

# S1

## 1.) DATA

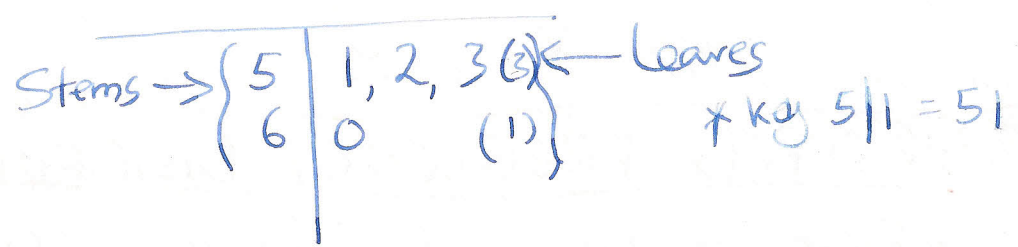
→ Continuous: Can take any value in a given range  
Eg. Time, Volume, Height, etc.

→ Discrete: Can only take particular values  
Eg. For a Die; you can get 1, 2, 3, 4, 5, 6 only  
[singular for Dice]

## ⇒ STEMS LEAF DIAGRAMS

\* A way of ordering & Presenting Data

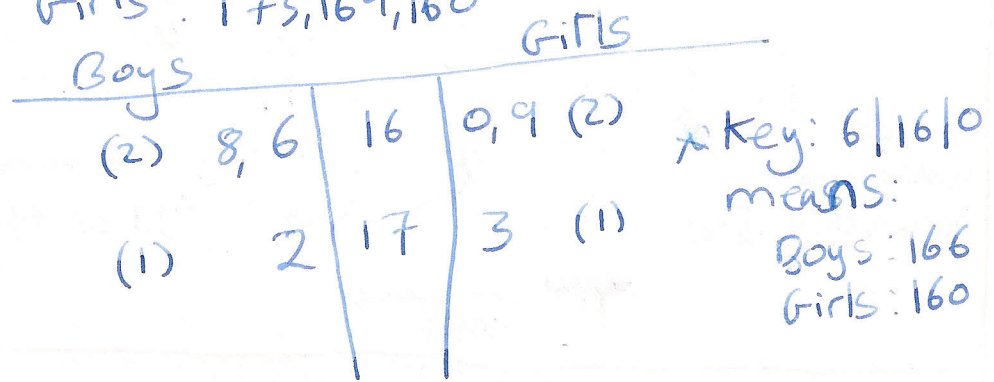
→ Plain Eg. 60, 51, 53, 42



\* make sure the leaves digits are in order  
\* Always give a key

## ⇒ Back-to-Back

Eg. Boys: 168, 166, 172  
Girls: 173, 169, 160



## → MEAN, MODE, MEDIAN & RANGE

(2)

\* MODE: is that value of a variable which occurs most frequently.

• modal class → has the highest frequency

\* MEDIAN: is the middle value of an ordered set of data.

