

# **Statistical Learning**

## **Measures of central tendency, dispersion and correlation**

# Outline

1. Raw Data
2. Frequency Distribution - Histograms
3. Cumulative Frequency Distribution
4. Measures of Central Tendency
5. Mean, Median, Mode
6. Measures of Dispersion
7. Range, IQR, Standard Deviation, coefficient of variation
8. Normal distribution, Chebyshev Rule.
9. Five number summary, boxplots, QQ plots, Quantile plot, scatter plot.
10. Visualization: scatter plot matrix, parallel coordinates.
11. Correlation analysis

# Data versus Information

When managers are bewildered by plethora of data, which do not make any sense on the surface of it, they are looking for methods to classify data that would convey meaning. The idea here is to help them draw the right conclusion. This session provides the nitty-gritty of arranging data into **information**.

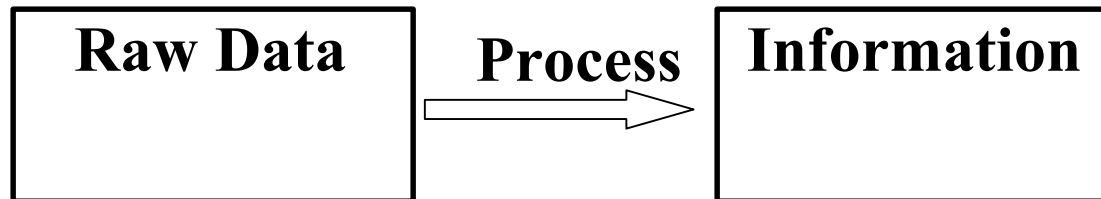
# Raw Data

## Meaning of Raw Data:

Raw Data represent numbers and facts in the original format in which the data have been Collected. You need to convert the raw data into information for managerial decision Making.

# Information is Key

Large and massive raw data tend to bewilder you so much that the overall patterns are obscured. You cannot see the wood for the trees. This implies that the raw data must be processed to give you useful information.

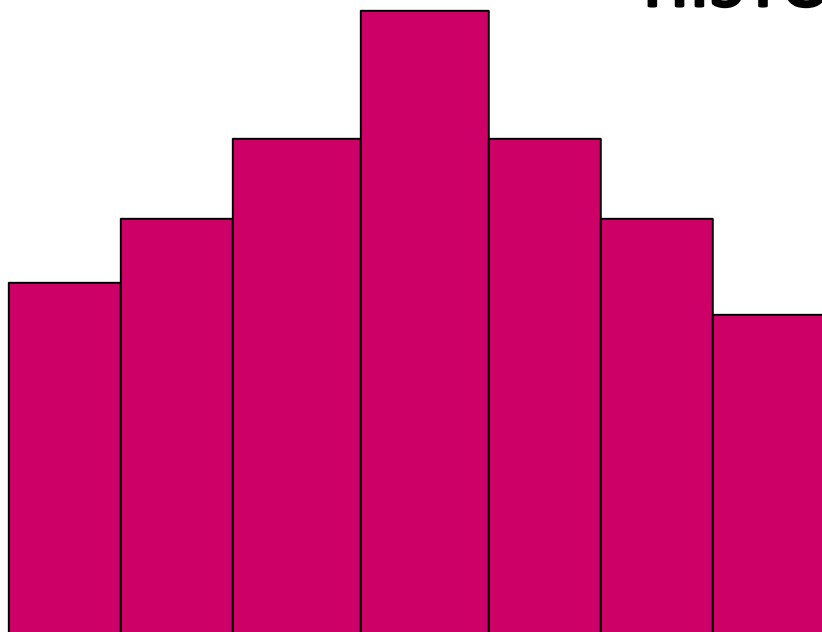


# Frequency Distribution

In simple terms, frequency distribution is a summarized table in which raw data are arranged into classes and frequencies.

Frequency distribution focuses on classifying raw data into information. It is the most widely used data reduction technique in descriptive statistics.

# HISTOGRAM



**Histogram** (also known as frequency histogram) is a snap shot of the frequency distribution.

Histogram is a graphical representation of the frequency distribution in which the X-axis represents the classes and the Y-axis represents the frequencies in bars

Histogram depicts the pattern of the distribution emerging from the characteristic being measured.

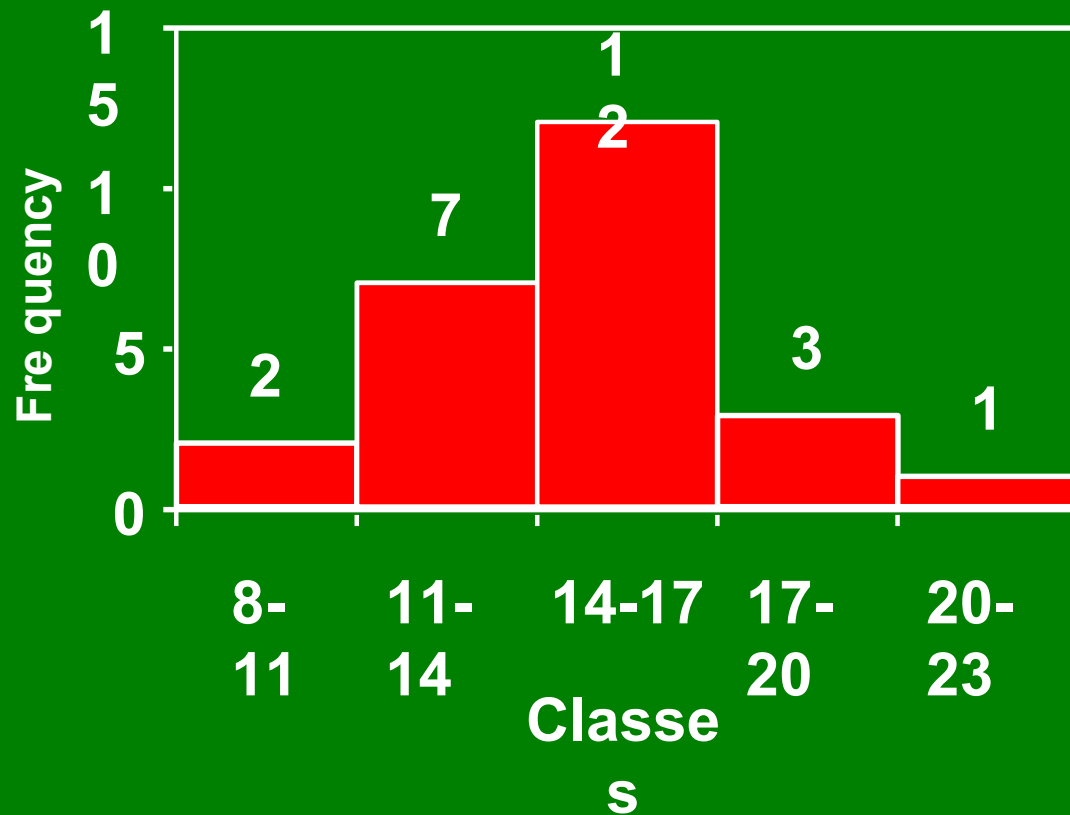
# Histogram- Example

The inspection records of a hose assembly operation revealed a high level of rejection. An analysis of the records showed that the "leaks" were a major contributing factor to the problem. It was decided to investigate the hose clamping operation. The hose clamping force (torque) was measured on twenty five assemblies. (Figures in foot-pounds). The data are given below: Draw the frequency histogram and comment.

8	13	15	10	16
11	14	11	14	20
15	16	12	15	13
12	13	16	17	17
14	14	14	18	15



# Histogram Example Solution



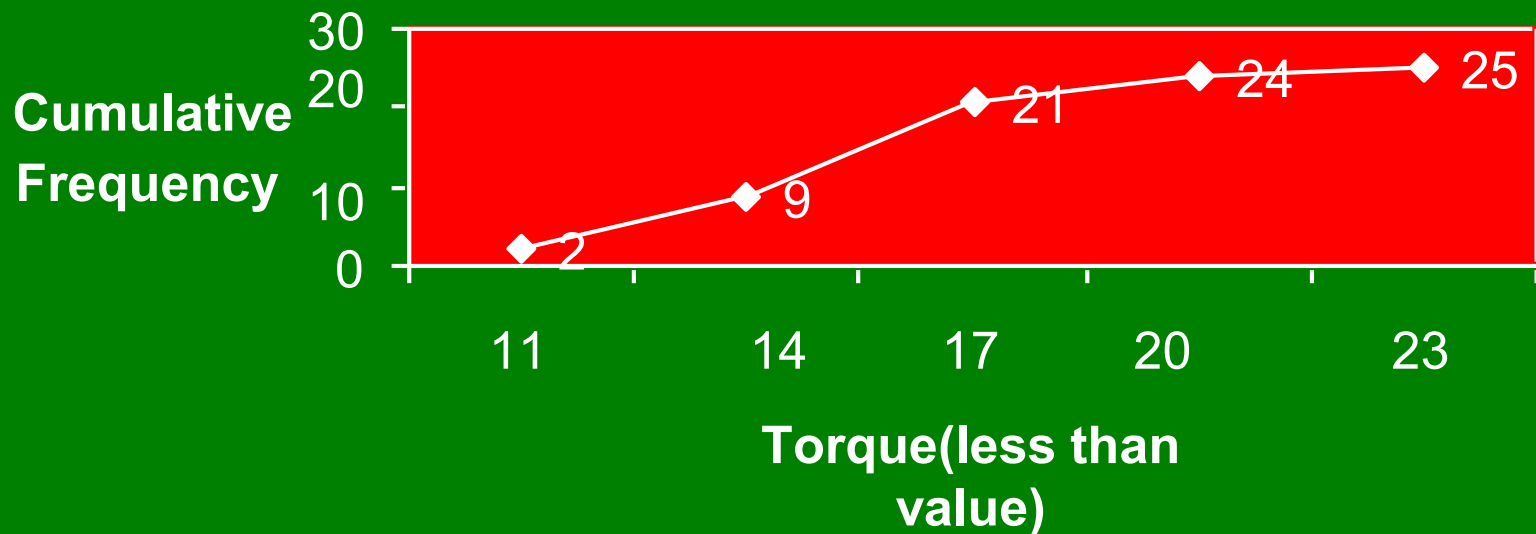
# Cumulative Frequency Distribution

A type of frequency distribution that shows how many observations are above or below the lower boundaries of the classes. You can formulate the following from the previous example of hose clamping

