

02-08-2025

Agenda:

- Bayes theorem
- PDF, PMF, CDF
- Skewness
- Correlation and covariance

A = patient has the disease

B = patient test positive

Traditional prob.

Q: what is the prob of getting a positive test if a person has the disease?

$$P(B|A)$$

→ 99% of people with the disease test positive. (data-ground truth)

$P(B|A) = 0.99$ (if someone has disease, the test will be positive 99% of time)

→ Bayesian view: $P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$

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

- \rightarrow 1% of population has the disease $\rightarrow P(A) = 0.01$
 $\rightarrow P(B|A) = 0.99$ (test detect disease 99% of time)] TP
 $\rightarrow P(B|A') = 0.01$ (1% false positive rate)] FP

Q. The patient tested positive, what is the prob that they actually have the disease?

test outcome \rightarrow Disease likelihood (Bayes)

Disease likelihood \rightarrow test positive (traditional)

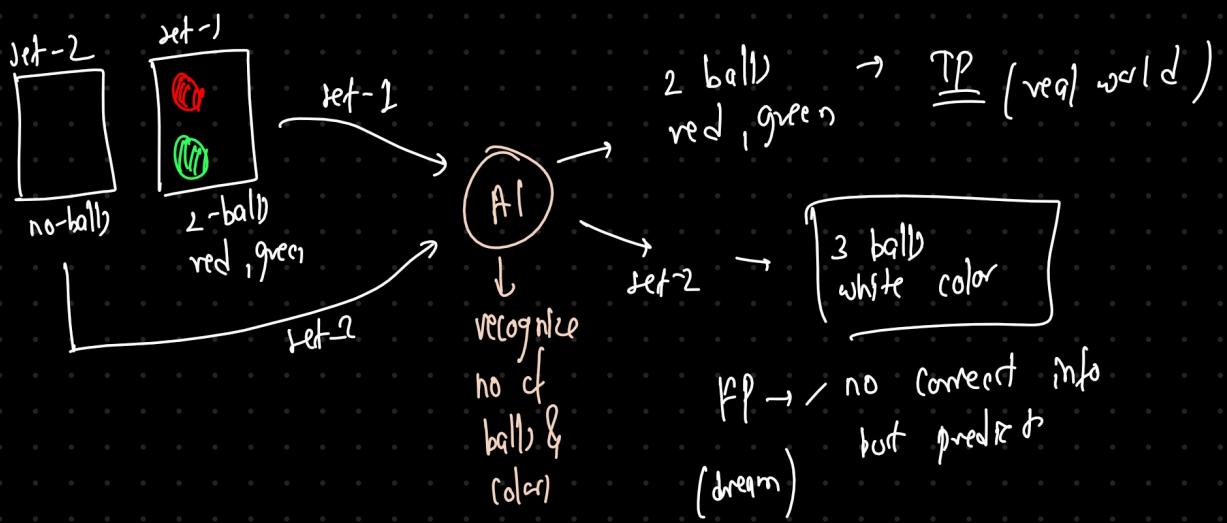
$$P(A|B) = \frac{0.99 \cdot 0.01}{0.0198} = \frac{0.0099}{0.0198} = 0.5$$

Even the test 99%, there only 50% chance the person actually has the disease, even after testing positive.

Recap:

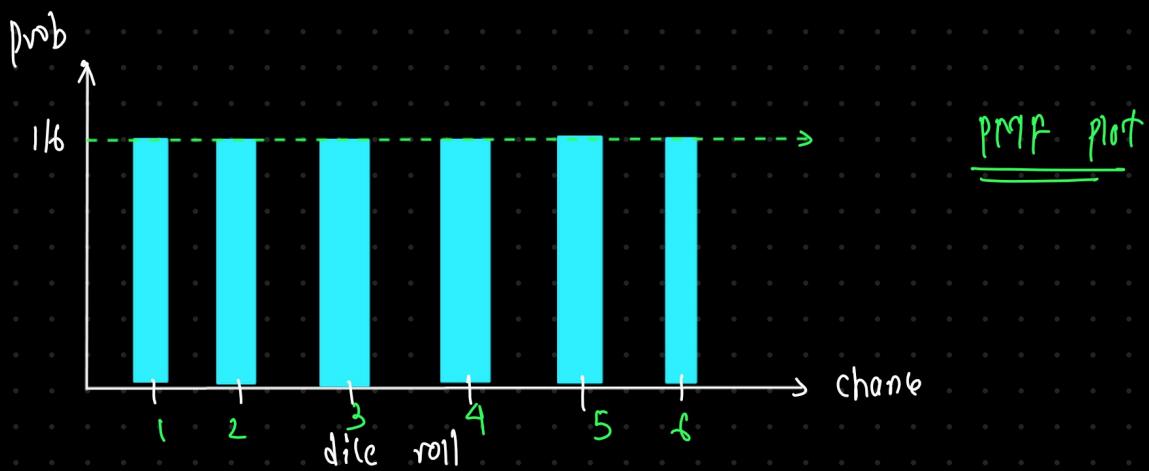
Traditional: — what is the chance of positive test if person is sick?

Bayes: — Given a positive test, how likely is the person sick?



PMF (Probability Mass Function)

- It is for discrete data (like dice rolls, coin flip, count)
- gives the prob. of each possible value in a dataset.



PMF is a function, Not just a plot

sum of prob. = 1

why: what's the chance of X happening exactly?

↳ 1, 2, 3, 4 → discrete data

1-2, 2-4 → cannot be range

→ PMF is not a plot, it's a function.

→ pmf function is the math rule behind the plot, the plot just visualize it,

marks of C14
↓
↓
↓
C14

Concept → pmf

CH 01

From → Data

enum

probably with their age

$$\begin{array}{l} \text{40: } 20\% \\ \text{60: } 40 \\ \text{80: } 30\% \end{array}$$

↓

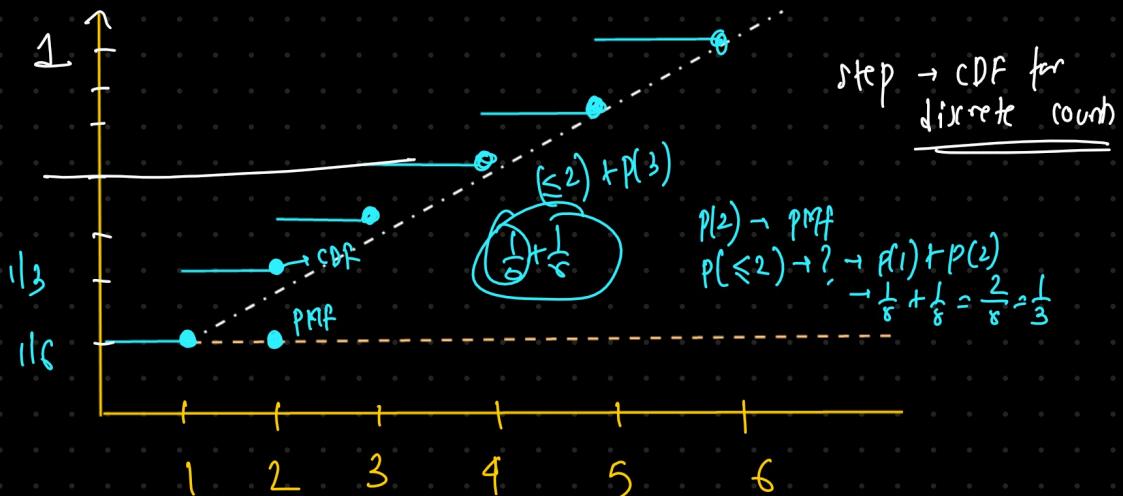
CDF (Cumulative Distribution function)

