

Central limit Theorem

DATE / /

The central limit theorem can also be explained as the "distribution of sample mean" which approximated the Normal distribution when this is apply, when sample size larger.

here is the assumptions all the samples are similar and shape of population distribution could anything.

e.g - let assume that there are ten teams of Cricket in your school, every team has a total 100 students in it.

If we want to measure the average height of all the students in the ~~stop~~ sports team, then that would be humongous task.

1- One way, go to the about this task and find of all the students height and divide by total no of student

$$\text{Mean} = \frac{\text{All Students height}}{\text{Total No. of student}}$$

but if 2000 student then?

every time not possible to find the student height of All students

The Central Limit theorem come into the picture

— Take a random 30 sample of every team-team.

— Take all the sample and try to find the "Mean" for every single sample

— Once done, then find the mean of Means of the samples

— The value recient approx height of All the student

20 Team, every Team has ~~100~~ 100 students
 $20 \times 100 = 2000$ students

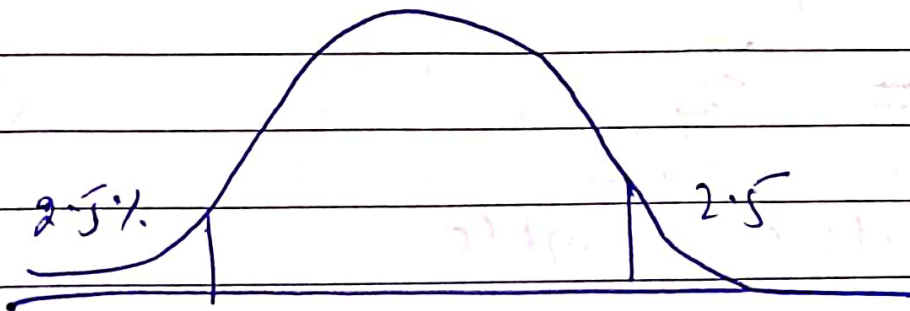
— Now — Take a sample of

every Team at list 30 or more
at this status we take 30
(Random Sample)

$$30 \times 20(\text{Team}) \\ = 600 (\text{sample})$$

$$= \frac{600}{2000} \quad \left(\begin{array}{l} 600 = \text{Total Sample} \\ 2000 = \text{Total No. of Sample} \end{array} \right)$$

$$= .3 = 30\%$$



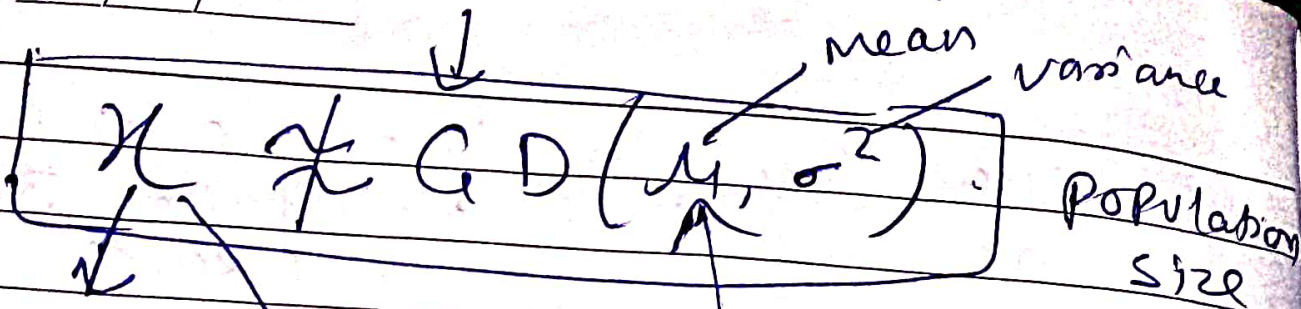
Mean \rightarrow Average value

Standard deviation \rightarrow is a measure how spread out of the values are

SD increase \rightarrow Normal distribution curve wider

DATE / /

Gaussian distribution



Random variable

We took some ~~set~~ Random Sample

$$S_1 = x_1 \dots x_{30}$$

$$= \bar{x}_1$$

Sample size ≥ 30
($n \geq 30$)

