

Ex:- Suppose we have the following values for salary (in Thousands of dollars), shown in increasing order:

30, 36, 47, 50, 52, 56, 60, 63, 70, 70, 110, we have

$$\begin{aligned}\text{Mean } \bar{x} &= 30 + 36 + 47 + 50 + 52 + 56 + 60 \\ &\quad + 63 + 70 + 70 \\ &\quad + 110 / 12 \\ &= \frac{656}{12} = 58\end{aligned}$$

Thus, The Mean salary is \$58,000.

\* Suppose that we had only the first 11 value in The list. Given an odd number of values, The Median is The middlemost value. This is The sixth value in This list, which has a value of 52,000\$

\* There is an even number of observations (i.e. 12); Therefore, The Median is not unique; it can be any value within The Two middlemost value of 52 and 56 (within The sixth and seventh value in the list).

By convention, we assign The average of The Two middlemost values as The Median; i.e.  $\frac{52+56}{2} = \frac{108}{2} = 54$   
Thus The Median is 54,000\$

\* Mode can be determined for qualitative and quantitative attributes. It is possible for greatest frequency to correspond to several different values, which results in more than one mode. Data sets with one, two, or three modes are respectively called unimodal, bimodal, and trimodal. In general, data set with two or more modes is multimodal. At the other extreme, if each data value occurs only once, then there is no mode.

\* The data in this example are bimodal. The two modes are \$52,000 and \$70,000.

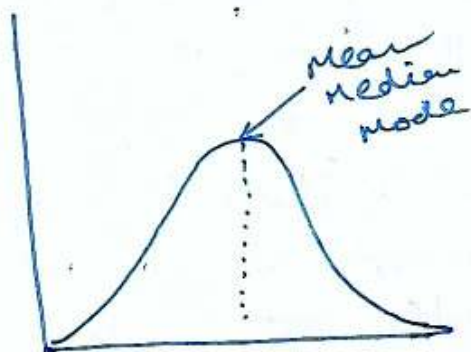
\* Midrange of the data of same example is

$$\frac{30,000 + 10,000}{2}$$

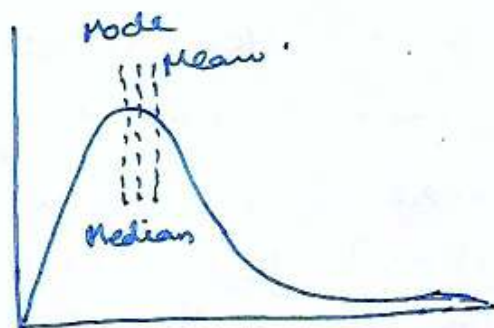
$$= \$70,000$$

In a unimodal frequency curve with perfect symmetric data distribution, the mean, median, and mode are all at the same center value, as shown in fig (a)

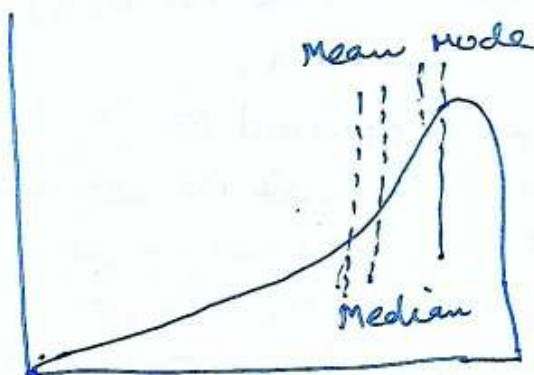
Data in most real applications are not symmetric. They may instead be either positively skewed, where the mode occurs at a value that is ~~also~~ smaller than the median, or negatively skewed, where the mode occurs at a value greater than the median.



(a) Symmetric data



(b) Positively skewed data



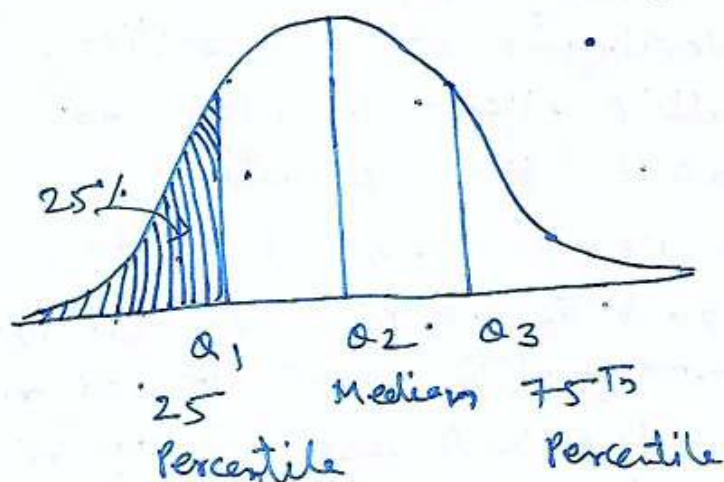
(c) Negatively skewed data



\* Suppose that the data for attribute  $x$  are sorted in increasing numeric order. Imagine that we can pick certain data points so as to split the data distribution into equal-size consecutive sets, as in fig. These data points are called quantiles. Quantiles are points taken at regular intervals of a data distribution, dividing it into essentially equal size consecutive sets.

