

TOPIC V

FULL NOTES ⇒ STATISTICS

- PROBABILITY

5.1 - PRESENTATION OF DATA

→ DISCRETE DATA ⇒ Exact number valves. Can be counted.

EG - 1 ⇒ Number of people, test scores, shoe sizes,

→ CONTINUOUS DATA ⇒ Takes numerical valves within a range, usually measured,

EG - 1 ⇒ Height, weight, temparature, ...

→ PRESENTATION METHODS ⇒

→ TABLES ⇒ We start with a list :

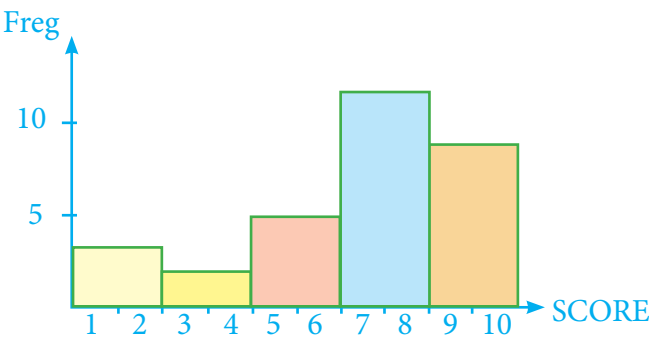
This data can be shown with a frequency table (night), and a groyed frequency table w / equal class size :

| SCORE | TALLY | FREQ | RELATIVE FREQ % |
|---------|-----------------------|------|--------------------|
| 1 - 2 | III | 3 | 10 % |
| 3 - 4 | II | 2 | 6.7 % |
| 5 - 6 | III | 5 | 16.7 % |
| 7 - 8 | III III II | 8 | 40 % |
| 9 - 10 | III III | 5 | 26.7 % |
| TOTAL : | | 30 | 100 % |

EG - 1 ⇒ Test scores, out of 10 : (n = 30)
8 7 8 4 1 9 9 10 6 8 6 7 8 2 10
10 9 8 6 8 2 8 7 3 7 9 5 9 6 8

| SCORE | TALLY | FREQUENCY |
|---------|---------|-----------|
| 1 | I | 1 |
| 2 | II | 2 |
| 3 | I | 1 |
| 4 | I | 1 |
| 5 | I | 1 |
| 6 | III | 4 |
| 7 | III | 4 |
| 8 | III III | 8 |
| 9 | III | 5 |
| 10 | III | 3 |
| TOTAL : | | 30 |

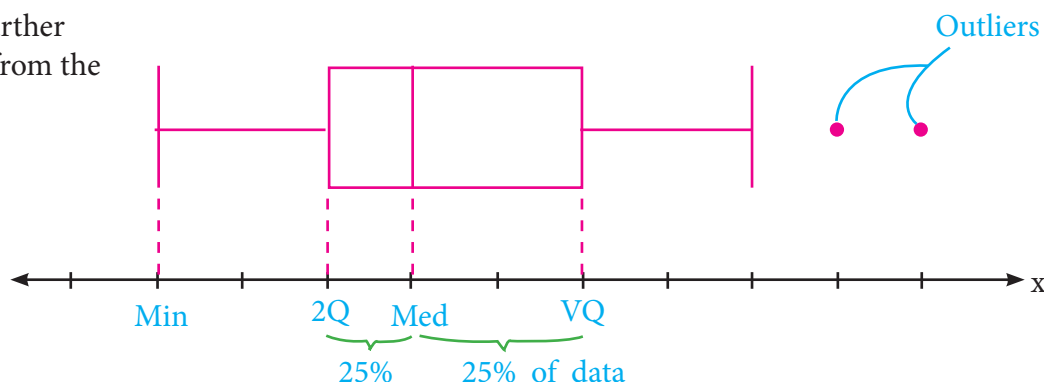
The next step here would be to display the data in a HISTOGRAM
↳ We can say here that the 'modal class' is a score of 7 to 8.



→ **BOX & WHISKER PLOTS** ⇒ This is another way of showing your data. It shows :

the median, upper / lower quartiles (see 5.2) and 'outliers'

Outliers : These are data points that lie further than $1.5 \times \text{IQR}$ from the median.



