

* EDA and Probability Theory :

① Mean : average, $\mu = \frac{1}{n} \sum_{i=1}^n x_i = \frac{x_1 + x_2 + \dots + x_n}{n}$

② Standard deviation, $\sigma = \sqrt{\text{variance}} = \sqrt{\sigma^2}$

variance $\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 = \frac{(x_1 - \mu)^2 + (x_2 - \mu)^2 + \dots + (x_n - \mu)^2}{n}$

③ Median : Given n observations $x_1, x_2, x_3, \dots, x_n$

sort them $\Rightarrow x_1 < x_2 < x_3 < \dots < x_n$

n is even

n is odd

Median = $\frac{\left(\frac{n}{2}\right)^{\text{th}} \text{ term} + \left(\frac{n}{2} + 1\right)^{\text{th}} \text{ term}}{2}$

Median = $\left(\frac{n+1}{2}\right)^{\text{th}} \text{ term}$

$x_1, x_2, x_3, x_4, x_5, x_6$

$x_1, x_2, \boxed{x_3}, x_4, x_5$

$\left(\frac{5+1}{2}\right) = 3^{\text{rd}} \text{ term}$

Median = $\frac{x_3 + x_4}{2}$

Median = x_3

④ Percentile

Quantiles

IQR

↓
divide given range
in 100 parts

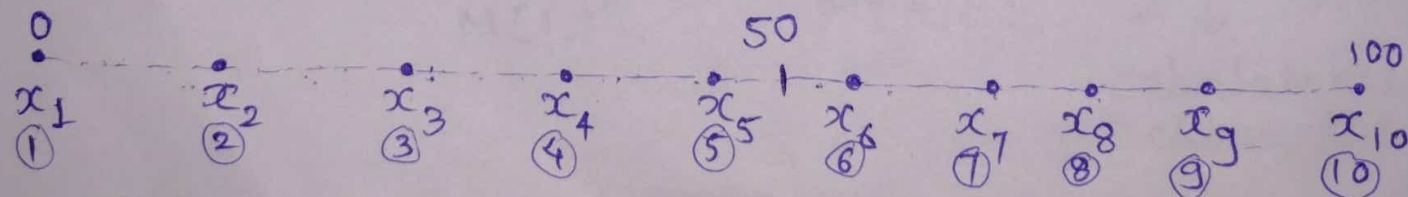
↓
25th, 50th, 75th
percentiles

↓
75th percentile
- 25th percentile

\Rightarrow 0th percentile is first term

\Rightarrow 50th percentile is Median

\Rightarrow 100th percentile is Last term



\Rightarrow p^{th} percentile term = $\frac{n \cdot p}{100} + 1 = \text{int} \frac{n \cdot p}{100} + \text{fraction}$

p^{th} percentile = $x_k + f(x_{k+1} - x_k)$ (Linear interpolation)

$$\left. \begin{aligned} Q_1 &= P(25) = 2.5 \\ Q_2 &= P(50) = \text{Median} = 5.5 \\ Q_3 &= P(75) = 7.5 \end{aligned} \right\} \text{Quantiles}$$

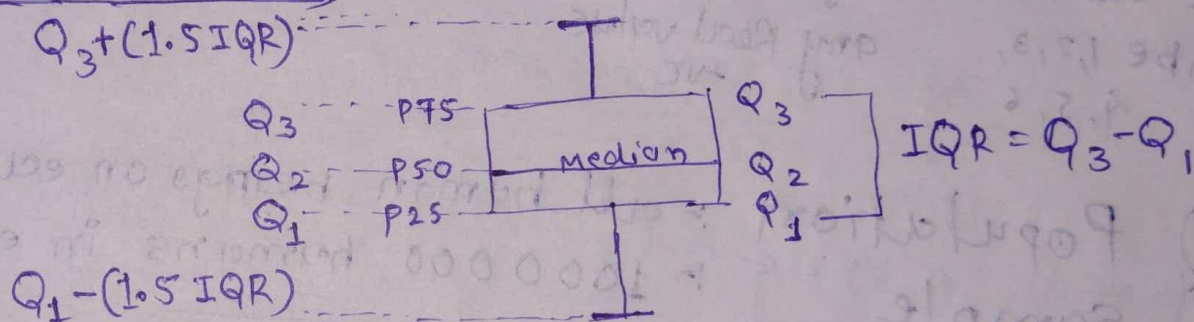
$$\Rightarrow \text{IQR} = Q_3 - Q_1 = 7.5 - 2.5 = 5$$