$$= \frac{1}{2}(n+1)\text{th observation}$$

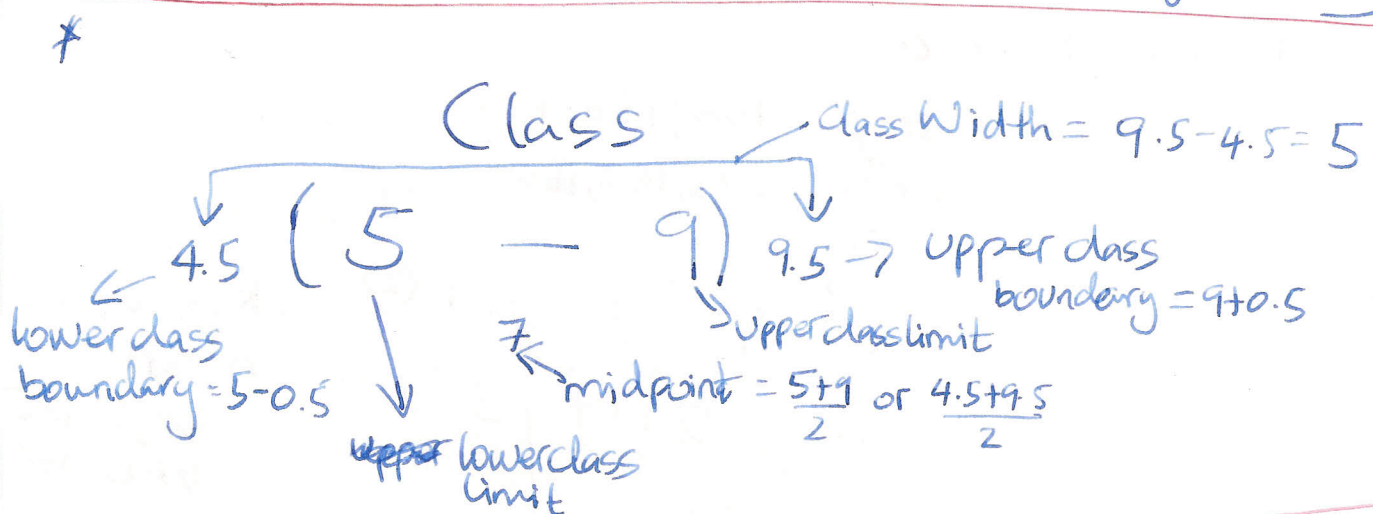
\* RANGE = Highest Value - Lowest Value

\* MEAN = Average =  $\frac{\text{Sum of Observations}}{\text{Total Observations}} = \frac{\sum x \text{ or } \sum fx}{n \quad \sum f}$

Signs =  $\bar{x}$ ;  $\bar{y}$ ;  $\mu$

## → GROUPED FREQUENCY DISTRIBUTION

\* First get class width =  $\left[ \frac{\text{largest value} - \text{Smallest value}}{\text{Number of groups}} \right]$



\* Cumulative Frequency = Running Total

# → HISTOGRAMS

(3)

Eg.

	Class	Number of children	Class Width	frequency Density
(129.5-134.5)	130-134	10	5	2.0
(134.5-139.5)	135-139	22	5	4.4
(139.5-144.5)	140-144	38	5	7.6

$$* \text{ Frequency Density} = \frac{\text{Frequency}}{\text{Class Width}}$$

\* Upper & lower class bound can be found by taking the end of one class + beginning of next class  
2

$$\text{Eg. } \frac{135 + 134}{2} = 134.5$$

$$\frac{140 + 139}{2} = 139.5$$

$$* \text{ AREA} = \text{FREQUENCY}$$
$$\text{FREQUENCY} = \text{Height of bar} \times \text{Class Width}$$

$$* \text{ Height of Bar} = \frac{\text{FREQUENCY}}{\text{Class Width}} ; \quad h = \frac{f}{w}$$

$$* \text{ TOTAL AREA} \propto \text{TOTAL FREQUENCY}$$

↑  
Proportional.

## → MEDIAN CLASS

(4)

Eg.

Visits	f	cf
0-4	32	32
5-9	71	103
10-14	20	123

$$\text{Median} = b + \left( \frac{\frac{1}{2}n - f}{f_m} \right) \times w$$

\*  $b$  = lower class boundary

$f$  = Sum of all frequencies below  $b$

$f_m$  = frequency of median class

$w$  = Class Interval/width

• Step 1: Calculate median value

$$\rightarrow \frac{1}{2}(n+1) = \frac{1}{2}(123+1) = 62$$

\* This is found in class (5-9)

• Step 2: \* Lower class boundary = 4.5

\*  $f = 32 \rightarrow$  cf before lower class boundary

• Step 3: \* Class width =  $9.5 - 4.5 = 5$

\*  $f_m = 71 \Rightarrow$  frequency of median class

$$\begin{aligned} \text{Median} &= 4.5 + \left[ \frac{\frac{1}{2}(123) - 32}{71} \right] \times 5 \\ &= \underline{\underline{6.58}} \end{aligned}$$

→ Mean

(3)

Visits	midclass Value $x$	$f$	$fx$
0-4	2.25	32	<del>88</del> 72
5-9	7	71	497
10-14	<del>11</del> 12	20	<del>500</del> 240
Totals		<u>123</u>	<u><del>1077</del> 809</u>

\* To find Mid class Value ( $x$ )

$$= \frac{\text{Lower class boundary} + \text{Upper class boundary}}{2}$$

\* Mean =  $\frac{\sum fx}{\sum f} = \frac{809}{123} = 6.58$  → DO NOT FORGET TO ADD UNITS (to 2d.p.)

