

❖ Probability and statistics .

★ → Random variable :- variable which take values as outcome of experiment

Ex:- dice $\{1, 2, 3, 4, 5, 6\}$

R.V $X = \{1, 2, 3, 4, 5, 6\}$

tossing of a coin is a Random experiment.

R.V $Y = \{H, T\}$

dice roll :

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

$$P(X \text{ is an even No}) = \frac{1}{2}$$

dice roll :

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

★ height of a randomly picked student

$$Y = 162.45$$

$$\text{or } Y = 132.62$$

or Y can be any real number.

X is discrete RV Y is continuous RV.

A discrete RV takes discrete values from a finite set of values.

A continuous RV can take any real value.

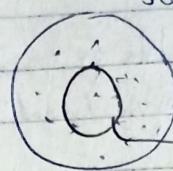
* Outlier: Y : height of student

$$\{122.2, 146.4, 132.5, \dots, \underbrace{12.26, 156.23}_{\text{outlier}}\}$$

* Population and sample:

estimate average human height

$$\bar{H} = \frac{1}{7B} \sum_{i=1}^{7B} h_i$$



Set of all the people in world.

It is very hard to measure height of each 7B+ humans.

But what we can do is we can take a sample out of population and calculate average height.

$$\bar{h} = \frac{1}{1000} * \sum_{i=1}^{1000} h_i$$

↳ sample size.

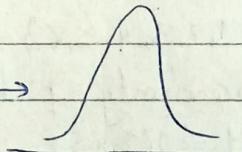
As sample size increases,

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

↳ sample mean. ↳ population-mean.

* Gaussian Distribution:

• PDF of a gaussian dist R.V \rightarrow



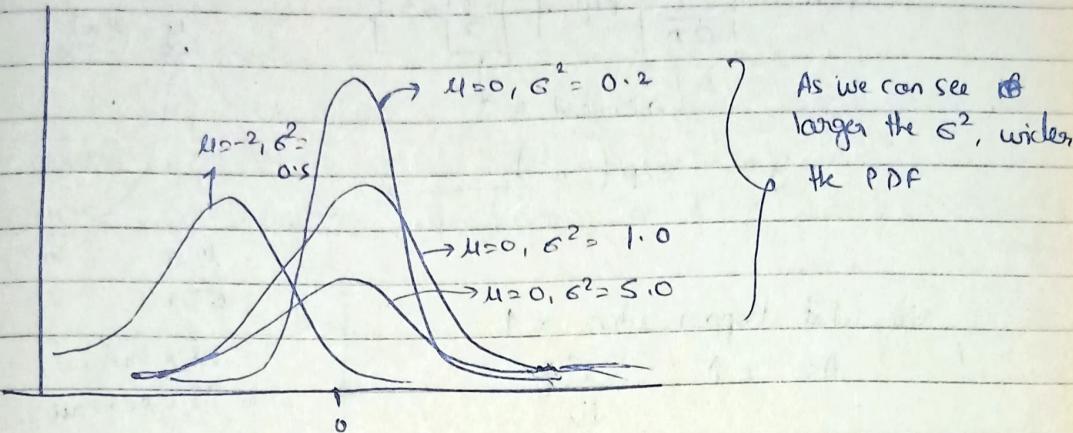
X: continuous R.V

Lot of ^{natural} things follows gaussian distribution, hence it becomes very important for us to learn gaussian distribution.

These distributions are very simple model, hence if we know that some feature follows a distribution we can predict lot of things about that feature.

mean and variance are parameters of a distribution because if we know ~~that~~ mean and variance we know which distribution is followed then we can draw PDF for those values.

~~so~~ But if we don't know the distribution followed by the values then mean and var are of no use.
Hence, to draw PDF we need all the three
 ① mean ② variance ③ type of distribution.



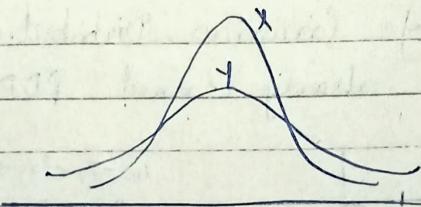
Parameters of Gaussian distribution : μ, σ^2

Let assume X is an R.V.,

$$X \sim N(\mu, \sigma^2)$$

X follows Normal/Gaussian distn. parameters.

Ex: $X \sim N(0, 2)$
 $Y \sim N(0, 4)$



~~$X \sim N(\mu, \sigma^2)$~~

$X \sim N(\mu, \sigma^2)$

$$P(X=x) = \frac{1}{\sqrt{2\pi} \sigma} \exp \left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\}$$

Let assume $\mu=0, \sigma^2=1$

$\sigma > 1$

$$P(x) = \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} \frac{x^2}{\sigma^2} \right)$$

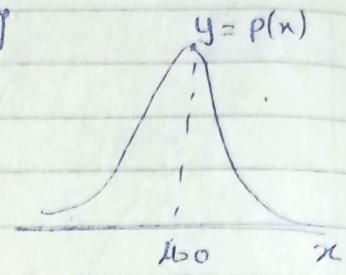
$$f(x) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}x^2\right\}$$

Let's try to plot the above eq.

$$f(x) = \left[\frac{1}{\sqrt{2\pi}} \right] \exp\left\{-\frac{1}{2}x^2\right\}$$