§ MAD : Median Absolute deviation

Let M be the Median of x_1, x_2, \dots, x_n

$$\begin{aligned} \text{MAD} &= \text{Median}(|x_i - M|) \\ &= \text{Median}(|x_1 - M|, |x_2 - M|, \dots, |x_n - M|) \end{aligned}$$

* Boxplot and whiskers



$$5^{\text{th}} \text{ per} = \frac{10-1}{\left(\frac{100}{50}\right)} + 1 = \frac{9}{2} + 1 = 5.5^{\text{th}} \text{ term} = 5^{\text{th}} + 0.5$$

$$\begin{aligned} \therefore p(50) &= x_5 + 0.5(x_6 - x_5) \\ &= 5 + 0.5(6 - 5) \\ &= 5.5 \end{aligned}$$

$$\Rightarrow p(25) = \frac{10-1}{\left(\frac{100}{25}\right)} + 1 = \frac{9}{4} + 1 = 3.25 = 3^{\text{rd}} + 0.25$$

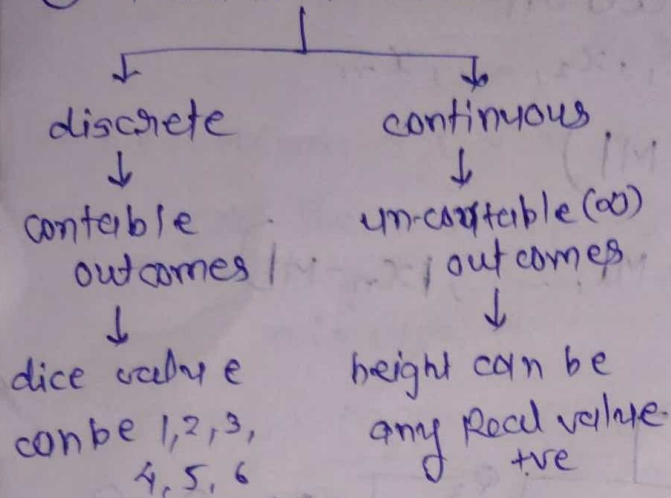
$$\begin{aligned} p(25) &= x_3 + 0.25(x_4 - x_3) \\ &= 3 + (0.25)(4 - 3) \\ &= 3.25 \end{aligned}$$

* Probability :-

① Experiment \rightarrow Rolling dice, Measure height, coin toss

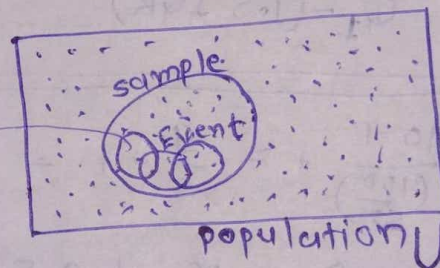
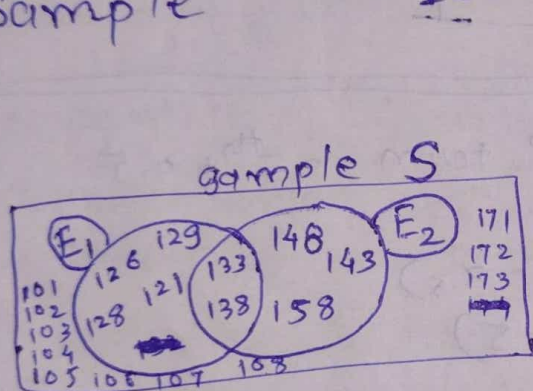
② Sample space $\rightarrow \{1, 2, 3, 4, 5, 6\}$, $R^+ - \{0\}$, $\{H, T\}$
 $\{1, 1.1, 2, 3, 26, 62.9, \dots\}$

