



DATA ANALYTICS

Unit 1:Data Exploration

Mamatha.H.R and Bharathi .R

Department of Computer Science and
Engineering

DATA ANALYTICS

Unit 1:Data Exploration

Mamatha H R

Department of Computer Science and Engineering

What is data exploration?

A preliminary exploration of the data to better understand its characteristics.

- Key motivations of data exploration include
 - Helping to select the right tool for preprocessing or analysis
 - Making use of humans' abilities to recognize patterns
 - People can recognize patterns not captured by data analysis tools

- In EDA, as originally defined by Tukey
 - The focus was on visualization
 - Clustering and anomaly detection were viewed as exploratory techniques
- In our discussion of data exploration, we focus on
 - Summary statistics
 - Visualization

- Summary statistics are numbers that summarize properties of the data
 - Summarized properties include frequency, location and spread
 - Examples: location - mean
spread - standard deviation
- Most summary statistics can be calculated in a single pass through the data

- **Cross-Sectional Data:** A data collected on many variables of interest at the same time or duration of time is called cross-sectional data.
- **Time Series Data:** A data collected for a single variable such as demand for smartphones collected over several time intervals (weekly, monthly, etc.) is called a time series data.
- **Panel Data:** Data collected on several variables (multiple dimensions) over several time intervals is called panel data (also known as longitudinal data).

DATA ANALYTICS

Data Type: Cross-Sectional Data

Cross-sectional data are usually data gathered from approximately the same period of time from a population.

Example: Responses from a questionnaire concerning the president's environmental policies



	A	B	C	D	E	F	G
1	Person	Age	Gender	State	Children	Salary	Opinion
2	1	35	Male	Minnesota	1	\$65,400	5
3	2	61	Female	Texas	2	\$62,000	1
4	3	35	Male	Ohio	0	\$63,200	3
5	4	37	Male	Florida	2	\$52,000	5
6	5	32	Female	California	3	\$81,400	1
7	6	33	Female	New York	3	\$46,300	5
28	27	27	Male	Illinois	3	\$45,400	2
29	28	63	Male	Michigan	2	\$53,900	1
30	29	52	Male	California	1	\$44,100	3
31	30	48	Female	New York	2	\$31,000	4

DATA ANALYTICS

Data Type: Time Series Data

- Time series data are data collected over time.

	A	B
1	Quarter	Revenue
2	Q1-2010	1026
3	Q2-2010	1056
4	Q3-2010	1182
5	Q4-2010	2861
6	Q1-2011	1172
7	Q2-2011	1249
8	Q3-2011	1346
9	Q4-2011	3402
10	Q1-2012	1286
11	Q2-2012	1317
12	Q3-2012	1449
13	Q4-2012	3893
14	Q1-2013	1462
15	Q2-2013	1452
16	Q3-2013	1631
17	Q4-2013	4200

- **Panel Data:** Data collected on several variables (multiple dimensions) over several time intervals is called panel data (also known as longitudinal data).

	F	G	H	I	J	K	L
	All the data in %						
Country	Year	Unempl	Remit	FDI	GrosCapForm	DomCredProvF inSec	
Armenia	2004	11.216	22.01857	6.909902	24.87835868	6.684198722	
Azerbaijan	2005	8	2.345269	54.36527	57.99045755	10.92778878	
Belarus	2006	9.272	1.034378	0.707732	28.66730144	21.21007845	
Estonia	2007	10.248	1.383825	9.009741	34.53914678	60.38037809	
Georgia	2008	12.62	6.996385	9.613621	31.90818318	18.94640069	
Kazakhstan	2009	8.4	0.132967	13.01286	26.31108968	29.05023214	
Kyrgyz Republic	2010	8.53	8.096525	7.933804	14.48838801	8.38417521	
Latvia	2011	11.708	1.519205	4.111913	33.03493107		
Lithuania	2012	10.684	2.547673	3.51495	22.68148798		
Moldova	2013	8.17	26.99413	5.81203	26.35933373	32.00577996	
Tajikistan	2014	13.412	12.13796	13.10239	10.37748137	6.923721381	
Ukraine	2015	8.59	2.889554	2.6458	21.13423719	31.68865316	
Uzbekistan	2016	8.058		1.467994	20.7		

A sample windowing and cross-sectional data extraction from the time series dataset.