\downarrow constant

$$\therefore y = \exp(-x^2)$$



So, what happens when $x \uparrow$

As $x \uparrow$, $-x^2 \downarrow$
 \downarrow

$\exp(-x^2) \downarrow$

$$6^2 = 1$$

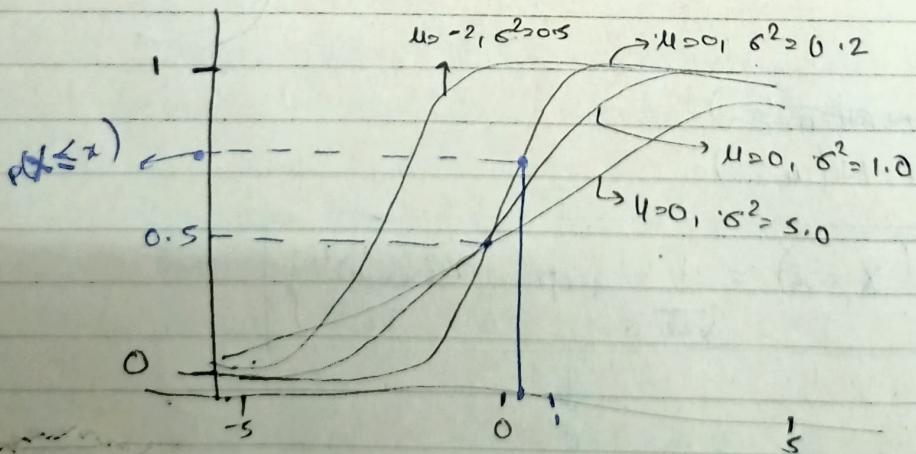
As x increases
y decreases.

Conclusion:

- (i) As x moves away from 0, y ↓
- (ii) y is a symmetric funcn.
- (iii) As x move away from 0, y reduces as $\exp(-x^2)$

→ CDF of Gaussian Distribution:

As we already learned PDF of Gaussian Distribution.



Height of CDF says $P(X \leq x)$

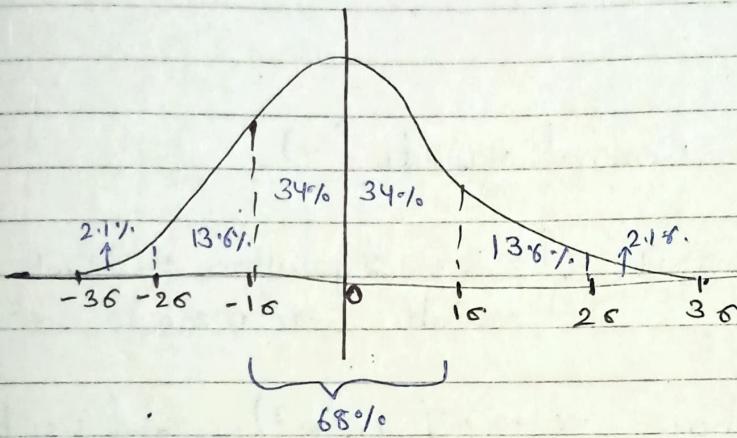
$$CDF = \frac{1}{2} \left[1 + \operatorname{erf} \left(\frac{x-\mu}{\sigma\sqrt{2}} \right) \right]$$

$$P(X \leq x) = \frac{1}{2} \left[1 + \operatorname{erf} \left(\frac{x-\mu}{\sigma\sqrt{2}} \right) \right]$$

* 68-95-99.7 rule.

(let we are given $X \sim N(\mu, \sigma^2)$ where $\mu=0, \sigma^2=4$
~~base~~ $\sigma^2=4 \rightarrow \sigma=2$

We can say lot of thing about X using 68-95-99.7 rule



B/w -1σ to 1σ we have 68% of our point

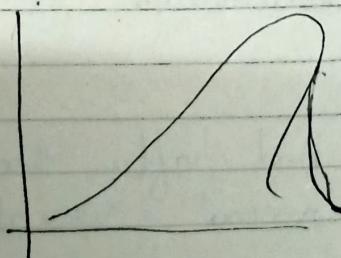
B/w -2σ to 2σ we have 95% " "

B/w -3σ to 3σ we have 99.7% " "

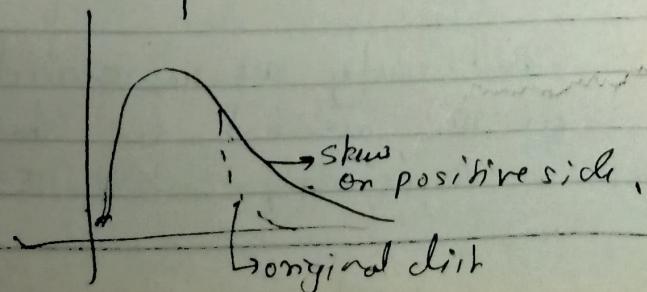
→ Symmetric distribution

→ Non-symmetric distribution.

→ Negative skew

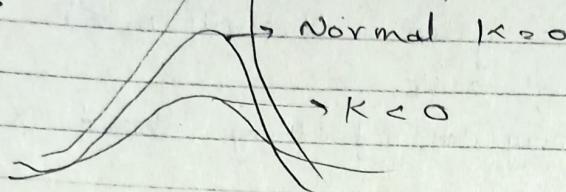


→ Positive skew



→ Kurtosis:- says how sharp the peak is in distribution. For a Normal/Gau distribution $K = 3$

Lesser the kurtosis flatter / or less peaked in distribution. $\rightarrow K > 3$



→ Standard Normal variate (z)

① $z \sim N(0, 1)$ { z is a Gaussian distribution variable with mean $= 0$, $\sigma^2 = 1$

② Let we have $x \sim N(\mu, \sigma^2)$, and we have sampled $[x_1, x_2, \dots, x_{50}]$

Given these observations of R.V x we can standardise these values. We can apply standardisation.

Standardisation: For x in X , define x'
Set $x'_i = \frac{x_i - \mu}{\sigma}$

And after that we get x' which is standard normal variate.

But why are we standardising?
At the moment we standardise our data we can apply 68-95-99.7 rule.

