

Chapter 3:

Measures of central Tendency and Dispersion.

A data set contain many variables and observations. We are not always interested in each of the measured values but rather in summary which interprets the data. Statistical functions fulfill the purpose of summarizing data in a meaningful yet concise way.

Now, we focus on the most important statistical concepts to summarize data: these are measures of central tendency and variability. The applications of each measure depends on the scale of the variables of interest.

* Measures of Central Tendency:

A natural human tendency is to make comparisons with the average.

e.g: average score, mean temperature in April in Pune, average income, etc.

- Various statistical concepts refer to the average of the data, but the right choice depends upon the nature and scale of the data as well as the objective of the study. We call statistical functions which describe the average or centre of the data location parameters or measures of central tendency.

Arithmetic Mean:

The arithmetic mean is one of the most intuitive measures of central tendency. Suppose a variable of size n consists of the values x_1, x_2, \dots, x_n . The arithmetic mean of this data is defined as

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i. \quad \text{--- (A)}$$

In informal language, we often speak of the average or just mean when using the formula (A).

Arithmetic mean for grouped data:

Class intervals a_j	$a_1 = e_0 - e_1$	$a_2 = e_1 - e_2$	\dots	$a_k = e_{k-1} - e_k$
Absolute freq. n_j	n_1	n_2	\dots	n_k
Relative freq. f_j	f_1	f_2	\dots	f_k

Note that a_1, a_2, \dots, a_k are the k class intervals & each interval a_j ($j=1, 2, \dots, k$) contains n_j observations with $\sum_{j=1}^k n_j = n$. The relative frequency of the j th class is $f_j = n_j/n$ and $\sum_{j=1}^k f_j = 1$. The mid-value of

the j th class interval is defined as $m_j = \frac{e_{j-1} + e_j}{2}$,

which is the mean of the lower and upper limits of the interval. The weighted arithmetic mean for grouped data is defined as

$$\bar{x} = \frac{1}{n} \sum_{j=1}^k n_j m_j = \sum_{j=1}^k f_j m_j.$$

The results of the mean and the weighted mean differ because we use the middle of each class as an approximation of the mean within the class. The implication is that we assume that the values are uniformly distributed within each interval. This assumption is obviously not met. If we had knowledge about the mean in each class, like in this example, we would obtain the correct result.

However, the weighted mean is meant to estimate the arithmetic mean in those situations where only grouped data is available. It is therefore typically used to obtain an approximation of the true mean.

* Properties of the Arithmetic mean:

① The sum of the deviation of each variable around the arithmetic mean is zero.

$$\sum_{i=1}^n (\bar{x}_i - \bar{x}) = \sum_{i=1}^n \bar{x}_i - n\bar{x} = n\bar{x} - n\bar{x} = 0.$$

② If the data is linearly transformed as $y_i = a + b\bar{x}_i$; where a and b are known constants, it holds that

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n (a + b\bar{x}_i) = \frac{1}{n} \sum_{i=1}^n a + \frac{b}{n} \sum_{i=1}^n \bar{x}_i$$

temperature ${}^{\circ}\text{C}$ we want of $\boxed{{}^{\circ}\text{F} = 32 + 1.8 \cdot {}^{\circ}\text{C.}}$

Median and Quantiles:

The median is the value which divides the observations into two equal parts such that at least 50% of values are greater than or equal to the median and at least 50% of the values are less than or equal to the median. The median is denoted by $\tilde{x}_{0.5}$; then, in terms of the empirical cumulative distribution function, the condition $F(\tilde{x}_{0.5}) = 0.5$ is satisfied.

Consider the n observations x_1, x_2, \dots, x_n which can be ordered as $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$. The calculations of the median depends on whether the number of observations n is odd or even. When n is odd, then $\tilde{x}_{0.5}$ is the middle ordered value. When n is even, then $\tilde{x}_{0.5}$ is the arithmetic mean of the two middle ordered values.

$$\tilde{x}_{0.5} = \begin{cases} x_{(n+1)/2} & \text{if } n \text{ is odd} \\ \frac{1}{2}(x_{(n/2)} + x_{(n/2+1)}) & \text{if } n \text{ is even.} \end{cases}$$

e.g.: $x_1 \dots x_{10}$, arrange in order $x_{(1)}, \dots, x_{(10)}$
find median.