CLT

Take another Random Sample

$$S_2 = x_1, x_2, \dots, x_{30} = \bar{x}_2$$

$$S_3 = x_1, x_3 \dots x_{30} = \bar{x}_3$$

⋮

$$S_{100} = x_1 \dots x_{30} = \bar{x}_{100}$$

100 Samples

Mean of every Random Sample

Now

$$X \approx GD\left(\mu, \frac{\sigma^2}{n}\right)$$

Where the sample mean = population mean

$$\bar{x}_1 + \bar{x}_2 + \bar{x}_3 + \dots + \bar{x}_{100} = \mu$$

So

$GD(\mu \text{ is same})$

it mean whatever random sample mean would be same sample for population too.

Sample mean = population
Variance is different = σ^2

$\frac{\sigma^2}{n}$ (Total
no. of
sample
size)

Imp

once this plot in histogram
then it would convert in
Normal distribution like
"bell curve"

Cumulative Distribution function

- PMF (probability distribution function) is one way to describe the distribution of "discrete random variable"
- PMF can not defined for "continuous random variables"
- CMF (Cumulative Distribution function)
This is another method describe "distribution of random variable"
- Advantage of CMF, it can defined any kind of random variable (discrete, continuous, and mixed)

CDF formula \Rightarrow

CDF is real value random variable
 X is function given by

$$F_X(x) = P(X \leq x)$$

where $F_X(x)$ = function of x

X = real value variable

P = probability that X will have a value less than or equal to

$$F_X(x) = \int_{-\infty}^x f_X(t) dt$$

in case of random variable X which has distribution having a discrete component at a value b

$$P(X=b) = F_X(b) - \lim_{z \rightarrow b^-} F_X(z)$$

where f_X is continuous at b .

DATE

Probability Distribution function

Discrete

probability mass function
PMF

Continuous

probability density function
PDF

Cumulative distribution function
CDF

PDF (probability distribution function)

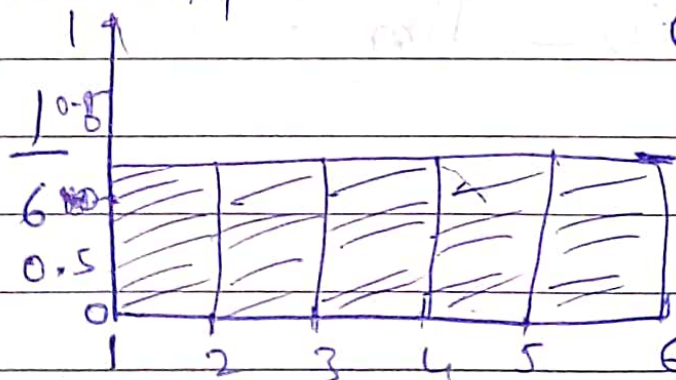
1 - probability mass function

2 - probability density function

3 - Cumulative distribution function

① Probability mass function \rightarrow
(PMF)