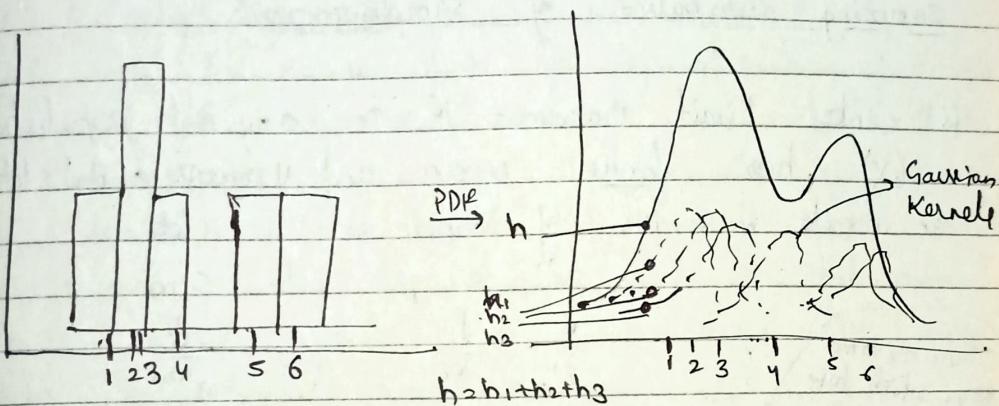
i.e. 68% of our x_i lie b/w -1 and 1,
($\sigma = 1$) here.

95% of x_i lie b/w -2 and 2

* Kernel density estimation

As we already know, that given a histogram we can get PDF from that histogram. We use Kernel density estimation to construct PDF using histogram.

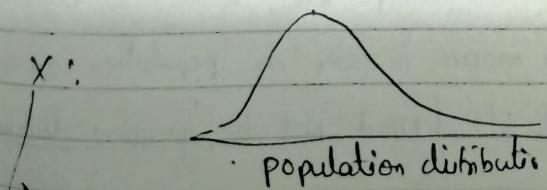
Ex:-



for each p-data point we draw Gaussian Kernel function. At any point we sum up all the Kernel heights at that point to get the height of PDF at that point variance of Kernel is called bandwidth.

* Sampling distribution:-

Assume we have distribution of R.V X (~~not necessarily gaussian distn~~)



X is distribution of incomes of all people in the world.

Now suppose we pick a sample S_1 from the population, suppose sample size is n .

let say $n=30$, Sample S_1 = sample of 1000 observations
let again sample S_2 of size 30

Let we keep doing the random sampling of 30 people, and got in samples, and computing means of each sample.

$$S_1 \rightarrow \bar{x}_1$$

$$S_2 \rightarrow \bar{x}_2$$

$$S_3 \rightarrow \bar{x}_3$$

$$\vdots$$

$$S_m \rightarrow \bar{x}_m$$

So we have $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_m$, there are ~~less~~ values also have a distribution, ~~so~~ this distribution is called sampling distribution of sample mean.

* Centre Limit theorem: If our original population (X) has finite mean and variance. And let's assume we create m samples of sample size n , $n \rightarrow \infty$.

There are some dist. which have infinite mean & variance
like Pareto dist.

$$\therefore S_1, S_2, \dots, S_m$$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$\bar{x}_1, \bar{x}_2, \bar{x}_m$$

Sample mean.

\Rightarrow dist. of \bar{x}_i = sampling dist. of sample mean

so, CLT says:

\bar{x}_i is distributed with a gaussian distribution with mean μ & variance $\frac{\sigma^2}{n}, n \rightarrow \infty$

Variance of population

$$\boxed{\bar{x}_i \sim N\left(\mu, \frac{\sigma^2}{n}\right) n \rightarrow \infty}$$

sample size

This mean is same as population mean

Ex: X : income of people: (Need Not be gaussian distibuted)

$$\rightarrow \mu, \sigma^2$$

$$\rightarrow S_1, S_2, S_3, \dots, S_m$$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$\bar{x}_1, \bar{x}_2, \bar{x}_3$$

$$\vdots$$

$$\bar{x}_m$$

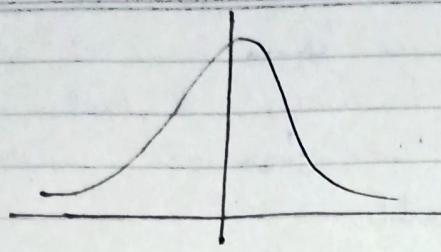
$$(|S| = 30)$$

$$m = 1000$$

Ans T

can mean income in each sample.

\downarrow plot a dist. of means :



\bar{x}_i mean of \bar{x}_i = population mean.

Hence we can predict the mean of population using mean of sample mean.

Hence, for any distribution, to get their mean & variance we just need to ensure that the population has finite mean & variance and we are done, we can calculate the mean & pop. var of a very large popule. by just using CLT in which we just take small samples, calculate their means and the distri of that ~~means~~ sample mean follows a Normal/Normal dist with mean = population mean & vari = $\frac{\sigma^2}{n}$.

* Quantile - Quantile (Q-Q) plot:-

Let we have R.V $x: x_1 x_2 x_3 \dots x_{500}$, How can we know if x is gaussian distribut. To Ans this question we have various techniques, Q-Q plot is one of them., Q-Q plot is graphical methode.

→ How to plot Q-Q plot:

(1) Sort x_i & compute percentiles.

$$x_1, x_2, \dots, x_{500} \quad \downarrow \text{Sort}$$

$$x'_1, x'_2, \dots, x'_{500} \quad (\text{s.t. } x'_1 \leq x'_2 \leq x'_3 \dots)$$

$$\textcircled{1} \quad \textcircled{2} \quad \textcircled{3} \quad \downarrow \text{Percentile}$$

$$x'_5, x'_{10}, x'_{15}, \dots, x'_{500} \quad 100$$

$$(x^{(1)}), (x^{(2)}), (x^{(3)}), \dots, (x^{(100)})$$

(1) Create a R.V $y \sim N(0, 1)$, create 1000 obserm from y : $y_1, y_2, \dots, y_{1000}$

$y_1, y_2, y_3, \dots, y_{1000}$
↓ sort

$y'_1, y'_2, y'_3, \dots, y'_{1000}$
↓ Percentile

$y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(100)}$

③ plot Q-Q plot using $x^{(1)}, x^{(2)}, \dots, x^{(100)}$
 $y^{(1)}, y^{(2)}, \dots, y^{(100)}$

we plot each pair

$(y^{(1)}, x^{(1)})$

$(y^{(2)}, x^{(2)})$

⋮

$y^{(100)}, x^{(100)}$

$x^{(5)}$

$(y^{(5)}, x^{(5)})$

$y^{(5)}$

$y \sim N(0, 1)$

* If the plot is roughly create a straight line then x has a gaussian distribution

* If plot is not straight line the x has not gaussian distribution.