5.2 STATISTICAL MEASURES

→ **CENTRAL TENDANCY** :

MODE ⇒ Most frequent value / class :

MEDIAN ⇒ Order data, then observe the $\left(\frac{n+1}{2}\right)^{\text{th}}$ value :

EG - ⇒

For : 1, 2, 1, 6, 5

①

②

$$\frac{15}{5} = ③$$

For : 18, 17, 14, 15
15, 11, 10, 4

⑮

Halfway between

14 & 15 = ⑭.5

$$\frac{104}{8} = ⑬$$

$$\text{MEAN} \Rightarrow \frac{\text{TOTAL}}{\text{FREQUENCY}} = \frac{\sum x}{\sum f}$$

↳ **For grouped frequency** ⇒ You can only 'estimate' the total by multiplying midpoints by Freq. and summing these

$$\frac{\sum (f.x)}{\sum f}$$

EG - 1 ⇒

↗ Add column

| SCORE | TALLY | FREQUENCY |
|---------|-------|-----------|
| 1 - 2 | 3 | 4.5 |
| 3 - 4 | 2 | 7 |
| 5 - 6 | 5 | 27.5 |
| 7 - 8 | 12 | 90 |
| 9 - 10 | 8 | 76 |
| TOTAL : | 30 | 205 |

Modal class

$$\text{MEAN} : \frac{205}{30} = 6.83$$

→ **QUARTILES** ⇒ Similar to how the median splits the data into 50% part, quartiles are 4 equal parts of 25%, after ordering

LOWER QUARTILE (LQ / Q1) ⇒ Point that is greater than 25% of the data

UPPER QUARTILE (UQ / Q3) ⇒ Point that is greater than 75% of the data

→ **PERCENTILES** ⇒ An extension of this breakdown of the data is splitting it into 1% sections, called 'percentiles'

DISPERSION

→ **RANGE** ⇒ **HIGHEST VALUE - LOWEST VALUE = RANGE**

→ **INTERQUARTILE RANGE** ⇒ **IQR = Q₃ - Q₁**

→ **VARIANCE** ⇒ Another measure of spread of data. It is essentially the average of the differences between each value and the mean.

$$(\text{Var}) = S_n^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

→ **STANDARD DEVIATION** ⇒ Related measure, just calculated by doing the square root of the variance.
S.D. is often used more frequently



CALCULATOR :

TI - nspire ⇒ Enter values in 'List & Spreadsheets', then on a calculator page, press MEND → 6 :
Statistik → 1 1 : One - var stats . S. D. is Qx

TI - 84 ⇒ Press STAT → 1 : EDIT → ENTER, then fill column with values. Then press
STAT → CALC → 1 : 1 var stats ENTER. S. D. is Qx

5.3 CUMULATIVE FREQUENCY

⇒ It is hard to see what proportion of the data is above or below a certain value. Adding on the frequency of each successive class, we build a 'cumulative freq' column. We plot this against the UB of each class, then analyse :

EG - 1 ⇒

| SCORE | FREQ. | Add column | |
|-------------------|-------|------------|------|
| | | U.B | C.F. |
| $10 \leq x < 20$ | 2 | 20 | 2 |
| $20 \leq x < 30$ | 5 | 30 | 7 |
| $30 \leq x < 40$ | 7 | 40 | 14 |
| $40 \leq x < 50$ | 21 | 50 | 35 |
| $50 \leq x < 60$ | 36 | 60 | 71 |
| $60 \leq x < 70$ | 40 | 70 | 111 |
| $80 \leq x < 80$ | 27 | 80 | 138 |
| $80 \leq x < 90$ | 9 | 90 | 147 |
| $90 \leq x < 100$ | 3 | 100 | 150 |

→ ANALYSIS ⇒

FINDING MEDIAN

↳ Start at $\left(\frac{n+1}{2}\right)^{th}$ position on y - axis, trace across to the curve, then down to the x - axis, then read your median. (See blue line)

LOWER QUARTILE

↳ Start at $\left(\frac{n+1}{4}\right)^{th}$ position then follow the same process as above.