Probability



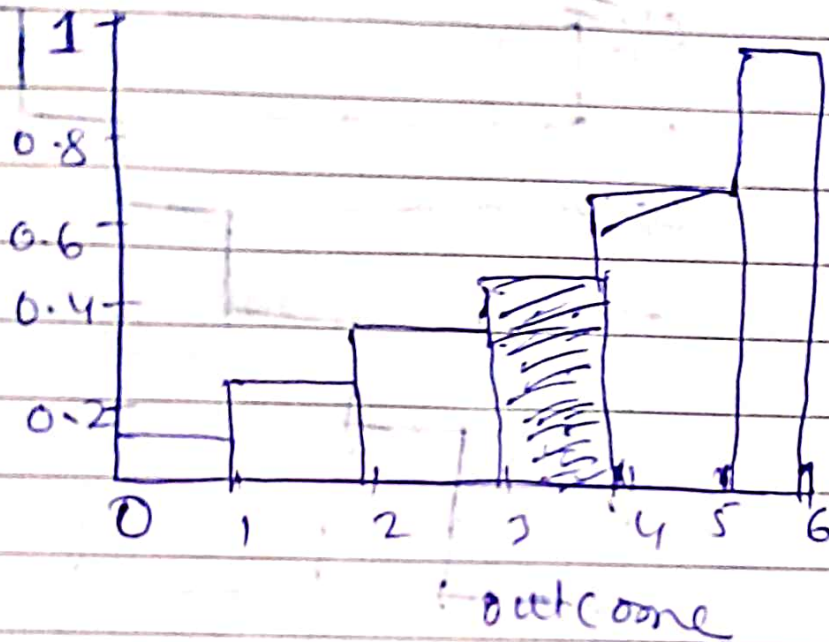
dice each $\frac{1}{6} = 0.167$
(equal chance) (Approx)

discrete

because why discrete? it outcome is 1, 2, 3

CDF

CDF →



$$P(X \leq 4) = P(X=1) + P(X=2) + P(X=3) + P(X=4)$$

~~Rolling Probability~~

Why doing PDF and CDF?

both the part of probability function.