So, we can check for any distribution, we change y distribution and check if x follow y 's distribution.

Code:

$\text{std_normal} = \text{np.random.normal}(\mu^{\text{mean}}, \sigma^2)$

$\text{std_normal} = \text{np.random.normal}(\text{loc} = 0, \text{scale} = 1, \text{size} = 1000)$

$\text{measurements} = \text{np.random.normal}(\text{loc} = 20, \text{scale} = 5, \text{size} = 100)$

`stats.probplot(measurements, dist = "norm", plot = pylab)
pylab.show()`:

* Limitation of Q-Q: If no. of sample obsrns in R.V is small then it becomes hard to interpret if R.V really follow Gaussian distn or not.

Using Q-Q we can check any 2 R.V if they have same distribution or not.

* Till now we have learned, R.V, PDF, CDF, Gaussian distn.
But question is How/where to use these distn in real world
These concepts are used in Exploratory Data Analysis.

* Chebyshev's Inequality:-

As we already know, if X is Normal dist and has mean μ and std. dev. σ

$$\text{P}(\mu - 2\sigma \leq X \leq \mu + 2\sigma) = 95\%$$

$$\text{P}(\mu - \sigma \leq X \leq \mu + \sigma) = 68\%$$

$$\text{P}(\mu - 3\sigma \leq X \leq \mu + 3\sigma) = 99.7\%$$

Q) what if, we don't know the distribution, but we know mean = μ , and std. dev. σ , and we also know that mean is finite and $\sigma \neq 0$ and finite.

In this case we can use Chebyshev's Inequality:

Chebyshev's Inequality: If X is a R.V with finite mean μ and a non-zero and finite std. dev. σ . We don't know the distribution of X . Then.

$$\boxed{\text{P}(|X - \mu| \geq k\sigma)} \leq \frac{1}{k^2}$$

OR

$$\boxed{\text{P}(\mu - k\sigma < X < \mu + k\sigma) \geq 1 - \frac{1}{k^2}}$$

A Uniform distribution

There are 2 types of UD

(i) Discrete uniform distn

(ii) Continuous " "

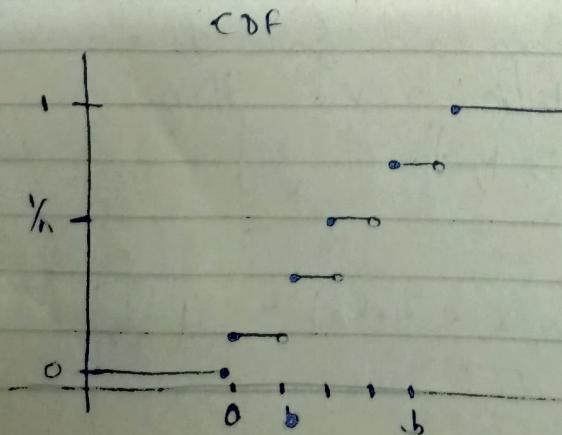
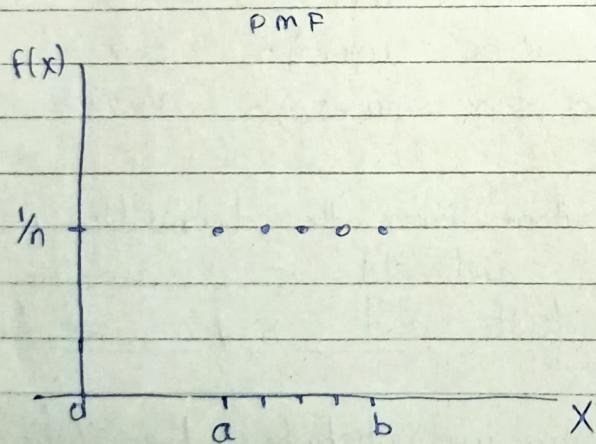
* Similar to continuous distribution, discrete distribution have Probability Mass function (PMF). It is similar to PDF, but PDF is for continuous R.V!

There are 2 parameters of Uniform dist, a and b .
 a and b are integers s.t. $b \geq a$.

$$n = b - a + 1$$

Notation $\rightarrow \text{unif}(a, b)$

In a uniform distribution all the values are equally probable., i.e. each value have prob = $\frac{1}{n}$



Notation : $\text{U}(a, b)$ or $\text{unif}(a, b)$

Parameter : $a \in \{-\dots, -2, -1, 0, 1, 2, \dots\}$

$b \in \{-\dots, -2, -1, 0, 1, 2, \dots\}, b > a$

$$n = b - a + 1$$

Support : $\{a, a+1, a+2, \dots, b-1, b\}$

PMF = $\frac{1}{n}$

CDF = $\frac{\lfloor x \rfloor - a + 1}{n}$

mean = $(a+b)/2$

Media = $(a+b)/2$

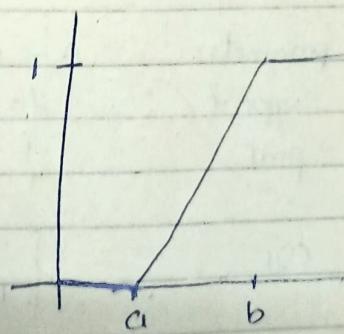
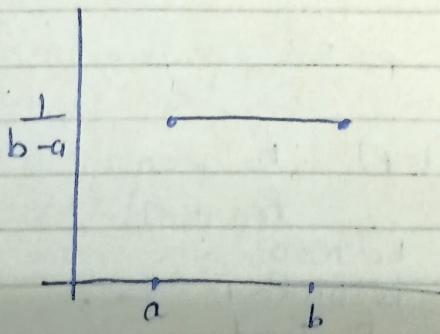
Mode = N/A

Variance = $\frac{(b-a+1)^2 - 1}{12}$

① Continuous UD: the continuous uniform distri or rectangular distribution is a family of symmetric probability distribution s.t for each member of the family, all intervals of same length are equal probable.

parameters : a, b

As it has have continuous R.V, It have PDF



PDF

Notation: $\text{U}(a, b)$ or $\text{unif}(a, b)$

Mean: $\frac{1}{2}(a+b)$

Parameter: $-\infty < a < b < \infty$

Media: $\frac{1}{2}(a+b)$

Support: $x \in [a, b]$

Variance: $\frac{1}{12}(b-a)^2$

PDF : $\begin{cases} \frac{1}{b-a} & \text{for } x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$

CDF : $\begin{cases} 0 & \text{for } x < a \\ \frac{x-a}{b-a} & \text{for } x \in [a, b] \\ 1 & \text{for } x > b \end{cases}$

* Random Number generator : UD

C++ , `rand()` or ~~or~~ in any language `rand()` function uses uniform distribution to generate a random numbers.