UPPER QUARTILE

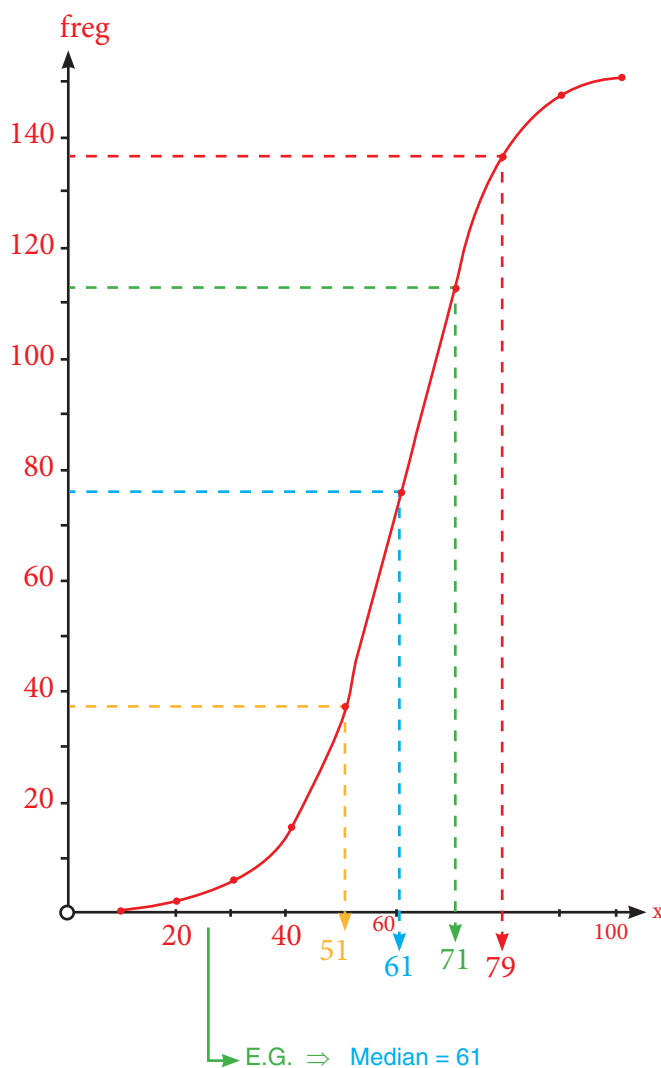
↳ Start at $\left(\frac{3(n+1)}{4}\right)^{th}$ position then follow the same process as above.

FINDING PERCENTILES

↳ Similar process, but with $\left(\frac{p(n+1)}{100}\right)$ if you are trying to find the pth percentile

↳ E.G.(above) : 90 th percentile :

$$\frac{90(151)}{100} = 1.35.9$$



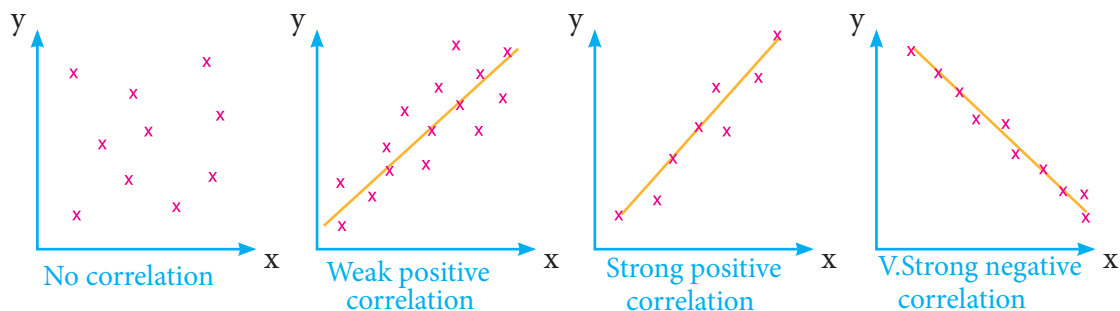
$$\left. \begin{array}{l} \text{LQ} = 51 \\ \text{UQ} = 71 \end{array} \right\} \text{IQR} = 71 - 51 = 20$$

90th percentile = 79

5.4 CORRELATION

→ **CORRELATION** ⇒ Testing the correlation between two sets of data, x and y, means that you are testing whether a change in x causes a similar change in y, and to what extent.