Date	inputYt	Window id	inputYt + 1 (horizon)	inputYt - 5	inputYt - 4	inputYt - 3	inputYt - 2	inputYt - 1	inputYt - 0
Jan 1, 2009	0.709	0	1.169	0.709	1.886	1.293	0.822	-0.173	0.552
Feb 1, 2009	1.886	1	1.604	1.886	1.293	0.822	-0.173	0.552	1.169
Mar 1, 2009	1.293	2	0.949	1.293	0.822	-0.173	0.552	1.169	1.604
Apr 1, 2009	0.822	3	0.080	0.822	-0.173	0.552	1.169	1.604	0.949
May 1, 2009	-0.173	4	-0.040	-0.173	0.552	1.169	1.604	0.949	0.080
Jun 1, 2009	0.552	5	1.381	0.552	1.169	1.604	0.949	0.080	-0.040
Jul 1, 2009	1.169	6	0.761	1.169	1.604	0.949	0.080	-0.040	1.381
Aug 1, 2009	1.604	7	2.312	1.604	0.949	0.080	-0.040	1.381	0.761
Sep 1, 2009	0.949	8	1.795	0.949	0.080	-0.040	1.381	0.761	2.312
Oct 1, 2009	0.080	9	0.586	0.080	-0.040	1.381	0.761	2.312	1.795
Nov 1, 2009	-0.040	10	-0.077	-0.040	1.381	0.761	2.312	1.795	0.586
Dec 1, 2009	1.381	11	0.613	1.381	0.761	2.312	1.795	0.586	-0.077
Jan 1, 2010	0.761	12	1.845	0.761	2.312	1.795	0.586	-0.077	0.613
Feb 1, 2010	2.312	13	1.984	2.312	1.795	0.586	-0.077	0.613	1.845
Mar 1, 2010	1.795	14	1.861	1.795	0.586	-0.077	0.613	1.845	1.984
Apr 1, 2010	0.586	15	0.661	0.586	-0.077	0.613	1.845	1.984	1.861

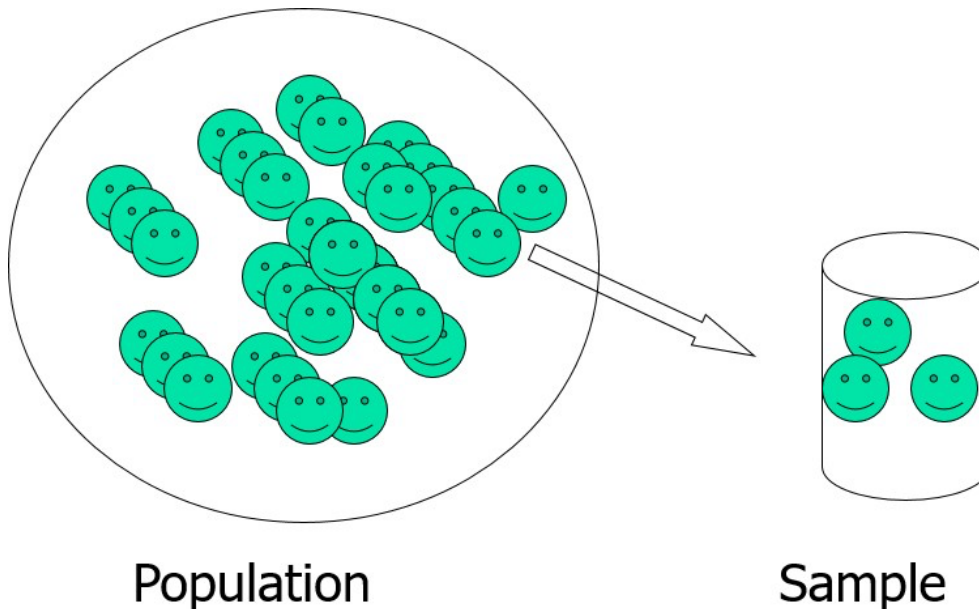
(A) Time series data set

(B) Cross-sectional data set

- **Nominal scale** refers to variables that are basically names (qualitative data) and also known as categorical variables.
- **Ordinal scale** is a variable in which the value of the data is captured from an ordered set, which is recorded in the order of magnitude.
- **Interval scale** corresponds to a variable in which the value is chosen from an interval set. Variable such as temperature measured in centigrade) or intelligence quotient (IQ) score are examples of interval scale
- Any variable for which the ratios can be computed and are meaningful is called **ratio scale**.