Class (torque lb)	Frequency	Relative Frequency	Cumulative Frequency	Cumulative Relative Frequency
8-11	2	0.08	2	0.08
11-14	7	0.28	9	0.36
14-17	12	0.48	21	0.84
17-20	3	0.12	24	0.96
20-23	1	0.04	25	1.00
<b>Total</b>	<b>25</b>	<b>1.00</b>		

# Ogive (Cumulative Frequency Distribution)

## Cumulative Distribution(Ogive Curve) for the Example



# What is Central Tendency?

Whenever you measure things of the same kind, a fairly large number of such measurements will tend to cluster around the middle value. Such a value is called a measure of "Central Tendency". The other terms that are used synonymously are "Measures of Location", or "Statistical Averages".

# Measures of Central Tendency

As a manager, You need the summary measures of central tendency to draw conclusions in your functional area of operation. The most widely used measures of central tendency are **Arithmetic Mean**, **Median**, and **Mode**.

# Arithmetic Mean

Arithmetic Mean (called mean) is defined as the sum of all observations in a data set divided by the total number of observations. For example, consider a data set containing the following observations:

In symbolic form mean is given by  $\bar{X} = \frac{\sum X}{n}$

$\bar{X}$   
Mean = Arithmetic

$\sum$   
 $X$  = Indicates sum all X values in the data set

$n$  = Total number of observations(Sample Size)

# Arithmetic Mean -Example

The inner diameter of a particular grade of tire based on 5 sample measurements are as follows: (figures in millimeters)

565, 570, 572, 568, 585

Applying the formula  $\bar{X} = \frac{\sum X}{n}$

We get mean =  $(565+570+572+568+585)/5 = 572$

Caution: Arithmetic Mean is affected by extreme values or fluctuations in sampling. It is not the best average to use when the data set contains extreme values (Very high or very low values).

# Median

Median is the middle most observation when you arrange data in ascending order of magnitude. Median is such that 50% of the observations are above the median and 50% of the observations are below the median.

Median is a very useful measure for ranked data in the context of consumer preferences and rating. It is not affected by extreme values (greater resistance to outliers)

$$\text{Median} = \frac{n+1}{2} \text{ th value of ranked data}$$

$n$  = Number of observations in the sample



## Median -

### Example

Marks obtained by 7 students in Computer Science Exam are given below: Compute the median.

45      40      60      80      90      65      55

Arranging the data after ranking gives

90      80      65      60      55      45      40

Median =  $(n+1)/2$  th value in this set =  $(7+1)/2$  th observation = 4<sup>th</sup> observation = 60

Hence Median = 60 for this problem.

# Mode

Mode is that value which occurs most often. It has the maximum frequency of occurrence. Mode also has resistance to outliers.

Mode is a very useful measure when you want to keep in the inventory, the most popular shirt in terms of collar size during festival season.

# Mode -Example

The life in number of hours of 10 flashlight batteries are

as follows: Find the mode.  
340 350 340 340 320 340 330 330  
340 350

340 occurs five times. Hence, mode=340.

# Comparison of Mean, Median, Mode

Mean	Median	Mode
Defined as the arithmetic average of all observations in the data set.	Defined as the middle value in the data set arranged in ascending or descending order.	Defined as the most frequently occurring value in the distribution; it has the largest frequency.
Requires measurement on all observations.	Does not require measurement on all observations	Does not require measurement on all observations
Uniquely and comprehensively defined.	Cannot be determined under all conditions.	Not uniquely defined for multi-modal situations.
Proprietary content. ©Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited		

## Comparison of Mean, Median, Mode Cont.

Mean	Median	Mode
Affected by extreme values.  Can be treated algebraically. That is, Means of several groups can be combined.	Not affected by extreme values.  Cannot be treated algebraically. That is, Medians of several groups cannot be combined.	Not affected by extreme values.  Cannot be treated algebraically. That is, Modes of several groups cannot be combined.

# Measures of Dispersion

In simple terms, measures of dispersion indicate how large the spread of the distribution is around the central tendency. It answers unambiguously the question "What is the magnitude of departure from the average value for different groups having identical averages?".

# Range

Range is the simplest of all measures of dispersion. It is calculated as the difference between maximum and minimum value in the data set.

$$\text{Range} = X_{\text{Maximum}} - X_{\text{Minimum}}$$

# Range-Example

## Example for Computing Range

The following data represent the percentage return on investment for 10 mutual funds per annum. Calculate Range.

12, 14, 11, 18, 10.5, 11.3, 12, 14, 11, 9

$$\text{Range} = X_{\text{Maximum}} - X_{\text{Minimum}} = 18 - 9 = 9$$

**Caution:** If one of the components of range namely the maximum value or minimum value becomes an extreme value, then range should not be used.



# Inter-Quartile Range(IQR)

IQR= Range computed on middle 50% of the observations after eliminating the highest and lowest 25% of observations in a data set that is arranged in ascending order. IQR is less affected by outliers.

$$\text{IQR} = Q_3 - Q_1$$

## Interquartile Range-Example

The following data represent the return on percentage investment for 9 mutual funds per annum. Calculate interquartile range.

Data Set: 12, 14, 11, 18, 10.5, 12, 14, 11, 9

Arranging in ascending order, the data set becomes  
9, 10.5, 11, 11, 12, 12, 14, 14, 18

$$\text{IQR} = Q_3 - Q_1 = 14 - 10.75 = 3.25$$

# Standard Deviation

Standard deviation forms the cornerstone for Inferential Statistics.

To define standard deviation, you need to define another term called variance. In simple terms, standard deviation is the square root of variance.

# Key Formulas

## Important Terms with Notations

Sample Variance

$$s^2 = \frac{\sum (X - \bar{X})^2}{n - 1}$$

Sample Standard Deviation

$$s = \sqrt{\frac{\sum (X - \bar{X})^2}{n - 1}}$$

Population Variance

$$\sigma^2 = \frac{\sum (X - \mu)^2}{N}$$

Population Standard

$$\text{Deviation } \sigma = \sqrt{\frac{\sum (X - \mu)^2}{N}}$$

Where  $\bar{X} = \frac{\sum X}{n}$  (Sample

Mean) and

$$\mu = \frac{\sum X}{N} \text{ (Population Mean)}$$

n = Number of observations in the sample (Sample size)

N = Number of observations in the Population (Population Size)

## Remarks

1.  $\frac{(X - \bar{X})^2}{n - 1}$  is an unbiased estimator of  $\sigma^2 = \frac{\sum (X - \mu)^2}{N}$
2.  $\bar{X} = \frac{\sum X}{n}$  is an unbiased estimator of  $\mu = \frac{\sum X}{N}$
3. The divisor  $n - 1$  is always used while calculating sample variance for ensuring property of being unbiased
4. Standard deviation is always the square root of variance

## Example for Standard Deviation

The following data represent the percentage return on investment for 10 mutual funds per annum. Calculate the sample standard deviation.

12, 14, 11, 18, 10.5, 11.3, 12, 14, 11, 9

# Solution for the Example

A	B	C	D
1			
2	X	$X - \bar{X}$	$(X - \bar{X})^2$
3	12	-0.28	0.08
4	14	1.72	2.96
5	11	-1.28	1.64
6	18	5.72	32.72
7	10.5	-1.78	3.17
8	11.3	-0.98	0.96
9	12	-0.28	0.08
10	14	1.72	2.96
11	11	-1.28	1.64
12	9	-3.28	10.76
13	Mean =		56.96
14	12.28	Variance=	6.33
15		Standard Deviation=	2.52

# Coefficient of Variation (Relative Dispersion)

Coefficient of Variation is defined as the ratio of (CV) Standard Deviation to Mean.

In symbolic form

$CV = \frac{S}{\bar{X}}$  for the sample data and  $= \frac{\sigma}{\mu}$  for the population data.

# Coefficient of Variation

## Example

Consider two Sales Persons in the territory. The sales performance of these two in the context of selling PCs are given below. Comment on the results.