→ **GRAPHS** ⇒ This is one way of showing what the strengths of correlation mean in reality :



→ **CORRELATIONS COEFFICIENT** ⇒ An accurate measure of the relationship between two variables 'r' measured on a scale from -1 (perfect negative corr.) to +1 (perfect positive correlation)

→ **BY HAND** ⇒
$$r = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2 \sum (y - \bar{y})^2}}$$
 (You will use a calculator in exams)

→ **BY HAND** ⇒

TI - 84

- **STAT** | 1 : Edit **ENTER**
- Put x values in L1, y in L2
- **MODE** | Check diagnostics are ON
- **2nd** **QUIT**
- **STAT** | **CALC** | 4 : Lin Reg (ax + b)
- **ENTER** until info shown

TI - ASPIRE

- Add 'Lists & Spreadsheet'
- Name columns : 'x' & 'y'
- Type data
- **MENU** | 4 : Stats | 1 : Stat Calc. | 3 : Lin. Reg (mx + b)
- Choose 'x' for X list, 'y' for Y list
- **ENTER** , then observe

→ **CORRELATION** ⇒ On the same calculator screen as the r value, you will see a value for ' a ' & ' b '. This creates an equation in the form $y = ax + b$.

This is a more advanced version of what you may have seen as a 'line of best fit', which is more of an estimate.

You may well be asked to use this equation to estimate y values given x values

5.5 PROBABILITY

→ **DEFINITIONS** ⇒ '**TRIAL**' → Each time an experiment is repeated.

· '**OUTCOMES**' → Possible results of one trial.

'**SAMPLE SPACE (U)**' → Set of all possible outcomes in an experiment

→ **PROB. RULES** ⇒ If A is a set of results from an experiment with all equally likely results, then :

· → Prob. Of A occurring = $P(A) = \frac{\text{EVENTS IN } A}{\text{EVENTS IN } U} = \frac{n(A)}{n(U)}$

EG - ⇒ Probability of rolling a multiple of 3 on a fair dice :

↳ **SOL :** $P(\text{Mult. of } 3) = \frac{2}{6} = \frac{1}{3}$

⇒ Two events are **complementary** if exactly one of the two events must occur, so :

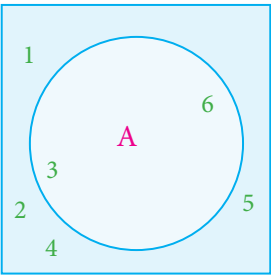
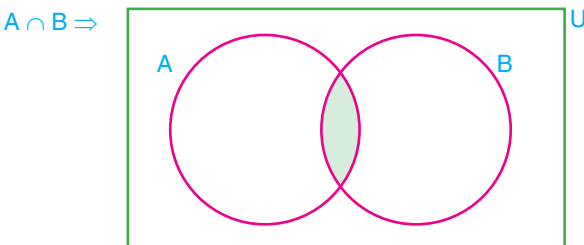
→ $P(A) + P(A') = 1$

EG - ⇒ Rolling a 6 and not rolling 6 are **complementary**, as the probabilities add to 1 :

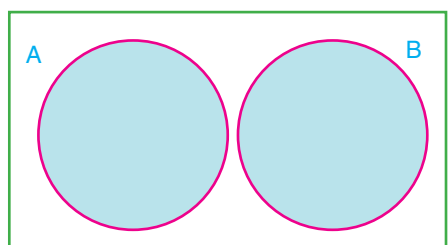
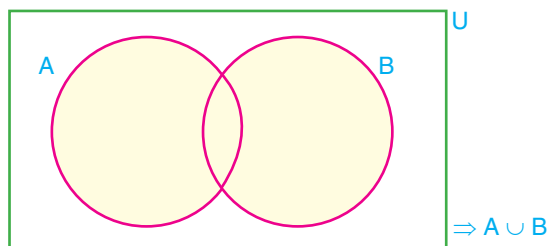
↳ **SOL :** $\frac{1}{6} + \frac{5}{6} = 1$

→ **VENN DIAGRAMS** ⇒ The venn diagram on the left represents the earlier multiples of 3 problem. You can see where the 2 out of 6 probability is derived from :

→ **INTERSECTION** ⇒ The venn to the right represents the 'intersection' between sets A & B , denoted $A \cap B$. These are the elements common to both A and B :