The 2-quantile is the data point dividing the lower and upper halves of the data distribution. It corresponds to the median. The 4-quantiles are the three data points that split the data distribution into four equal parts; each part represents one-fourth of the data distribution. They are more commonly referred to as quartiles.

The 100-quantiles are most commonly referred to as Percentiles; they divide the data distribution into 100 equal-sized consecutive sets. The median, quartiles and Percentiles are the most widely used forms of quantiles.



Quartiles give an indication of a distribution's center, spread and ~~center~~ shape. The first quartile, denoted by  $Q_1$ , is the 25th percentile: it cuts off the lowest 25% of the data. The third quartile, denoted by  $Q_3$ , is the 75th percentile - it cuts off the lowest 75% (or highest 25%) of the data. The second quartile is the 50th percentile. As the median, it



\* The distance between the first and third quartiles is a simple measure of spread that gives the range covered by the middle half of the data. This distance is called the IQR and is defined as

$$IQR = Q_3 - Q_1$$

\* is our data values;  $Q_1 = \$47000$  and  $Q_3 = \$63000$   
Thus the interquartiles range is  $IQR = 63 - 47 = 16000 \$$   
and Median is  $52000 \$$ .

Five number summary, Box Plots, and outliers: →

In the symmetric distribution, the median (and other M.S of Central) splits the data into equal-size halves. This does not occur for skewed distributions.

Therefore, it is more informative to also provide the two quartiles  $Q_1$  and  $Q_3$ , along with median. A common rule of thumb for identifying suspected outliers is to single out values falling at least  $1.5 \times IQR$  above the third quartile or below the first quartile.

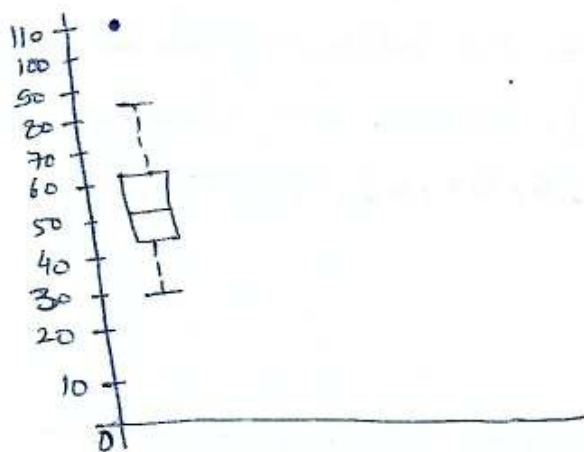
Because  $Q_1$ , the median ( $Q_2$ ), and  $Q_3$  together contain no information about the end points (or tails) of the data, a fuller summary of the shape of a distribution can be obtained by providing the lowest and highest data values as well. This is known as the five number summary. The five number summary of a distribution consists of the median ( $Q_2$ ), the quartiles  $Q_1$  and  $Q_3$ , and the smallest and largest individual observations, written in the order of minimum,  $Q_1$ , median,  $Q_3$ , maximum.

Box plots are a popular way of visualizing a distribution. A box plot incorporates the five number summary as follows:

- \* Typically the ends of the box are the quartiles so that the box length is the interquartile range.

- \* The Median is marked by a line within the box.

- \* Two lines (called whiskers) outside the box extend to the smallest (minimum) and largest (maximum) observations.



Variance and standard deviation: → Variance and

standard deviation are measures of data dispersion.

They indicate how spread out a data distribution is.

A low standard deviation means that the data observations tend to be very close to mean, while a high standard deviation indicates that the data are spread out over a large ~~number~~ range of values.

The variance of  $N$  observations,  $x_1, x_2, \dots, x_N$ , for a numeric attribute  $x$  is

$$\sigma^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 = \left( \frac{1}{N-1} \sum_{i=1}^N x_i^2 \right) - \bar{x}^2$$



The standard deviation,  $\sigma$ , of the observations is the square root of the variance,  $\sigma^2$ .

Properties of  $\sigma$  :-

$\sigma$  measures spread about the mean and should be considered only when the mean is chosen as the measure of center.

$\sigma = 0$  only when there is no spread, that is, when all observations have the same value, otherwise  $\sigma > 0$ .

Example:- Suppose we have the following values for salary (in thousands of dollars), shown in increasing order:  
30, 36, 47, 50, 52, 52, 56, 60, 63, 70, 70, 110

Sol:- we have  $\bar{x} = 58$

$$\sigma^2 = \frac{1}{12-1} (30^2 + 36^2 + 47^2 + \dots + 110^2) - 58^2$$

$$= \frac{1}{11} (44918) - 3364$$

$$\Rightarrow 4083.45 - 3364$$

$$\Rightarrow 719.454$$

$$\sigma = 26.82 \quad \underline{\text{Ans.}}$$