→ Quartiles:

$$Q_1 = \text{lower quartile} = b + \left[ \frac{\frac{1}{4}n - f}{f_m} \right] \times w$$

\* Similar to median formula

$$Q_2 = \text{median} = b + \left[ \frac{\frac{1}{2}n - f}{f_m} \right] \times w$$

$$Q_3 = \text{Upper Quartile} = b + \left[ \frac{\frac{3}{4}n - f}{f_m} \right] \times w$$

→ Interquartile range =  $(Q_3 - Q_1)$

→ Semi-Interquartile Range (S.I.Q.R) =  $\frac{1}{2}(Q_3 - Q_1)$

$$\rightarrow \text{Percentile} = \frac{l}{100} n$$

(6)

Eg. Percentile 85th

$$\therefore P_{85} = b + \left[ \frac{\left( \frac{85}{100} n - f \right)}{f_m} \right] \times w$$

$$\rightarrow \text{Decile} = \frac{l}{10} n$$

Eg. 3rd Decile

$$\therefore D_3 = b + \left[ \frac{\left( \frac{3}{10} n - f \right)}{f_m} \right] \times w$$

\* Make sure you calculate the value of  $b, f, f_m, w, \text{etc.}$  for EACH QUANTILE, DECILE & PERCENTILE as they are NOT the same for each of them.

\* Follow the same steps mention on Pg. 4.

\* Remember: To identify the relevant classes, use:

$$* Q_1 = \frac{1}{4} n$$

$$* Q_2 = \frac{1}{2} (n+1)$$

$$* Q_3 = \frac{3}{4} n$$

$$* \text{Percentile} = \frac{l}{100} n$$

$$* \text{Decile} = \frac{l}{10} n$$

# → Stem & Leaf Diagrams

For  $Q_1 = \frac{1}{4}n$

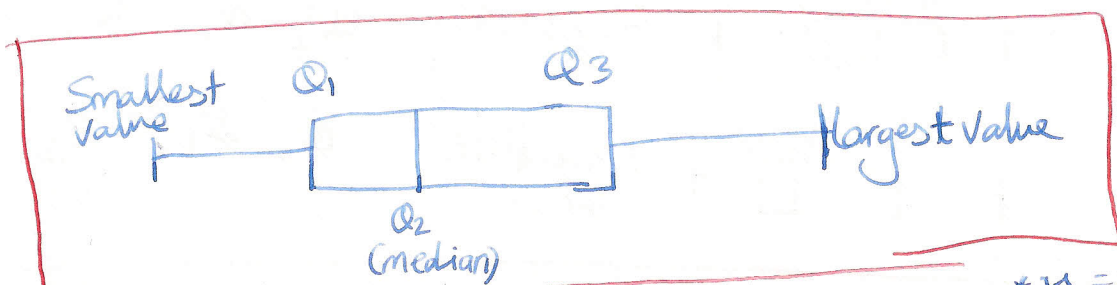
For  $Q_2 = \frac{1}{2}(n+1)$

For  $Q_3 = \frac{3}{4}n$

This Tells you which observation to take eg.  $Q_1 = \frac{51}{4} = 13\text{th observation}$

$\therefore Q_1 = 32$   
↑  
taken from the Stem & Leaf Diagram

# → Box Plots [or Box & Whisker Diagram]



# → VARIANCE & STANDARD DEVIATION ( $\sigma^2$ )

\*  $\mu = \bar{x} = \text{mean} = \frac{\sum x}{\sum f}$

\* VARIANCE =  $\sigma^2 = \frac{\sum (x - \mu)^2}{n} = \frac{\sum fx^2}{\sum f} - \mu^2$

\* STANDARD DEVIATION =  $\sqrt{\sigma^2} = \sqrt{\left(\frac{\sum fx^2}{\sum f} - \mu^2\right)}$

# → CODING

\*  $x$  = Original value  
\*  $y$  = new 'coding' value

$y = \frac{x - b}{a}$

$\bar{y} = \frac{\bar{x} - b}{a}$

$\sigma_y = \frac{\sigma_x}{a}$

\* All you are doing is simplifying longer values, in order to save time.

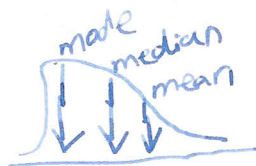
\* to get the  $y$  values, minus any thing common eg 1 million or the smallest value, & divide by something that is divisible for all original values.

\* Find mean & Standard Deviation of  $y$ , then replace to find of  $x$

# ⇒ SKEWNESS & OUTLIERS

8

\* methods



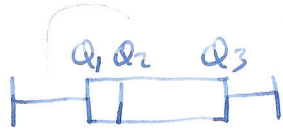
\*  $\angle$  means 'less than'

① Positive Skew: mode  $\angle$  median  $\angle$  mean.

