

5/1/2021

* Random Experiment:

→ An experiment in which all possible outcomes are known in advance but we don't know which possible outcome occurs.

Eg 1: Tossing a coin

$$\{H, T\}$$

2: Throwing a die

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

3: Drawing a card from pack of cards

52 cards

Ques: * Sample Space

* → Set of all possible outcomes of a Random Experiment.

Eg: 1 Tossing a coin

$$S = \{H, T\}$$

2: Tossing a pair of coins

$$S = \{HH, HT, TH, TT\}$$

Note:

→ If we toss 'n' coins number of possible outcomes = 2^n

3: Tossing 3 coins

$$2^3 = 8 \text{ Possibilities}$$

$$\frac{3}{2} = 8 \quad H \quad H \quad H \quad \checkmark \quad 1^{\text{st}} \text{ Sample point } 2$$

$$\frac{8}{2} = 4 \quad H \quad H \quad T$$

$$H \quad T \quad H$$

$$H \quad T \quad T$$

$$\frac{4}{2} = 2 \quad T \quad H \quad H$$

$$T \quad H \quad T$$

$$\frac{2}{2} = 1 \quad T \quad T \quad H$$

$$T \quad T \quad T - 8^{\text{th}}$$

∴ Total 8 possibilities

8 possible outcomes

4) 4 coins

$$2^4 = \frac{16}{2} = 8H \quad 8T$$

$$\frac{8}{2} = 4H \quad 4T$$

$$\frac{4}{2} = 2H \quad 2T$$

$$\frac{2}{2} = 1H \quad 1T$$

H H H H \checkmark 1st Possibility

H H H T

H H T H

H H T T

H T H H

H T H T

H T T H

H T T T

T H H H

T H H T

T H T H

T H T T

T T H H

T T H T

T T T H

T T T T

$\checkmark 16^{\text{th}}$

S = {HHHH, ..., TT TT}

3

II) Throwing a dice → Random experiment

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

Throwing a pair of dice

$$S = \{(1,1), (1,2), (1,3), \dots, (1,6)\}$$

$$(2,1), (2,2), \dots, \dots, \dots$$

..... (6,6) } → 36 Possibilities.

note:

→ If I throw n dice then number of possible outcomes 6^n

Q1 A class consists of 6 girls and 10 boys. If a committee of 3 is chosen at random from the class. find the probability that

- (i) 3 boys are Selected
- (ii) Exactly 2 girls are Selected.

Sol: Total number of Students = 16

6 G, 10 B

Prob: $\frac{\text{(i) Number of favourable cases}}{\text{Total number of cases}}$

Selection



Combination

arrangement



permutation → Batting order

$$\frac{{}^{10}C_3}{{}^{16}C_3} \rightarrow {}^nC_r = \frac{n!}{(n-r)!r!}$$

$$\frac{10 \times 9 \times 8}{1 \times 2 \times 3} = 120$$

Eg: Cricket

$$16C_3 \rightarrow \frac{16 \times 15 \times 14}{1 \times 2 \times 3}$$

(ii) Exactly 2 girls

$$\frac{6C_2 \times 10C_1}{16C_3}$$

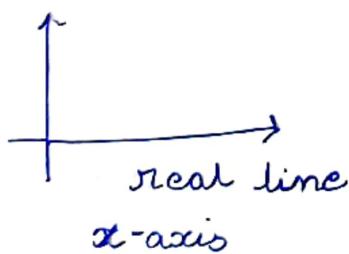
3 → Committee

$$[\because nC_1 \rightarrow n]$$

7/1/2021

* Random Variable:

✓ A real Variable X whose Value is determined by the output of a random Experiment is called a random Variable.



Eg:

0, 10, 10.5, 9.6 real numbers

Types of Random Variable

- a) Discrete } Random Variable
- b) Continuous } Random Variable

* Discrete RV

Eg: 1, 2, 3, 4, ...

* number of students in a class

* Continuous Random Variable:

Eg: 1.1, 1.2, ...

✓ weight, height of a person.

① Let X denote the number of heads in a single toss of 4 fair coins. Determine

(i) $P(X < 2)$

(ii) $P(1 < X \leq 3)$

Sol: Possible values of X is 4 $\rightarrow \{0, 1, 2, 3, 4\}$

(\because Given

Single toss of 4
fair coins)

↓
Discrete Random
Variable

(RV \rightarrow capital
letters)

Probability Distribution

X	0	1	2	3	4
$P(X)$	$\frac{1}{16}$	$\frac{4}{16}$	$\frac{6}{16}$	$\frac{4}{16}$	$\frac{1}{16}$

0 heads \rightarrow All tails

$$\Rightarrow \sum P(X) = 1$$

$$\text{I} \Rightarrow P(X < 2) = P(X = 0) + P(X = 1)$$

$$= \frac{1}{16} + \frac{4}{16}$$

$$= \frac{5}{16}$$

H H H H
H H H T
H H T H
H H T T
H T H H
H T H T
H T T H
H T T T
T H H H
T H H T
T H T H
T H T T
T T H H
T T H T
T T T H
T T T T

iii) $P(1 < x \leq 3)$

$$P(x=2) + P(x=3)$$

$$\frac{6}{16} + \frac{4}{16} = \frac{10}{16} = \frac{5}{8}$$

U
R
V

R
1
2
3
4

Q 2 Two dice are thrown. Let x assigns to each point (a, b) in S . The maximum of its numbers i.e. $x(a, b) = \max_{\min}(a, b)$. find the probability distribution. X is a random variable with $X(S) = \{1, 2, 3, 4, 5, 6\}$. Also find mean, variance.

Sol:

$$\left\{ \begin{array}{l} (1,1) (1,2) (1,3) (1,4) (1,5) (1,6) \\ (2,1) (2,2) (2,3) (2,4) (2,5) (2,6) \\ (3,1) (3,2) (3,3) (3,4) (3,5) (3,6) \\ (4,1) (4,2) (4,3) (4,4) (4,5) (4,6) \\ (5,1) (5,2) (5,3) (5,4) (5,5) (5,6) \\ (6,1) (6,2) (6,3) (6,4) (6,5) (6,6) \end{array} \right\}$$

$$X[(1,1)] = \max(1,1) = 1$$

$$P(x=1) = P(1) = \frac{1}{36}$$

$$P[x=2] \quad x = \{(1,2) (2,1) (2,2)\}$$

$$\frac{3}{36}$$

$$P[x=3] = \{(1,3) (3,1) (3,2) (3,3) (2,3) (3,1)\}$$

$$= \frac{5}{36}$$

$X = x$ 1 2 3 4 5 6

$P(X)$ $\frac{1}{36}$ $\frac{3}{36}$ $\frac{5}{36}$ $\frac{7}{36}$ $\frac{9}{36}$ $\frac{11}{36}$

1	2	3	4	5	6
1	2	3	4	5	6

③ A Random variable X has following probability distribution

x	0	1	2	3	4	5	6	7
$P(x)$	0	K	$2K$	$2K$	$3K$	K^2	$2K^2$	$7K^2 + K$

(i) Determine K

$$\sum P(x=a) = 1$$

$P(x)$

$$\sum P(X=a) = 1$$

↓
0 to any
value

$$0 + K + 2K + 2K + 3K + K^2 + 2K^2 + 7K^2 + K = 1$$

$$10K^2 + 9K - 1 = 0$$

$$K = \frac{1}{10}, -1$$

-1 X

Rule 2: $0 \leq P(E) \leq 1$

$$K = \frac{1}{10}$$

(ii) find $P(x < 6)$

$$= P(x=0) + P(x=1) + P(x=2) + \dots + P(x=5)$$

$$\therefore \{P(x < 6) = P(x=5)\}$$

$$= 0 + \frac{1}{10} + 2 \cdot \frac{1}{10} + 2 \cdot \frac{1}{10} + 3 \cdot \frac{1}{10}$$

$$= 0.81$$

$$\text{iii) } P(x \geq 6) = P(x=6) + P(x=7)$$

$$P(x \geq 6) = 1 - P(x < 6)$$

$$\text{iv) } P(0 \leq x \leq 4)$$

$$P(x=0) + P(x=1) + P(x=2) + P(x=3) + P(x=4)$$

+ 0

iv) $P(x \leq \underline{x}) > \frac{1}{2}$ → find the minimum value of $\underline{x} \approx$

$$P(x \leq 0) > \frac{1}{2}$$

$$0 > 0.5 \times$$

$$\underbrace{P(x \leq 3)}_{x=3} > \frac{1}{2}$$

$$\underbrace{x=1}$$

$$0 + K + 2K + 2K$$

$$P(x \leq \underline{x}_1) > \frac{1}{2}$$

$$5K > \frac{1}{2}$$

$$0.1 > 0.5 \times \text{ false}$$

$$5 \cdot \frac{1}{10} > \frac{1}{2}$$

$$\underbrace{x=2}$$

$$\frac{5}{10} > \frac{1}{2}$$

$$P(x \leq 2) > \frac{1}{2}$$

$$0.5 > 0.5 \times$$

$$\frac{1}{10} + 2 \cdot \frac{1}{10} > \frac{1}{2}$$

false

$$\frac{3}{10} > \frac{1}{2}$$

$$\underbrace{x=4}$$

$$0.3 > 0.5 \times \text{ false}$$

$$P(x \leq 4) > \frac{1}{2}$$

x satisfies 4, 5, 6, 7

$$8K > \frac{1}{2}$$

$0.8 > 0.5 \checkmark$
True

$\left\{ \begin{array}{l} \text{min value} \\ \therefore \text{ } \end{array} \right. \therefore x \text{ value is } 4 \left. \begin{array}{l} \checkmark \\ \end{array} \right.$

* Distribution function

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

Denoted by $F(x)$

x	$F(x)$
0	0
1	$1/10$
2	$3/10$
3	$5/10$
4	$4/10$
5	0.81
6	0.83
7	1

$$\underline{x=5}$$

$$P(X \leq 5) > \frac{1}{2}$$

$$8k^2 + k^2$$

$$8\left(\frac{1}{10}\right) + \left(\frac{1}{10}\right)^2 = 0.81$$

$$\underline{x=6}$$

$$P(X \leq 6) > \frac{1}{2}$$

$$8\left(\frac{1}{10}\right)^2 + 3\left(\frac{1}{10}\right)^2 = 0.83$$

$$\underline{x=7}$$

$$P(X \leq 7) > \frac{1}{2}$$

$$8k + 10k^2 + k = 9k + 10k^2 = \frac{10}{10} = 1$$

* Cumulative Distribution function (CDF)

* Probability mass function (PMF) →
↓
discrete

* Probability density function (PDF) →
↳ continuous.