Population And Sample

- ▶ **Population** is the set of all possible observations (often called cases, records, subjects or data points) for a given context of the problem.
- ▶ **Sample** is the subset taken from a population.



An Illustration: Which Group is Smarter?

Class A--IQs of 13 Students

102	115
128	109
131	89
98	106
140	119
93	97
110	

Class B--IQs of 13 Students

127	162
131	103
96	111
80	109
93	87
120	105
109	

Each individual may be different. If you try to understand a group by remembering the qualities of each member, you become overwhelmed and fail to understand the group.

DATA ANALYTICS

Descriptive Statistics



Which group is smarter now?

Class A--Average IQ

110.54

Class B--Average IQ

110.23

They're roughly the same!

With a summary descriptive statistic, it is much easier to answer our question.

Types of descriptive statistics:

- Organize Data

 - Tables

 - Graphs

- Summarize Data

 - Central Tendency

 - Variation

Types of descriptive statistics:

Organize Data

Tables

- Frequency Distributions
- Relative Frequency Distributions

Graphs

- Bar Chart or Histogram
- Stem and Leaf Plot
- Frequency Polygon

Summarizing Data:

- Central Tendency (or Groups' "Middle Values")
 - Mean
 - Median
 - Mode

- Variation (or Summary of Differences Within Groups)
 - Range
 - Interquartile Range
 - Variance
 - Standard Deviation

- Mean (or Average) Value

Mean is the arithmetical average value of the data and is one of the most frequently used measures of central tendency.

$$\text{Mean} = \bar{x} = \frac{x_1 + x_2 + \cdots + x_n}{n} = \sum_{i=1}^n \frac{x_i}{n}$$

Symbol \bar{x} is frequently used to represent the estimated value of the mean from a sample. If the entire population is available and if we calculate mean based on the entire population, then we have the population mean which is denoted by μ (population mean).

Property of Mean

An important property of mean is that the summation of deviation of observations from the mean is zero, that is

$$\sum_{i=1}^n \left(X_i - \bar{X} \right) = 0$$

- Median is the value that divides the data into two equal parts, that is, the proportion of observations below median and above median will be 50%.
- Easiest way to find the median value is by arranging the data in the increasing order and the median is the value at position $(n + 1)/2$ when n is odd. When n is even, the median is the average value of $(n/2)^{\text{th}}$ and $(n + 2)/2^{\text{th}}$ observation after arranging the data in the increasing order.

- **Mode** is the most frequently occurring value in the dataset
- Mode is the only measure of central tendency which is valid for qualitative (nominal) data since the mean and median for nominal data are meaningless.
- For example, assume that a customer data with a retailer has the marital status of customer, namely, (a) Married, (b) Unmarried, (c) Divorced Male, and (d) Divorced Female. Mean and median are meaningless when we try to use them on a qualitative data such as marital status. On the other hand, mode will capture the customer type in terms of marital status that occurs most frequently in the database

- **Percentile**, decile and quartile are frequently used to identify the position of the observation in the dataset.
- Percentile, denoted as P_x , is the value of the data at which x percentage of the data lie below that value

Position corresponding to $P_x \approx x(n+1)/100$

- P_x is the position in the data calculated, where n is the number of observations in the data.

- ▶ **Decile** corresponds to special values of percentile that divide the data into 10 equal parts. First decile contains first 10% of the data and second decile contains first 20% of the data and so on.
- ▶ **Quartile** divides the data into 4 equal parts. The first quartile (Q_1) contains first 25% of the data, Q_2 contains 50% of the data and is also the median. Quartile 3 (Q_3) accounts for 75% of the data

DATA ANALYTICS

Example

Time between failures (in hours) of a wire cutter used in a cookie manufacturing oven is given in table below. The function of the wire-cut is to cut the dough into