- If we do

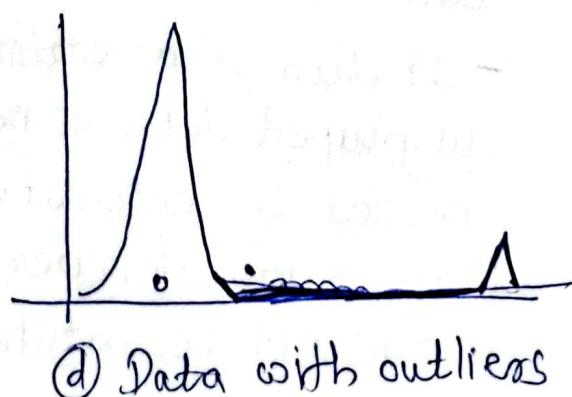
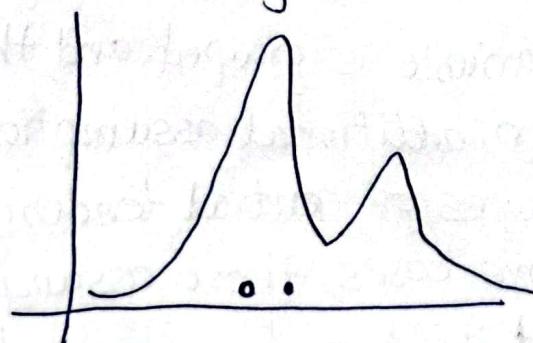
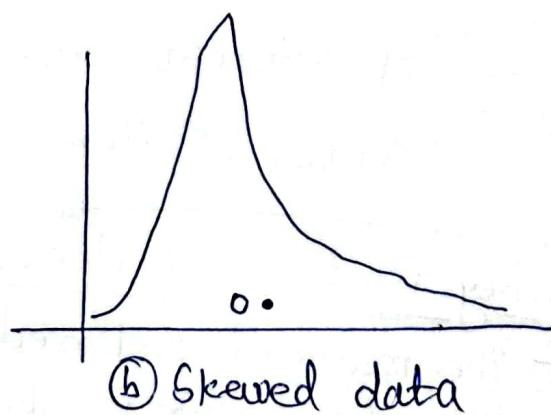
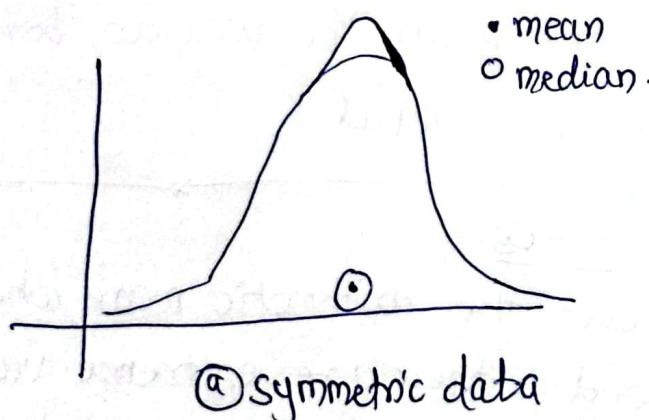
When we deal with grouped data, we can calculate the median under the assumption that the values within each class are equally distributed. Let K_1, K_2, \dots, K_k be k classes with observations of size n_1, n_2, \dots, n_k , respectively. First, we need to determine which class is the median class, i.e., the class that includes the median. We define the median class as the class K_m for which

$$\sum_{j=1}^{m-1} f_j < 0.5 \quad \text{and} \quad \sum_{j=1}^m f_j \geq 0.5 \quad \text{hold.}$$

Then, we can determine the median as

$$\tilde{x}_{0.5} = l_{m-1} + \frac{0.5 - \sum_{j=1}^{m-1} f_j}{f_m} d_m,$$

where l_{m-1} denotes the lower limit of the interval K_m and d_m is the width of the interval K_m .