→ **INTERSECTION** ⇒ Denoted $A \cap B$, represents elements in A or B or both :

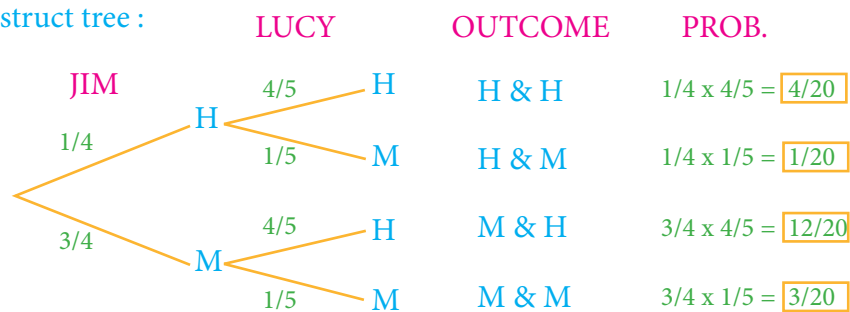


- * These sets are disjoint
- * We say $A \cap B = \emptyset$, where \emptyset is the empty set
- * A & B are mutually exclusive - meaning you can have both occurring at the same time

→ **PROBABILITY TREES** ⇒ When you have two or more trials, and the possible outcomes are not too numerous, we can use 'probability trees' :

EG - ⇒ In archery, Jim hits a target $\frac{1}{4}$ of the time, and Lucy hits it $\frac{4}{9}$ of the time.

i) Construct tree :



ii) Prob of at least one hit ? 1st , 2nd & 3rd options : $4/20 + 1/20 + 12/20 = \boxed{17/20}$

5.6 COMBINED PROB :

- **UNION PROB.** ⇒ $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- **INTERSECTION** ⇒ $P(A \cap B) = P(A) \times P(B)$ [if independent]
- **CONDITIONAL** ⇒ $P(A / B) = \frac{P(A \cap B)}{P(B)}$ → where $P(A / B)$ means : the prob. of a occurring, given that B has occurred.
- **INDEPENDENT** ⇒ Two event are independent if one event occurring does not affect the probability of the other event occurring.
It can also be shown formally by checking whether :
 $P(A / B) = P(A)$ holds
OR
 $P(A \cap B) = P(A) \times P(B)$ holds.

5.7 DISCRETE RANDOM VARIABLES

- **RANDOM VARIABLE** ⇒ This represents, in number form, the possible outcomes which could occur for some random experiment.
- **DISCRETE** ⇒ A set of distinct possible values, which you can count.
- **CONTINUOUS** ⇒ Values are measured between a certain range.
- **PROBABILITY DISTRIBUTIONS** ⇒ For any random variable, there is a prob. dist. which describes the probability of each value occurring
- ↳ **NOTATION** ⇒ The prob. that the variable X takes value x is denoted : $P(x = x)$
- EG - ⇒ Tossing 2 coins, counting how many ' heads ' occur
- ↳ $P(x = 0) = 1/4$, $P(x = 1) = \frac{2}{4} = \frac{1}{2}$, $P(x = 2) = 1/4$
- **RULE** ⇒ For something to be a valid probability distribution function, the sum of the probabilities needs to equal 1.
i. e : $\sum P(x) = 1$
- EG - ⇒ (From above ↑) : $\frac{1}{4} + \frac{1}{2} + \frac{1}{4} = 1$

EG - \Rightarrow Find k : $\frac{x}{P(x=x)} \begin{array}{c|c|c|c} 0 & 1 & 2 \\ \hline 0,3 & k & 0,5 \end{array} \rightarrow 1 - 0,3 - 0,5 = 0,2, \quad k = 0,2$

\rightarrow EXPECTED VALUE \Rightarrow Take n trials of an experiment, where in each of the trials the event has prob. of p of occurring, then the number of times we expect the event to occur is $n \times p$
The expected outcome for the random variable X is the mean result, μ

$$E(X) = \mu = \sum_{i=1}^a x_i q_i$$

EG - \Rightarrow In a magazine store, 23% of customers purchased 1 magazine, 38% bought 2, 21% bought 3, 13% bought 4, and 5% bought 5. Calculate the expected number of magazines bought