2	22	32	39	46	56	76	79	88	93
3	24	33	44	46	66	77	79	89	99
5	24	34	45	47	67	77	86	89	99
9	26	37	45	55	67	78	86	89	99
21	31	39	46	56	75	78	87	90	102

Time between failures of wire-cut (in hours)

2	22	32	39	46	56	76	79	88	93
3	24	33	44	46	66	77	79	89	99
5	24	34	45	47	67	77	86	89	99
9	26	37	45	55	67	78	86	89	99
21	31	39	46	56	75	78	87	90	102

1. Calculate the mean, median, and mode of time between failures of wire-cuts
2. The company would like to know by what time 10% (ten percentile or P_{10}) and 90% (ninety percentile or P_{90}) of the wire-cuts will fail?
3. Calculate the values of P_{25} and P_{75} .

- a) Mean = 57.64, median = 56, and mode = 46
- a) Note that the data in Table is arranged in increasing order in columns. The position of $P_{10} = 10 \times (51)/100 = 5.1$. We can round off 5.1 to its nearest integer which is 5. The corresponding value from table is 21 (10 percentage of observations in Table have a value of less than or equal to 21). That is, by 21 hours, 10% of the wire-cuts will fail. In asset management (and reliability theory), this value is called P_{10} life.

Instead of rounding the value obtained from Eq, we can use the following approximation:

$$P_{10} = 10 \times (51)/100 = 5.1$$

Value at 5th position is 21. Value at position 5.1 is approximated as $21 + 0.1 \times (\text{value at 6}^{\text{th}} \text{ position} - \text{value at 5}^{\text{th}} \text{ position}) = 21 + 0.1(1) = 21.1$

$$P_{90} = 90 \times 51/100 = 45.9$$

The value at position 45 is 90 and at position 45.9 is

$$90 + 0.9 \times (3) = 92.7$$

That is, **90%** of the wire-cuts will fail by **92.7 hours**

P_{25} (1st Quartile or Q_1) = $25 \times 51/100 = 12.75$,
Value at 12th position is 33, so

$P_{25} = 33 + 0.75$ (value at 13th position – value at 12th position) = $33 + 0.75 (1) = 33.75$

P_{75} (3rd Quartile or Q_3) = $75 \times 51/100 = 38.25$

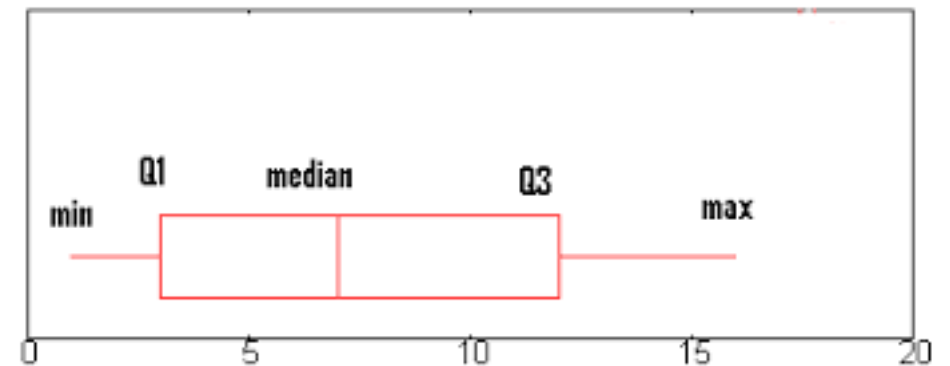
Value at 38th position is 86, so

$P_{75} = 86 + 0.25$ (value at 39th position – value at 38th position) = $86 + 0.25 (0) = 86$

- Predictive analytics techniques such as regression attempt to explain variation in the outcome variable (Y) using predictor variables (X)
- Variability in the data is measured using the following measures:
 - Range
 - Inter-Quartile Distance (IQD)
 - Variance
 - Standard Deviation

- **Range** is the difference between maximum and minimum value of the data. It captures the data spread.
- **Inter-quartile distance** (IQD), also called inter-quartile range (IQR) is a measure of the distance between Quartile 1 (Q_1) and Quartile 3 (Q_3)
- **Variance** is a measure of variability in the data from the mean value. Variance for population, σ^2 , is calculated using

$$\text{Variance} = \sigma^2 = \sum_{i=1}^n \frac{(X_i - \mu)^2}{n}$$



- In case of a sample, the Sample Variance (S^2) is calculated using

$$S^2 = \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{n - 1}$$

While calculating sample variance S^2 , the sum of squared

deviation $\sum_{i=1}^n (X_i - \bar{X})^2$ is divided by $(n-1)$, this is known as **Bessel's correction**.