Comparing mean and median:

In ①, the raw data is summarized by using ticks at the bottom of the graph and by using a kernel density estimator. The mean and the median are similar here because the distribution of the observations is symmetric around the centre. If we have skewed data ②, then the mean and the median may differ. If the data has more than one centre, such as in ③ neither the median nor the mean has meaningful interpretations. If we have outliers ④, then it is wise to use the median because the median is not sensitive to outliers.

- These examples show that depending on the situation of interest either the mean, the median, both or neither of them can be useful.

* key points & further issues:

- The median is preferred over the arithmetic mean when the data distribution is skewed or there are extreme values.
- If data of a continuous variable is grouped, and the original ungrouped data is not known, additional assumptions are needed to calculate measures of central tendency and dispersion. However, in some cases, these assumptions may not be satisfied, and the formulae provided may give imprecise results.

Quantiles:

Quantiles are a generalization of the idea of the median. The median is the value which splits the data into two equal parts. Similarly, a quantile partitions the data into other proportions.

e.g., a 25%-quantile splits the data into two parts such that at least 25% of the values are less than or equal to the quantile and at least 75% of the values are greater than or equal to the quantile.

In general, let α be a number between zero and one. The $(\alpha \times 100)\%$ -quantile, denoted as \tilde{x}_α , is defined as the value which divides the data in proportions of $(\alpha \times 100)\%$ and $(1-\alpha) \times 100\%$ such that at least $\alpha \times 100\%$ of the values are less than or equal to the quantile and at least $(1-\alpha) \times 100\%$ of the values are greater than or equal to the quantile.

- In terms of the empirical cumulative distribution function, we can write $F(\tilde{x}_\alpha) = \alpha$.
- It follows immediately that for n observations, at least $n\alpha$ values are less than or equal to \tilde{x}_α and at least $n(1-\alpha)$ observations are greater than or equal to \tilde{x}_α .

- The median is the 50% quantiles $\hat{x}_{0.5}$.
- If α takes the values $0.1, 0.2, \dots, 0.9$, the quantiles are called deciles.
- If $\alpha \cdot 100$ is an integer number (e.g. $\alpha \cdot 100 = 95$), the quantiles are called percentiles, i.e., the data is divided into 100 equal parts.
- If α takes the value $0.2, 0.4, 0.6 \& 0.8$, the quantiles are known as quartiles and they divide the data into five equal parts.
- If α takes the values $0.25, 0.5$ and 0.75 , the quantiles are called quartiles.

Consider n ordered observations $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$.

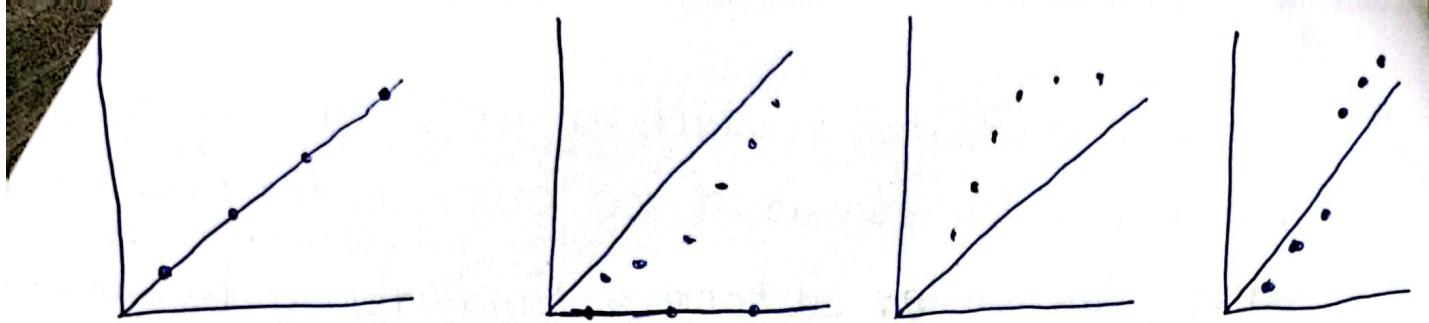
The $\alpha \cdot 100\%$ -quantile \hat{x}_α is calculated as

$$\hat{x}_\alpha = \begin{cases} x_{(k)} & \text{if } n\alpha \text{ is not an integer number,} \\ & \text{choose } k \text{ as the smallest integer} > n\alpha. \\ \frac{1}{2}(x_{(n\alpha)} + x_{(n\alpha+1)}) & \text{if } n\alpha \text{ is an integer} \end{cases}$$

* Quantile - Quantile Plots (QQ Plots).