③ Random variable $\rightarrow X = \text{Event} \rightarrow$ subset of S



getting dice value > 3
 getting height value $> 125\text{cm}$
 getting H two times $\rightarrow \{HH, HT, TH, TT\} = S$
 $X = \{HH\}$

④ Population : all human beings on earth
 ⑤ Sample : 1000000 humans in experiment



$$E_i \subseteq S \subseteq U$$

E_1 : height is between 120 to 140 , $\{120 \leq x \leq 140\}$

E_2 : height is between 130 to 160 , $\{130 \leq x \leq 160\}$

$$E_1 = \{121, 129, 126, 128, 133, 138\} \quad n(E_1) = 6$$

$$E_2 = \{133, 148, 138, 158, 143\} \quad n(E_2) = 5$$

Assume we have sample S of size 20 , $n(S) = 20 = n$

union $E_1 \cup E_2 = \{121, 129, 126, 128, 133, 138, 148, 158, 143\}$

$$n(E_1 \cup E_2) = 9$$

inter section $E_1 \cap E_2 = \{133, 138\}$

$$n(E_1 \cap E_2) = 2$$

$$n(E_1 \cup E_2) = n(E_1) + n(E_2) - n(E_1 \cap E_2)$$

$$E_1^c = S - E_1 = \{101, 102, 103, 104, 105, 106, 107, 108, 148, 143, 158, 171, 172, 173\}$$

$$n(E_1^c) = n(S) - n(E_1) = 20 - 6 = 14$$

→ Frequentist approach

1) Experiment is performed n times

2) $n(E)$: number of times event E occurred

$$\therefore P(E) = \lim_{n \rightarrow \infty} \frac{n(E)}{n}$$

In general $P(E) = \frac{n(E)}{n(S)}$

classical $\leftarrow \frac{\text{favourable outcome}}{\text{total outcomes}}$

$$\Rightarrow P(E_1) = P(120 \leq x \leq 140) = \frac{n(E_1)}{n(S)} = \frac{6}{20} = 0.3 = 30\%$$

$$\Rightarrow P(E_2) = P(130 \leq x \leq 160) = \frac{n(E_2)}{n(S)} = \frac{5}{20} = 0.25 = 25\%$$

$$\Rightarrow P(E_1 \cup E_2) = P(120 \leq x \leq 160) = \frac{n(E_1 \cup E_2)}{n(S)} = \frac{9}{20} = 45\%$$

$$= P(E_1) + P(E_2) - P(E_1 \cap E_2)$$

$P(E_1 \cap E_2)$ is also written as $P(E_1 E_2)$

$$= \frac{6}{20} + \frac{5}{20} - \frac{2}{20} = 45\%$$

$$\Rightarrow P(E_1^c) = 1 - P(E_1) = 1 - \frac{6}{20} = \frac{14}{20} = 70\%$$

* $P(E_1 \cup E_2) \leq P(E_1) + P(E_2)$ $\rightarrow \geq 0$ +ve value

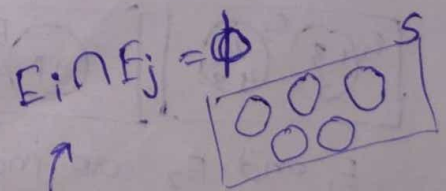
* $P(E_1 \cup E_2 \cup E_3) = P(E_1) + P(E_2) + P(E_3) - P(E_1 E_2) - P(E_2 E_3) - P(E_1 E_3) + P(E_1 E_2 E_3)$