- The population standard deviation (σ) and sample standard deviation (S) are given by

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

$$S = \sqrt{\sum_{i=1}^n \frac{(X_i - \bar{X})^2}{n-1}}$$

- **Degrees of freedom** is equal to the number of independent variables in the model (Trochim, 2005). For example, we can create any sample of size n with mean value of \bar{x} by randomly selecting $(n - 1)$ values. We need to fix just one out of n values. Thus the number of independent variables in this case is $(n - 1)$
- Degrees of freedom is defined as the difference between the number of observations in the sample and number of parameters estimated (Walker 1940, Toothaker and Miller, 1996). If there are n observations in the sample and k parameters are estimated from the sample, then the degrees of freedom is $(n - k)$.

- **Chebyshev's theorem** (also known as Chebyshev's inequality) is an empirical rule that allows us to predict proportion of observations that is likely to lie between an interval defined using mean and standard deviation. Probability of finding a randomly selected value in an interval defined by $\mu \pm k\sigma$ is $1 - \frac{1}{k^2}$ that is

$$P(\mu - k\sigma \leq X \leq \mu + k\sigma) \geq 1 - \frac{1}{k^2}$$

Example

- Amount spent per month by a segment of credit card users of a bank has a mean value of 12000 and standard deviation of 2000. Calculate the proportion of customers who are spending between 8000 and 16000?

- **Solution:**

$$P(8000 \leq X \leq 16000) = P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) \geq 1 - \frac{1}{2^2} = 0.75$$

That is, the proportion of customers spending between 8000 and 16000 is at least 0.75 (or 75%)

- **Skewness** is a measure of symmetry or lack of symmetry. A dataset is symmetrical when the proportion of data at equal distance (measured in terms of standard deviation) from mean (or median) is equal. That is, the proportion of data between μ and $\mu - k\sigma$ is same as μ and $\mu + k\sigma$, where k is some positive constant.
- **Pearson's moment coefficient of skewness** for a dataset with n observations is given by

