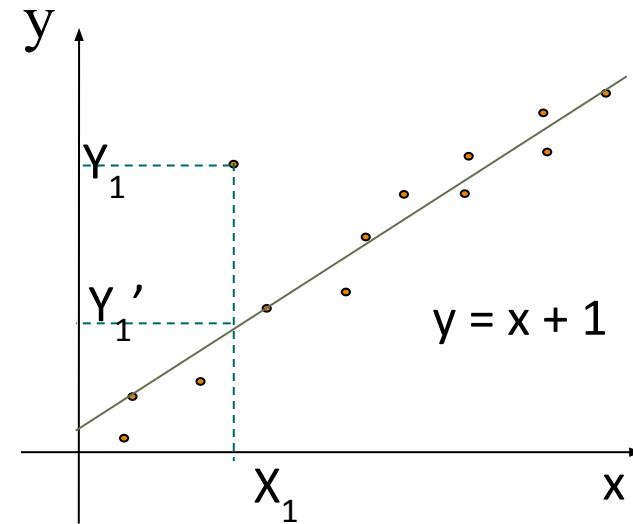


Clustering on the Raw Data

Stratified Sampling

Parametric Data Reduction: Regression Analysis

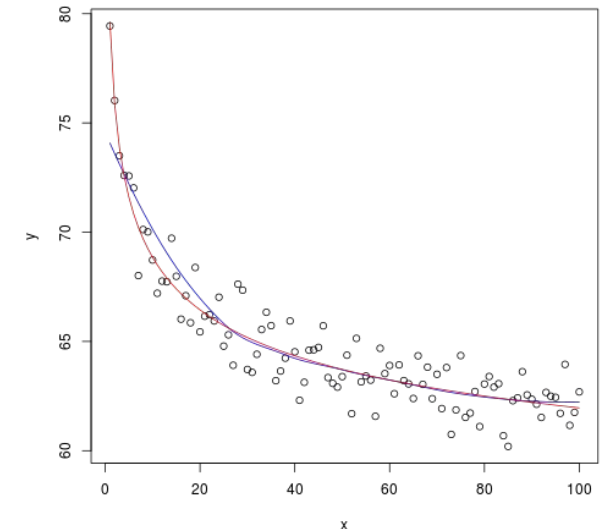
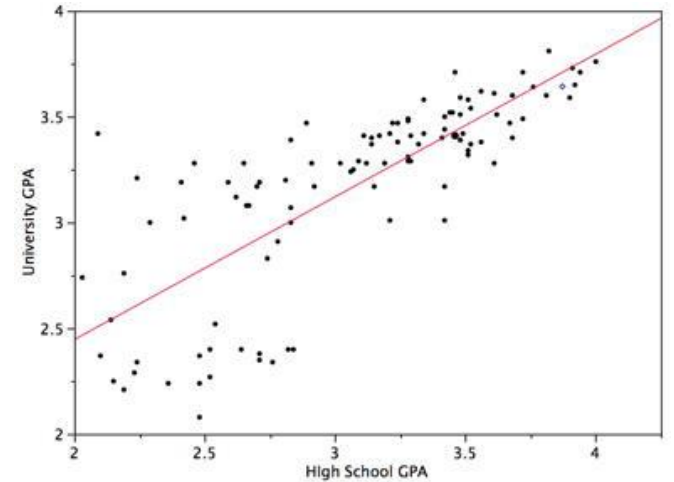
- Regression analysis: A collective name for techniques for the modeling and analysis of numerical data consisting of values of a **dependent variable** (also called **response variable** or *measurement*) and of one or more **independent variables** (also known as **explanatory variables** or **predictors**)
- The parameters are estimated so as to give a "**best fit**" of the data
- Most commonly the best fit is evaluated by using the **least squares method**, but other criteria have also been used



- Used for prediction (including forecasting of time-series data), inference, hypothesis testing, and modeling of causal relationships