### Sales Person 1

Mean Sales (One year average) 50

units  
Standard Deviation 5  
units

### Sales Person 2

Mean Sales (One year average) 75 units

Standard deviation  
25 units



## Interpretation for the Example

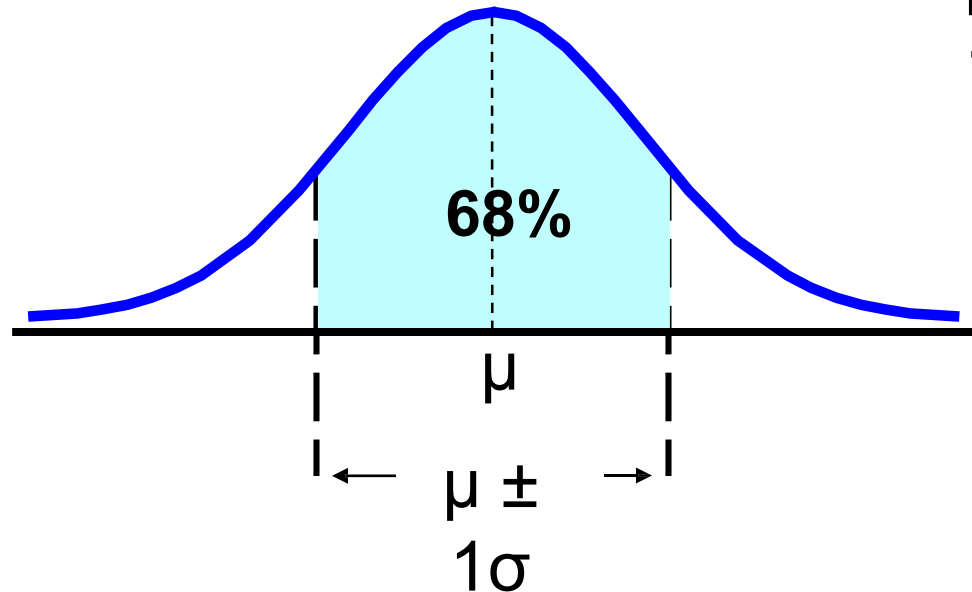
The CV is  $5/50 = 0.10$  or 10% for the Sales Person1 and  $25/75 = 0.33$  or 33% for sales Person2.

The moral of the story is "don't get carried away by absolute number". Look at the scatter. Even though, Sales Person2 has achieved a higher average, his performance is not consistent seems erratic, and

# The Empirical Rule

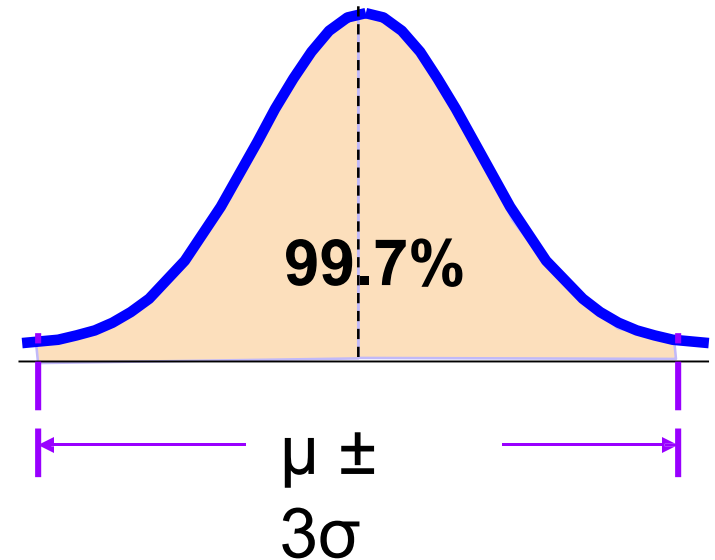
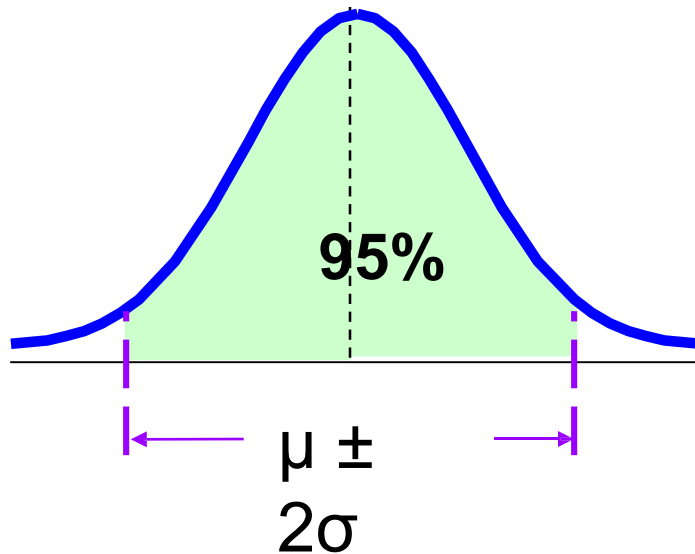
- The empirical rule approximates the variation of data in a bell-shaped distribution
- Approximately 68% of the data in a bell shaped distribution is within 1 standard deviation of the mean or

$$\mu \pm 1\sigma$$



# The Empirical Rule **greatlearning**

- Approximately 95% of the data in a bell-shaped distribution lies within two standard deviations of the mean, or  $\mu \pm 2\sigma$
- Approximately 99.7% of the data in a bell-shaped distribution lies within three standard deviations of the mean, or  $\mu \pm 3\sigma$



## Chebyshev Rule

- Regardless of how the data are distributed, at least  $(1 - 1/k^2) \times 100\%$  of the values will fall within  $k$  standard deviations of the mean (for  $k > 1$ )
- For Example, when  $k=2$ , at least 75% of the values of any data set will be within  $\mu \pm 2\sigma$

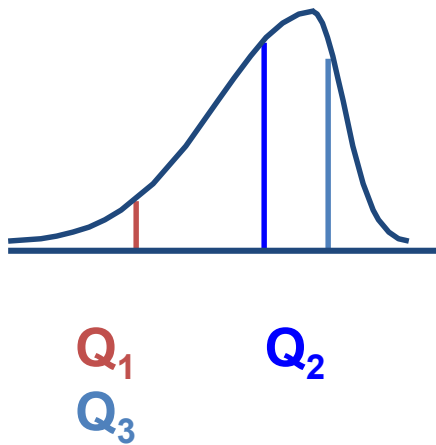
# The Five Number Summary

The five numbers that help describe the center, spread and shape of data are:

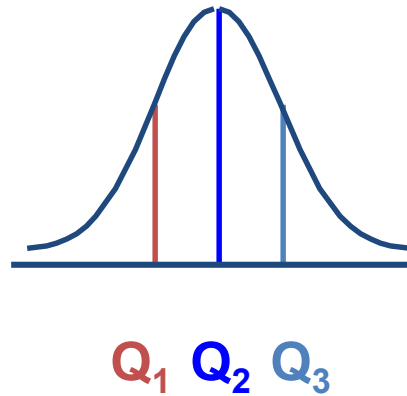
- $X_{\text{smallest}}$
- First Quartile ( $Q_1$ )
- Median ( $Q_2$ )
- Third Quartile ( $Q_3$ )
- $X_{\text{largest}}$