If we plot the quantiles of two variables against each other, we obtain a Quantile-Quantile plot (QQ-Plot).

This provides a simple summary of whether the distributions of the two variables are similar with respect to their location or not.



(a) setting 1

(b) Setting 2

(c) Setting 3

(d) Setting 4

As a summary, let us consider four important patterns:

- (a) If all the pairs of quantiles lie (nearly) on a straight line at an angle of 45° from the α -axis, then the two samples have similar distributions.
- (b) If the y -quantiles are lower than the α -quantiles, then y -values have a tendency to be lower than α -values.
- (c) If the α -quantiles are lower than the y -quantiles, then the α -values have a tendency to be lower than the y -values.
- (d) If QQ plot is like (d), it indicates that there is a break point upto which the y -quantiles are lower than the α -quantiles and after that point, the y -quantiles are higher than the α -quantiles.

Mode:

Consider a situation in which an ice cream shop owner wants to know which flavour of ice cream is the most popular among his ~~custom~~ customers. Similarly, a footwear shop owner may like to find out what design and size of shoes are in highest demand. To answer this type of question, one can use the mode which is another measure of central tendency.

- The mode \bar{x}_m of n observations $x_1, x_2, x_3, \dots, x_n$ is the value which occurs the most compared with all other values, i.e., which has maximum absolute frequency. It may happen that two or more values occur with the same frequency in which case the mode is not uniquely defined.
 - A formal definition of the mode is
- $$\bar{x}_m = a_j \Leftrightarrow n_j = \max\{n_1, n_2, \dots, n_k\}.$$

The mode is typically applied to any type of variable for which the number of different values is not too large. If continuous data is summarized in groups, then the mode can be used as well.

Geometric Mean:

Consider n observations x_1, x_2, \dots, x_n which are all positive and collected on a quantitative variable. The geometric mean \bar{x}_G of this data is defined as

$$\bar{x}_G = \sqrt[n]{\prod_{i=1}^n x_i} = \left(\prod_{i=1}^n x_i\right)^{1/n}$$

Time t	Inventory B_t	Growth factor x_t	Growth rate r_t
0	B_0		$((\bar{x}_1 - 1) \cdot 100)\%$
1	B_1	$x_1 = \bar{B}_1 / B_0$	$((x_1 - 1) \cdot 100)\%$
2	B_2	$x_2 = B_2 / B_1$	$((x_2 - 1) \cdot 100)\%$
\vdots			
T	B_T	$x_T = \frac{B_T}{B_{T-1}}$	$((x_T - 1) \cdot 100)\%$

The ratio of B_t and B_{t-1} , $x_t = \frac{B_t}{B_{t-1}}$ is called the t th growth ~~rate~~ factor.

The growth rate r_t is defined as $r_t = ((x_t - 1) \cdot 100)\%$? It gives us an idea about the growth or decline of our value at time t .

Geometric mean $\bar{x}_G = \sqrt[x_1 x_2 \cdots x_T]{}$

Thus, B_t at time t can be calculated as

$$B_t = B_0 \cdot \bar{x}_G^t$$

Harmonic Mean:

The harmonic mean is typically used whenever different α_i contributes to the mean with a different weight i.e., when we implicitly assume that the weight of each α_i is not one. It can be calculated as

$$\bar{\alpha}_H = \frac{\omega_1 + \omega_2 + \dots + \omega_k}{\frac{\omega_1}{\alpha_1} + \frac{\omega_2}{\alpha_2} + \dots + \frac{\omega_k}{\alpha_k}} = \frac{\sum_{i=1}^k \omega_i}{\sum_{i=1}^k \frac{\omega_i}{\alpha_i}}$$

e.g: When calculating the average speed, each weight relates to the relative distance travelled, n_i/n , with speed α_i . Using $\omega_i = n_i/n$ & $\sum_i \omega_i = \sum_i \frac{n_i}{n} = 1$, the harmonic mean can be written as

$$\bar{\alpha}_H = \frac{1}{\sum_{i=1}^k \frac{\omega_i}{\alpha_i}}.$$

Measures of Dispersion:

Measures of central tendency, as introduced earlier, give us an idea about the location where most of the data is concentrated. However, two different data sets may have the same value for the measure of central tendency, say the same arithmetic means, but they may have different concentrations around the mean.