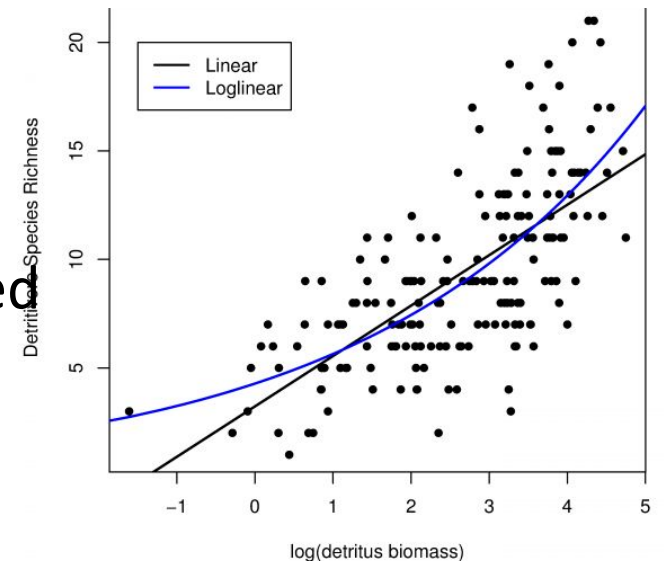
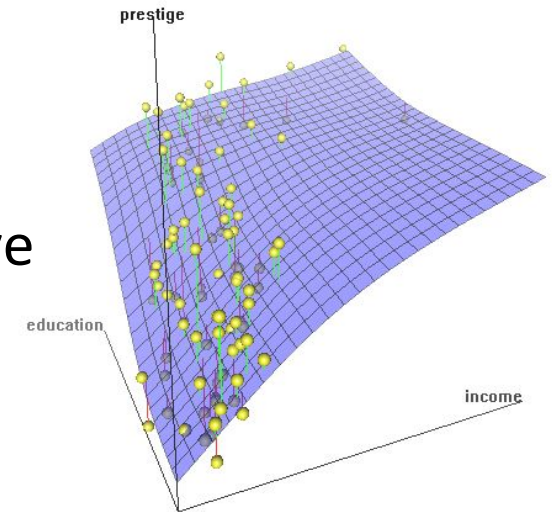
Linear and Multiple Regression

- Linear regression: $Y = wX + b$
 - Data modeled to fit a straight line
 - Often uses the least-square method to fit the line
 - Two regression coefficients, w and b , specify the line and are to be estimated by using the data at hand
 - Using the least squares criterion to the known values of $Y_1, Y_2, \dots, X_1, X_2, \dots$
- Nonlinear regression:
 - Data are modeled by a function which is a nonlinear combination of the model parameters and depends on one or more independent variables
 - The data are fitted by a method of successive approximations



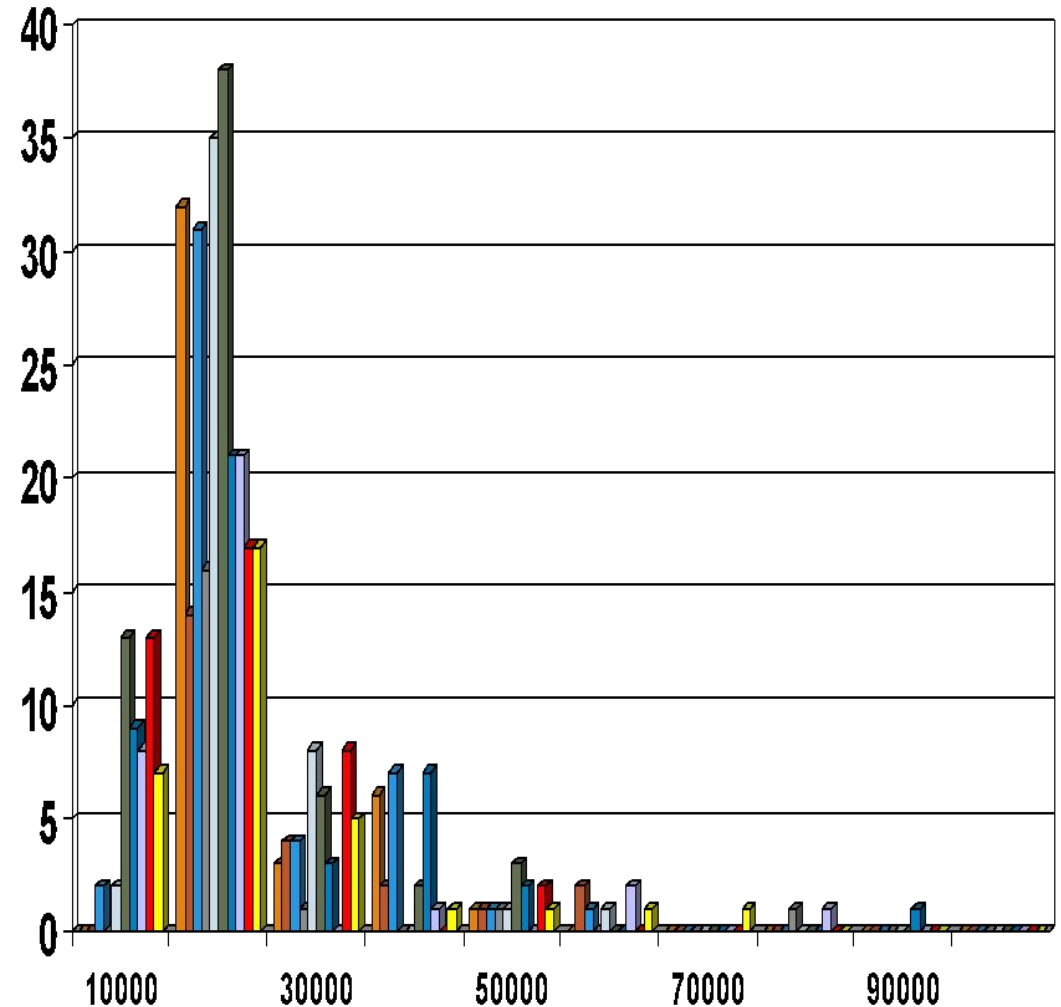
Multiple Regression and Log-Linear Models