# Distribution Shape

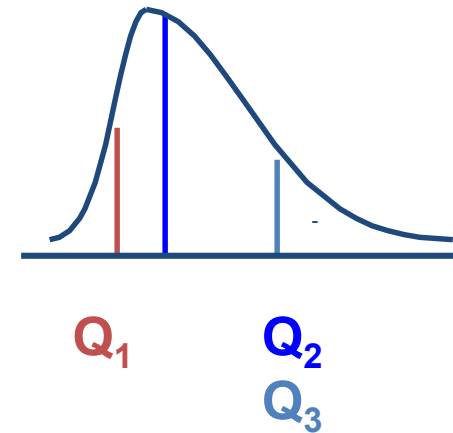
## Left-Skewed



## Symmetric



## Right-Skewed

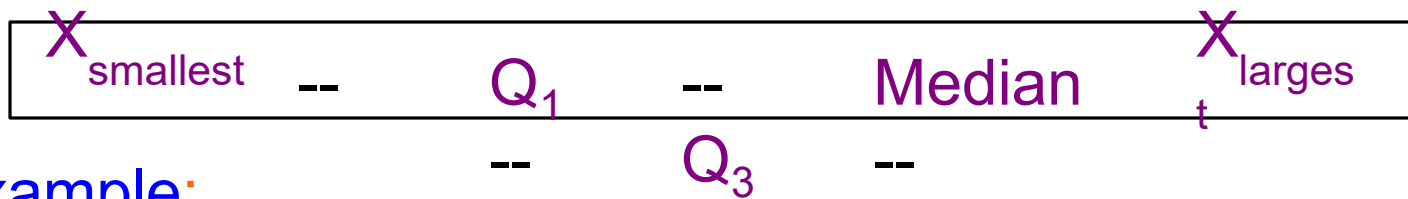


# Relationships among the five-number summary and distribution shape

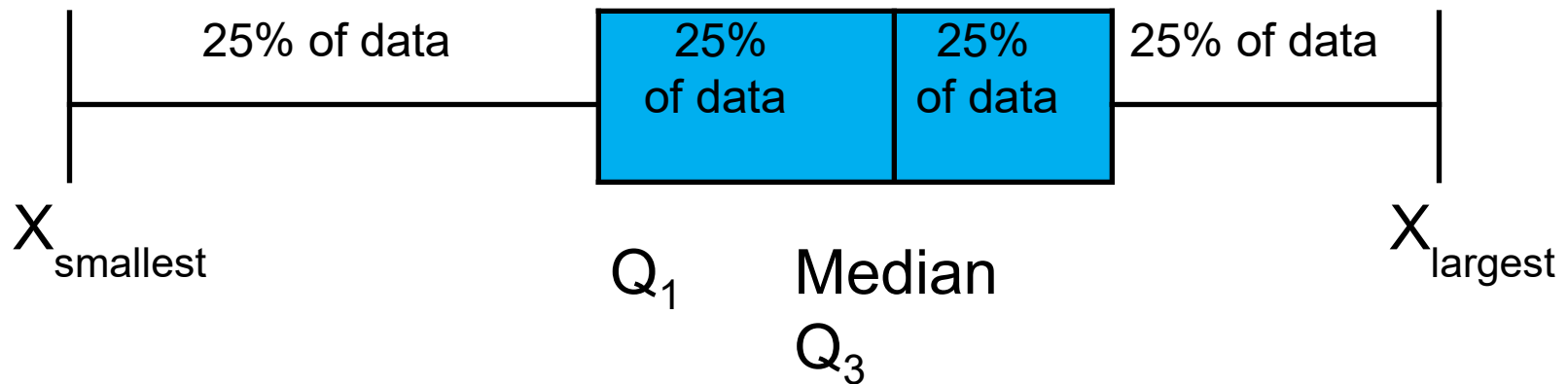
Left-Skewed	Symmetric	Right-Skewed
Median – $X_{\text{smallest}}$ >	Median – $X_{\text{smallest}}$ $\approx$	Median – $X_{\text{smallest}}$ <
$X_{\text{largest}}$ – Median $Q_1 - X_{\text{smallest}}$ >	$X_{\text{largest}}$ – Median $Q_1 - X_{\text{smallest}}$ $\approx$	$X_{\text{largest}}$ – Median $Q_1 - X_{\text{smallest}}$ <
$X_{\text{largest}}$ – $Q_3$	$X_{\text{largest}}$ – $Q_3$	$X_{\text{largest}}$ – $Q_3$
Median – $Q_1$ >	Median – $Q_1$ $\approx$	Median – $Q_1$ <
$Q_3$ – Median	$Q_3$ – Median	$Q_3$ – Median

# Five Number Summary and The Boxplot

- **The Boxplot:** A Graphical display of the data based on the five-number summary:



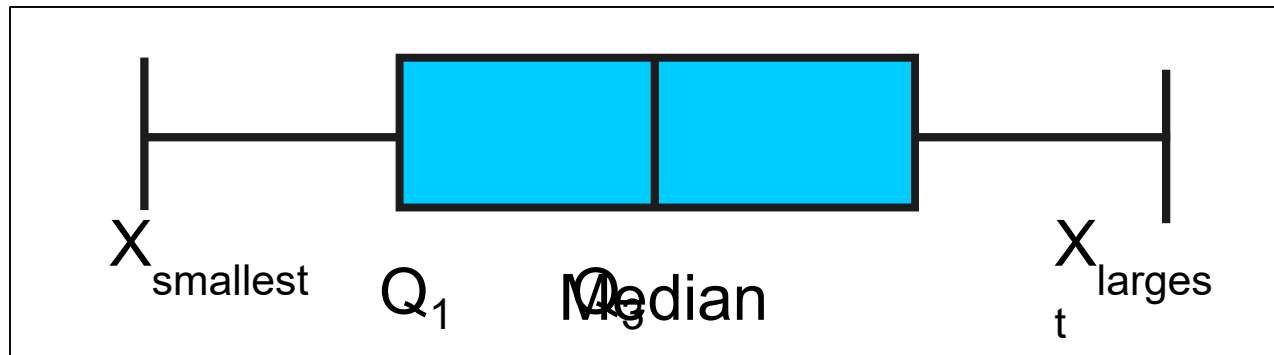
**Example:**





## Five Number Summary: Shape of Boxplots

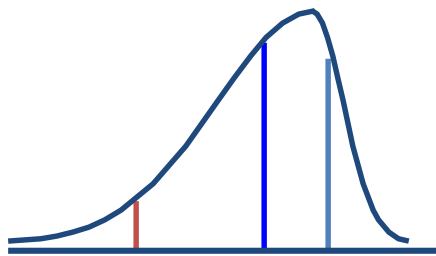
- If data are symmetric around the median then the box and central line are centered between the endpoints



- A Boxplot can be shown in either a vertical or horizontal orientation

# Distribution Shape and The Boxplot

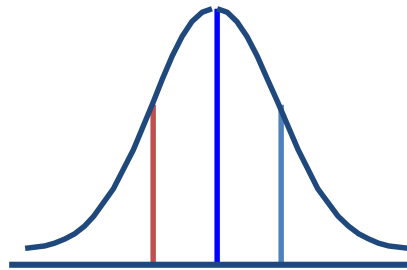
## Left-Skewed



$Q_1$   $Q_2$   
 $Q_3$



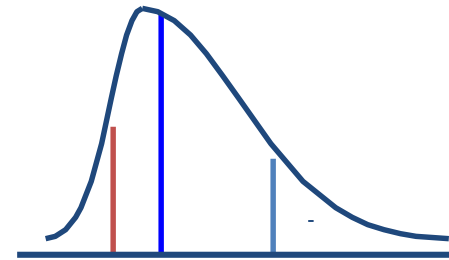
## Symmetric



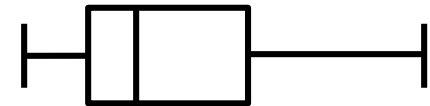
$Q_1$   $Q_2$   $Q_3$



## Right-Skewed



$Q_1$   $Q_2$   
 $Q_3$



# Boxplot Example

Below is a Boxplot for the following data:

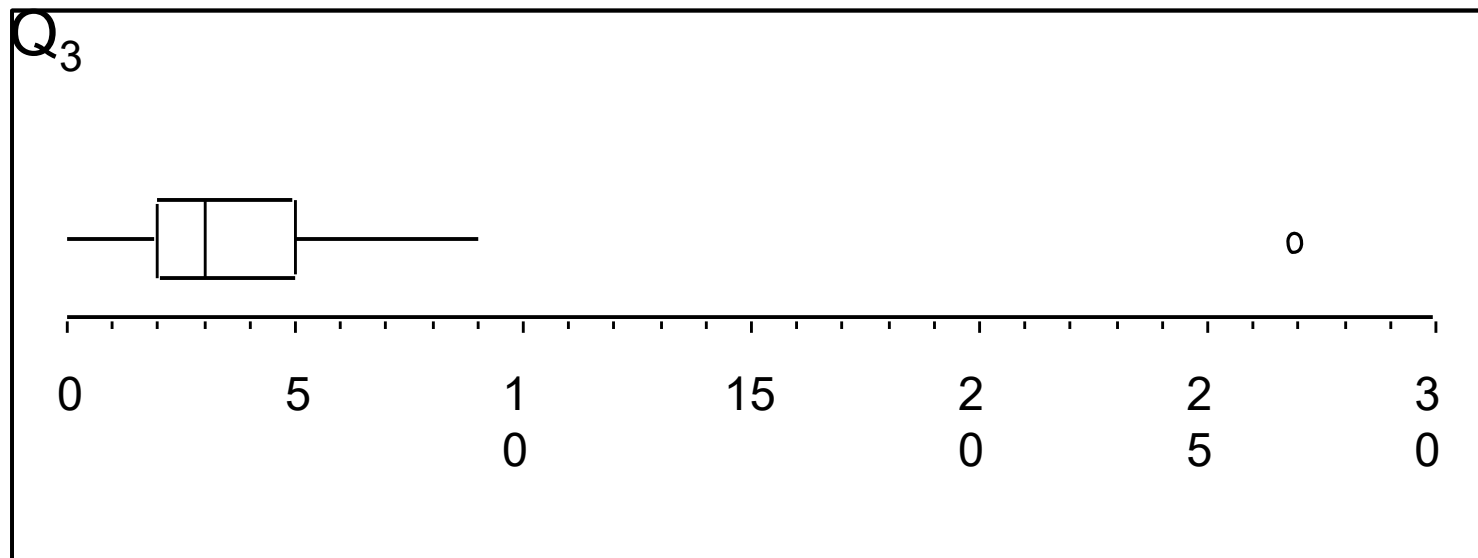
$X_{\text{smallest}}$	$Q_1$	$Q_2$	$Q_3$	$X_{\text{largest}}$
0	2	2	3	27
		4	5	
2		9		