UD can also be used to uniformly sample the data.
Suppose our iris data set which has 150 samples. From 150 we want a sample of 30 data.

$$n = 150$$

$$m = 30$$

$$p = m/n$$

$$\text{sampled_data} = []$$

for i in range(0, n):

if `random.random() <= p:`

`sampled_data.append(d[i, :])`

`rand()` used UD and gives a no. b/w (0-1) \rightarrow treat that no. as prob.

and check if prob $\leq p$

then append

to `sampled_data`

* Bernoulli distributions

Bernoulli distribution is used when R.V can take only 2 values. Like a coin toss experiment we have 2 outcome probability of one outcome is p and other outcome is $1-p$.

parameters : $0 < p < 1$, $p \neq 0.5$

Support : {0, 1}

pmf :
$$\begin{cases} q = (1-p) & \text{for } k=0 \\ p & \text{for } k=1 \end{cases}$$

cdf :
$$\begin{cases} 0 & \text{for } k<0 \\ 1-p & \text{for } 0 \leq k < 1 \\ 1 & \text{for } k \geq 1 \end{cases}$$

mean : p

media :
$$\begin{cases} 0 & \text{if } q > p \\ 0.5 & \text{if } q = p \\ 1 & \text{if } q < p \end{cases}$$

variance : pq

* Binomial Distribution

Binomial Distribution with parameters n and p is the discrete probability distribution of the no. of successes in a seq. of n independent experiments, each asking a yes-no question.
When $n=1$, Binomial distribution = Bernoulli distribution.

Notation : $B(n, p)$

parameters : $n \in \mathbb{N}_0$ - no. of trials.

$p \in [0, 1]$ - success probability in each trial.

support : $K = \{0, 1, \dots, n\}$

pmf : $n \binom{k}{n} p^k \cdot (1-p)^{n-k}$

cdf : $I_{1-p}(n-k, 1+k)$

Mean : np

Variance : $np(1-p)$

* Log-normal distribution

- continuous probab. distn. of a R.V whose logarithm is normally distributed. Thus if the R.V x is log-normally distributed, $Y = \log(x)$ has a normal distn.

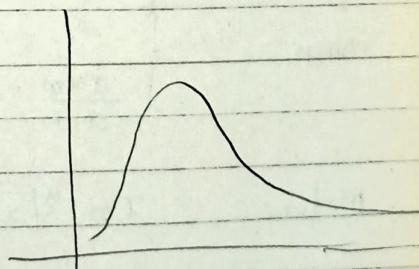
Notation : $\text{lognormal}(\mu, \sigma^2)$

Parameter : $\mu \in (-\infty, +\infty)$

$\sigma > 0$

Support : $x \in (0, +\infty)$

PDF : $\frac{1}{x \sigma \sqrt{2\pi}} \cdot e^{-\frac{(\log x - \mu)^2}{2\sigma^2}}$



mean : $\exp(\mu + \sigma^2/2)$

median : $\exp(\mu)$

Variance : $[\exp(\sigma^2) - 1] \exp(2\mu + \sigma^2)$

* The length of comments posted in Internet discussion forums follows a log-normal distribution.

* The user's dwell time on online articles follow a log-normal distn.

* Income of 97% - 99% population is distributed log-normally.

* Power law distribution:-

Assume we have 2 variables, a power law is a functional relationship b/w two quantities, where a relative change in one quantity results in a proportional change relative change in other quantity,

independent of the initial size of those quantities : one quantity varies as a power of another.

For ex:- area of square in terms of length of its side if the length is doubled, area is multiplied by factor of 4

→ Pareto distribution :- when a distribution follows power law it is called pareto distribution

Parameters : $x_m > 0$ scale (real)

$\alpha > 0$ shape (real)

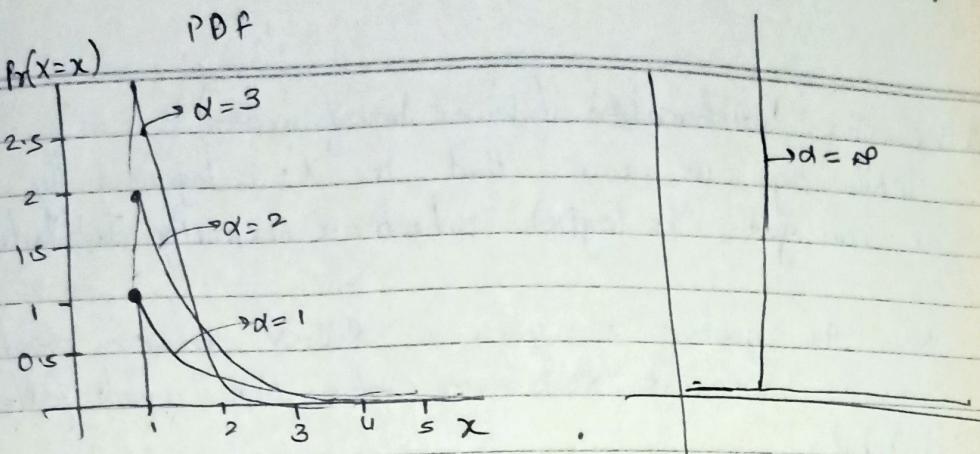
Support : $x \in [x_m, +\infty]$

PDF : $\frac{\alpha x_m^\alpha}{x^{\alpha+1}}$ for $x \geq x_m$

CDF : $1 - \left(\frac{x_m}{x}\right)^\alpha$ for $x \geq x_m$

mean : $\begin{cases} \frac{\alpha}{\alpha-1} x_m & \text{for } \alpha > 1 \\ \frac{\alpha x_m}{\alpha-1} & \text{for } \alpha \leq 1 \end{cases}$