Source: towards Data Science

$$x=7 \text{ [contd]}$$

$$\therefore 9k + 10k^2$$

$$\frac{9}{10} + 10\left(\frac{1}{10}\right)$$

$$= \frac{10}{10}$$

$$= 1$$

12/1/2021

Probability distributions:

★ ★ ★ ★

v.Gmp

v.Gmp

→ machine Learning

→ Data Science

areas → Research

Machine Learning
Data Science
Research

Probability Distributions

Discrete distribution

Continuous distribution

→ Binomial

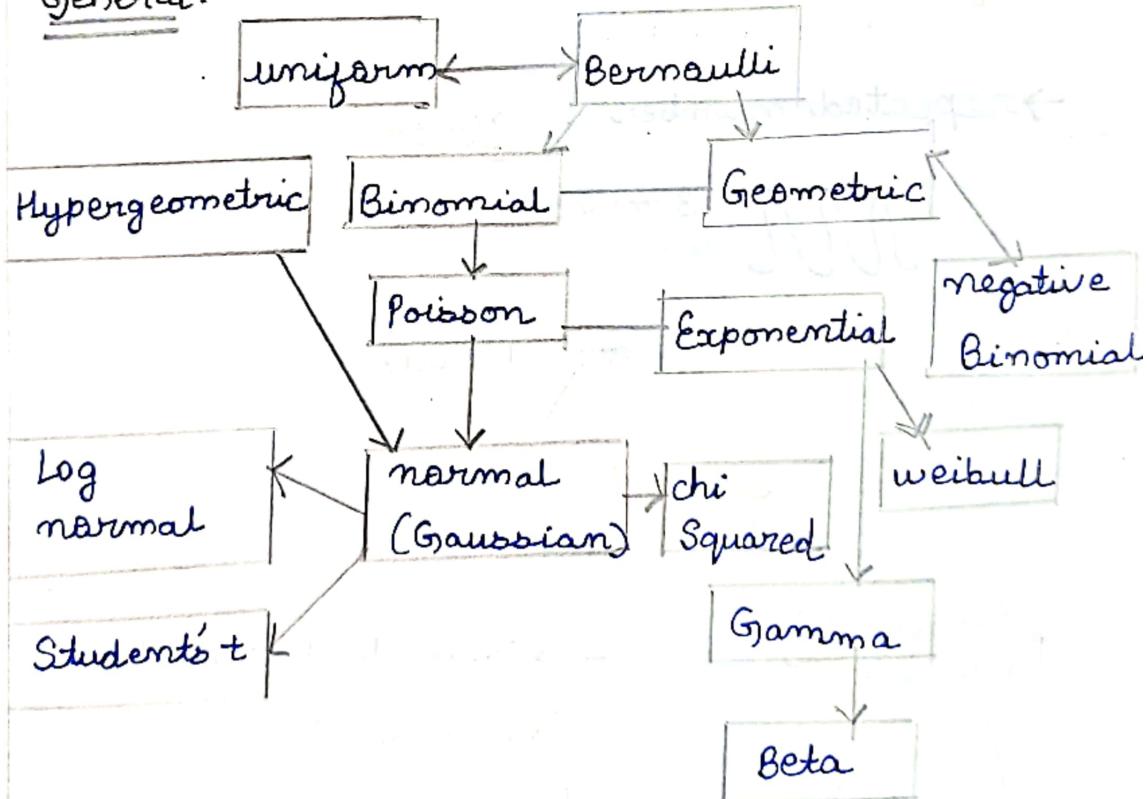
→ Poisson

distributions

→ Normal (Gaussian)

→ Exponential distribution

General:



* * * * * NORMAL DISTRIBUTIONS:

✓ Mother of all distributions is normal distribution.

✓ Gaussian distribution (Normal)

✓ Biologists use ~~this~~ information of the normal Distribution to Study patterns in nature.

Graph To Study patterns in nature.

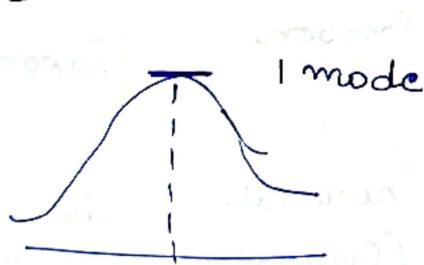
Gauss is

✓ Father of first Predictive algorithm.

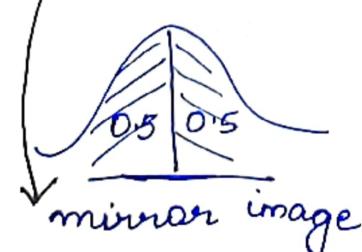
✓ Linear Regression is base for most of machine learning and data Science algorithms.

→ repeated number is mode

III 3 modes



★ Symmetric



→ Normal distribution

↳ 1. Symmetric

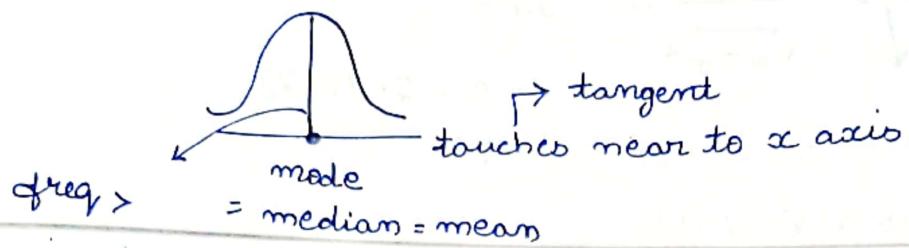
↳ 2. unimodal curve

↳ 3. Asymptotic

Asymptote: Tangent at ∞

↳ Touches / meets at Single Point.

SC axis is tangent for normal curve

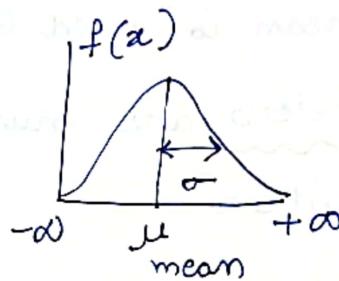


★ ★ ★ $\text{mean} = \text{median} = \text{mode}$ [In normal distribution]

* Properties

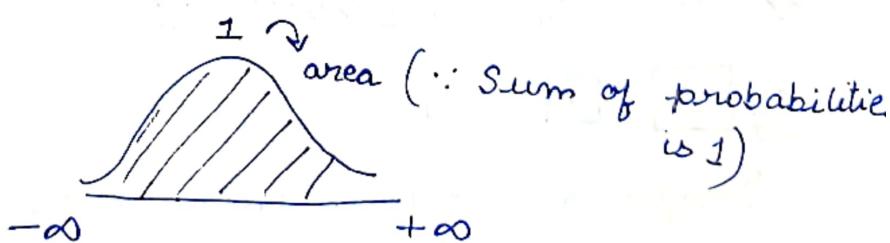
1. Bell Shaped
2. Symmetrical
3. mean, median, mode are equal
4. Location is characterized by mean μ
5. spread is characterized by Standard deviation σ

★ ★ information/
Variance



$$\sigma^2 = \text{Variance}$$

Random Variable has an infinite theoretical range $-\infty$ to ∞



V. Gmp

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

μ, σ are parameters

$f(x)$ = Probability density function of normal distribution
 x Values $-\infty$ to $+\infty$ x is variable

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

values $-\infty$ to $+\infty$ $\pi =$

x is a variable

μ, σ are parameters

[\therefore mean = median = mode in normal distribution
mean is used Generally]

parameters are numerical descriptors of a distribution.

* μ, σ decides shape, structure in normal distribution.

15/1/2021:

General Maths
lecture

✓ Linear Algebra
✓ Calculus

Statistics

✓ Descriptive { Statistics
✓ Inferential

→ 64K images of $\{0, 1, \dots, 9\}$ MNIST data Set

e.g.

1

1

1

Images are represented by matrices

Each image to represent we require 784
 $28 \times 28 = 784$ Pixels

✓ Each image is feature

$$D = \begin{bmatrix} 1 & f_1 & f_2 & f_3 & \dots & f_{784} \\ 2 & & & & & \\ 3 & & & & & \\ \vdots & & & & & \\ 64K & & & & & \end{bmatrix}$$

64000×784

→ dimensions → columns

784 $\xrightarrow{\text{Convert}}$ 2D or 3D to visualize data

(Linear Algebra) plays crucial role
[∴ ease]

Dimensionality Reduction $\left\{ \begin{array}{l} \text{Data Visualization (2D or 3D)} \\ \text{To apply machine learning algorithm} \end{array} \right.$ even 10D

→ To convert to 2D

we need Eigen Values; Consider top two Eigen values λ_1, λ_2

\downarrow λ_1 λ_2 Each Eigen Value \downarrow has Eigen Vector

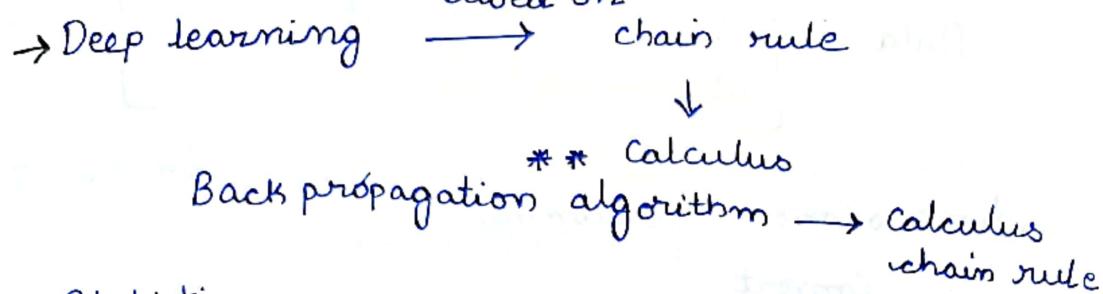
with principle component analysis (\uparrow 1933)

$D = \begin{bmatrix} f_1, f_2 \\ ? \\ 64K \end{bmatrix}$ 2D visualization.