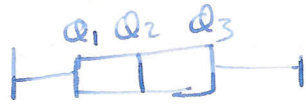
Negative Skew: mean  $\angle$  median  $\angle$  mode.

Symmetrical: mean = mode = median

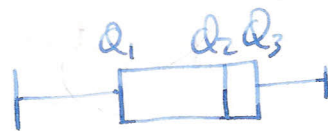
②



; Positive skew =  $Q_2 - Q_1 < Q_3 - Q_2$



Symmetry =  $Q_2 - Q_1 = Q_3 - Q_2$



Negative skew =  $Q_2 - Q_1 > Q_3 - Q_2$

③

$$\text{Pearson's Coefficient of Skewness} = \frac{\text{mean} - \text{mode}}{\text{Standard deviation}} = \frac{3(\text{mean} - \text{median})}{\text{Standard deviation}}$$

⇒ Outlier = is found outside the rest of the readings



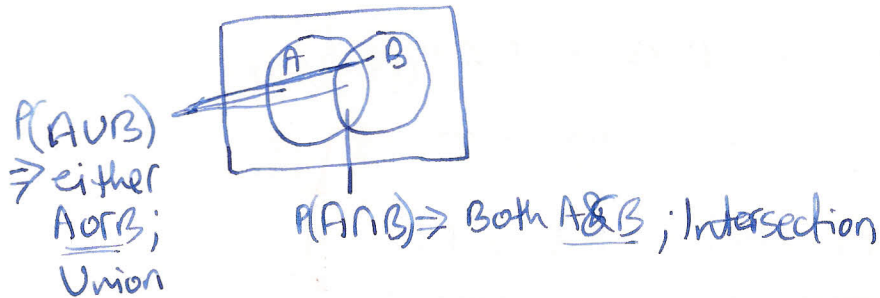
## 2.) PROBABILITY

(9)

$$* P(\text{event}) = \frac{\text{No. of outcomes where event happens}}{\text{Total number of possible outcomes}}$$

\* Sample Space = set of all possible outcomes

\* Venn Diagram = shows which outcome corresponds to which events.



→ ADDITIONAL RULE

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

→ Complementary probability:  $P(A') = 1 - P(A)$  or  $P(A') + P(A) = 1$   
not event A      event A      event A

→ CONDITIONAL PROBABILITY

~~$P(A|B)$~~   $P(m|B) = P(m \text{ given } B \text{ has already happened})$

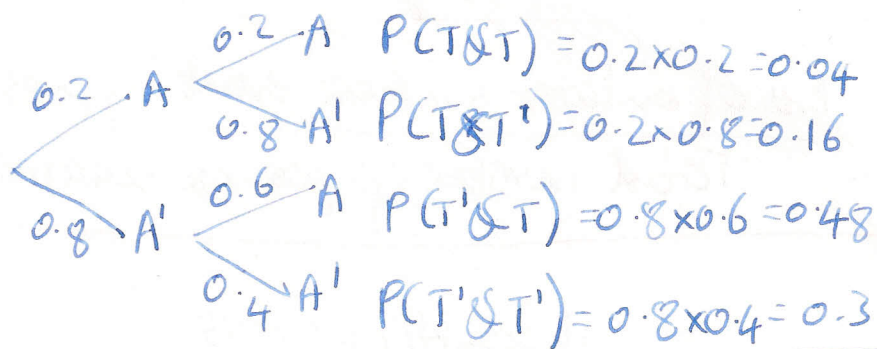
$$P(m|B) = \frac{P(m \cap B)}{P(B)}$$

\* MULTIPLICATION RULE =  $P(m|B) \times P(B) = P(m \cap B)$

→ Tree diagrams: Show Probabilities of 2 or more events

(10)

Eg.



1.00 → All probabilities  
Should add up  
to ONE.

→ Independent Events have no effect on Each other

$$P(A \cap B) = P(A) P(B)$$

Since  $P(A|B) = \frac{P(A \cap B)}{P(B)}$  &  $P(A|B) = P(A)$

→ Mutually Exclusive: Events have no overlap

$\Rightarrow P(A \cap B) = 0$

Since  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ , &  $P(A \cap B) = 0$

$$\therefore P(A \cup B) = P(A) + P(B)$$



# → Scatter Diagrams.

(12)



\* one variable gets bigger, so does the other

\* one variable gets smaller, the other becomes bigger.



no Correlation

\* The two variables are not linked.