- Multiple regression: $Y = b_0 + b_1 X_1 + b_2 X_2$
 - Allows a response variable Y to be modeled as a linear function of multidimensional feature vector
 - Many nonlinear functions can be transformed into the above
- Log-linear model:
 - A math model that takes the form of a function whose logarithm is a linear combination of the parameters of the model, which makes it possible to apply (possibly multivariate) linear regression
 - Estimate the probability of each point (tuple) in a multi-dimen. space for a set of discretized attributes, based on a smaller subset of dimensional combinations
 - Useful for dimensionality reduction and data smoothing



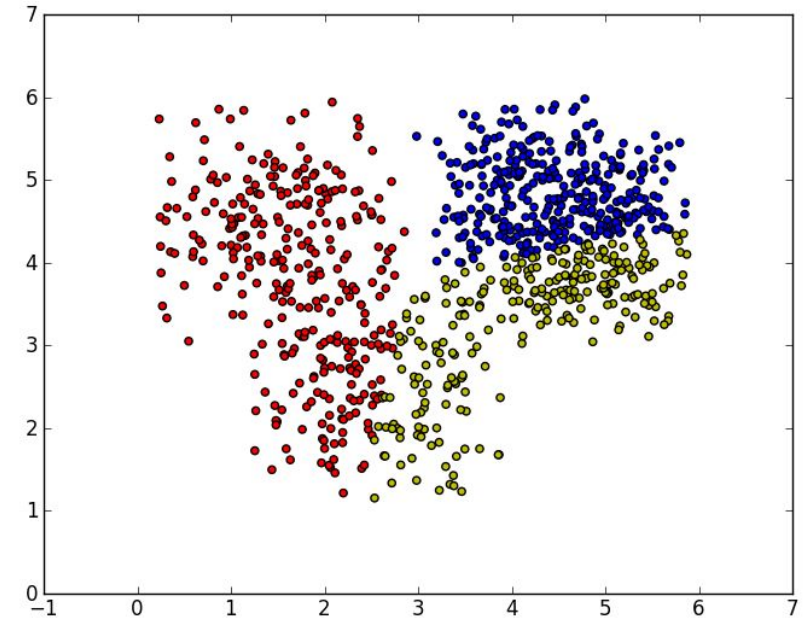
Histogram Analysis

- Divide data into buckets and store average (sum) for each bucket
- Partitioning rules:
 - Equal-width: equal bucket range
 - Equal-frequency (or equal-depth)



Clustering

- ❑ Partition data set into clusters based on similarity, and store cluster representation (e.g., centroid and diameter) only
- ❑ Can be very effective if data is clustered but not if data is “smeared”
- ❑ Can have hierarchical clustering and be stored in multi-dimensional index tree structures
- ❑ There are many choices of clustering definitions and clustering algorithms
- ❑ Cluster analysis will be studied in depth in Chapter 10



Dimensionality Reduction

- What Is Dimensionality Reduction?
- Dimensionality Reduction Methods
 - Principal Component Analysis
 - Attribute Subset Selection
 - Nonlinear Dimensionality Reduction Methods

What Is Dimensionality Reduction?