Eg: To identify dog, cat in an image we need linear algebra

Calculus:

✓ Geoffrey Hinton is father of modern deep Learning



Statistics:

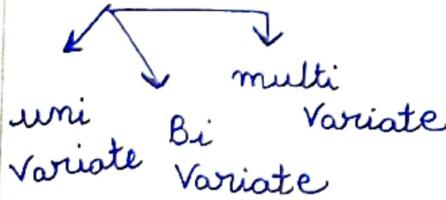
* Decision making

- ↳ fact based
- ↳ not by Emotions
- Backed up data [information]

→ Branch of Mathematics is Statistics

Collect, organize, interpret data

Descriptive Statistics



To Solve inferential
depends on
descriptive statistics

→ Summarize data

*** v. Simple

Inferential Statistics:

↳ Research

↳ academics

↳ Industry

Eg: Covid

↳ medical.

↳ Engineering

Hypothesis testing → Inferential Statistics: model fitting → conforms which model fits best.
 ↳ regression
 Logistic
 multi linear } model

→ Model is relation between dependent Variable and independent Variable

Descriptive Statistics:

- Mean
- median
- mode
- Standard deviation
- Sample variance
- Kurtosis
- Skewness
- Range
- maximum
- minimum
- Sum
- Count
- Geometric mean
- Harmonic mean
- AAD
- MAD
- IQR

$$X = [\quad] \quad 64K \times 784$$

$$X^T = [\quad] \quad 784 \times 64K$$

$$C = X^T X [\quad] \quad 784 \times 784$$

Covariance matrix

↓ under
multivariate
descriptive
matrix

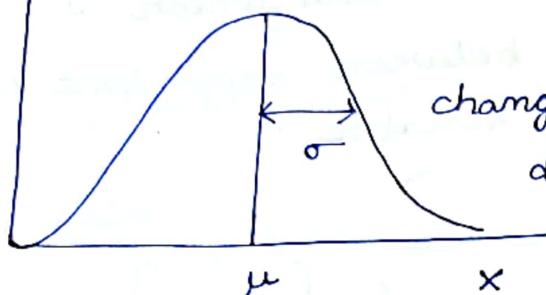
16/1/2021:

Normal distribution [Contd]

Shape

$f(x)$

changing μ shifts the distribution left or right



changing σ increases or decreases the spread.

$\rightarrow \mu \rightarrow$ we can shift mean

$\rightarrow \sigma \rightarrow$ spread of location

* Converting $X \rightarrow Z$

↳ Normal distribution (X)

~~Imp~~ Standardized Normal distribution (Z)

① Needs to be transform X units into Z units

② mean is 0 and standard deviation is 1

③ $\mu = 0 \quad \sigma = 1$ (in Z)

* Standardized Normal distribution:

$$Z = \frac{X - \mu}{\sigma}$$

$Z \rightarrow$ Standard normal variate variable

$X \rightarrow$ normal distribution

$\mu \rightarrow$ mean

$\sigma \rightarrow$ Standard deviation

$$f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} \rightarrow \begin{array}{l} \text{Pdf of Standard} \\ \text{normal deviation} \\ \text{distribution} \end{array}$$

$$z = \frac{x-\mu}{\sigma} \quad [\because \mu=0, \sigma=1] \text{ conditions.}$$

$$z = \frac{x-0}{1} = x$$

z is increasing; $f(z)$ is decreasing

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2} \quad -\infty \leq x \leq \infty$$

$$\begin{aligned} &= \frac{1}{1 \times \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-0}{1}\right)^2} \quad -\infty \leq x \leq \infty \\ &= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} z^2} \quad -\infty \leq z \leq \infty \\ &\quad \left(\frac{x-0}{1} = z \right) \end{aligned}$$

mean is 0 and Standard deviation is 1

$$f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$

Probability density function
of Standardized normal
distribution.

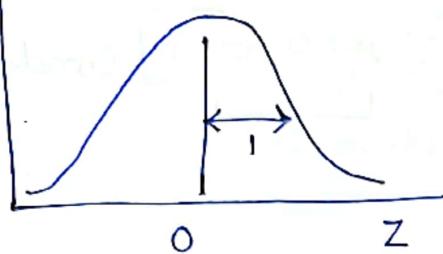
where e = the mathematical constant approximated by 2.71828

π = the mathematical constant approximated by 3.14159

z = any value of the Standardized normal distribution

Standardized Normal distribution:

$$f(z)$$



68
34 34
95%. $\mu=20 \quad \sigma=20$

→ Z distribution

→ mean 0

→ Standard deviation 1

→ positive Z values

→ Negative Z values

Eg:

Supposely

$$f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$

$$f(z) = e^{-z^2}$$

1 SD
34

$$z = 0$$

$$f(0) = 1$$

2 SD
2.35

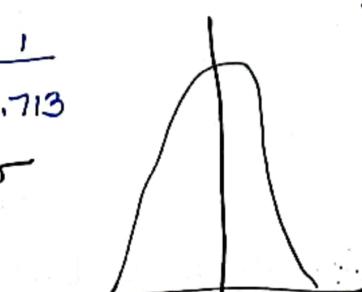
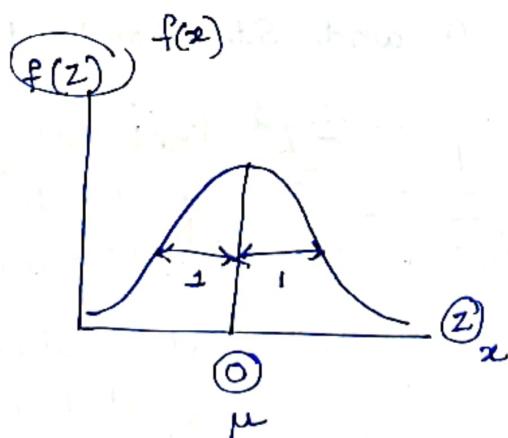
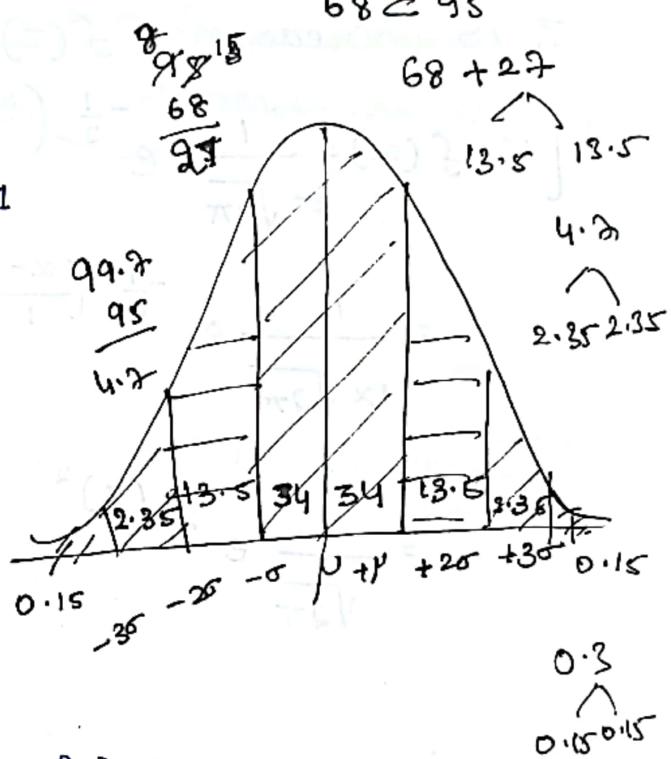
$$z = 1$$

$$f(1) = e^{-1}$$

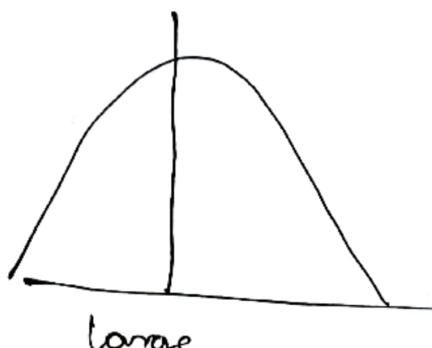
3 SD
0.15

$$= \frac{1}{e}$$

0.15
0.2
0.35
0.1
0.05
Role of σ
0.05
0.00
3.6 5.0



small SD few points
from mean



large



* Calculate the height of all individuals in world. 68% of people will be within 1 SD. avg in terms of ht.

① If x is distributed normally with mean of 100 and σ of 50, the Z value for $x=200$ is

Sol:

$$Z = \frac{x - \mu}{\sigma}$$

; σ is standard deviation.

$$\Rightarrow \frac{200 - 100}{50} \Rightarrow \frac{100}{50} = 2$$

* applicable for forecasting when actual data pred. is

* Empirical Rule for Normal Distribution:
 ↓ also known as (applicable)
 * Rule holds good only for (68-95-99.7 rule) Normal distribution

Symmetrical curve. Normal Bell shaped.

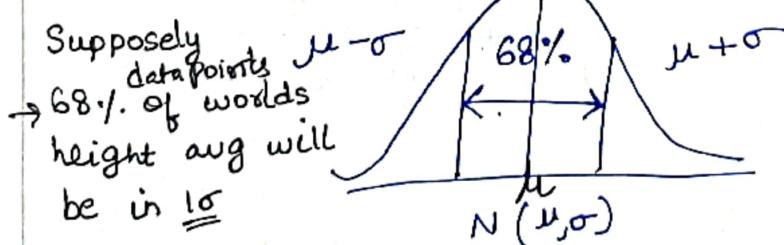
Empirical rule is also called as

↳ 68-95-99.7 rule

↳ Three Sigma rule

↳ 3 Standard deviations of a mean

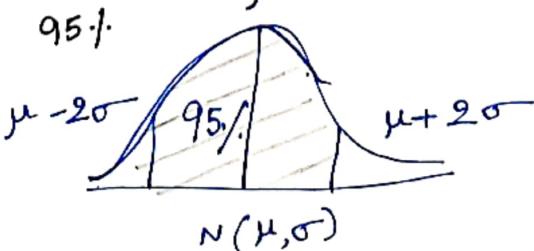
* one Standard deviation distance within μ 1 SD.



many points closer to average.

* Two Standard deviation

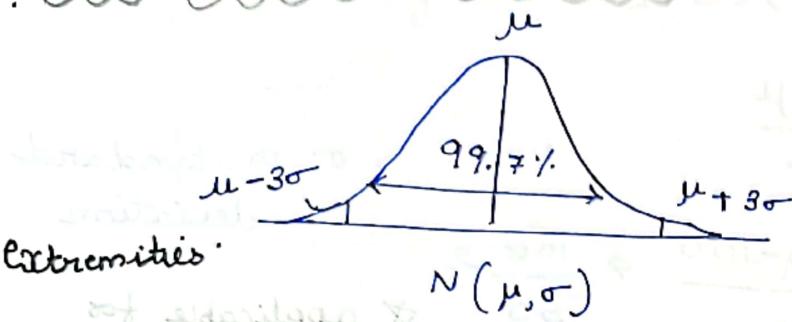
extreme values exists.