AXIOMS OF PROBABILITY

① $0 \leq P(E_i) \leq 1$

② $P(S) = 1$

③ if E_1, E_2, \dots, E_k are mutually exclusive events then $P(E_1 \cup E_2 \cup E_3 \cup \dots \cup E_k) = \sum_i P(E_i) = P(E_1) + P(E_2) + \dots + P(E_k)$



* Conditional probability & Baye's Theorem

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

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

$$P(B|A) \cdot P(A)$$

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} \neq 0$$

OR

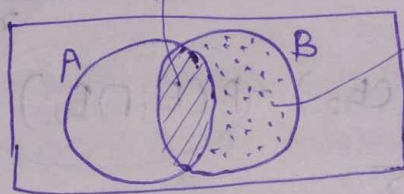
$$P(B|A) = \frac{P(A|B) \cdot P(B)}{P(A)} \neq 0$$

$$P(A|B) \neq P(B|A)$$

B has already occurred
Now what is $P(A)$?

$$\text{Now } P(B) = P(B \cap A) + P(B \cap A^c)$$

$$\begin{aligned} (B \cap A) \cup (B \cap A^c) &= B \cap (A \cup A^c) \\ &= B \cap \Omega \\ &= B \end{aligned}$$



$$\begin{aligned} \text{OR} \\ &= P(B) + P(A \cap A^c) \\ &= P(B) + 0 \\ &= P(B) \end{aligned}$$

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B \cap A) + P(B \cap A^c)}$$

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B|A) \cdot P(A) + P(B|A^c) \cdot P(A^c)}$$

No need to know $P(B)$

* Independent event

* Mutually exclusive event

Exp 1 \Rightarrow toss a coin

Exp 2 \Rightarrow throw a dice



E_1 and E_2 are independent

E_3 and E_2 are independent

E_1 : you get even number on dice

E_2 : you get H on coin

E_3 : you get odd number on dice

E_1 and E_3 are mutually exclusive events

$$P(E_1 \cap E_3) = 0$$

$$P(E_1 \cap E_2) = P(E_1) \cdot P(E_2)$$

$$P(E_1|E_2) = P(E_1)$$

* Random variable: Quantity/Specify event outcome

↳ Mapping from an event to R .

X : getting 4 or less (dice throw)

X : amount of chocolate eaten by students

X : Time spent of website

This random variable can take values

$\{x_1, x_2, x_3, \dots\}$ denoted by small x

Random variable

discrete

set of outcomes are countable

PMF

probability Mass Function

continuous

set of outcomes are not countable

PDF

probability density function

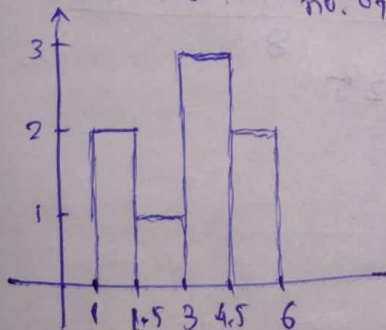
X : Number of pizza sold in a day (discrete X)

$X: \{1, 1, 2, 3, 4, 4, 5, 6\}$

X : diameter of Mango in inch (continuous X ; it can be $2.1, 2.2223, 4.599, \dots, \infty$)

Histogram: counting results in bin range

$$\text{bin range} = \frac{\text{max} - \text{min}}{\text{no. of bins}} = \frac{6 - 1}{4} = 1.5$$



$$1 \leq x < 1.5 \Rightarrow 2 \quad 2/8$$

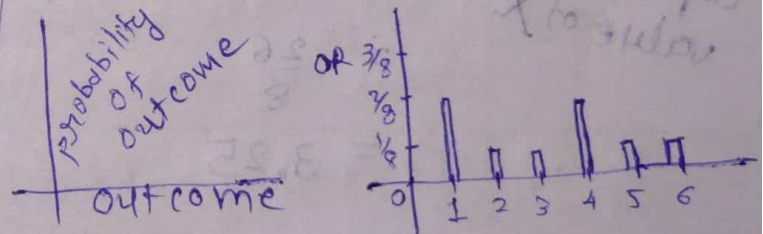
$$1.5 \leq x < 3 \Rightarrow 1 \quad 1/8$$