\hookrightarrow

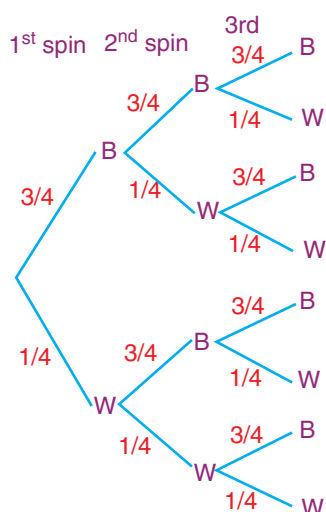
| | | | | | |
|-------|------|------|------|------|------|
| x_i | 1 | 2 | 3 | 4 | 5 |
| p_i | 0.23 | 0.38 | 0.21 | 0.13 | 0.05 |

$$\begin{aligned} \mu &= \sum x_i p_i \\ &= 1(0.23) + 2(0.38) + 3(0.21) + 4(0.13) + 5(0.05) \\ &= 2.39 \text{ magazines} \end{aligned}$$

5.8 BINOMIAL DISTRIBUTION

\rightarrow BINOMIAL EXPERIMENTS \Rightarrow These are experiments where there are just 2 possible results : **success** or **failure** (event occurring or not). This is then repeated in a number of independent trials.
For each trial, prob. of success is p , and prob. of failure is therefore $1 - p$

\rightarrow OPENING PROBLEM \Rightarrow A 'spinner' has three blue edges and one white edge so it has $3/4$ prob. of getting blue, and $1/4$ of getting white.
If we call spinning a blue a 'success', so $p = 3/4$, then we can analyse the prob. of getting 0, 1, 2, 3 success from 3 spins.



PROBABILITIES \Rightarrow

$$P(3 \text{ blues}) = P(x = 3) = \left(\frac{3}{4}\right)^3 = 0.4219$$

$$P(0 \text{ blues}) = P(x = 0) = \left(\frac{1}{4}\right)^3 = 0.0156$$

Those two are easy to calculate, but 1 blue and 2 blues both have three different paths leading to that total. So we have a 'multiplier' of $\times 3$:

$$P(1 \text{ blues}) = P(x = 1) = \left(\frac{3}{4}\right) \left(\frac{1}{4}\right)^2 \times 3 = 0.1406$$

\nwarrow
BWW, WBW or WWB

$$P(2 \text{ blues}) = P(x = 2) = \left(\frac{3}{4}\right)^2 \left(\frac{1}{4}\right) \times 3 = 0.4219$$

Notice that $0.0156 + 0.1406 + 0.4219 + 0.4219 = 1$

→ **FUNCTION** ⇒ Above is just an example for 3 trials. We need to know how calculate probabilities
 . for any n. The multiplier effect is illustrated in the $\binom{n}{r}$ part of the following function :

For n trials, the probability that there are r successes & n - r failures is :

$$P(x = r) = \binom{n}{r} p^r (1 - p)^{n-r} \quad \text{for } r = 0, 1, \dots, n$$

EG - ⇒ 8 rolls of a dice, prob. of rolling 5 sixes ?

$$n = 8$$

$$r = 5$$

$$p = 1/6$$

$$P(x = 5) = \binom{8}{5} \left(\frac{1}{6}\right)^5 \left(\frac{5}{6}\right)^3 = 56 \cdot \left(\frac{1}{6}\right)^5 \left(\frac{5}{6}\right)^3 = 0,00417$$

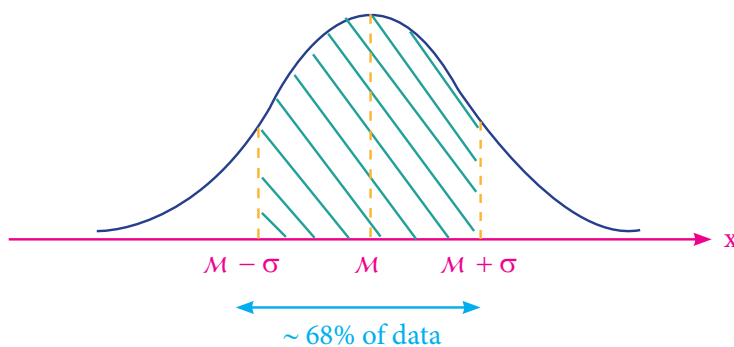
See **1.3**

5.9 NORMAL DISTRIBUTION

When you have a large single set of values, a histogram / bar chart becomes unfeasible. So an alternative is to approximate this into one bell - shaped curve