- In this case, the location measures may not be adequate enough to describe the distribution of the data.
- The concentration or dispersion of observations around any particular value is another property which characterizes the data and its distribution.

Now, we introduce statistical methods which describe the variability or dispersion of data.

* Range & Interquartile Range:

Consider a variable X with n observations x_1, x_2, \dots, x_n . Order these n observations as $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$.

Range: The range is a measure of dispersion defined as the difference between the maximum and minimum

value of the data as

$$R = \bar{x}_{(n)} - \bar{x}_{(1)}.$$

Interquartile Range:

The interquartile range is defined as the difference between the 75th and 25th quartiles as

$$dq = \bar{x}_{0.75} - \bar{x}_{0.25}$$

It covers the centre of distribution and contains 50% of the observations.

* Absolute Deviation, Variance and Standard Deviation:

Another measure of dispersion is the variance. The variance is one of the most important measures in statistics & is needed throughout this course.

- We use the idea of absolute deviation to give some more background and motivation for understanding the variance as a measure of dispersion.

Consider the deviation of n observations around a certain value A and combine them together, for instance, via the arithmetic mean of all the deviations :

$$D = \frac{1}{n} \sum_{i=1}^n (\bar{x}_i - A) \quad \text{--- } \textcircled{A}$$

This measure has the drawback that the deviations $(\bar{x}_i - A)$, $i=1, 2, \dots, n$, can be either positive or negative and consequently their sum can potentially

very small or even zero. Thus, using D as a measure of variability is therefore not a good idea since D may be small even for a large variability in the data.

- Using absolute values of the deviations solves this problem and we introduce the following measure of dispersion.

$$D(A) = \frac{1}{n} \sum_{i=1}^n |\bar{x}_i - A|. \quad \text{--- } \textcircled{B}$$

- It can be shown that the absolute deviation attains its minimum when A corresponds to the median of the data

$$D(\tilde{x}_{0.5}) = \frac{1}{n} \sum_{i=1}^n |\bar{x}_i - \tilde{x}_{0.5}|. \quad \text{--- } \textcircled{C}$$

* We call $D(\tilde{x}_{0.5})$ the absolute median deviation.

* When $A = \bar{x}$, we called absolute mean deviation

given by $D(\bar{x}) = \frac{1}{n} \sum_{i=1}^n |\bar{x}_i - \bar{x}|. \quad \text{--- } \textcircled{D}$

- = Another solution to avoid the positive and negative signs of deviation in \textcircled{A} is to consider the squares of deviations $\bar{x}_i - A$, rather than using the absolute value. This provides another measure of dispersion

as

$$s^2(A) = \frac{1}{n} \sum_{i=1}^n (\bar{x}_i - A)^2 \text{ which is known as the mean squared error (MSE) w.r.t. } A.$$

- The MSE (mean squared error) is another important measure in statistics, see Unit IV.

- It can be shown that $s^2(A)$ attains its minimum when $A = \bar{x}$. This is the (sample) variance

$$\tilde{s}^2 = \frac{1}{n} \sum_{i=1}^n (\bar{x}_i - \bar{x})^2.$$

After expanding, we get

$$\tilde{s}^2 = \frac{1}{n} \sum_{i=1}^n \bar{x}_i^2 - \bar{x}^2.$$

The positive square root of the variance is called the (sample) standard deviation, defined as

$$\tilde{s} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\bar{x}_i - \bar{x})^2}$$

- The standard deviation has the same unit of measurement as the data whereas the unit of the variance is the square of the units of the observations.