$$g_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})^3 / n}{\sigma^3}$$

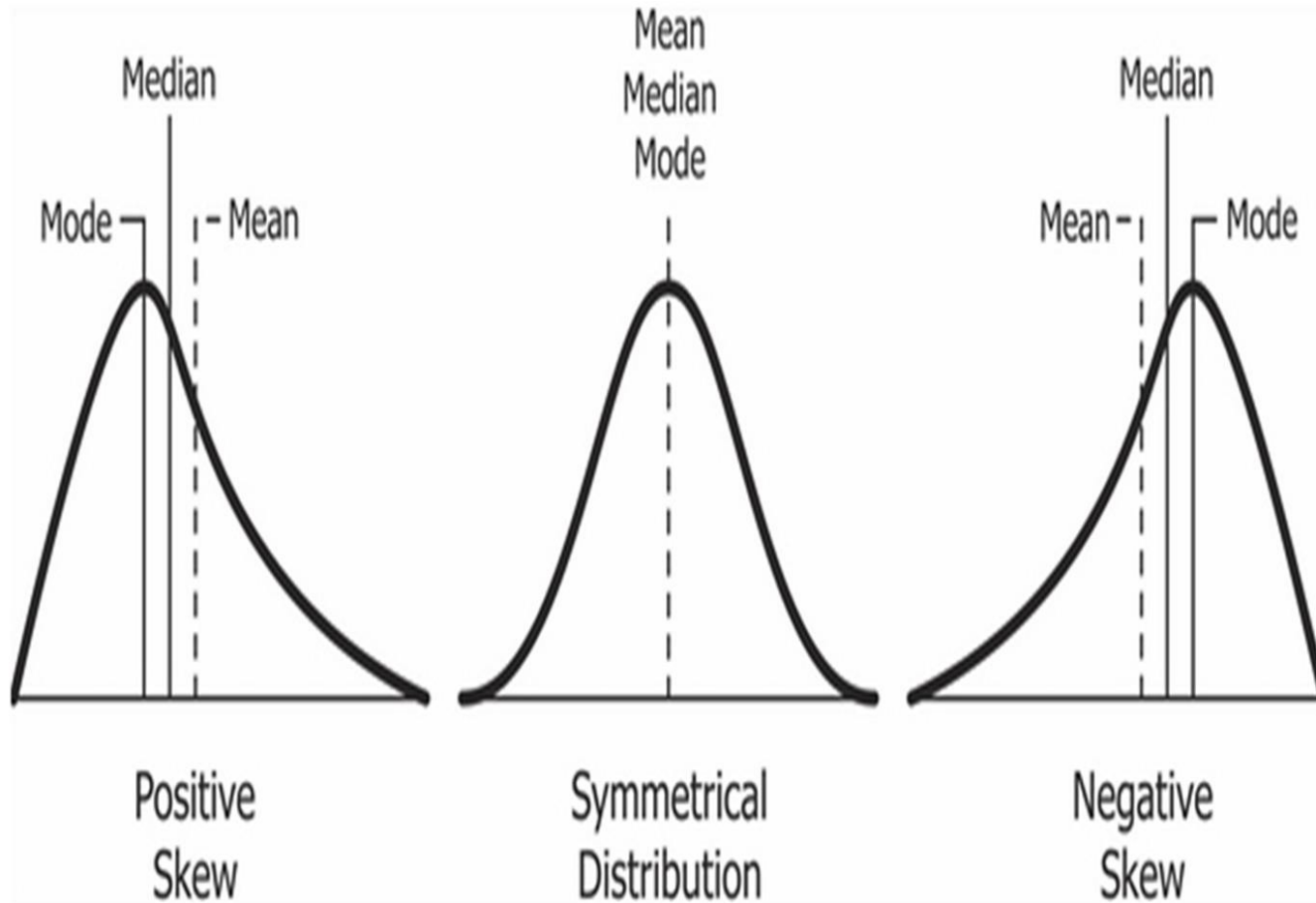
- The value of g_1 will be close to 0 when the data is symmetrical. A positive value of g_1 indicates a positive skewness and a negative value indicates **negative skewness**.

- The following formula is used usually for a sample with n observations (Joanes and Gill, 1998):

$$G_1 = \frac{\sqrt{n(n-1)}}{n-2} g_1$$

- The value of $\frac{\sqrt{n(n-1)}}{n-2}$ will converge to 1 as the value of n increases.