$$3 \leq x < 4.5 \Rightarrow 3 \quad 3/8$$

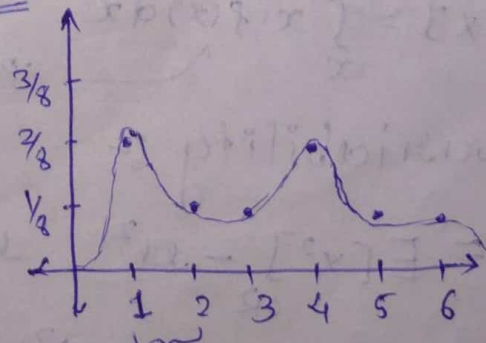
$$4.5 \leq x \leq 6 \Rightarrow 2 \quad 2/8$$

count probability

PMF: discrete



PDF: continuous



value between 1 and 2 is also possible for continuous RV X :

OR

x	1	2	3	4	5	6
$P(x)$	$2/8$	$1/8$	$1/8$	$2/8$	$1/8$	$1/8$

PDF or PMF = $f(x=a) = p(a)$ OR $p(x) = p(x)$

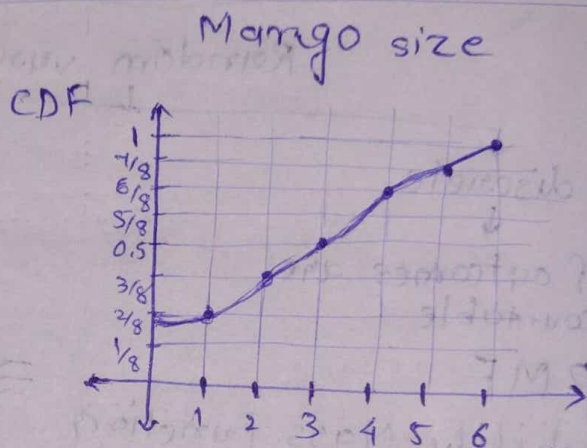
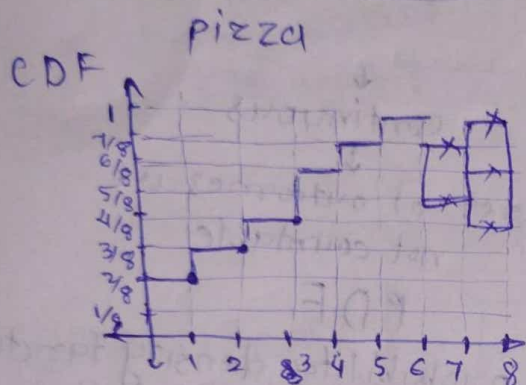
CDF
cumulative
distribution
function

$$F_x(x \leq a) = F(a) = \sum_{x \leq a} p(x=x) = \sum_{x \leq a} p(x)$$

$x=a$	1	2	3	4	5	6
$n(x \leq a)$	2	3	4	6	7	8
$p(x)$	$\frac{2}{8}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{2}{8}$	$\frac{1}{8}$	$\frac{1}{8}$
$F(x)$	$\frac{2}{8}$	$\frac{3}{8}$	$\frac{4}{8}$	$\frac{6}{8}$	$\frac{7}{8}$	$\frac{8}{8}$