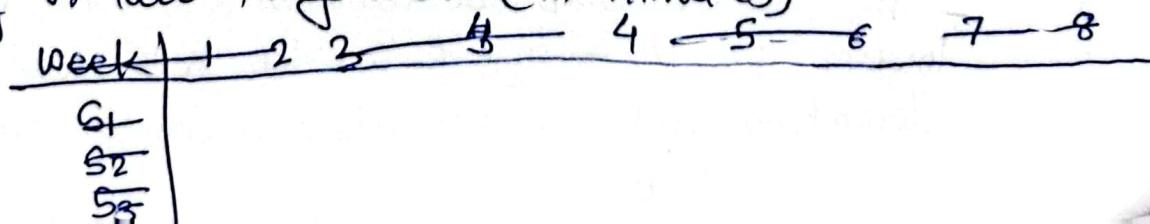
e.g: if X is distance, measured in km, then \bar{x} & \tilde{s} are also measured in km, whereas \tilde{s}^2 is measured in km² (which may be difficult to interpret).

* The variance is a measure which we use in next units to obtain measures of association between variables and to draw conclusions from a sample about a population of interest.

The standard deviation is typically preferred for a descriptive summary of the dispersion of data.

- The standard deviation measures how much the observations vary or how they are dispersed around the arithmetic mean.
- A low value of the standard deviation indicates that the values are highly concentrated around the mean. A high value of the standard deviation indicates lower concentration of the observations around the mean, and some of the observed values may be far away from the mean.
- If there are extreme values or outliers in the data, then the arithmetic mean is more sensitive to outliers than the median. In such a case, the absolute median deviation (MAD) may be preferred over the standard deviation.

exa: Suppose three students S₁, S₂, S₃ arrive at different times in the class to attend their lectures. Let us look at their arrival time in the class after or before the starting time of lecture, i.e., let us look how early or late they were (in minutes)



week	1	2	3	4	5	6	7	8	9
M1	0	0	0	0	0	0	0	0	0
M2	-10	+10	-10	+10	-10	+10	-10	+10	-10
M3	3	5	6	2	4	6	8	4	5

We calculate the variance & absolute median deviation:

$$\tilde{S}_{M1} = \frac{1}{10} \sum_{i=1}^{10} (\bar{x}_i - \bar{x})^2 = \frac{1}{10} ((0-0)^2 + \dots + (0-0)^2) = 0$$

$$\tilde{S}_{M2} = \frac{1}{10} \sum_{i=1}^{10} (\bar{x}_i - \bar{x})^2 = \frac{1}{10} ((-10-0)^2 + \dots + (10-0)^2) \approx 111.1$$

$$S_{M3} = \frac{1}{10} \sum_{i=1}^{10} (\bar{x}_i - \bar{x})^2 \approx 3.3 \quad (\because \bar{x} = 5)$$

$$D(\tilde{x}_{0.5}, M_1) = \frac{1}{10} \sum_{i=1}^{10} |\bar{x}_i - \tilde{x}_{0.5}| = \frac{1}{10} (|0-0| + \dots + |0-0|) = 0$$

$$D(\tilde{x}_{0.5}, M_2) = \frac{1}{10} \sum_{i=1}^{10} |\bar{x}_i - \tilde{x}_{0.5}| = \frac{1}{10} (|-10-0| + \dots + |10-0|)$$

$$D(\tilde{x}_{0.5}, M_3) = \frac{1}{10} \sum_{i=1}^{10} |\bar{x}_i - \tilde{x}_{0.5}| = (|3-5| + \dots + |7-5|) \frac{1}{10} = 1.4$$

— Using the arithmetic mean, we concluded that both M_1, M_2 , M_3 arrive on time, whereas M_3 is always late.

But we saw the variation of arrival times differs substantially among the three students.

— We observe that the variation/dispersion/variability is the lowest for M_1 & highest for M_2 . Both median absolute deviation and variance allow a comparison between the two students.

A hiking enthusiast has a new app for his smartphone which summarizes his hikes by using a GPS device. Let us look at the distance hiked (in km) and maximum altitude (in m) for the last 10 hikes.