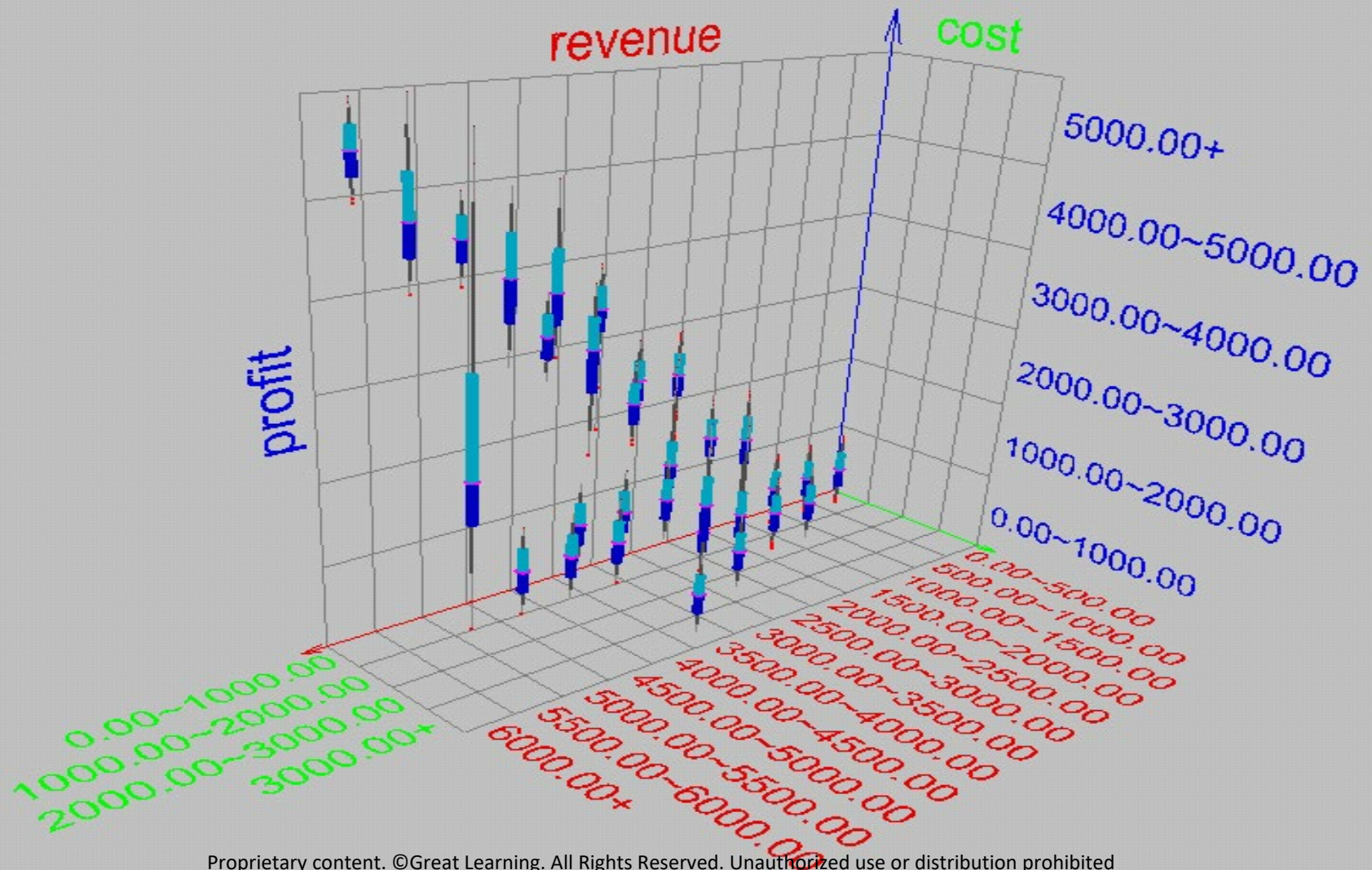
The data are right skewed, as the plot depicts

# Box plot example showing an outlier

- The boxplot below of the same data shows the outlier value of 27 plotted separately
- A value is considered an outlier if it is more than 1.5 times the interquartile range below  $Q_1$  or above



# Visualization of Data Dispersion: 3-D Boxplots



# Graphic Displays of Basic Statistical Descriptions

**Boxplot:** graphic display of five-number summary

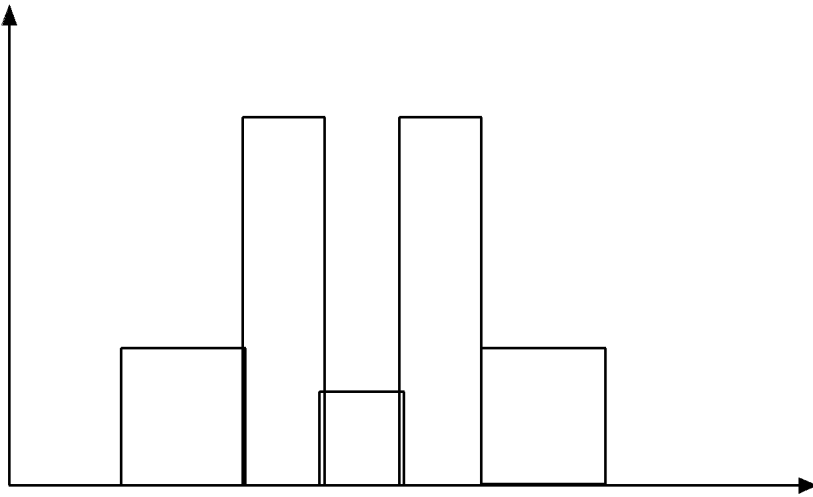
**Histogram:** x-axis are values, y-axis repres. frequencies

**Quantile plot:** each value  $x_i$  is paired with  $f_i$  indicating that approximately  $100 f_i \%$  of data are  $\leq x_i$

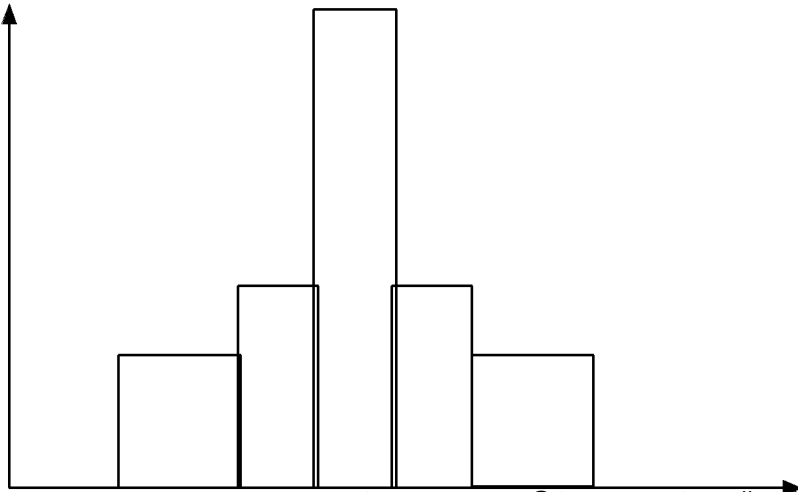
**Quantile-quantile (q-q) plot:** graphs the quantiles of one univariant distribution against the corresponding quantiles of another