Median : $x_m \sqrt[{\alpha}]{2}$



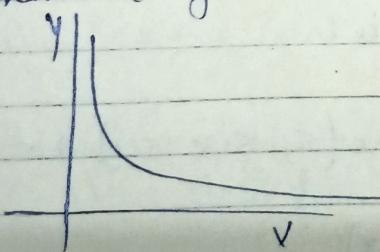
$\alpha = \infty$: give a ~~discrete~~ Dirac delta func, when $\alpha = \infty$ there is only one value and give a peak value.

In pareto distri we can see there is a peak value and then distribution starts falling.

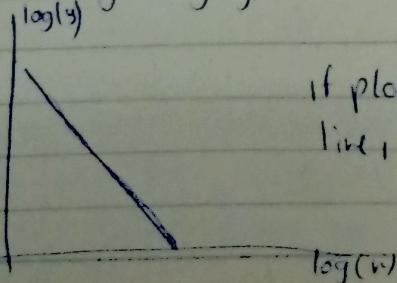
- file size dist of Internet traffic uses the TCP proto follower pareto distri
- Hard disk drive error rates follow pareto distribution.

* log-log plot is used to determine if the given plot is power law plot or not.

Let assume we have 2 vari: x & y ; when we plot them we get something like this



To check if the above plot is power law or not, plot another plot of $\log(y)$ and $\log(n)$



If plot comes out like this a straight line, then the two variables follow power law.

* Box-Cox transformation (Power transformation)

while ago we saw that if X is log-normally distributed we can get $Y = \log(x)$ which is normally distributed.

Now the question is given a R.V X having pareto distribution can we convert X to some other R.V which follows gaussian distribution.

Ans is Yes, we can do that using Box-Cox distribution.

Pareto $x : x_1, x_2, \dots, x_n \rightarrow$ conversion.

Gaussian $y : y_1, y_2, \dots, y_n \rightarrow$

$$\textcircled{1} \quad \text{boxcox}(x) = \lambda \text{lambda}(\lambda)$$

$x_1, x_2, x_3, \dots, x_n$

$$\textcircled{2} \quad y_i = \begin{cases} \frac{x_i^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \log(x_i) & \text{if } \lambda = 0 \end{cases} \quad i : 1 \rightarrow n$$

This case suggest that X is lognormal when $\lambda = 0$

code:

```
from scipy import stats
```

$x = \text{stats.loggamma.rvs}(5, size=500) + 5$ R.V
of
gamma
dist.
It is one
of pareto
dist

$xt, - = \text{stats.boxcox}(x)$

* Application of non-gaussian distribution

① Uniform \rightarrow Random generate

② Bernoulli, Binomial

③ Log-normal, pareto

and many more non-gaussian dist are there.

* Covariance :- lets assume we have 2 R.V 'x' and 'y'

x = heights of student

y = weight of student

Q: is there a relationship b/w x & y

There are 3 measure to check if 2 R.V are related

① Covariance

② Pearson correlation coeffs

③ Spearman rank corr. coeffs.

W's understand co-variance :-

$$\text{cov}(x,y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

(mean of x) (mean of y)

$\text{cov}(x,y) = +ve$, if $x \uparrow$ then $y \uparrow$

$\text{cov}(x,y) = -ve$ if $x \uparrow$ then $y \downarrow$

* One problem with cov is that for same data set if we change the unit the cov may not be same. suppose we find value $\text{cov}(x, y)$ then

$\text{cov}(x, y)$ both of them will not be equal. in cm kg

(ii) Pearson correlation coeff (PCC)

$$\text{PCC} = \rho_{x,y} = \frac{\text{Cov}(x,y)}{(\sigma_x \cdot \sigma_y)}$$

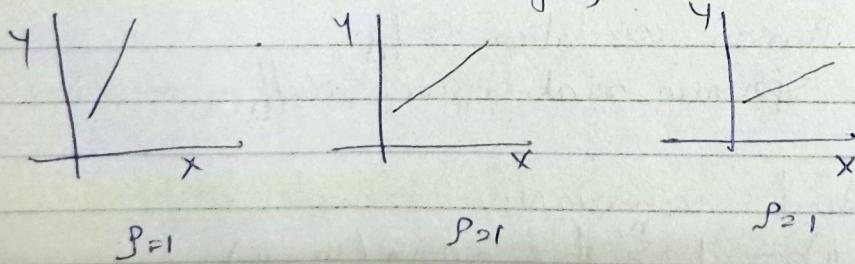
problem with cov was if was not taking into account of variability that std.

In cov we knew that if $x \uparrow$ and $y \uparrow$, the cov is +ve but we didn't know how much positive.

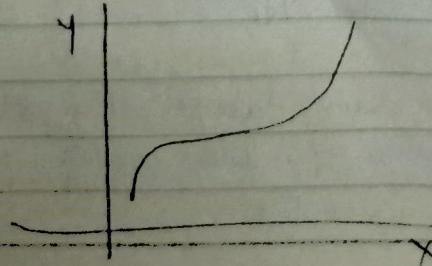
value of ρ is always $-1 \leq \rho \leq 1$
 when $\rho = -1$ that means as $x \uparrow$ $y \downarrow$ perfectly
 when $\rho = 1$ " " " " $x \uparrow$ $y \uparrow$ perfectly
 on straight line.
 $\rho = 0$ if there is no relation b/w X and Y

But ~~Pearson~~ PCC also fails, PCC works well only if X and Y have linear relationship.

The other thing is if X and Y are lying in straight line then slope of line is not considered that means it does not tell anything about how the R.V are changing.