→ works for both discrete & continuous data,

- It gives the probability of $X \leq$ a value.

what is the chance of X being less than or equal to 9,

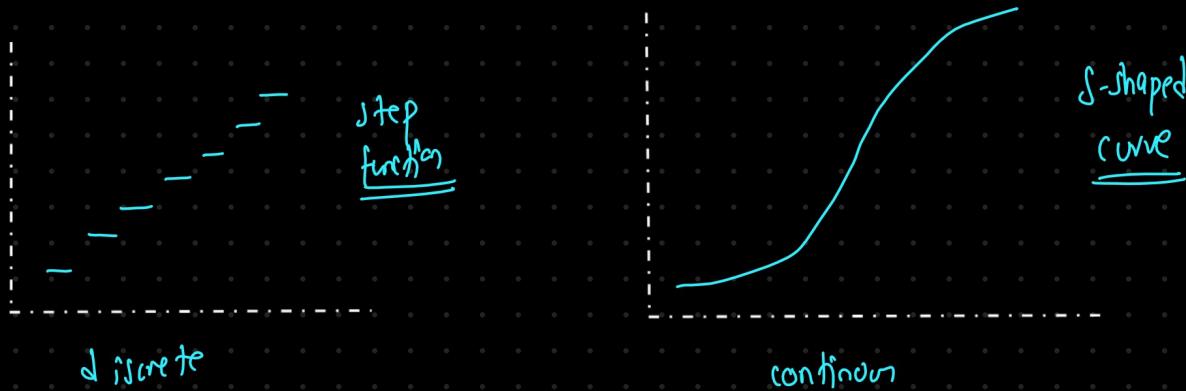
$\rightarrow P(\leq 4)$ on a die!



0 :	
1 :	
3 :	
4 :	20
5 :	10

0: 10 + Nothing \rightarrow 10
 1: 20 + 10 \rightarrow 30
 3: 10 + 20 + 10 \rightarrow 40
 4: 20 + 10 + 20 + 10 \rightarrow 60
 5: 60 + 10 \rightarrow 70

CDF always starts with 0 and ends at 1



PMF \rightarrow probability of outcome

CDF \rightarrow cumulative probability

PDF (Probability Density function)

- PDF is for continuous data

- shows the relative likelihood of a range of value (not exact point) \curvearrowright why?

why?

\rightarrow what is the chance of X falling between a & b ?

\rightarrow what's the prob. of a student height in between 160cm & 170cm?

why PDF over range (Not exact value)

→ continuous data has infinite precision.

1 180.00 cm ?, prob

160,000 | 160,00002 | 160,0000000009 | 160,X
↓ | ↓ | possibilities

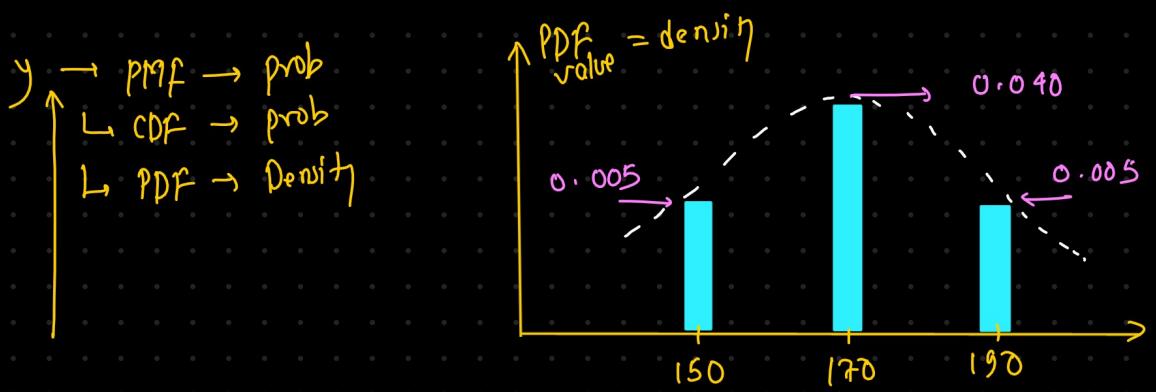
$$P(X) = \frac{\text{No. of favourable outcome}}{\text{Total possible outcome}} = \frac{1}{\infty} = 0$$

range → for continuous data

continuous prob are only meaningful over interval (a, b)