**Scatter plot:** each pair of values is a pair of coordinates and plotted as points in the plane

# 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

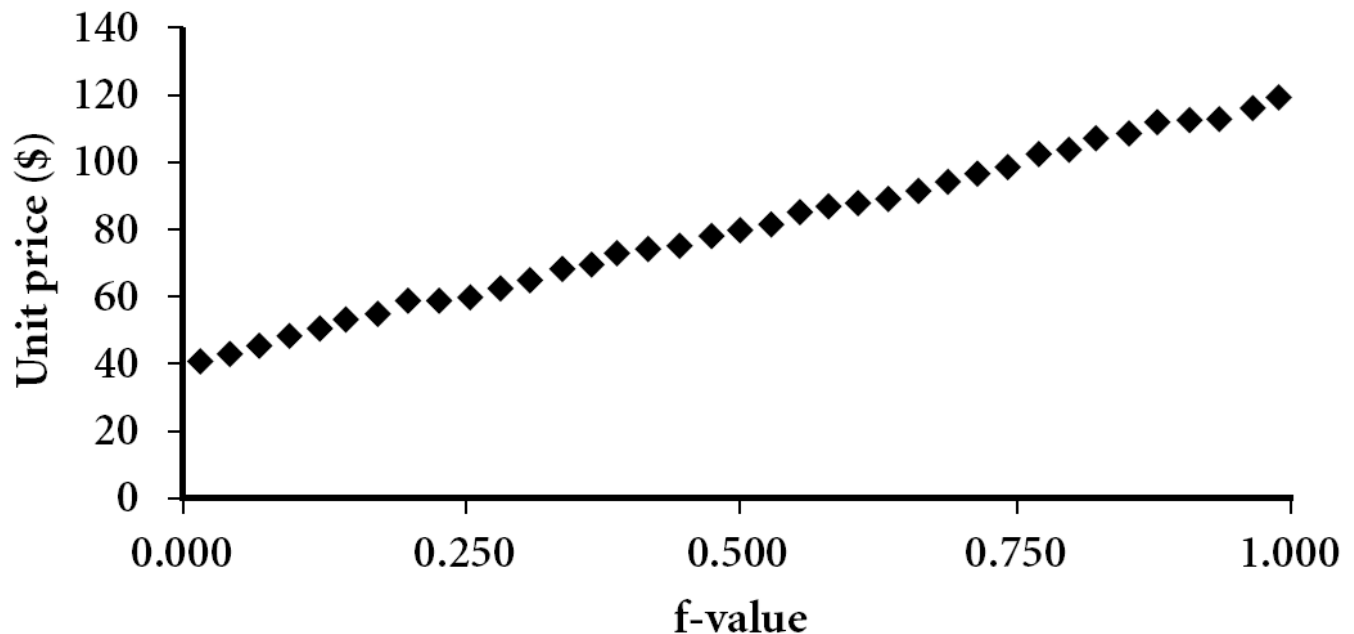


# Quantile Plot

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  $x_i$



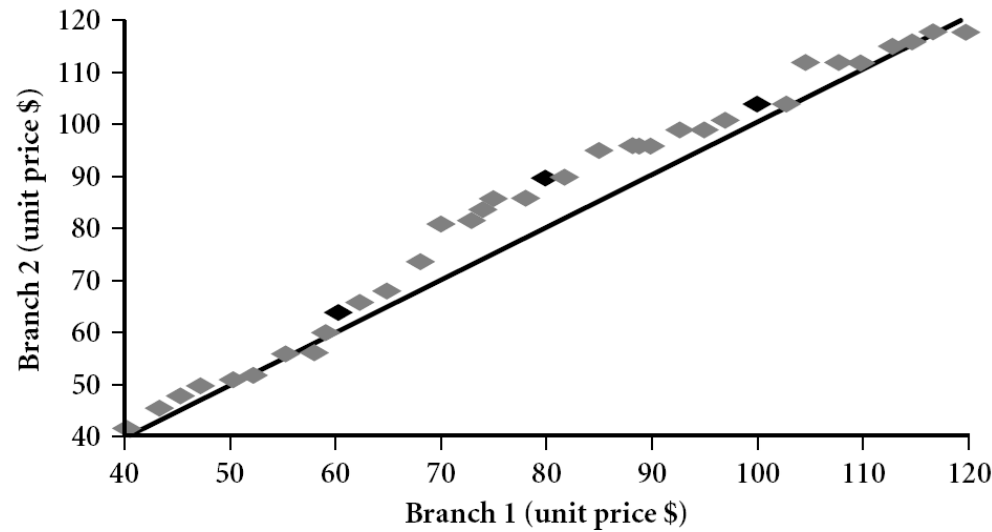


# Quantile-Quantile (Q-Q) Plot greatlearning

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?

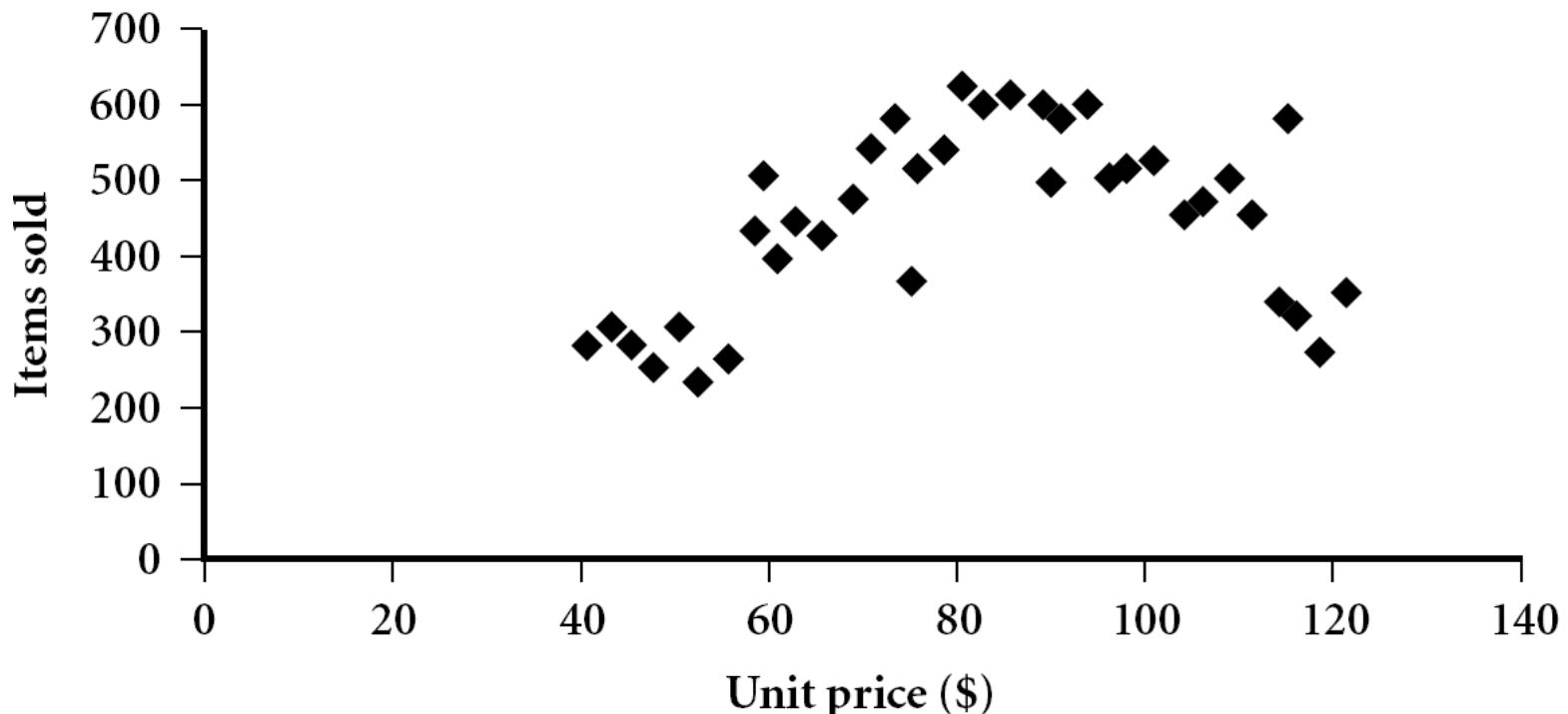
Example shows unit price of items sold at Branch 1 vs. Branch 2 for each quantile. Unit prices of items sold at Branch 1 tend to be lower than those at Branch 2.



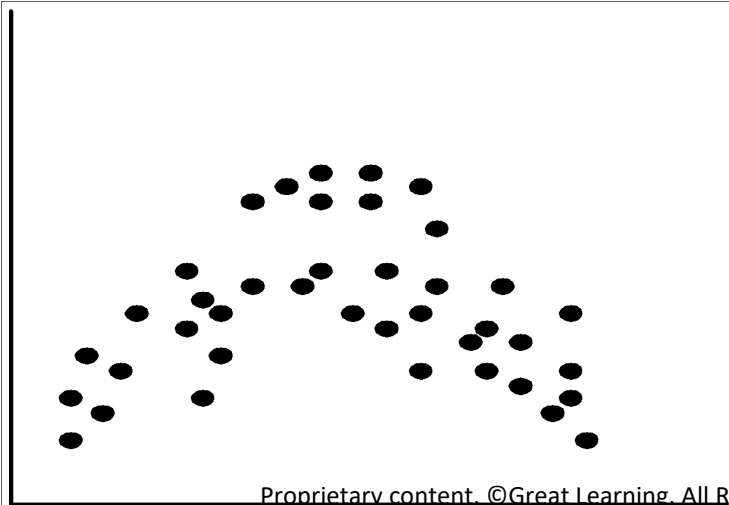
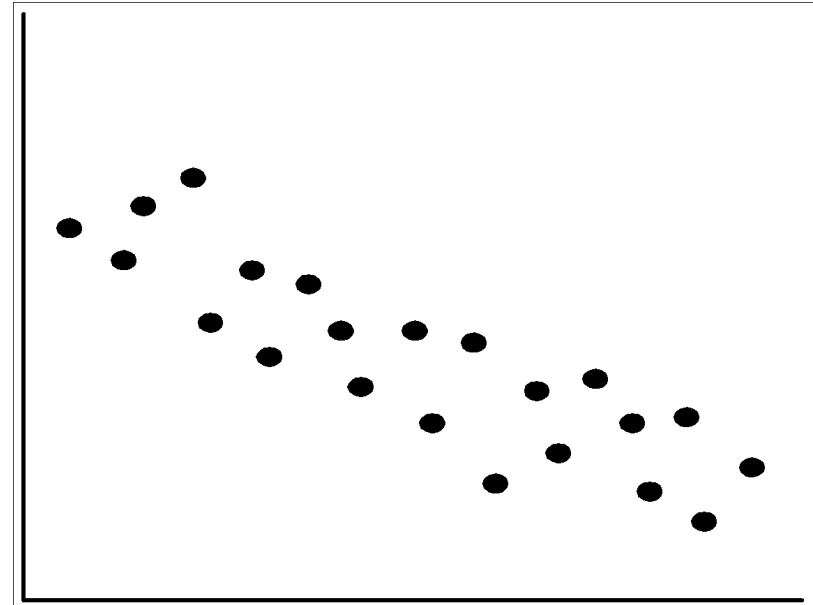
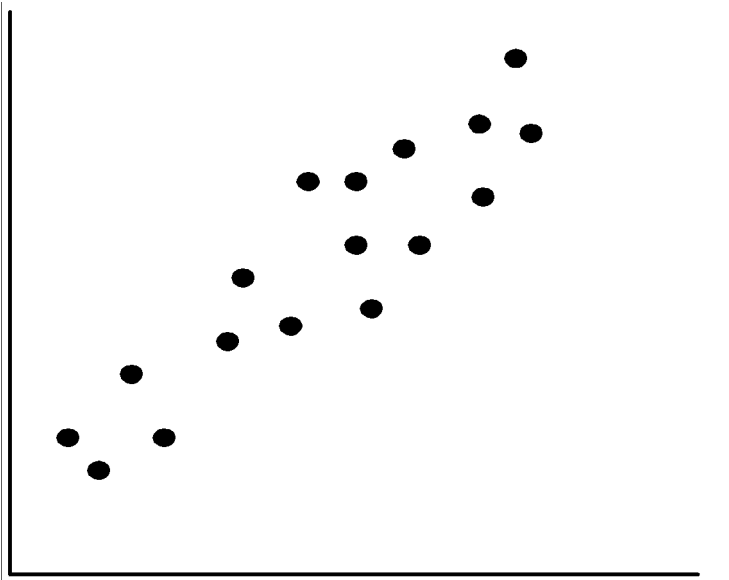
# Scatter plot