(iii) Spearman rank corr coeff : This works for some slightly linear relationships.



$$\rho_{x,y} = 0.88$$

$$\text{Spearman} = 1$$

	$X = \text{height}$	$Y = \text{weight}$	r_{xy}	r_y
S1	160	52	4	3
S2	150	66	2	4
S3	170	68	5	5
S4	140	46	1	1
S5	158	51	3	2

$$r = r_{xy}$$

Rank-cor : doesn't care about linearity, it just check if $x \uparrow y \uparrow$ or $x \uparrow y \downarrow$ using ranks of x & y . \therefore

→ Correlation does not imply causation.

Correlation says only if $x \uparrow y \uparrow$ but we cannot conclude that y happens due to x . Just because 2 f.v are correlated that doesn't mean that x causes y or y causes x .

Application of Correlation:

* C.I for mean(μ) of a R.V X any distn not nec. gaussian.
 let a R.V $X \sim N(\mu, \sigma^2)$ with pop-mean of μ & std-dev of σ .
 $\{x_1, x_2, \dots, x_{10}\} \rightarrow$ sample of size, $n=10$

Q) What is the 95% confidence Interval (C.I.) of μ .

Case 1:- Assume somebody told us the popule std-deviation

* Confidence Interval :- Assume we have R.V X we don't know distribution

Sample: $\{x_1, x_2, \dots, x_{10}\} \rightarrow$ sample of size 10 from X
 we want to estimate population mean of X -
 one way is we can say populat. mean will approx eq. to sample mean.

$$\mu \approx \bar{x}$$

pop mean ↳ sample mean

But This is not but approximation:-

Suppose $\{x_1, x_2, \dots, x_{10}\}$

heights of population
 $\{180, 162, 158, 172, 168, 150, 171, 183, 165, 176\}$

$$\text{POINT ESTIMATE of } \mu = \frac{1}{10} \sum_{i=1}^{10} x_i = 168.5 \text{ cm}$$

other thing we can say that

$\mu \in [162.1, 174.9]$ with 95% probability.

↳ Confidence Interval

* Computing C.I given the underlying distribution:-

$$X \sim N(\mu, \sigma^2); \mu = 168 \text{ cm}$$

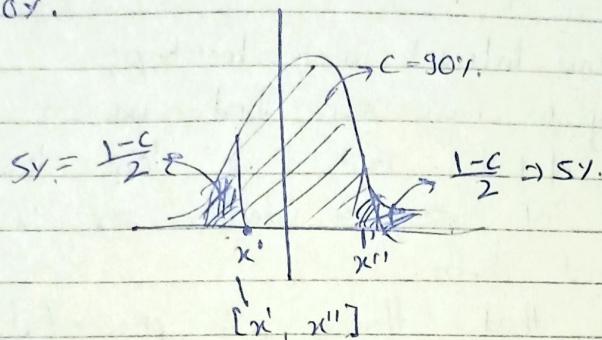
heights. ↳ Gaussian distn. $\sigma = 5 \text{ cm}$

Q) C.I. of height?

for a gaussian distn. we know that 95% of data lie in $[\mu - 2\sigma, \mu + 2\sigma]$
 $\hookrightarrow [158, 178]$

hence this is C.I. with 95% probability that height lie in interval $[158, 178]$

Suppose someone ask a interval with confidence = 90% or 80% of 70%.



We can use Normal distn table and find the value of x', x'' .

Q1. for mean(μ) of a R.V

Let a R.V X with some ~~distn~~ (Not known) with pop-mean μ & std-dev σ .

we have a sample of size 10 of $x_1, x_2, x_3, \dots, x_{10}$
 $\{180, 162, 158, 172, 168, 150, 171, 183, 165, 176\}$

Q) what is the 95% C.I. of μ .

There are 2 cases for this

Case 1: Assume somebody told us the pop.pop std-devi:

$$\sigma = 5\text{cm}$$

Acc. to CLT: $\bar{x} = \text{sample mean} = \frac{1}{10} \sum_{i=1}^{10} x_i$
 $\bar{x} \sim N(\mu, \frac{\sigma}{\sqrt{n}})$

Now, we can say μ lies in $[\bar{x} - \frac{2\sigma}{\sqrt{n}}, \bar{x} + \frac{2\sigma}{\sqrt{n}}]$ with 95%.

Case II : If we don't know σ
we use t-distribution

Acc. to t-distr. $\bar{x} \sim t(n-1)$
sample mean. \hookrightarrow degrees of freedom

→ t-distribution :- This distr. was developed to calculate
C.I. when pop. std. mean is unknown.

* Confidence Interval using bootstrap:-

Previously we have seen that we can calculate C.I. for
 μ of R.V. But we don't know how to calculate
C.I. of σ , CI, media or 90th per centile or any of
statistical data.

For that there is a powerful technique called
bootstrap which can be used to calculate C.I.
for any of the above term (σ , median etc.)

→ $X \sim$ Any distr.

task : estimate C.I. of median of X .

Let's assume we have a sample of size n

$$\{x_1, x_2, \dots, x_n\}$$

using only this we need to compute C.I. of median
of x . We don't know type of distr.

① Generate new sample using give sample.

$$S_1 : x_1^{(1)}, x_2^{(1)}, x_3^{(1)}, \dots, x_m^{(1)} \text{ s.t } m \leq n$$

↑ \hookrightarrow random sample of size
generated from S .
Repetition is allowed
i.e. some value can be repeated.

One way to do this is generate a Uniform R.V. $U(1, n)$
and generate Random sample m times with U R.V.