160cm 170cm
90%

prob. of an exact point is 0 because there's no "width" the interval,



Higher PDF value = More likely
to find value nearby

Step 1: collect data

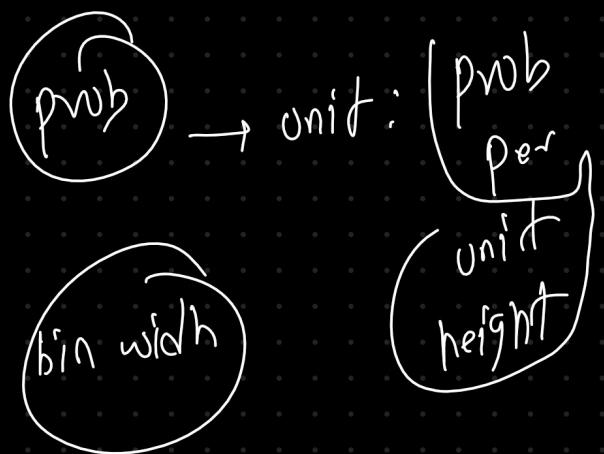
Height (in cm) of 20 people:

[162, 168, 171, 154, 175, 178, 165, 169, 172, 160, 167, 173, 170,
178, 169, 179, 166, 170, 174, 168] $\rightarrow \text{len}() \rightarrow 20$

Step 2: Create a histogram.

Height (in cm)	frequently	bin-width = 5
150 - 155	1	
155 - 160	0	
160 - 165	3	
165 - 170	6	
170 - 175	6	
175 - 180	4	

$$\text{PDF(Density)} = \frac{\text{frequency of bin}}{\text{Total observation (Q11)}} \xrightarrow{\text{Bin width}} \text{bin width}$$



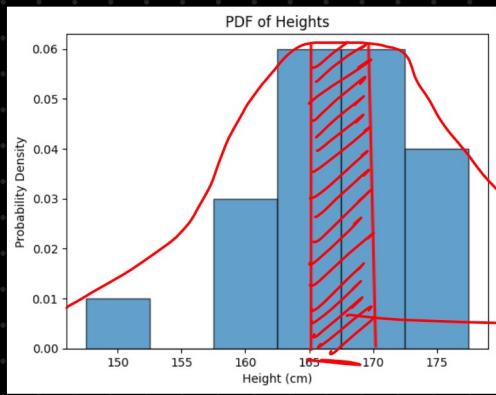
Step 3 : Normalize the Histogram

Height (in cm)	frequency	PDF
150 - 155	1	$(1/20 \times 5) = 0.01$
155 - 160	0	0
160 - 165	3	0.03
165 - 170	6	0.06
170 - 175	6	$(6/20 \times 5) = 0.06$
175 - 180	4	0.04

$$\text{PDF} = \frac{\text{Freq. of bin}}{\frac{\text{Total observation}}{\text{Bin width}}} \rightarrow \frac{\text{Freq. of bin}}{\text{Total observ} \times \text{Bin width}}$$

→ Dividing a number is equivalent to multiplying by its reciprocal.