# → Coding

$$X = \frac{x-a}{b} ; Y = \frac{y-c}{d}$$

\* Use codes when values are massive

\* Use the codes in the calculation. The PMCC(r) you get would be the same.

# → LINEAR REGRESSION

\* This is the 'Statistics' way of find equation of the best fit line.

⇒ Regression line (line of best fit) ⇒  $y = a + bx$

$$* a = \bar{y} - b\bar{x} \text{ or } \frac{\sum y}{n} - b \frac{\sum x}{n}$$

$$* b = \frac{\sum xy}{\sum x^2} \text{ or } \frac{\sum xy - \frac{(\sum x)(\sum y)}{n}}{\sum y^2 - \frac{(\sum y)^2}{n}}$$

\* Residual = Observed y-value - Estimated y-value

↳ The value taken from graph - value calculated using regression line equation.

# 4.) PROBABILITY DISTRIBUTION

\*  $X$  (upper case)  $\rightarrow$  name of random variable

$x$  (lower case)  $\rightarrow$  particular value

\* Random Variable  $\rightarrow$  no fixed value

\* Discrete Random Variable - only has a certain number of Random Variables.

Eg. Rolling a die with 4 faces.

$x:$	1	2	3	4
$P(X=x):$	$\frac{k}{1}$	$\frac{k}{2}$	$\frac{k}{3}$	$\frac{k}{4}$

$\swarrow$  possible results

$\nwarrow$  probabilities

$\uparrow$   
Probability function

$$\sum_{\text{all } x} P(X=x) = 1$$

$\rightarrow$  all probabilities add up to 1.

$\therefore \frac{k}{1} + \frac{k}{2} + \frac{k}{3} + \frac{k}{4} = 1 \Rightarrow k = \frac{12}{25}$

$\rightarrow$  mode = most likely value = highest probability

Eg. 2

$x:$	1	2	3	4
$P(X=x):$	$\frac{12}{25}$	$\frac{6}{25}$	$\frac{4}{25}$	$\frac{3}{25}$

(a)  $P(X=1) = \frac{12}{25}$

(b)  $P(X=0) = 0$

(c)  $P(X < 3) = \frac{6+12}{25} = \frac{18}{25}$   
 $\downarrow P(2)+P(1)$

(d)  $P(X \leq 3) = \frac{12+6+4}{25} = \frac{22}{25}$   
 $\downarrow P(2)+P(1)+P(3)$

(e)  $P(X \geq 1) = P(X \leq 4)$

(f)  $P(2 \leq X \leq 4) = \frac{6+4}{25} = \frac{10}{25}$   
 $\downarrow P(X=2)+P(X=3)$

(g) mode = 1 (it has the highest probability)

→ Distribution function = Cumulative Distribution Function

$$F(x_0) = P(X \leq x_0)$$

$$\therefore F(3) = P(X \leq 3)$$

Eg.

$$P(X=x) = \begin{cases} \frac{kx}{6} & \text{for } x=1,2,3 \\ \frac{k(7-x)}{6} & \text{for } x=4,5,6 \\ 0 & \text{otherwise} \end{cases}$$

\* Re-write this in the table format. It's always useful.

$x$	1	2	3	4	5	6
$P(X=x)$	$\frac{k}{6}$	$\frac{2k}{6}$	$\frac{3k}{6}$	$\frac{3k}{6}$	$\frac{2k}{6}$	$\frac{k}{6}$

$$\left( \underbrace{\frac{k}{6} + \frac{2k}{6} + \frac{3k}{6}}_{\hookrightarrow \frac{kx}{6}} + \underbrace{\frac{3k}{6} + \frac{2k}{6} + \frac{k}{6}}_{\hookrightarrow \frac{k(7-x)}{6}} \right) = 1 \Rightarrow \therefore k = \frac{1}{2}$$

$$* F(3) = P(X \leq 3) \Rightarrow P(X=1) + P(X=2) + P(X=3) = \frac{6}{12} = \frac{1}{2}$$

\* You can add a row for Cumulative distribution:

$x$	1	2	3	4	5	6
$F(x)$	$\frac{k}{6}$	$\frac{3k}{6}$	$\frac{6k}{6}$	$\frac{9k}{6}$	$\frac{11k}{6}$	$\frac{12k}{6}$

← This should be 1

\* It's just like Cumulative frequency.