- **PROPERTIES** ⇒ *
- * SYMMETRICAL BELL - SHAPED CURVE
 - * THE AREA UNDER THE CURVE IS EQUAL TO 1 (or 100%)
 - * IT IS DEFINED BY ITS MEAN (μ) & STANDARD DEV. (σ)
 - * THE MEAN IS IN THE CENTER
 - * ~ 68% OF THE DATA IS WITHIN 1σ OF THE MEAN
 - * ~ 95% IS + 2σ OF THE MEAN ~99% IS $\pm 3\sigma$

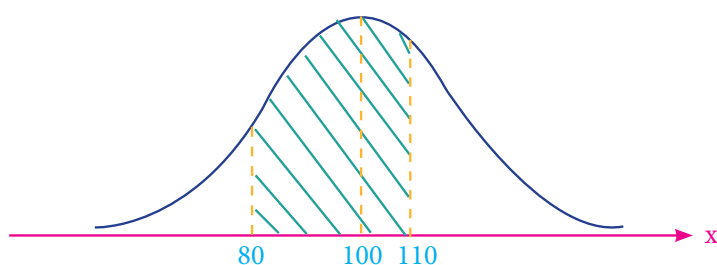
→ **GENERAL DIAGRAM** ⇒



→ **NORMAL PROB. CALCULATIONS**

⇒ When you are given the μ & σ of a curve, you can find the probability that a randomly selected valve would lie within a certain range of x :

EG - ⇒ A set of 2000 IQ scores is normally distributed with $\mu = 100$ & $\sigma = 10$. Find the probability if picking an IQ between 80 & 110 :



NOTE : If no US / LB given, you may have to use 10^{99} & -10^{99} respectively.

SOL : (TI - nspire)

→ **MENU** | 5 : Prob. | 5 : Dist. | 2 : Normal CDF
 → Enter L.B. = 80, U.B = 110, $\mu = 100, \sigma = 10$
 → Answer will be shown as decimal

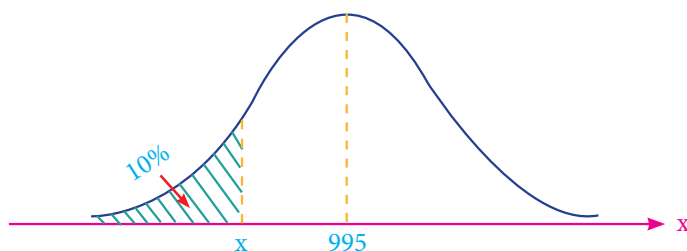
(TI - 84)

→ **2nd** | **DISTR** | 2 : normal cdf | **ENTER**
 → Rest is the same as above ...

→ **INVERSE NORM CALCULATIONS**

⇒ This style of question is the opposite. You will be given an area / prob., and asked to find a boundary :

EG - ⇒ The vol. of milk contains has : $\mu = 995$ mL & $\sigma = 5$ mL. 10% of contains are $< x$ mL. Find x :



SOL : (TI - nspire)

→ **MENU** | 5 : Prob. | 5 : Dist. | 3 : Inverse.
 → Enter Area = 0,1 $\mu = 995, \sigma = 5$
 → The Answer (998.6 mL) will be an upper boundary. To do find a Lower boundary do [1 - answer]

(TI - 84)

→ **2nd** | **DISTR** | 3 : Invnorm | **ENTER**
 → Rest is the same as above ...

→ **INVERSE NORM CALCULATIONS**

⇒ In certain cases, it is helpful to work with a normal distribution with $\mu = 0$, $\sigma = 1$. This is called the **Z - distribution**, or $Z \sim N(0,1)$

→ **METHOD** ⇒ To use the Z - dist, you must also transform each x - value to what we call a **z - score**

$$z = \frac{x - \mu}{\sigma} \text{ [This also represents how many S.D's from } \mu \text{ it is]}$$