Density = probability per unit height



→ plot

$$0.06 \times 5 = 0.30$$

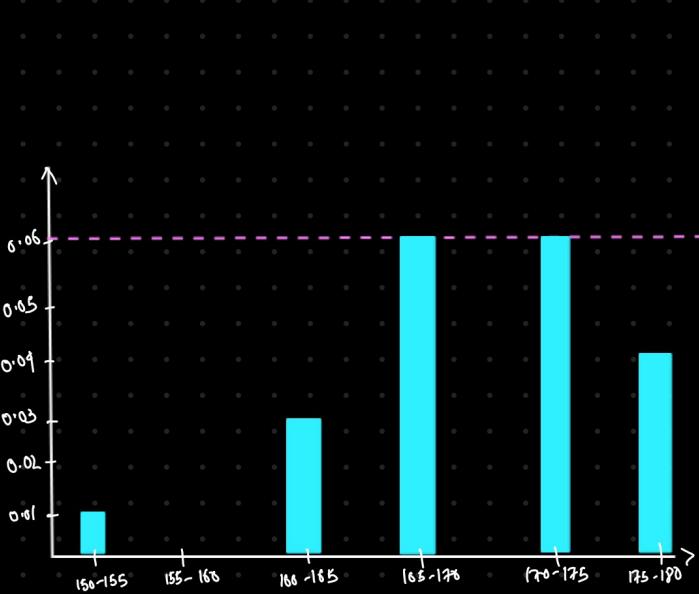
Step 4: calculate area under PDF curve: $\rightarrow \text{height} \geq 165 \& \leq 170$

e.g., find $(P(165 \leq \text{height} \leq 170))$

(1) Identify bin involved = 0.06 (Density)

$$\underline{\text{prob}} = \text{Area} = \text{PDF value} \times \text{Bin width} = 0.06 \times 5 = 0.30 \quad (\text{Ans. chart})$$

$$\underline{\text{Verification}} = 165 - 170 \rightarrow \frac{5}{20} = 0.30$$



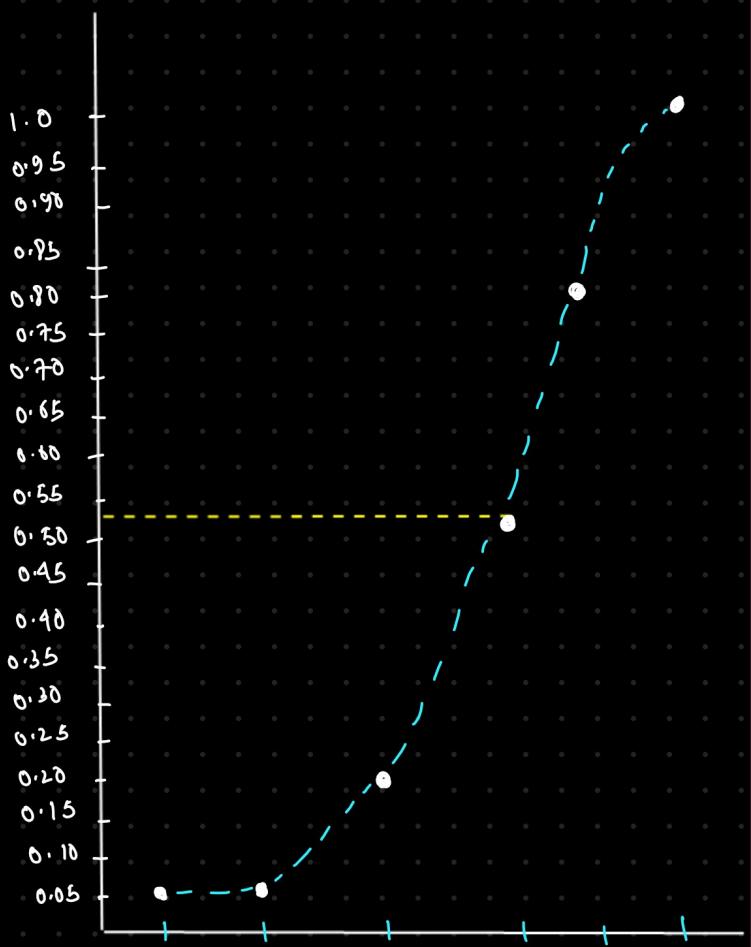
$$0.01 \times 5 = 0.05 \quad (\text{Prob of bin})$$