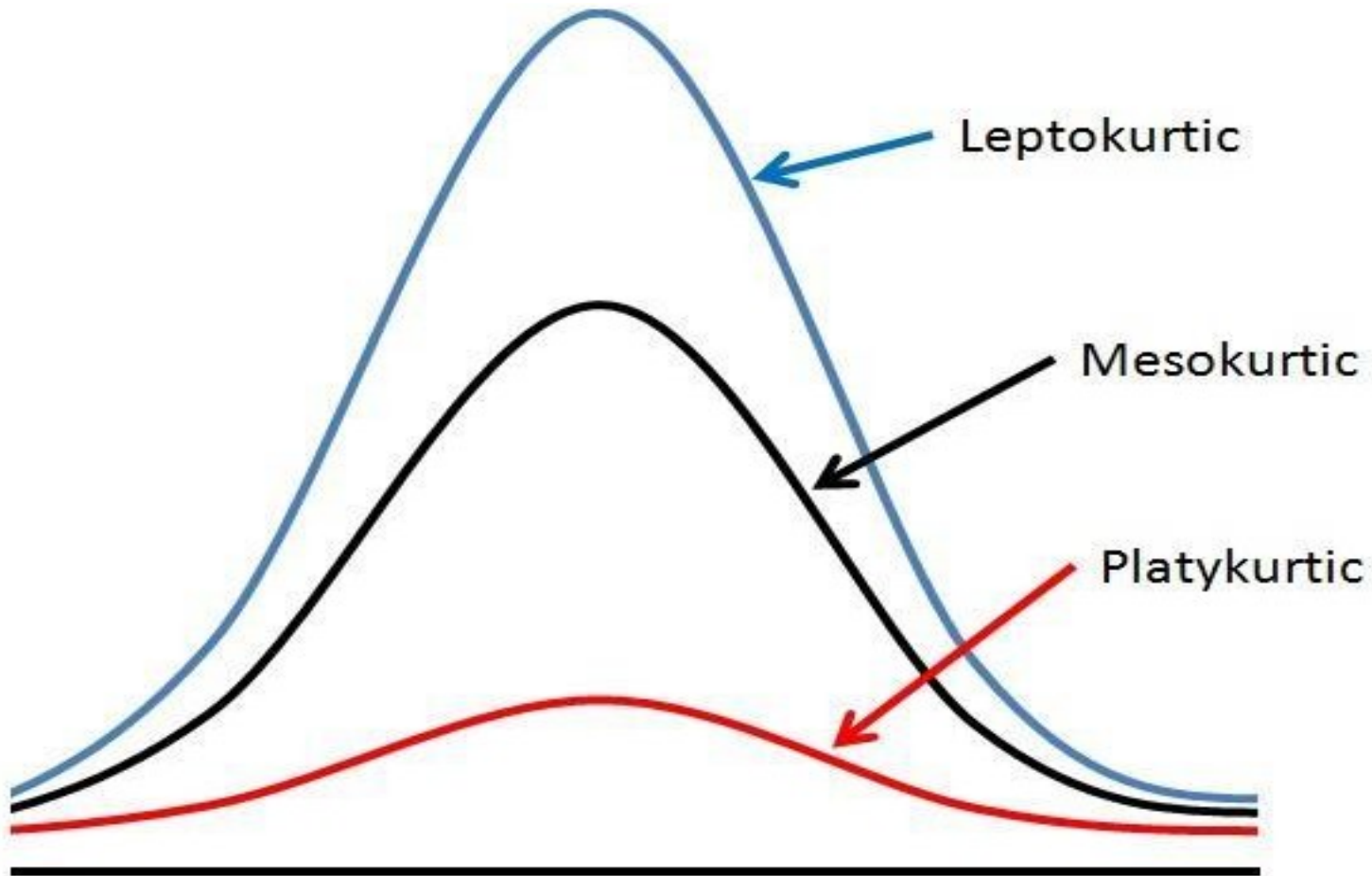
- Median, mean and mode of symmetric, positively and negatively skewed data



- **Kurtosis** is another measure of shape, aimed at shape of the tail, that is, whether the tail of the data distribution is heavy or light. Kurtosis is measured using the following equation:

$$\text{Kurtosis} = \frac{\sum_{i=1}^n \left(X_i - \bar{X} \right)^4 / n}{\sigma^4}$$

- Kurtosis value of less than 3 is called **platykurtic distribution** and greater than 3 is called **leptokurtic distribution**. The kurtosis value of 3 indicates standard normal distribution (also called **mesokurtic**)



- The excess kurtosis is a measure that captures deviation from kurtosis of a normal distribution and is given by:

$$\text{Excess Kurtosis} = \frac{\sum_{i=1}^n \left(X_i - \bar{X} \right)^4 / n}{\sigma^4} - 3$$

The daily football at a retail store in Bangalore over the last 30 days is shown in Table 1. calculate the Mean, Median, Mode and Standard Deviation.

Table 1. Footfall data

232	277	261	173	283	197	251	212	213	213
229	164	219	196	186	247	244	269	216	272
252	314	161	165	221	260	219	290	225	251

For the data in Table 1, calculate the skewness and kurtosis. what can you infer from the skewness and kurtosis of the football data?

For the data in Table 1, calculate the values of first quartile and third quartile. Are there any outliers in the data?

References

Text Book:

- [“Business Analytics, The Science of Data-Driven Decision Making”](#), U. Dinesh Kumar, Wiley 2017
- [Data Mining: Concepts and Techniques](#) by Jiawei Han, Micheline Kamber and Jian Pei, The Morgan Kaufmann Series in Data Management Systems, 3rd Edition.
- [Introduction to Data Mining](#) , Tan, Steinbach, Kumar, 2nd Edition



THANK YOU

Dr.Mamatha H R

Professor,Department of Computer Science

mamathahr@pes.edu

+91 80 2672 1983 Extn 834