❑ **Curse of dimensionality**

- ❑ When dimensionality increases, data becomes increasingly sparse
- ❑ Density and distance between points, which is critical to clustering, outlier analysis, becomes less meaningful
- ❑ The possible combinations of subspaces will grow exponentially

❑ **Dimensionality reduction**

- ❑ Reducing the number of random variables under consideration, via obtaining a set of principal variables

❑ **Advantages of dimensionality reduction**

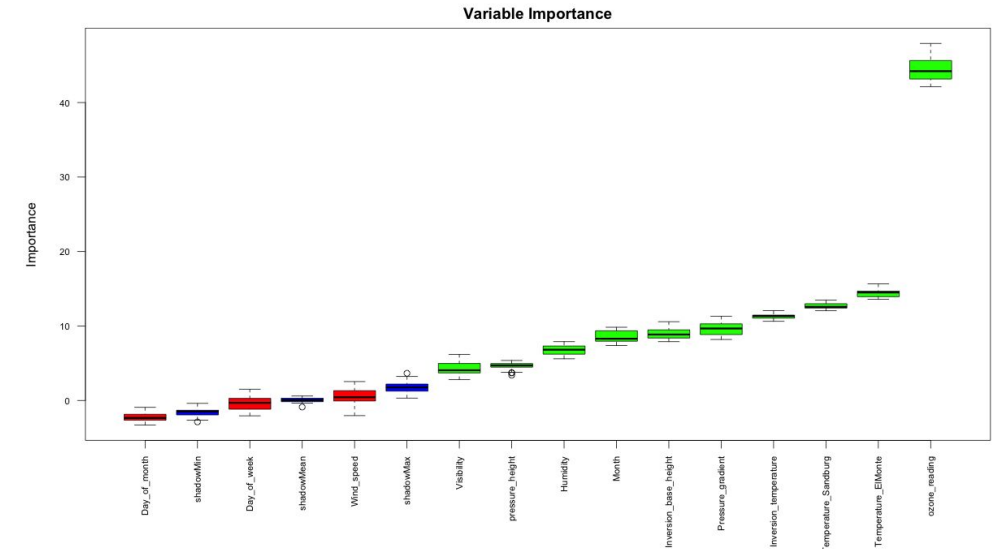
- ❑ Avoid the curse of dimensionality
- ❑ Help eliminate irrelevant features and reduce noise
- ❑ Reduce time and space required in data mining
- ❑ Allow easier visualization

Dimensionality Reduction Methods

- Dimensionality reduction methodologies
 - **Feature selection:** Find a subset of the original variables (or features, attributes)
 - **Feature extraction:** Transform the data in the high-dimensional space to a space of fewer dimensions
- Some typical dimensionality reduction methods
 - Principal Component Analysis
 - Attribute Subset Selection
 - Nonlinear Dimensionality Reduction

Attribute Subset Selection

- Another way to reduce dimensionality of data
- Redundant attributes
 - Duplicate much or all of the information contained in one or more other attributes
 - E.g., purchase price of a product and the amount of sales tax paid
- Irrelevant attributes
 - Contain no information that is useful for the data mining task at hand
 - Ex. A student's ID is often irrelevant to the task of predicting his/her GPA



Heuristic Search in Attribute Selection

- ❑ There are 2^d possible attribute combinations of d attributes
- ❑ Typical heuristic attribute selection methods:
 - ❑ Best single attribute under the attribute independence assumption: choose by significance tests
 - ❑ Best step-wise feature selection:
 - ❑ The best single-attribute is picked first
 - ❑ Then next best attribute condition to the first, ...
 - ❑ Step-wise attribute elimination:
 - ❑ Repeatedly eliminate the worst attribute
 - ❑ Best combined attribute selection and elimination
 - ❑ Optimal branch and bound:
 - ❑ Use attribute elimination and backtracking