$$0 \times 5 = 0$$

$$0.03 \times 5 = 0.15$$

$$0.06 \times 5 = 0.3 \quad , \quad 0.06 \times 5 = 0.3$$

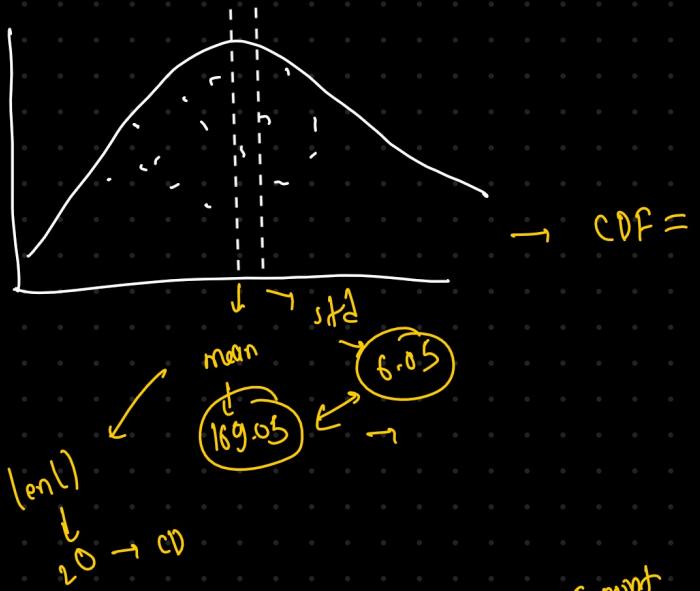
$$0.04 \times 5 = 0.20$$



CDF \rightarrow Area \times width