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



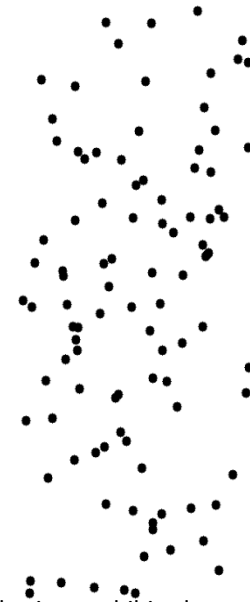
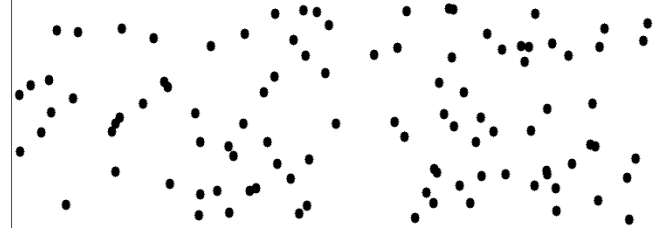
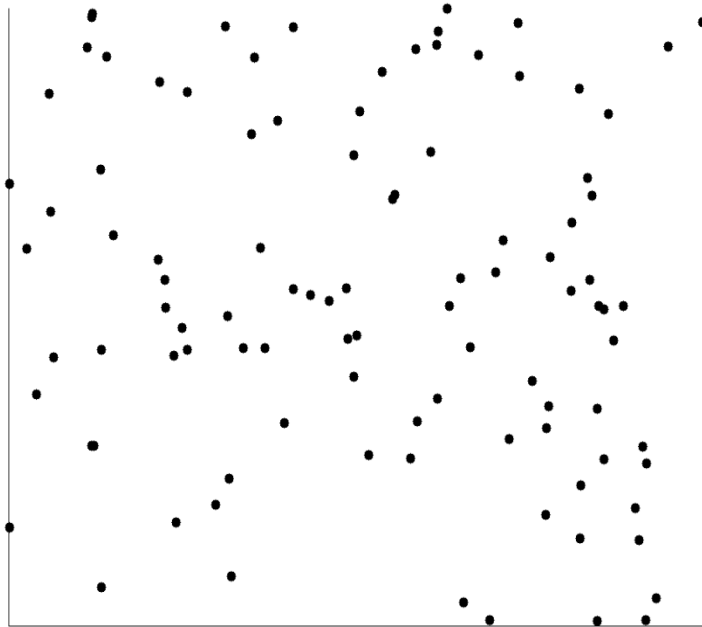
# Positively and Negatively Correlated Data



The left half fragment is positively correlated

The right half is negative correlated

# Uncorrelated Data



# 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

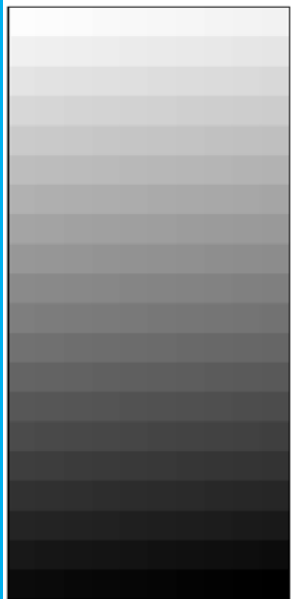
- Visualizing complex data and relations

# Pixel-Oriented Visualization Techniques

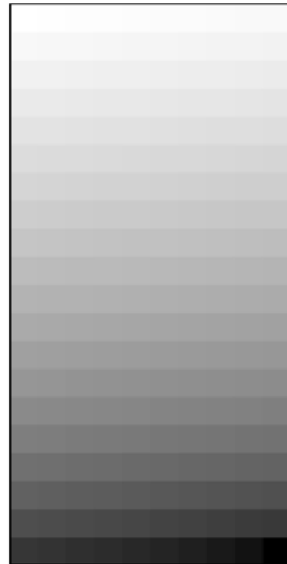
For a data set of  $m$  dimensions, create  $m$  windows on the screen, one for each dimension