↑ ↓ outliers.

Empirical Rule

* Three Standard deviation:



Extremities

of standard deviation

area within 3 standard deviations

of mean

- ② Assume that the mean of 1 year-old girls in the US is normally distributed with mean of about $\mu = 9.5 \text{ kg}$ with standard deviation of $\sigma = 1.1 \text{ kg}$. Without using calculator, estimate the % of 1 year-old girl in the US that meet the following conditions:

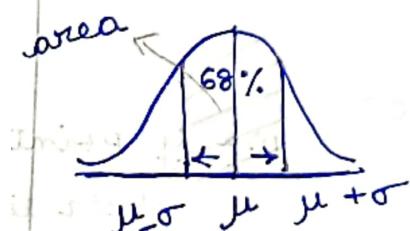
1. Less than 8.4 kg

2. Between 7.3 kg and 11.7 kg

3. more than 12.8 kg

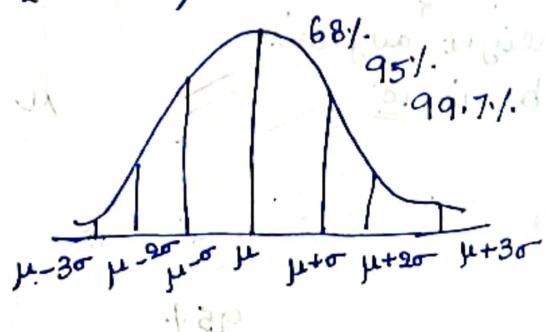
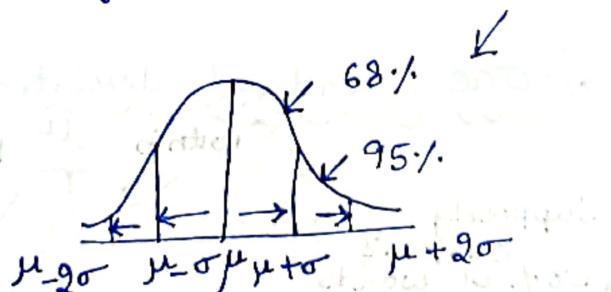
Given $\mu = 9.5 \text{ kg}$, $\sigma = 1.1 \text{ kg}$

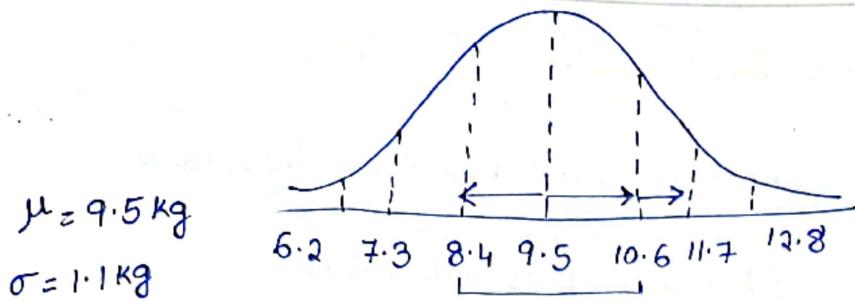
rule * $68-95-99.7$



68% is area between

$\mu - \sigma$ and $\mu + \sigma$



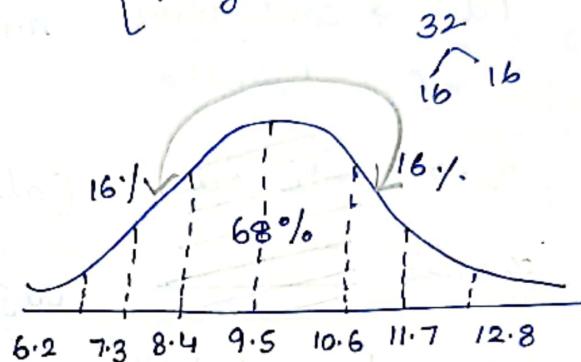
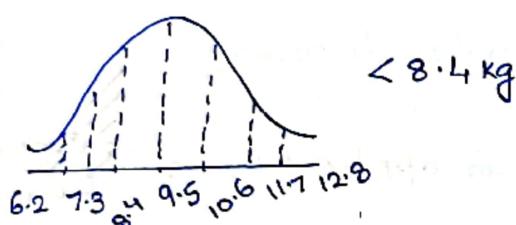


$$\mu - \sigma = 9.5 - 1.1 = 8.4$$

$$\mu + \sigma = 9.5 + 1.1 = 10.6$$

\therefore Symmetric distribution

(i) Less than 8.4 kg.

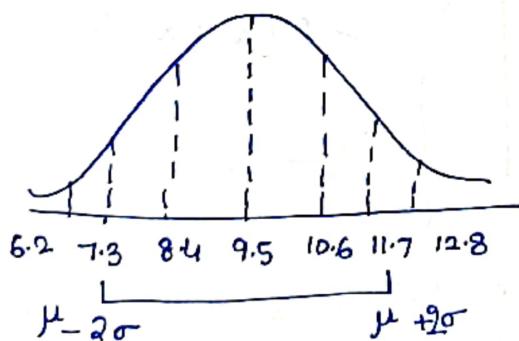


$$100\% \rightarrow 100 - 68 = 32$$

\therefore Less than 8.4 kg = 16%.

$\frac{16}{100}$ is Probability

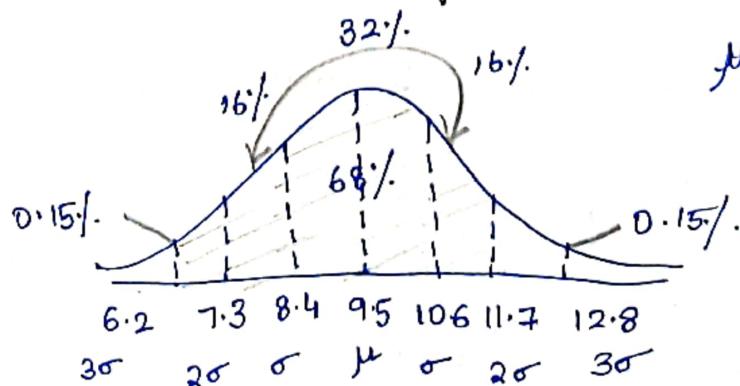
(ii) Between 7.3 Kg and 11.7 Kg = $\frac{95}{100} = 0.95$



area between $\mu - 2\sigma$ and
 $\mu + 2\sigma$ is 95%
 $100 - 95 = 5\%$

$\frac{95}{100}$ probability

(iii) more than 12.8 kg = 0.15%. (iv) below than $\frac{6.2}{0.15}$ = 0.15%.



$\mu - 3\sigma$ and $\mu + 3\sigma$

99.7%

$100 - 99.7 = 0.3\%$

0.15 0.15

Cumulative Distribution Function

CDF of random variable X is defined as

$$F_X(x) = P(X \leq x), \text{ for all } x \in \mathbb{R}$$

Cdf \rightarrow describes distribution of a

Pdf \rightarrow continuous random variable

Pmf \rightarrow discrete

$\mathbb{R} \rightarrow$ continuous (also for discrete)

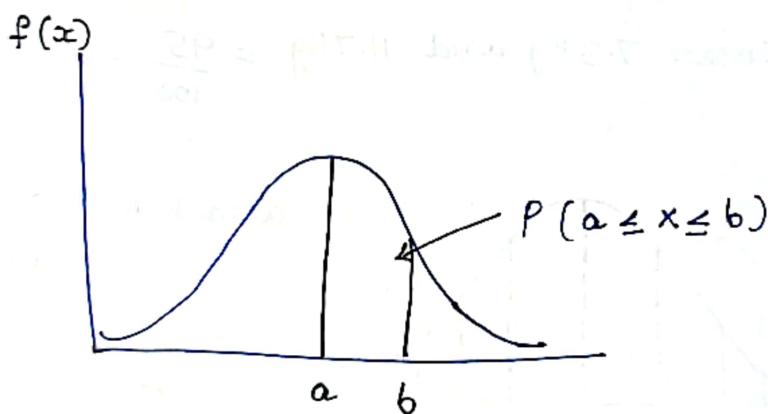
x axis

x axis

Cdf is applicable for continuous
discrete

Normal probabilities:

\rightarrow Probability is measured by the area under the curve



Note that the probability of any zero individual value is zero.

$$\int f(x)dx = 1 (\because \text{area})$$

$-\infty$ $\left[\because \text{In } X \text{ Normal Distribution} \right] \int \text{integration is complex}$

Convert to Z

no integration required.

Cumulative Distribution Function

CDF of random variable X is defined as

$$F_X(x) = P(X \leq x), \text{ for all } x \in \mathbb{R}$$

Cdf \rightarrow describes distribution of a
Pd f \rightarrow continuous random variable

Pmf \rightarrow discrete

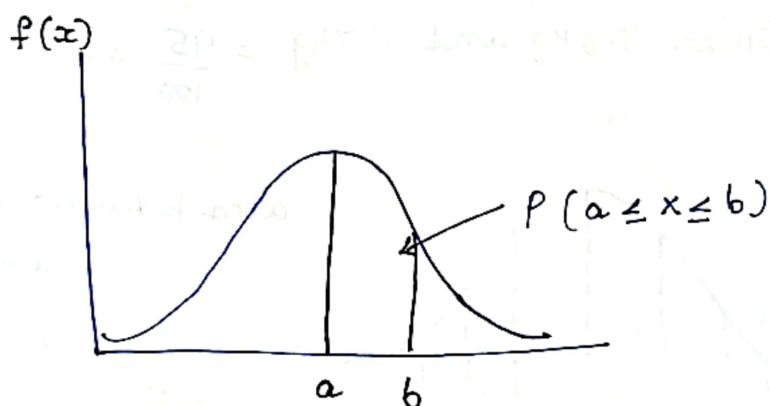
$\mathbb{R} \rightarrow$ continuous (also for discrete)

x axis ——

Cdf is applicable for continuous
/ discrete

Normal probabilities:

\rightarrow Probability is measured by the area
under the curve



Note that the probability of any zero individual value is zero.

$$\int f(x)dx = 1 \quad (\because \text{area})$$

$\int_{-\infty}^{\infty} \quad [\because \text{In } X \text{ Normal Distribution } \int \text{ integration is complex}]$

Convert to Z

no integration required.