CDF \rightarrow Area

PDF gives Density



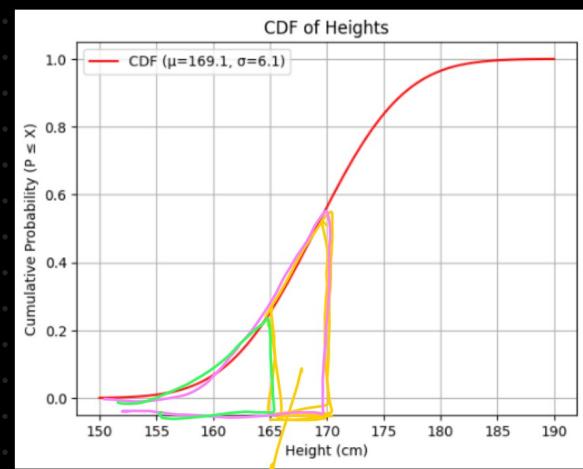
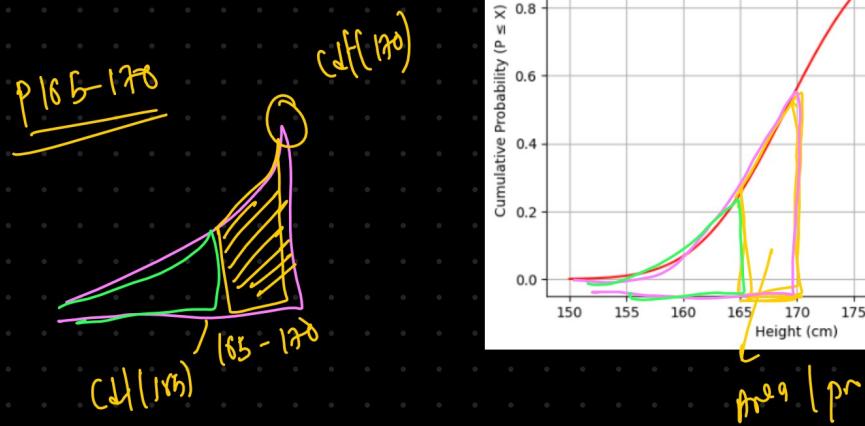
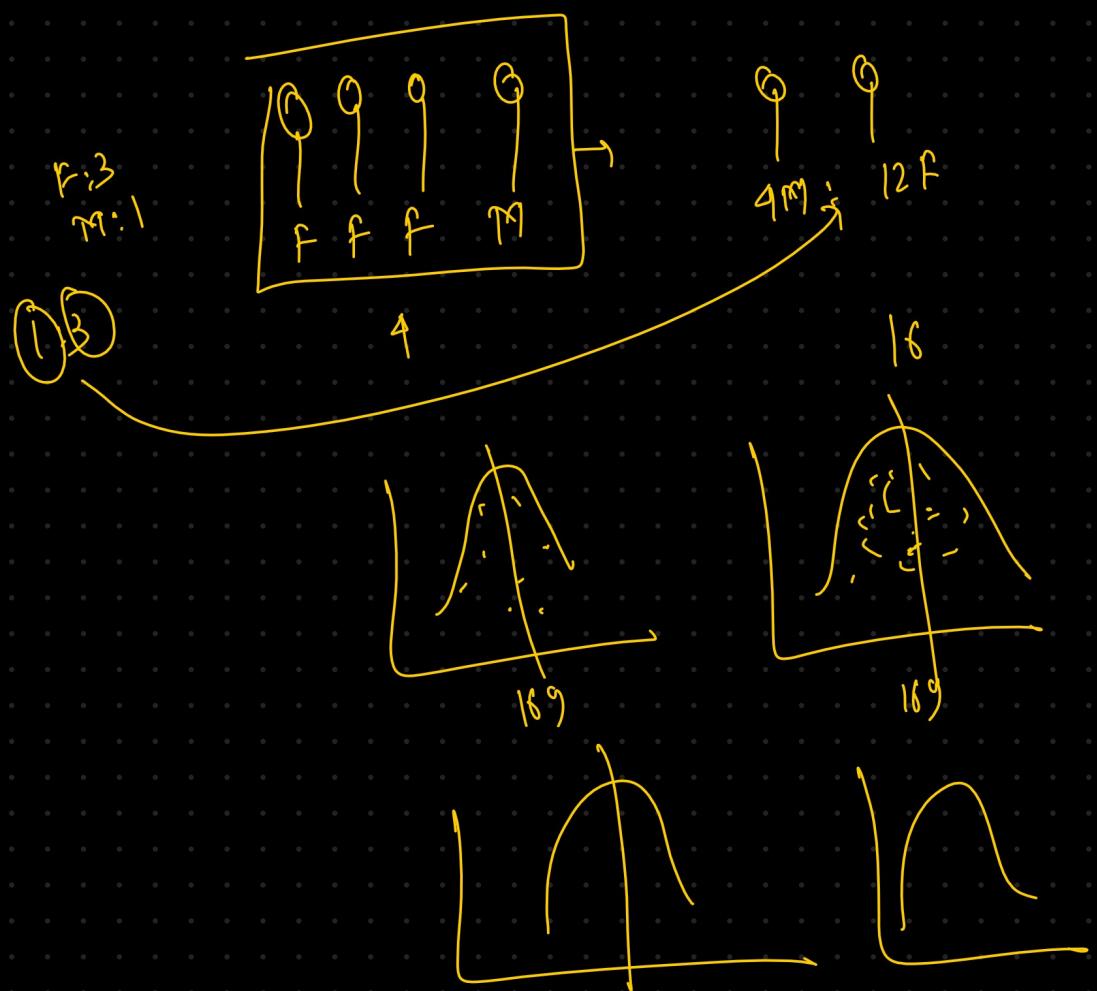
$M \rightarrow$ 5 point
 $M \rightarrow$