# \* Discrete Uniform Distribution:

$\mu = \text{Mean} = \left( \frac{a+b}{2} \right)$  ; Variance =  $\left[ \frac{(b-a+1)^2 - 1}{12} \right]$   
 also same as Expected Value  $= \left( \frac{(n+1)(n-1)}{12} \right)$



## \* EXPECTED VALUES, MEAN & VARIANCE

Mean = Expected Value =  $E(X) = \sum x P(X=x)$

Variance =  $E(X^2) - (E(X))^2 = \sum x^2 P(X=x) - \left[ \sum x P(X=x) \right]^2$

Eg.  $X :$

1	2	3	4
$\frac{12}{25}$	$\frac{6}{25}$	$\frac{4}{25}$	$\frac{3}{25}$

\* Expected Value =  $E(X) = \sum x P(X=x) =$   
 $\left( 1 \times \frac{12}{25} + 2 \times \frac{6}{25} + 3 \times \frac{4}{25} + 4 \times \frac{3}{25} \right)$   
 $= \frac{48}{25} = 1.92$

\* Variance =  $\sum x^2 P(X=x) - \left( \sum x P(X=x) \right)^2 = E(X^2) - (E(X))^2$   
 $\Rightarrow \text{Var}(X)$   
 $\Rightarrow \sigma^2 = \frac{24}{5} - \frac{2304}{625} = \frac{696}{625} = \underline{1.1136}$

\* Expected Value & Variance formulas for functions

\*  $E(ax) = aE(x)$

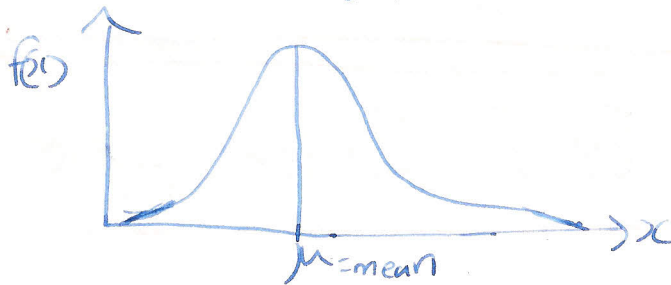
\*  $E(ax+b) = aE(x)+b$

\*  $Var(ax) = a^2 Var(x)$

\*  $Var(ax+b) = a^2 Var(x)$

## ⇒ THE NORMAL DISTRIBUTION

\* Continuous Distribution ⇒ Area = Probability  
↳ has no gaps

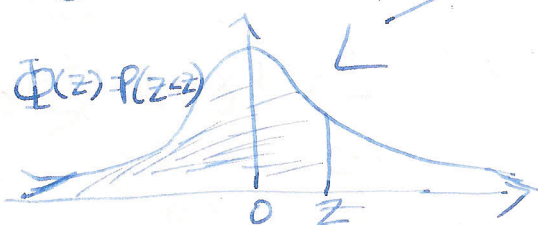


\* We 'Normalise' a variable to use the Z-tables.

$X \sim N(\mu, \sigma^2) \Rightarrow X$  is normally distributed with mean  $\mu$  & variance  $\sigma^2$

To normalise to Z

\*  $X \sim N(\mu, \sigma^2) \Rightarrow \frac{X-\mu}{\sigma} = Z, Z \sim N(0,1)$



Eg.  $X \sim N(5, 16)$

(a)  $P(X < 7)$

\* first 'Normalise'

$$N(0, 1^2) \sim \frac{X - \mu}{\sigma} = Z$$

$$\therefore P(X < 7) = P\left(Z < \frac{7-5}{4}\right) = P(Z < 0.5) = 0.6915$$

\* On the table, look for 0.50 under 'Z' and write the value opposite it, under ' $\Phi Z$ '

(b)  $P(X > 9)$

$$P(Z > a) = 1 - \Phi(a)$$

$$\therefore P(X > 9) = P\left(Z > \frac{9-5}{4}\right) = P(Z > 1) = 1 - P(Z < 1) = 1 - 0.8413 = 0.1587$$

\* Remember  $\Phi(1) = P(Z < 1)$

↳ Cumulative distribution function.

(c)  $P(5 < X < 11) = P\left(\frac{5-5}{4} < Z < \frac{11-5}{4}\right) = P(0 < Z < 1.5)$

$$P(a < Z < b) = \Phi(b) - \Phi(a) = P(Z < b) - P(Z < a)$$

$$= P(Z < 1.5) - P(Z < 0) = 0.9332 - 0.5 = 0.4332$$

$$* \Phi(-2) = 1 - \Phi(2)$$