Applications of Cumulative Distribution Function

- * Standard Normal table
- * Unit Normal Table } also called is applicable in
* Z table field of data Science.

\because we need ^{not} to find integration to find probability using above tables.

→ using PDF and CDF we can find probability of a Random Variable.

$$\left\{ \begin{array}{l} f(x) = \text{P.d.f} \\ F(x) = \text{C.d.f} \\ f(x) = \frac{d}{dx} [F(x)] \end{array} \right.$$

* Statistics in Understanding Data

* To understand data we require Statistics.

"There are two kinds of Statistics, the kind you look up and the kind you make up"
Rex Stout

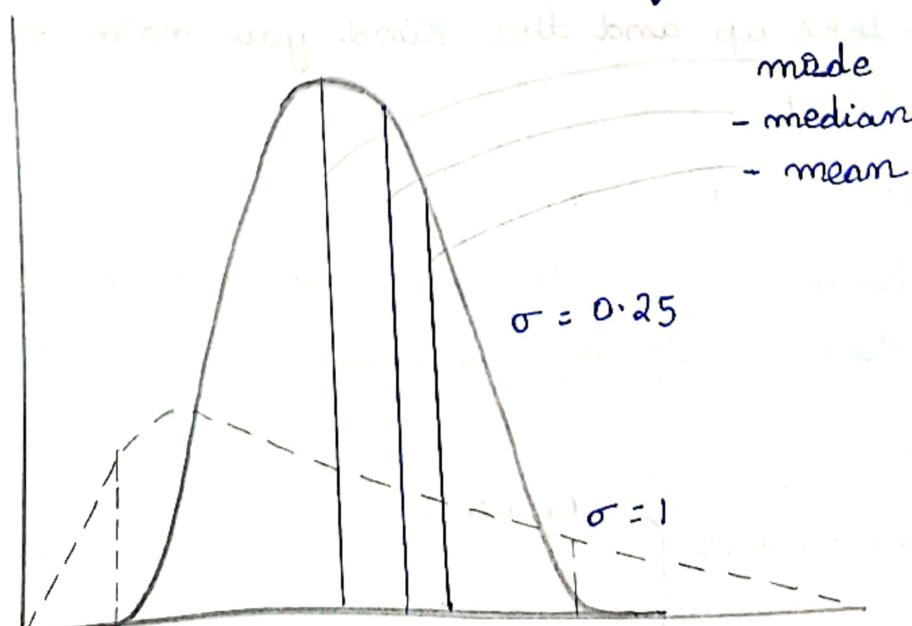
What is Statistics?

* A branch of mathematics that deals with collecting, analyzing and interpreting data.

✓ Descriptive }
✓ Inferential } Statistics

Descriptive Statistics:

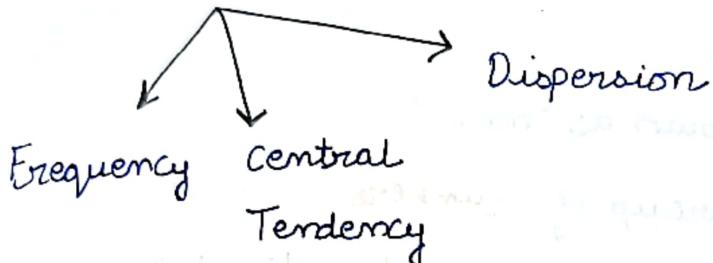
- ✓ Summarize, describes the data
- ✓ Identifies important elements in dataset without seeking to find explanations of those elements
- ✓ Depending on how many variables are involved descriptive statistics are of 3 types
- ✓ Summarize single variable → univariate
 - { Relation b/w 2 variables → BiVariate
 - Describe → multiple variables multivariate
- ✓ Summarize data as it is
- ✓ Do not posit any hypothesis about data
- ✓ Detect outliers
- ✓ Plan how to prepare data
- ✓ precursor to feature engineering



✓ Feature engineering involves extracting features from data to build model
↓ we require features
we require feature engineering

✓ Precursor to feature engineering.

univariate (Single variable)



Measures of frequency:

↳ Frequency Tables

↳ Histograms

Summary of Statistics:

Measures of Central Tendency:

Measures of Location

• Arithmetic mean

✓ weighted mean

✓ median

✓ mode

✓ Percentile

✓ Geometric

✓ Harmonic

{ means

↓

infrequently used

Measures of Dispersion

✓ Skewness

✓ Kurtosis

✓ Range

✓ Interquartile range

✓ Variance

✓ Standard score

✓ Co-efficient of

Variance

Arithmetic Mean

→ Data : {16, 17, 10, 13, 20, 18, 13, 14, 18}

$$\bar{x} = \frac{16 + 17 + 10 + 13 + 20 + 18 + 13 + 14 + 18}{9}$$

$$\bar{x} = \frac{139}{9}$$

$$\bar{x} = 15.444$$

★ Commonly known as 'mean'

★ Average of Group of numbers

★ Applicable for interval and ratio data

Eg:

Interval data

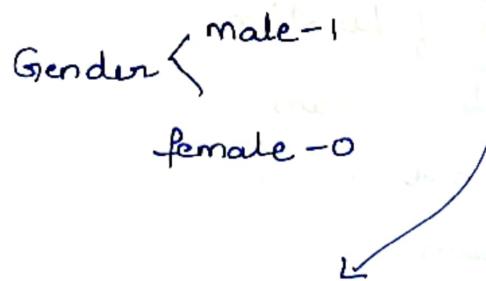
[2000, 2020]

Heights, weights

wts. are allowed

★ Not applicable for nominal or ordinal data

Ex:



1. Boy, AssnP, AssisP

2. Grades of Students

★ Affected by each value in the dataset, including extreme values

* Computed by summing all values in the dataset and dividing the sum by the number of values in the dataset.

Population Mean:

Data: 60 20 10 40 50 30

$$\mu = \frac{\sum x_i}{n} = \frac{60+20+10+40+50+30}{6} = 35$$

* Impact of outliers:

Data: 60 20 10 40 50 30 1000

$$\mu = \frac{\sum x_i}{n} = \frac{60+20+10+40+50+30+1000}{7}$$

$$\text{Mean} = 172.85$$

Population mean $\rightarrow \mu$

Sample mean $\rightarrow \bar{x}$

$$\bar{x} = \frac{\sum x}{n} = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n}$$

$$= \frac{57+86+42+38+90+66}{6}$$

$$= \frac{379}{6}$$

$$= 63.167$$

* Mean of Grouped Data

* weighted average of class midpoints

* class frequencies are the weights

$$\mu = \frac{\sum f M}{\sum f} \quad \begin{matrix} \text{frequency} \\ \text{class midpoint} \end{matrix}$$

$$= \frac{\sum f_m}{N}$$
$$f_1 M_1 + f_2 M_2 + f_3 M_3 + \dots + f_i M_i$$
$$= \frac{f_1 + f_2 + f_3 + \dots + f_i}{f_1 + f_2 + f_3 + \dots + f_i}$$

f_1 is frequency of first class

M_1 is class midpoint of first class

Calculation of Grouped Mean:

class Interval	Frequency (f)	class (M)	fM
20 - under 30	6	25	150
30 - under 40	18	35	630
40 - under 50	11	45	495
50 - under 60	11	55	605
60 - under 70	3	65	195
70 - under 80	1	75	75
			<u>2150</u>

$$\mu = \frac{\sum f M}{\sum f} = \frac{2150}{50}$$

$$= 43.0$$

20/1/2021

* Finding Normal probability procedure.

→ To find $P(a < X < b)$ when X is distributed normally:

→ Draw the normal curve for the problem in terms of X .

✓ Translate X -values to Z -values
use the Standardized Normal Table

→ Assessing Normality:

CLT

↳ Very Imp

many Analytical and Statistical tools

Data → Data cleaning → verify ND

* Demo on Normal Distribution:

```
from Scipy.stats import norm
```

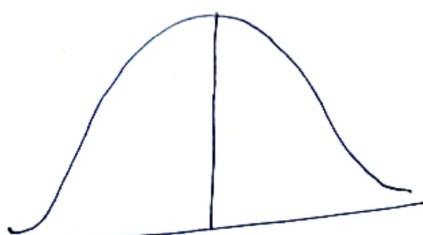
 $\text{val, } m, s = 68, 65.5, 2.5$
 $\text{print(norm.cdf(val, m, s))}$
 $\text{Val} = 68, \text{m} = 65.5, \text{s} = 2.5$
 $\downarrow \quad \downarrow \quad \downarrow$
 $x = 68 \quad \mu = 65.5 \quad \sigma = 2.5$

$$Z = \frac{x - \mu}{\sigma}$$

$$Z = \frac{68 - 65.5}{2.5}$$

$$= \frac{2.5}{2.5}$$

$$= 1$$



20/1/2021

gmp

Central Limit Theorem:

↳ Real time applications

↳ Super useful

↳ Statistics

↳ Mathematics

Developing a Sampling distribution:

Assume there is a population ...