\rightarrow 7 point

\rightarrow 100 point

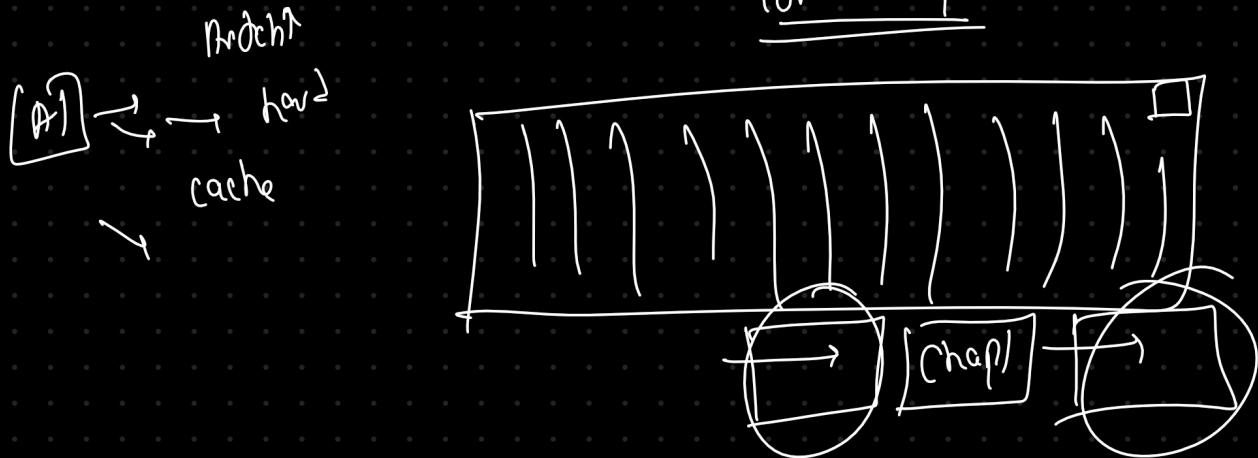
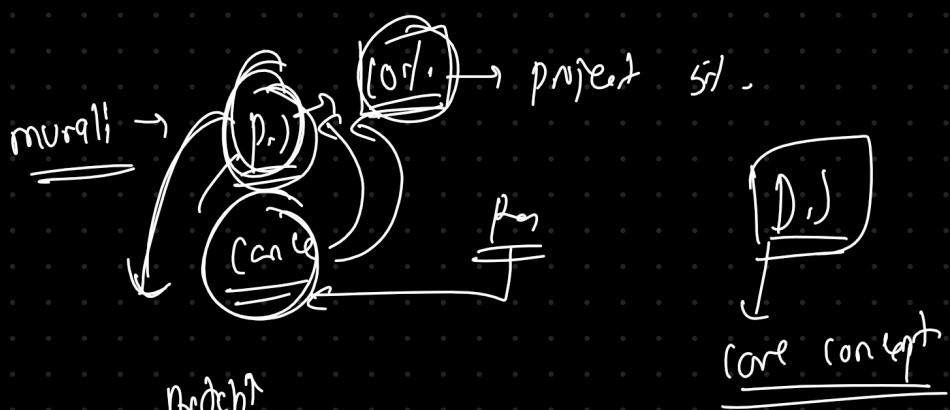
M







$\text{PMF} \rightarrow \text{discrete prob}$
 $\text{CDF} \rightarrow \text{cum. prob} \rightarrow \begin{matrix} \text{discrete} \\ \text{continuous} \end{matrix}$
 $\text{PDF} \rightarrow \text{prob of range of value} \rightarrow \begin{matrix} \text{CDF(PDF)} \\ \downarrow \\ \text{Area under curve} \end{matrix}$



1 1 1 1 1 1 1 [A]