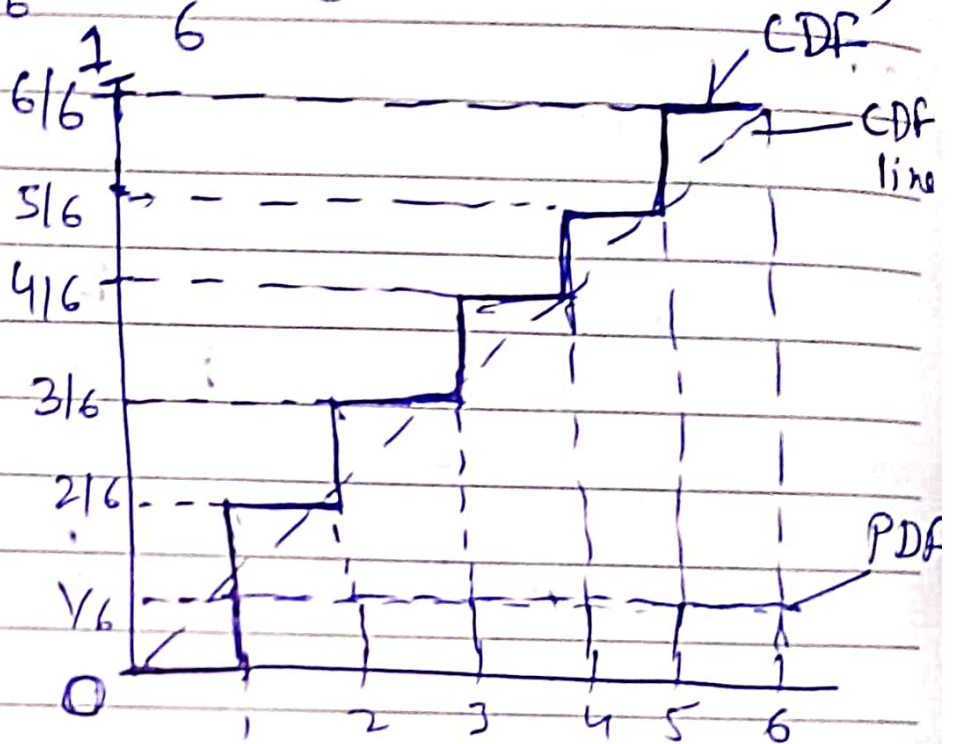
Probability distribution



DATE / / $F_x(x)$ & Capital F)CDF \rightarrow denoted by $F_x(x) = P(X \leq x)$

if we have dice,

$$P_{1,2,3,4,5,6} = \frac{1}{6} \quad (\text{probability occurs})$$

Calculate

$$F_x(0) = 0 \quad \neq \text{PDF is straight line}$$

$$F_x(1) = P(1)$$

$$F_x(2) = P(1) + P(2)$$

$$F_x(3) = P(1) + P(2) + P(3)$$

$$F_x(4) = P(1) + P(2) + P(3) + P(4)$$

$$F_x(6) = P(1) + P(2) + P(3) + P(4) + P(5) + P(6)$$

It is Cumulative Addition

$$F_X(0) = 0$$

$$F_X(1) = P(1) = \frac{1}{6}$$

$$F_X(2) = P(1) + P(2) = \frac{1}{6} + \frac{1}{6}$$

$$F_X(3) = P(1) + P(2) + P(3) \\ = \frac{1}{6} + \frac{1}{6} + \frac{1}{6}$$

$$F_X(3) = \frac{3}{6}$$

$$F_X(6) = P(1) + P(2) + P(3) + P(4) + P(5) + P(6) \\ = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6}$$

$$F_X(6) = \frac{6}{6}$$

Properties of CDF \rightarrow

1) CDF is Always greater than equal to zero

$$F_X(x) \geq 0$$

2) Range of CDF is

$$0 \leq F_X(x) \leq 1$$

CDF

is 0 to 1 (Range)

DATE / /

3) It is Always increasing function
min 0 to max 1.

PDF

denoted $f_x(x)$ (small f)

$$f_x(x) = P(x \leq x)$$

CDF is Cumulative ^{Addition} ~~Addition~~ of probability
PDF is Magnitude ^{Addition} ~~Addition~~ of probability

if we want PDF

$$f_x(x) = P(x \leq x)$$

$$P(5) = F_x(5) - F_x(4)$$

$$f_x(x) = \frac{dF_x(x)}{dx}$$

 Relation of PDF

$$F_x(x) = \int_{-\infty}^{+\infty} f_x(x) dx$$

 CDF

This
is the
Relation
of
CDF
&
PDF

~~what is~~Properties of PDF

1) PDF range Always 0 to 1

$$0 \leq f_x(x) \leq 1$$

2) Total Area of PDF = 1

$$1 = \int_{-\infty}^{\infty} f_x(x) dx$$

e.g. →

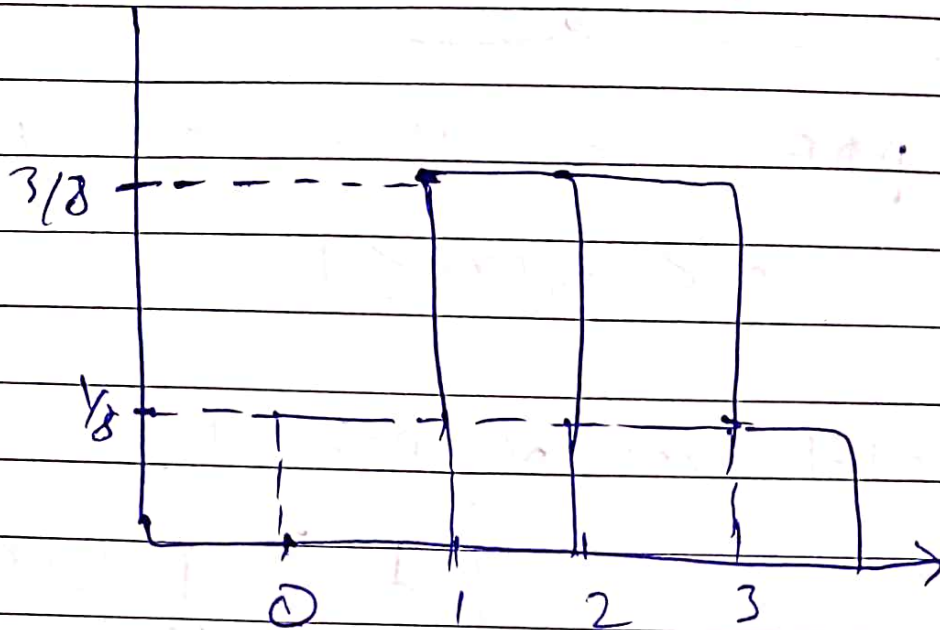
In an experiment, a trial consists of three successive tosses of a coin. If we define random variable x as the no. of heads appearing in trial, determine $f_x(x)$ and $F_x(x)$ (PDF and CDF)