Population Size $N=4$

Random variable x , is age of individuals

Values of x :

18, 20, 22, 24 (years)

$$\mu = \frac{\sum x_i}{N}$$

$$= \frac{18+20+22+24}{4} \Rightarrow 21$$

μ is Population mean

\downarrow
(by default)

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}$$

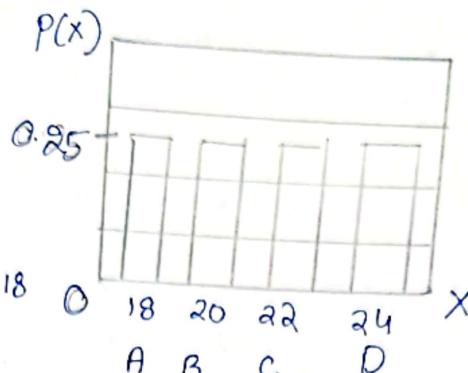
$$= (x_1 - \mu)^2 + (x_2 - \mu)^2 + (x_3 - \mu)^2 + (x_4 - \mu)^2$$

$$\begin{array}{cccc} 18 & 20 & 22 & 24 \\ | & | & | & | \\ x_1 & x_2 & x_3 & x_4 \end{array}$$

on Substitution

$$\sigma = 2.236$$

$$\therefore \mu = 21$$



\therefore Probability for 18

$$\frac{1}{4}$$

0 18 20 22 24 X

A B C D uniform distribution

Since 18, 20, 22, 24

$$\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4} \left\{ \begin{array}{l} \text{favourable cases} \\ \text{Total} \end{array} \right\}$$

Now consider all possible Samples of Size

$$n=2$$

1 st obs	2 nd observation			
	18	20	22	24
18	18, 18	18, 20	18, 22	18, 24
20	20, 18	20, 20	20, 22	20, 24
22	22, 18	22, 20	22, 22	22, 24
24	24, 18	24, 20	24, 22	24, 24

$N^n = 4^2 = 16$

note:

$N^n \xrightarrow{n} \text{without replacement}$

$N_c^n \xrightarrow{n=2} \text{with repetition}$

16 possible Samples

(Sampling with replacement)

$$\rightarrow N=4, n=2, N^n = 4^2 = 16$$

1 st obs	$n=2$				$= 16$	16 sample means
	18	20	22	24		
18	18, 18	18, 20	18, 22	18, 24	18	18
20	20, 18	20, 20	20, 22	20, 24	20	19
22	22, 18	22, 20	22, 22	22, 24	22	21
24	24, 18	24, 20	24, 22	24, 24	24	23

above 2nd table is obtained by averages

avg:

$$\frac{18+18}{2} = \frac{36}{2} = 18$$

$$\frac{18+18}{2} = 18, \frac{18+20}{2} = 19, \dots 30 \text{ on}$$

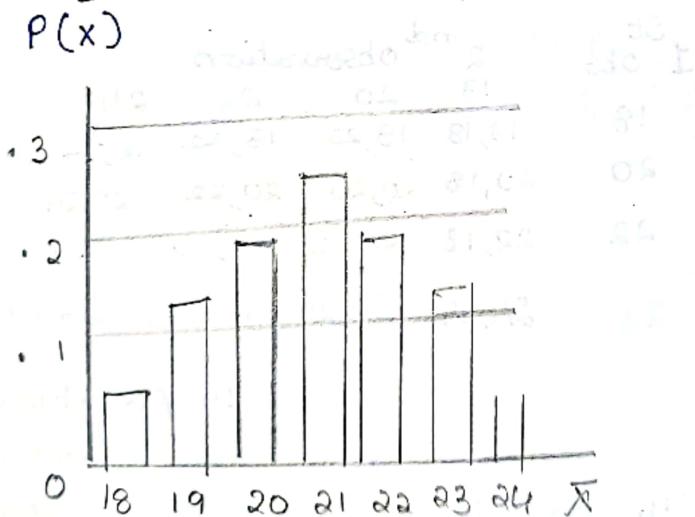
mean is called
first repeat mean

16 Sample means

CV		18	18	19	20	22	24	means of sample Size = 2
/	not applicable	20	19	20	21	22		
Pareto		22	20	21	22	23		
Cauchy	x	24	21	22	23	24		
all not used								

Sample means distribution

$P(\bar{x})$



no longer uniform

$\therefore \mu$ is always population mean

\bar{x} sample mean

Summary measures of this Sampling distribution

$$E(\bar{x}) = \frac{\sum \bar{x}_i}{N} = \frac{18+19+21+\dots+24}{16} = 21 = \mu$$

μ and $E(\bar{x})$ are always same (values)

$$\sigma_{\bar{x}} = \sqrt{\frac{\sum (\bar{x}_i - \mu)^2}{N}}$$

$$= \sqrt{\frac{(18-21)^2 + (19-21)^2 + \dots + (24-21)^2}{16}}$$

$$= 1.58$$

σ value 2.236

always $\sigma_{\bar{x}} < \sigma$ values

Comparing the population with its Sampling distribution

Population

$$N=4$$

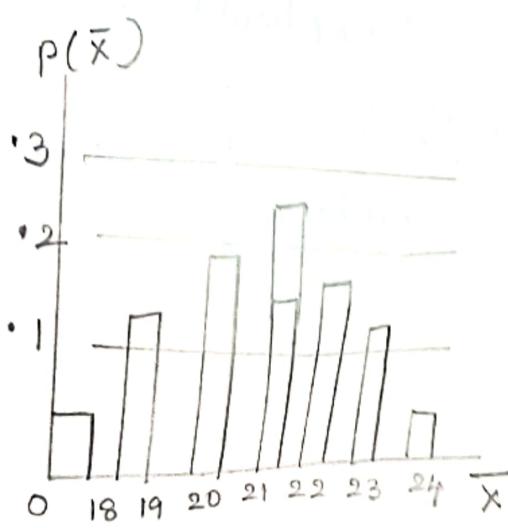
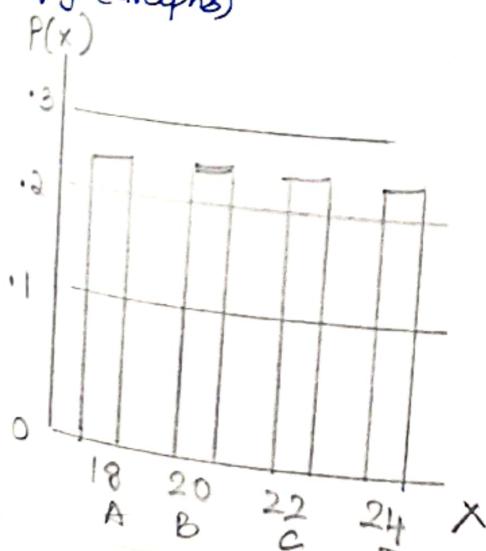
$$\mu=21, \sigma=2.236$$

Sample means distribution

$$n=2$$

$$\mu_{\bar{x}} = 21, \sigma_{\bar{x}} = 1.58$$

dig (Graphs)



Even if you're not

normal...

... the average...

... is normal!!!

* CLT [Central Limit Theorem]

Statement:

Let x be an independent identically distributed random variable with finite Population mean μ and finite population variance σ^2 then a random variable converges in distribution to the standard normal variables as n goes to infinity with mean μ and variance σ^2/n .

Here n is the size of sample from given distribution

{ no matter what is the shape of the population its distribution of sampling distribution will be normal if n is very large.

Limitations:

not applicable for

* Pareto { $\because \mu = \infty$
Cauchy } Distributions
 \therefore mean is \bar{x}

Let us assume x is a random variable with finite population mean and finite population variance $x(\mu, \sigma^2)$

I am collecting m samples

$s_1, s_2, s_3, \dots, s_m$

Size of each sample is $n=30$, my sample size

Let me assume I am collecting 1000 samples

$$n=30, m=1000$$

$$m \times n = 1000 \times 30$$

= 30k samples

For each of the sample, let us calculate mean

$$\begin{array}{c} s_1, s_2, s_3, \dots, s_m \\ \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ \bar{x}_1 \quad \bar{x}_2 \quad \bar{x}_3 \quad \dots \quad \bar{x}_m \end{array}$$

Distribution of $\bar{x}_i, i=1 \text{ to } m$

Distribution of \bar{x}_i = Sampling distribution of Sample mean

Now Central Limit theorem says

\bar{x}_i is distributed with Gaussian Normal distribution with mean μ and Variance $\frac{\sigma^2}{n}$

as $n \rightarrow \infty$

$$\text{i.e. } \bar{x}_i \xrightarrow{\text{as } n \rightarrow \infty} N\left(\mu, \frac{\sigma^2}{n}\right)$$

as $n \rightarrow \infty$

Sampling distribution Properties:

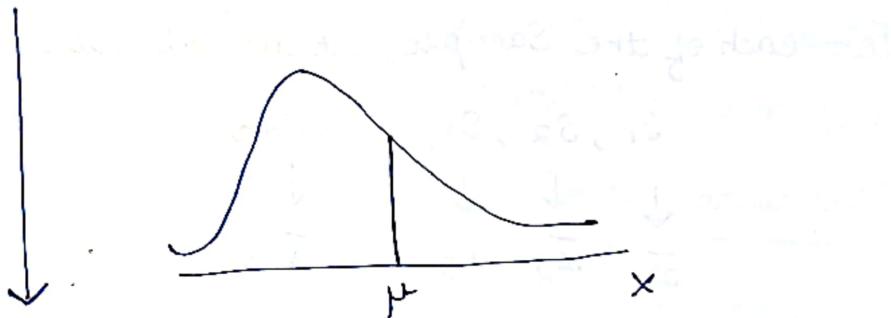
Central Tendency

$$\mu_{\bar{x}} = \mu$$

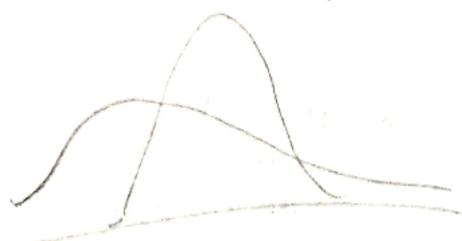
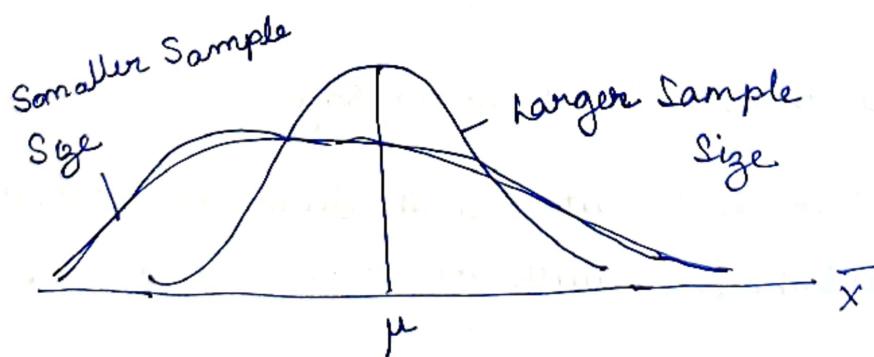
Variation

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

* Population distribution:



* Sampling distribution (becomes normal as n increases)



21/11/2021

① Suppose a large population has mean $\mu = 8$ and Standard deviation $\sigma = 3$. Suppose a random sample of size $n = 36$ is selected. What is the probability that the sample mean is between 7.8 and 8.2?