$\{1, 1, 2, 3, 4, 4, 5, 6\}$

← cumulative probability



* Expectation :- what is the expected number of pizza you will sell today? (average) (mean)

$$E[x] = \sum_x x \cdot p(x) \quad \text{PMF/PDF } p(x=x)$$

Expected
value of X

$$= 1 \cdot \left(\frac{2}{8}\right) + 2 \cdot \left(\frac{1}{8}\right) + 3 \cdot \left(\frac{1}{8}\right) + 4 \cdot \left(\frac{2}{8}\right) + 5 \cdot \left(\frac{1}{8}\right) + 6 \cdot \left(\frac{1}{8}\right)$$

$$= \frac{26}{8}$$

$$= 3.25$$

$$\text{OR } \mu = \frac{1+1+2+3+4+4+5+6}{8} = 3.25$$

in case of continuous variable

$$E[x] = \int_x x \cdot f(x) dx$$

PDF

* Variance, variability :-

$$\text{Var}(X) = E[x^2] - \mu^2$$

OR

μ is mean of $E[x]$

$$= E[x^2] - (E[x])^2$$

* Distributions :- distribution is nothing but to see and model how probability of an event is distributed i.e follows trend across given range.

Distribution	PDF/PMF $f(x)$	CDF $F(x) = P(X \leq x)$	Expectation	Variance
Bernoulli(p)	$\begin{cases} p, & x=1 \\ 1-p, & x=0 \end{cases}$	$\begin{cases} 1-p, & x \leq 0 \\ 1, & x > 1 \end{cases}$	p	$pq = p(1-p)$
Binomial(n, p)	$\binom{n}{k} p^k (1-p)^{n-k}$	$\sum_{i=0}^k \binom{n}{i} p^i (1-p)^{n-i}$	$n \cdot p$	$n \cdot p(1-p)$ $= npq$
UNIFORM(a, b)				
continuous	$\begin{cases} \frac{1}{b-a} & x \in [a, b] \\ 0 & \text{o.w.} \end{cases}$	$\begin{cases} 0 & x < a \\ \frac{x-a}{b-a} & x \in [a, b] \\ 1 & x > b \end{cases}$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$
discrete	$\begin{cases} \frac{1}{n} & x \in [a, b] \\ 0 & \text{o.w.} \end{cases}$	$\begin{cases} 0 & x < a \\ i/n & x \in [a, b] \\ 1 & x > b \end{cases}$	$\frac{a+b}{2}$	$\frac{(n^2-1)}{12}$
Poisson(λ)	$\frac{\lambda^k e^{-\lambda}}{k!}$	$e^{-\lambda} \sum_{i=0}^k \frac{\lambda^i}{i!}$	λ	λ
Exponential(λ) $\beta = 1/\lambda$	$\begin{cases} \lambda \cdot e^{-\lambda x}, & x \geq 0 \\ 0 & x < 0 \end{cases}$	$\begin{cases} 1 - e^{-\lambda x}, & x \geq 0 \\ 0 & x < 0 \end{cases}$	$\beta = 1/\lambda$	$\beta^2 = 1/\lambda^2$
	$P_n(X > s+t X > s) = P_n(X > t)$			
Normal/Gaussian (μ, σ^2)	$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$	$\int_0^x f(x) dx$	μ	σ^2
Log-Normal (μ, σ^2)	$\frac{1}{x\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{\ln x - \mu}{\sigma}\right)^2}$	$\int_0^x f(x) dx$	$e^{\left(\mu + \frac{\sigma^2}{2}\right)}$	$(e^{\sigma^2} - 1) \cdot e^{(2\mu + \sigma^2)}$
Pareto(x_m, α) scale \uparrow shape \uparrow	$\frac{\alpha \cdot x_m^\alpha}{x^{\alpha+1}}$	$1 - \left(\frac{x_m}{x}\right)^\alpha$	log normal $\xrightarrow{\log}$ normal Pareto $\xrightarrow{\text{Box-Cox}}$	

① For discrete RV $E(X) = x_1 p_1 + x_2 p_2 + \dots + x_n p_n$ $\sum p_i = 1$

② For continuous RV $E(X) = \int x \cdot f(x) dx$

p_i is prob.
 $\frac{f(x)}{\text{PDF}}$ is prob dist fun

$$V(X) = E[X^2] - E[X]^2$$

$$\sigma^2 = E[X^2] - \mu^2$$