R.V	$f(x)$		Toss to get head
0	$1/8$	TTT	0 0 0 = No heads = 0
1	$3/8$	TTH	0 0 1 = One head = 1
		THT	0 1 1 = Two heads = 2
2	$3/8$	TTH	1 1 1 = Three heads = 3
3	$1/8$	TTT	
		HTT	
		HHT	
		HHH	

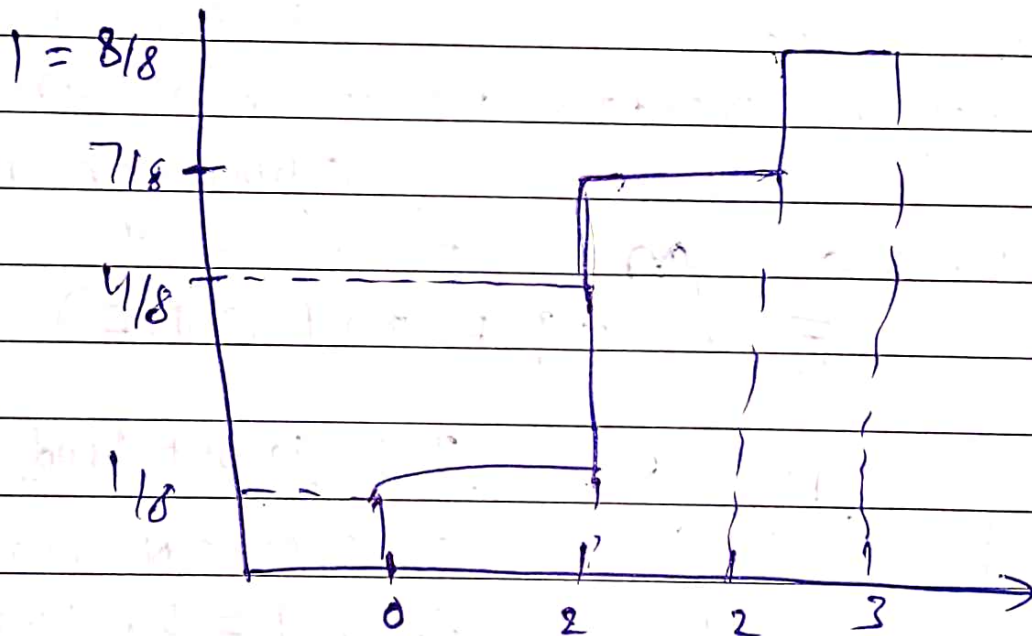
mean 4 Random Variable

DATE / /

Draw PDF



Draw CDF



Statistics Estimation

↓
Parametric



PDF, Normal distribution,
CDF

↓
Non-Parametric

↑
Kernel
density
function

Kernel density function

Formula

$$\hat{f}(x; h) = \frac{1}{nh} \sum_{i=1}^n K \left\{ \frac{(x - X_i)}{h} \right\}$$

where K = Kernel function (non-negative)

$h \rightarrow$ smoothing parameter (bandwidth)

X_i = given value of x

↓
Smoothing the
Curve

DATE / /

~~RDE & Histogram~~
~~and kernel~~

Bandwidth (Smoothing Parameter)

Smoothing of Curve

- 1) Cross validation
- 2) Normal optimal Smoothing

Kernel density estimates are closely
related to Histogram