Sol: $\because n = 36 \checkmark$ large size \rightarrow CLT is applicable $n \geq 25$

Even if the Population is not normal distributed, the central theorem can be used

$$\begin{array}{ccc} \mu_{\bar{x}} = \mu & & (n > 25) \\ \downarrow & \downarrow & \\ \text{Sample mean} & \text{population mean} & Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} \end{array}$$

So the Sampling distribution of \bar{x} is approximately normal

with mean $\mu_{\bar{x}} = 8$

and standard deviation

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{3}{\sqrt{36}} = 0.5$$

$$P(7.8 < \mu_{\bar{x}} < 8.2) = P\left(\frac{7.8 - 8}{3/\sqrt{36}} < \frac{\mu_{\bar{x}} - \mu}{\sigma/\sqrt{n}} < \frac{8.2 - 8}{3/\sqrt{36}}\right)$$

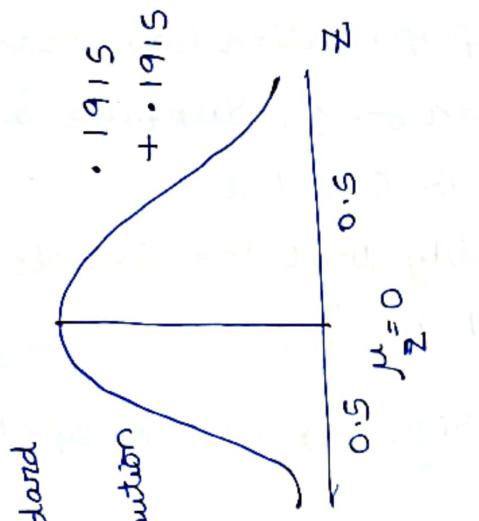
$$= P(-0.5 < Z < 0.5) = 0.3830$$

\therefore from Z table

Value of $0.5 + 0.1915$ is 0.1915

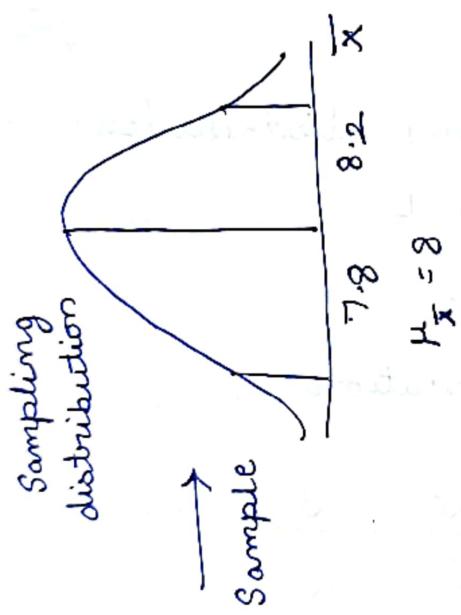


Scanned with OKEN Scanner



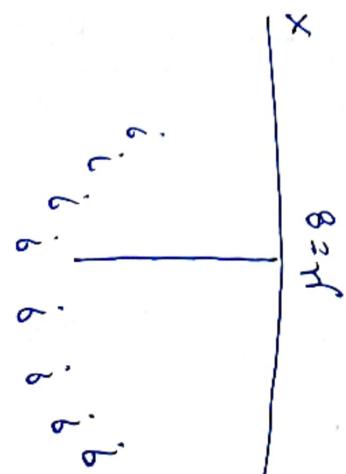
Standard
Normal
distribution

→ Standardize



Sampling
distribution

→ Sample



?

?

?

?

?

?

?

parametric family of distributions:



① Binomial distribution:

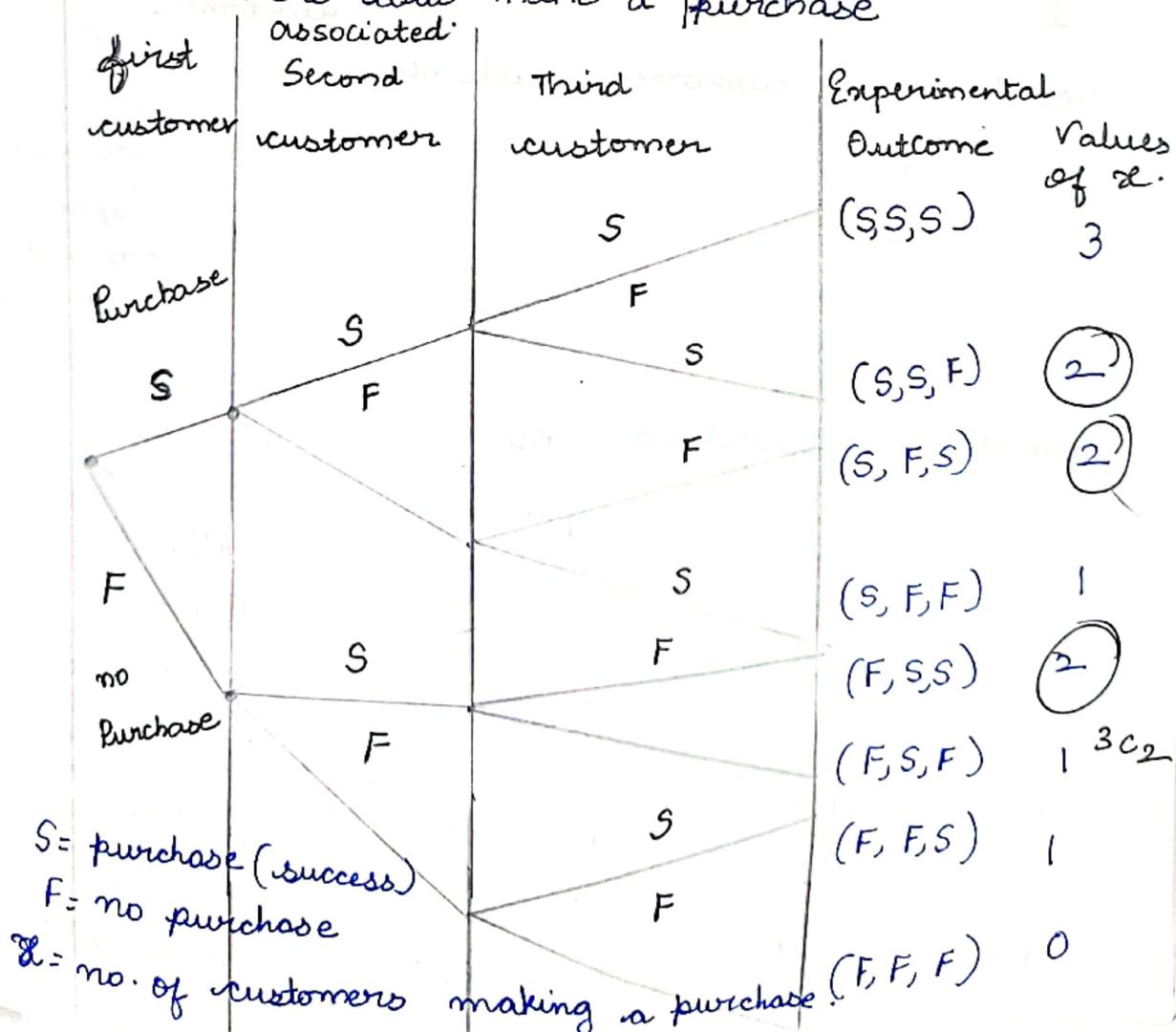
*** Martin clothing Store problem.

→ Let us consider the purchase decisions of the next 3 customers who enter a store.

→ On the basis of past experience, the store manager estimates the probability that any one customer will purchase is 0.30

$$\star \left\{ \begin{array}{ll} p = 0.30 & \text{purchasing } p^* \\ 1-p = 0.70 & \text{not purchasing } p^* \end{array} \right.$$

→ What is the probability that two of the next 3 customers will make a purchase



${}^3 C_2$ is no. of Combinations to purchase

2 customers out of 3 customers

2 2 2

$$3_{C_2} = 3$$

$$m_{C_{n_1}} = m_{C_{n-1}}$$

$$3_{C_0} = 3_{C_1} = 3$$

Trial outcomes :-

1st customer 2nd customer 3rd customer Experimental customer p ↓

purchase purchase m

purchase

Success P
Failure P)

Probability
of Experimental outcome.

$$P P(1-P) =$$

$$P^2(1-P)$$

$$= (0.30)^2 (0.70)$$

$$= 0.063$$

$$(S, F, S) = P(1-P)^P$$

$$= 0.063$$

$$(F, S, S) = (-P)(PP) = 0.063$$

Graphical representation of probability dist. for number of customers making a purchase

$$x \quad P(x)$$

$$0 \quad 0.7 \times 0.7 \times 0.7 = 0.343$$

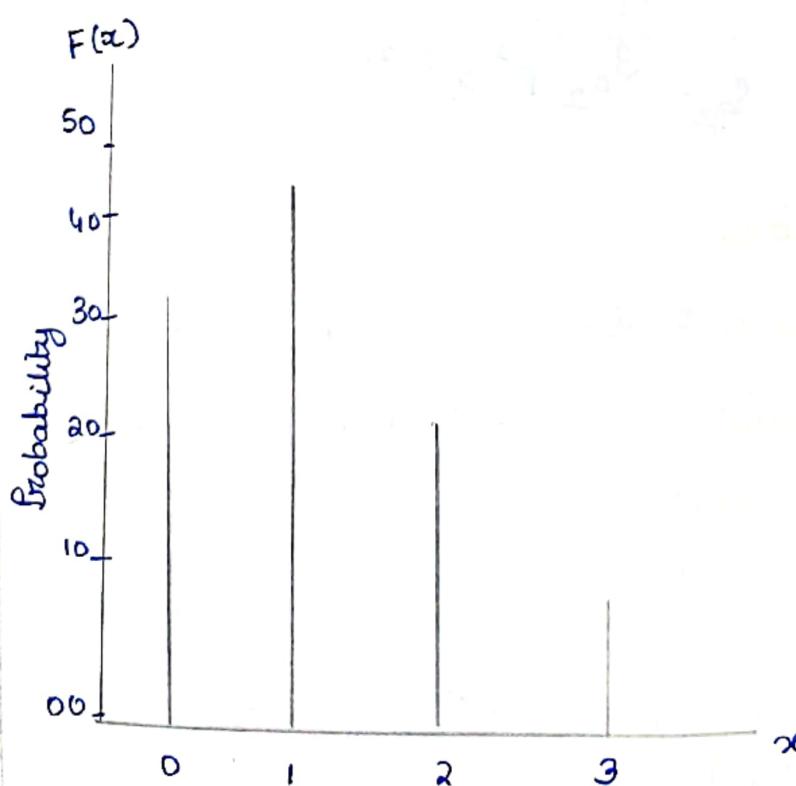
$$\begin{aligned} & 1 \quad SFF \\ & 0.3 \times 0.7 \times 0.7 + \\ & \quad \quad \quad FSF \\ & \quad \quad \quad 0.7 \times 0.3 \times 0.7 + \end{aligned}$$

$$0.7 \times 0.7 \times 0.3 = 0.441$$

FFS

$$2 \quad 0.189$$

$$3 \quad 0.027$$



Number of customers making a purchase

from trial outcomes.