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



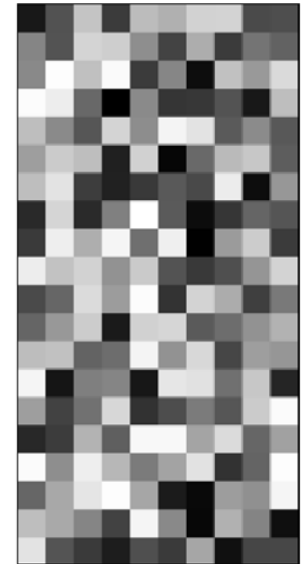
(a) Income



(b) Credit Limit



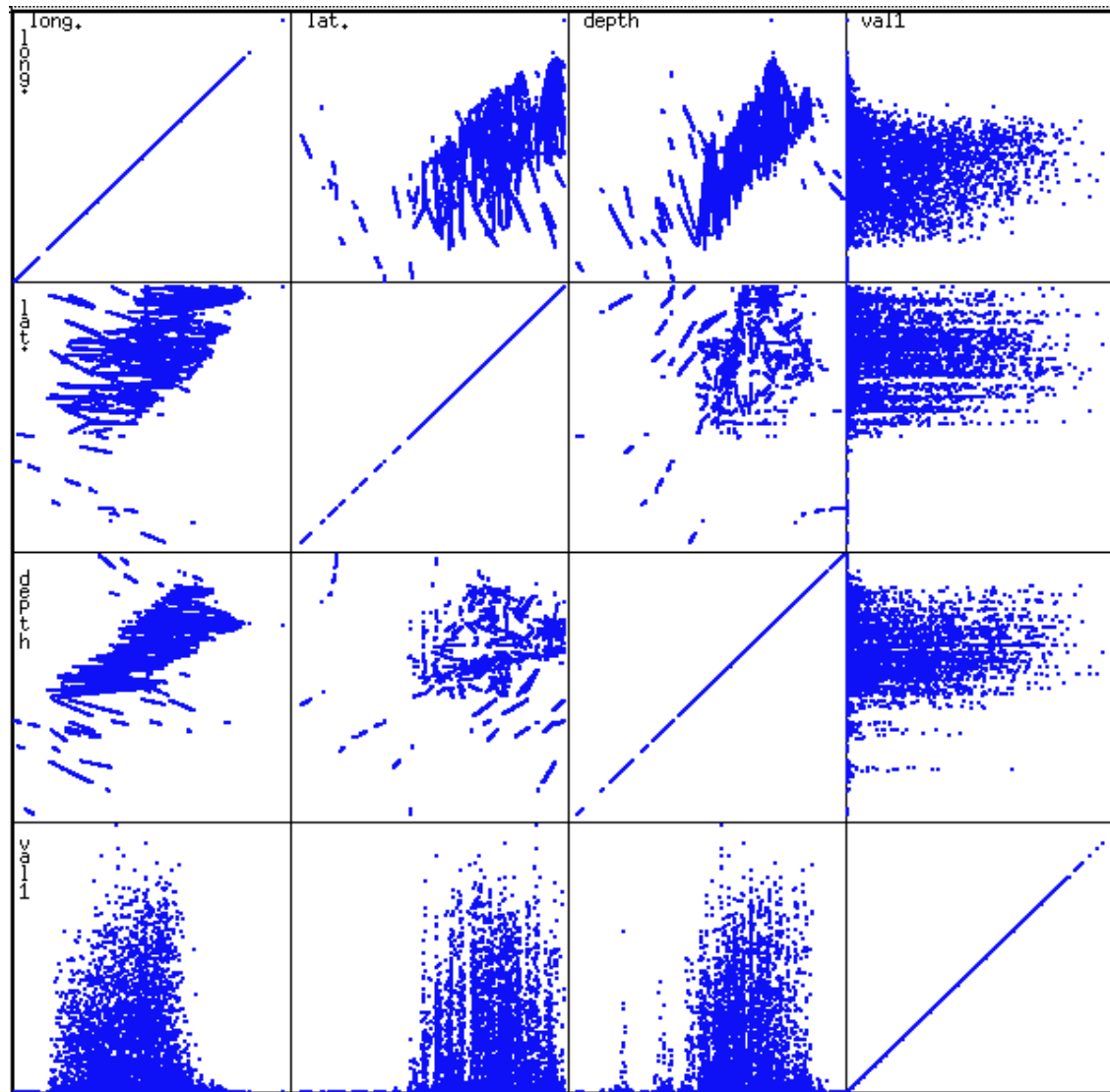
(c) transaction volume



(d) age

# Scatterplot Matrices

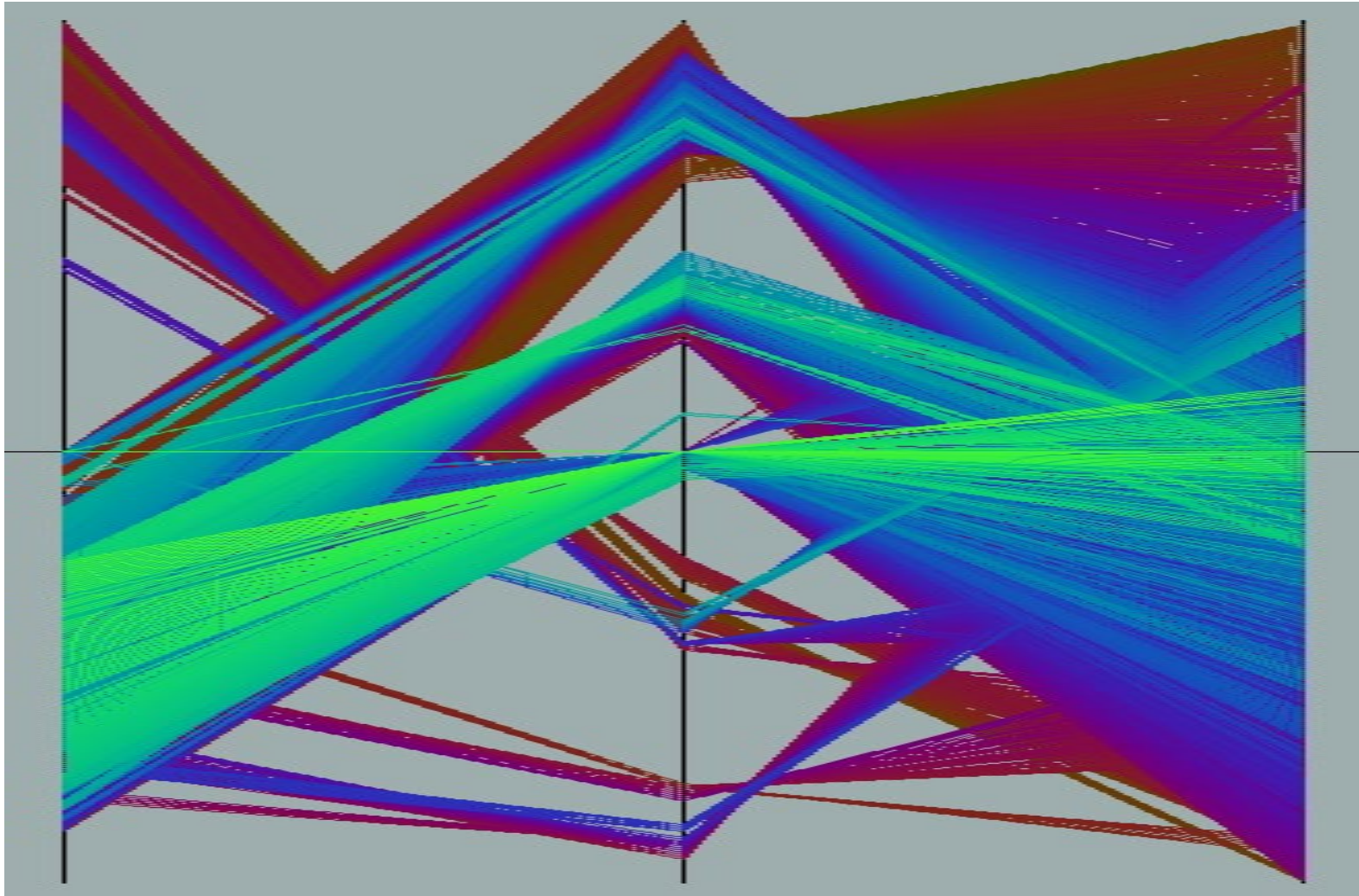
Used by permission of M. Ward, Worcester Polytechnic Institute



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

Proprietary content. ©Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited

# Parallel Coordinates of a Data Set



Proprietary content. ©Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited



# Correlation Analysis (Nominal Data): Chi-Square Test

	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

$$e_{ij} = \frac{\text{count}(A = a_i) \times \text{count}(B = b_j)}{n}$$

$$e_{11} = \frac{\text{count}(\text{male}) \times \text{count}(\text{fiction})}{n} = \frac{300 \times 450}{1500} = 90$$

For this  $2 \times 2$  table, the degrees of freedom are  $(2 - 1)(2 - 1) = 1$ . For 1 degree of freedom, the  $\chi^2$  value needed to reject the hypothesis at the 0.001 significance level is 10.828

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

$$\chi^2 = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}}$$

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

507.93 > 10.828 shows that like\_science\_fiction and play\_chess are correlated in the group

# Correlation Analysis (Numeric Data)

Correlation coefficient (also called **Pearson's product moment coefficient**)

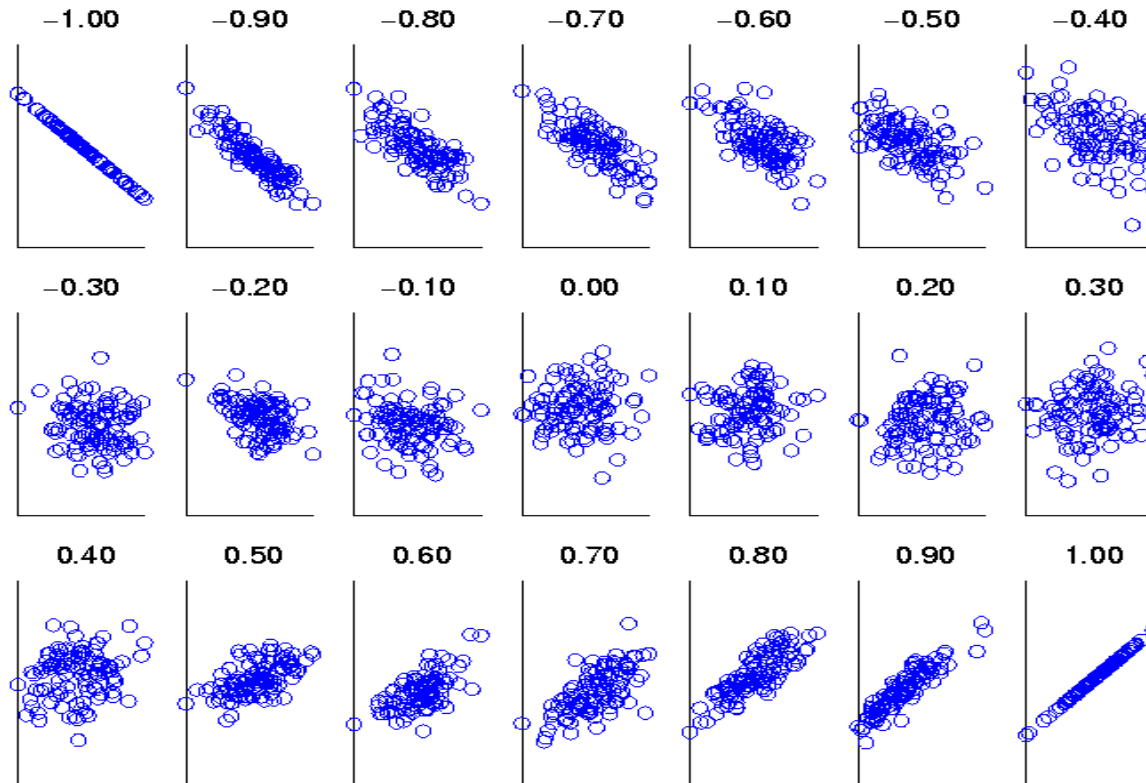
$$r_{A,B} = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{(n-1)\sigma_A\sigma_B} = \frac{\sum_{i=1}^n (a_i b_i) - n\bar{A}\bar{B}}{(n-1)\sigma_A\sigma_B}$$

where  $n$  is the number of tuples,  $\bar{A}$  and  $\bar{B}$  are the respective means of  $A$  and  $B$ ,  $\sigma_A$  and  $\sigma_B$  are the respective standard deviation of  $A$  and  $B$ , and  $\sum(a_i b_i)$  is the sum of the  $AB$  cross-product.

If  $r_{A,B} > 0$ ,  $A$  and  $B$  are positively correlated ( $A$ 's values increase as  $B$ 's). The higher, the stronger correlation.

$r_{A,B} = 0$ : independent;  $r_{AB} < 0$ : negatively correlated

# Visually Evaluating Correlation



Scatter plots showing the similarity from -1 to 1.

# Summary

- Histograms
- Measures of central tendency: mean, mode, median
- Measures of dispersion: range, IQR, variance, std deviation, coefficient of variation.
- Normal distribution, Chebyshev Rule.
- Five number summary, boxplots, QQ plots, Quantile plot, scatter plot.
- Visualization: scatter plot matrix, parallel coordinates.
- Correlation analysis.

- 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, 2<sup>nd</sup> ed., Graphics Press, 2001
- C. Yu et al., Visual data mining of multimedia data for social and behavioral studies, Information Visualization, 8(1), 2009