Distance	12.5	29.9	14.8	18.7	7.6	16.2	16.5	27.4	12.1	17.5
Altitude	342	1245	502	555	398	670	796	912	238	466

- a) Calculate the arithmetic mean and median for both distance and altitude. $\bar{x}_D = 17.32$ $\hat{x}_{0.5,D} = 16.35$
 $\bar{x}_A = 612.4$ $\hat{x}_{0.5,A} = 528.5$

- b) Determine the arithmetic mean & median for the distance and the altitude variables. (Discuss the shape of the distribution given the results of a & b).

- c) Calculate the interquartile range, absolute median, deviation & standard deviation for both variables.

- What is your conclusion about the variability of the data? $d_{0.5,A} = \text{approx} 398$, $d_{0.5,D} = 6$ = absolute median deviation
 $s_D = \sqrt{41.5} \Leftarrow s_D^2 \approx 41.5$, $s_A^2 = 82314$. $s_A = \sqrt{82314}$ $D_D(\hat{x}_{0.5}) = 4.68$
 $D_A(\hat{x}_{0.5}) = 223.2$

- d) One metre corresponds to approximately 3.28 ft.

What is the average altitude when measured in feet rather than in meters. $\bar{y} = 3.28 \times 528.5 = 1722.48$

- e) Assume distance is measured as only short (5-15km) moderate (15-20 km) & long (20-30 km). Summarize the data grouped data in a frequency table. Calculate the weighted arithmetic mean under the assumption that the raw data is not known.

Class intervals $(5; 15]$ $(15; 20]$ $(20; 30]$

n_j	4	4	2
f_j	$4/10$	$4/10$	$2/10$
$\sum f_j$	$4/10$	$8/10$	1

$$\bar{x} = \sum f_j m_j \approx 16$$

$$\tilde{x}_{0.5} = e_{m-1} + \frac{0.5 - \sum_{j=1}^{m-1} f_j}{f_m} d_m = 16.25.$$

- ② Suppose maximum temperature during the day (in degree Celsius) for ~~December~~ July 1-31 as follows.

22, 24, 21, 22, 25, 26, 25, 24, 23, 25, 25, 26, 27, 25, 26, 25,
26, 27, 27, 28, 29, 29, 29, 28, 30, 29, 30, 31, 30, 28, 29.

Summarize this data into three categories. mean = 26.48°C

Arithmetic
Class intervals $(0-25]$ $(25, 30]$ $(30, 35]$

Absolute freq. $n_1 = 12$ $n_2 = 18$ $n_3 = 1$

Relative freq. $f_1 = \frac{12}{31}$ $f_2 = \frac{18}{31}$ $f_3 = \frac{1}{31}$

Weighted arithmetic mean : $\bar{x} = \sum_{j=1}^k f_j m_j \approx 25.7$

$\tilde{x}_{0.5} = e_{m-1} + \frac{0.5 - \sum_{j=1}^{m-1} f_j}{f_m} d_m \approx 25.97$

for grouped data (median). / median = 26.

$\tilde{x}_{0.25} = 25$, $\tilde{x}_{0.5} = 26$, $\tilde{x}_{0.75} = 29$

The following frequency table gives the values obtained in 40 rolls of a die.

Value	frequency		Find	a) mean	3.05
1	9	9			
2	8	7			
3	5	4			
4	5	4			
5	6	4			
6	7	12			

$$\text{Mean} = \frac{1 \times 9 + 2 \times 7 + 3 \times 4 + 4 \times 4 + 5 \times 4 + 6 \times 12}{40}$$

$$= \frac{9 + 14 + 12 + 16 + 20 + 72}{40}$$

$$= \frac{143}{40}$$

$$= 3.575.$$

1-2	3-4	5-6
16	8	16

[1,2]	(2,4]	(4,6]
16	8	16

$$\frac{1.5 \times 16 + 3.5 \times 8 + 5.5 \times 16}{40} \quad \left| \begin{array}{l} 1.5 \times 16 + 3 \times 8 \\ + 5 \times 16 \end{array} \right.$$

$$\text{Weighted Mean} = \frac{140}{40} = \frac{14}{4} = \underline{\underline{3.5}} \quad \left| \begin{array}{l} = 128 \\ \frac{128}{40} \end{array} \right.$$

$$\text{Median} = \frac{3+4}{2} = \underline{\underline{3.5}}, \quad \text{Mode} = \underline{\underline{6}}$$