$$3C_2 p^2 (1-q)$$

$$p^2 (1-p) + p^2 (1-p) + p^2 (1-p)$$

$$= 3p^2 (1-p)$$

$$= 3C_2 p^2 (1-p)^{3-2}$$

$$nC_x p^x q^{n-x}$$

Binomial distribution

$$P(X=x) = nC_x p^x q^{n-x}$$

$$\text{Ex: } 3C_2 p^2 q^{3-2}$$

Assumptions:

- ★ Experiment involves n identical trials.
- ★ Each trial has exactly two possible outcomes success or failure.
- ★ Each trial is independent of the previous trial
- ★ p is the probability of a success on any one trial
- ★

- * $q = (1-p)$ is the probability of a failure on any one-trial.
- * p and q are constant throughout the experiment.

probability

function.

$$P(x) = \frac{n!}{x!(n-x)!} p^x q^{n-x} \text{ for } 0 \leq x \leq n$$

mean

Value

$$\mu = n \cdot p$$

Variance and

$$\sigma^2 = n \cdot p \cdot q$$

Standard

deviation.

$$\sigma = \sqrt{\sigma^2} = \sqrt{n \cdot p \cdot q}$$

$$P(X=x) = {}^n C_x p^x q^{n-x}$$

$${}^n C_x = \frac{n!}{x!(n-x)!}$$

$$n! (n-x)!$$

Mean and Variance

- * Suppose that for the next month the clothing store forecasts 1000 customers will enter the store.
- * What is the expected number of customers who will make a purchase?

* The answer is $\mu = np$ $n = 1000$

$$= (1000)(0.3) = 300$$

✓ for the next 1000 customers entering the store, the Variance and Standard deviation for the number of customers who will make a purchase are

$$\begin{aligned}\sigma^2 &= np(1-p) \\ &= 1000(0.3)(0.7) \\ &= 210 \\ \sigma &= \sqrt{210} \\ &= 14.49\end{aligned}$$

* $\mu = np$

$$\sigma^2 = npq$$

{
import scipy
import numpy as np
from Scipy.stats import binom

Poisson Distribution: Assumptions:

- * Statistical Distribution \rightarrow how many times event is occurred within specified period of time.
- * Discrete Event \nwarrow occurring not occurring
- * Describes rare events in large population mutation of cells in body: mutation acquisition.
- * Each occurrence is independent any other occurrences.
- (discrete occurrences over interval) resources
- * The expected number of occurrences must hold constant throughout the experiment

Example:

① Measurements of number of occurrences within a time period.

arrival patterns

Eg.1: Arrivals at queuing systems

- * airports - people, airplanes and baggage
- * Banks - people, automatic and loan application.
- * customers at a store.
- * computer file servers - read and write operations.

number of calls to a switch board.

The number of arrivals in any service facility like at an ATM, railway station, petrol pump.

Eg 2: Defects in manufactured goods

* number of defects per 1,000 feet of extruded

copper wire

* number of blemishes per square foot of

* painted surface

* number of errors per typed page.

Eg 3: Number of extreme weather events

Eg 4: The number of aircraft/road accidents

in any time interval.

* poisson distributions model the number of occurrences within a fixed time period.

* can be used to approximate binomial distributions for very small success probabilities.

Binomial:

$$P(X=x) = n c_x p^x q^{n-x}$$

$$\downarrow \quad n \rightarrow \infty, p \rightarrow 0$$

poisson distribution

$$\frac{e^{-\lambda} \lambda^x}{x!} = P(X=x)$$

poisson distribution formula:

$$\frac{e^{-\lambda} \lambda^x}{x!}$$

e = Euler's constant ≈ 2.718

λ = mean or expected value of the variable

x = number of successes for the event

! = factorial.

mean Variance

$$\lambda$$

$$\lambda$$

Standard deviation.

$$\lambda$$

$\lambda = 3.2$ customers/
4 minutes

$x = 10$ customers/
8 minutes

$$P(x) = \frac{\lambda^x e^{-\lambda}}{x!}$$

$$P(x=10) = \frac{6.4^{10} e^{-6.4}}{10!} \\ = 0.0528$$

$$P(x=x)$$

$$= \frac{\lambda^x e^{-\lambda}}{x!}$$

$\lambda = 3.2$ customers/
4 minutes

$x = 6$ customers/
8 minutes

Adjusted λ

6.4 customers/
8 minutes

$$P(x) = \frac{\lambda^x e^{-\lambda}}{x!}$$

$$P(x=6) \\ = \frac{6.4^6 e^{-6.4}}{6!} = 0.1856$$

=

* Sampling Theory:

* Population:
The totality of observations with which we are concerned, whether this number be finite or infinite, constitutes what we call Population.
(or)

Population (or universe) is the aggregate or totality of statistical data forming a subject of investigation.

Eg; Population of Heights of Indians.

2. Population of Standardized banks.

* Sampling:

→ Study of entire Population may not be possible to carry out and hence a Part alone is selected from Given Population.

↓
Sample.

←
we need Sample to understand about Population.

★ { Sample Size n }
 { Population Size N }

Large Sample:

If Size of the Sample (n) ≥ 30 , the Sample is said to be large Sample.

Small Sample: $n < 30$

Correction factor:

$$* \frac{N-n}{N-1}$$

(i) what is the value of correction factor if $n=5$, and $N=200$

$$\text{Sol: } \frac{200-5}{200-1} = \frac{195}{199} = \underline{\underline{0.98}}$$

(ii) How many different Samples of Size two can be chosen, from a finite Population of Size 25

$$n=2$$

Sol:

$$N=25$$

$$N_{C_n} = 25_{C_2} = \frac{25 \times 24}{1 \times 2}^1$$

without
replacement

= 300 Samples

$$\text{with } N^n = 25^2$$

replacement

$$= 625$$

Samples

③ A Population consists of 5 numbers 2, 3, 6, 8, 11. Consider all possible Samples of size two which can be drawn with replacement from this Population.

Sol: (i) find

6 (i) mean of the Population

3.29 (ii) Standard deviation of the Population

6 (iii) Mean of the Sampling distribution.

(iv) The Sample Standard deviation of Sampling distribution of means

2.32.

Sol: (i) mean = $\frac{2+3+6+8+11}{5}$

$$\mu = 6.0 \quad \sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n}$$

(i) $\sigma^2 = \frac{(2-6)^2 + (3-6)^2 + (6-6)^2 + (8-6)^2 + (11-6)^2}{5}$

$$\sigma^2 = 10.8$$

$$\sigma = 3.29$$

(iii) $n=2, N=5$

with replacement: $N^n = 5^2 = 25$ Samples

(2,2) (2,3) (2,6) (2,8) (2,11)

(3,2) (3,3) (3,6) (3,8) (3,11)

(6,2) (6,3) (6,6) (6,8) (6,11)

(8,2) (8,3) (8,6) (8,8) (8,11)

(11,2) (11,3) (11,6) (11,8) (11,11)

The Corresponding Samples are

$$\begin{array}{cccccc} 2.0 & 2.5 & 4.0 & 5.0 & 6.5 \\ 2.5 & 3.0 & 4.5 & 5.5 & 7.0 \\ 4.0 & 4.5 & 6.0 & 7.0 & 8.5 \\ 5.0 & 5.5 & 7.0 & 8.0 & 9.5 \\ 6.5 & 7.0 & 8.5 & 9.5 & 11.0 \end{array}$$

The mean distribution is

$$\bar{x} = \frac{\text{Sum of all Sample means (above)}}{25}$$
$$= \frac{150}{25}$$
$$= 6.0$$
$$\boxed{\mu = \bar{x}}$$

(iv) $\sigma_{\bar{x}} = \frac{(2-6)^2 + \dots + (11-6)^2}{25}$

$$= \frac{135}{25}$$

$$= 5.40$$

$$\sigma_{\bar{x}} = \sqrt{5.40}$$

$$= 2.32$$

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{3.29}{\sqrt{2}} = 2.32$$

8/1/2021:

* Moments:

- ✓ measure of turning or rotating in physics
- ✓ In Statistics these moments are used to measure the characteristics like dispersion, Skewness and Kurtosis of a frequency distribution.

Eg: find the first four moments for the set of numbers 2, 4, 6, 8

Sol: Here

$$N = 4$$

$$\text{mean} = \frac{2+4+6+8}{4} = 5$$

Numbers

	$d_i = x_i - 5$	d_i^2	d_i^3	d_i^4
x_1 2	-3	9	-27	81
x_2 4	-1	1	-1	1
x_3 6	1	1	-1	1
x_4 8	3	9	-27	81
$\sum x_i = 20$	$\sum d_i = 0$	$\sum d_i^2 = 20$	$\sum d_i^3 = 0$	$\sum d_i^4 = 164$

moments about mean

$$\mu_1 = \frac{\sum d_i}{N} = \frac{0}{4} = 0$$

$$\mu_2 = \sum \frac{d_i^2}{N} = \frac{20}{4} = 5; \quad \mu_3 = \sum \frac{d_i^3}{N} = \frac{0}{4} = 0$$

$$\sum \frac{d_i^4}{N} = \frac{164}{4} = 41$$

$$\beta_1 = \text{Skewness} = \frac{\mu_3^2}{\mu_2^{(3)2}} = 0$$

$\therefore \text{Skewness} = 0 \rightarrow \text{Normal distribution}$

$$\beta_2 = \text{Kurtosis} = \frac{\mu_4^2}{\mu_2^2} = \frac{41}{25} = 1.64$$

* Properties of central moments:

(i) The first moment about mean is always zero

$$\text{i.e } \mu_1 = 0$$

(ii) The second moment about mean measures

$$\text{Variance i.e } \mu_2 = \sigma^2 \text{ or } \sigma = \sqrt{\mu_2}$$

(iii) The third moment about mean measures Skewness of given distribution

(a) If $\mu_3 > 0$, the distribution is Positively Skewed.

(b) If $\mu_3 < 0$, the distribution is negatively Skewed

(c) If $\mu_3 = 0$, the distribution is Symmetrical