$S = \{n_1, n_2, n_3, \dots, n_n\}$
 Using sampling with replacement

$S_1 : x_1^{(1)}, x_2^{(1)}, x_3^{(1)}, \dots, x_m^{(1)}$
 $S_2 : x_1^{(2)}, x_2^{(2)}, x_3^{(2)}, \dots, x_m^{(2)}$
 $S_3 : x_1^{(3)}, x_2^{(3)}, \dots, x_m^{(3)}$
 \vdots
 $S_K : x_1^{(K)}, x_2^{(K)}, x_3^{(K)}, \dots, x_m^{(K)}$

} bootstrap samples

for each sample compute median

$S_1 \rightarrow m_1$

$S_2 \rightarrow m_2$

$S_3 \rightarrow m_3$

\vdots
 $S_K \rightarrow m_K$

let say $K=1000$

$m_1, m_2, m_3, \dots, m_{1000}$

↓ sort

$m_1, m_2, m_3, \dots, m_{1000}$ (increasing order)

↓

$[m_{25}, m_{75}]$

↓

950 value → 95% of 1000
 set of

This process can be done for anything i.e. $\mu, \sigma^2, \text{percentiles}$.

This called a non-parametric technique.

→ doesn't make any assumption about distribution.

Hypothesis testing:-

Assume we have 2 classes C₁ & C₂ and assume we collect the height of 50 students from both classes.

C ₁	C ₂
1 160	1 162
2 152	2 156
:	:
!	!
50 158	50 152

a) Is there a diff in heights of students in C₁ & C₂ students.

To ans this ques we use Hypothesis testing:-

① choosing a test-statistics :- choose a statistical parameter

$$\mu_2 - \mu_1 \\ \text{mean of } C_2 \leftarrow \text{mean of } C_1$$

② Null hypothesis (H₀)

H₀: no difference in μ_1 & μ_2

Alternative hypothesis (H₁):- Inver of H₀

H₁: there is diff in μ_1 & μ_2 .

Assume H₀ is true.

③ compute p-value:- p-value says what is the probability of observing your test statistic if ~~not~~ H₀ is true. Assume we got $\mu_2 - \mu_1 = 10$, then p-value says what is the prob to get $\mu_2 - \mu_1 = 10$ if H₀ is true.

If p-value = 0.9 it means that probability of $\mu_2 - \mu_1$ is 0.9 if H₀ is true.

If p-value is High then we accept H_0 is true.

Ex:- Given a coin, determine if the coin is biased towards heads or not.

~~biased towards~~.

We can answer questions like this ~~using~~ using Hypothesis testing.

Initially, we try to prove that ~~is~~ using basic probability,
→ Using Basic prob :-

→ design a expt: Flip a coin 5 times and count No. of heads.

Let No. of Head = $x \leftarrow$ Test statistic

→ perform the experiment :- flip, flip, flip, flip, flip

H H H H H

$x = 5 \leftarrow$ observation by experimenter.

$P(x=5 \mid \text{coin is not biased towards head}) =$

H_0

$$P(x=5 \mid H_0) \Rightarrow \frac{1}{(2)^5} = (0.5)^5 = 3\%$$

There is 3% chance of getting 5 Heads in 5 tosses if the coin is not biased towards Head.

Hypothesis testing works like this:-

$$P(\text{obs by exp/assumption}) = 3\% \quad \text{if small}$$

p-value.

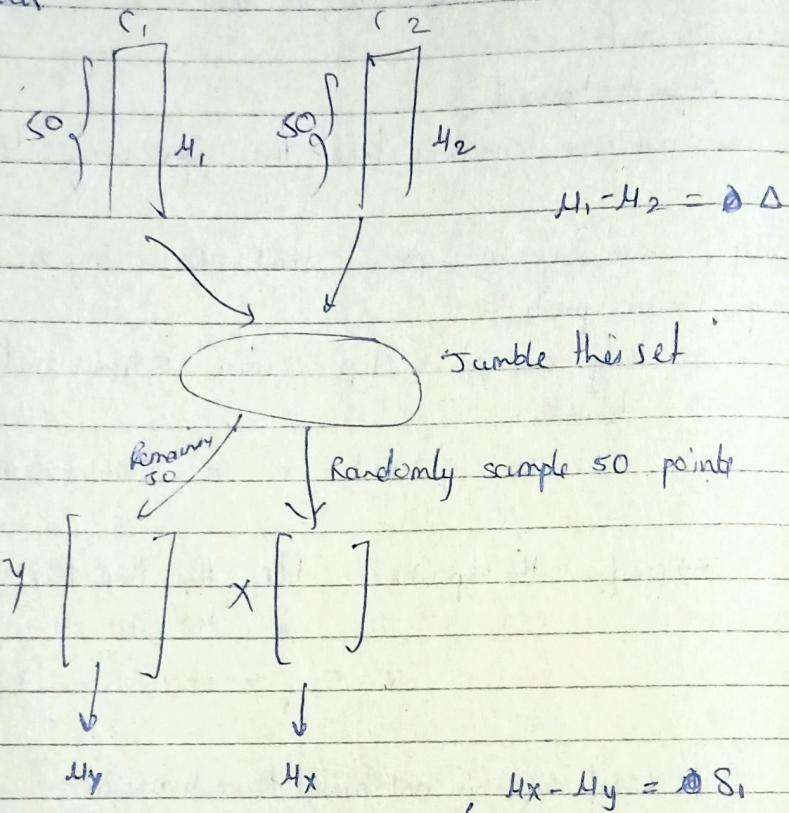
If $P(\text{Obs}/H_0) < 5\%$. Then H_0 is not true.

Hence coin is biased.

We need to be careful about
① Sampling or ~~design~~ experiment size
② Defining Null hypothesis.
③ Design X.

* p-value calculation (permutation testing):

Assume we have C_1 and C_2 and so height from both class



Again resample the jumbled (mixed) set and again calculate $H_y \times H_x$ and Repeat this process say 10K time.

We get 10K Δ_s 's

$\Delta_1, \Delta_2, \dots, \Delta_{10K}$

($\Delta_1, \Delta_2, \dots, \Delta_{10K}$) \curvearrowright sort

$\Delta_1', \Delta_2', \Delta_3', \dots, \Delta_{10K}'$

Find correct position of Δ in sorted Δ 's

If there are 5% value greater than Δ in sorted Δ 's then p-value = 0.05%.

If Δ value greater than Δ in sorted Δ 's
then p-value $\propto 1/\alpha$

Kolmogorov-Smirnov

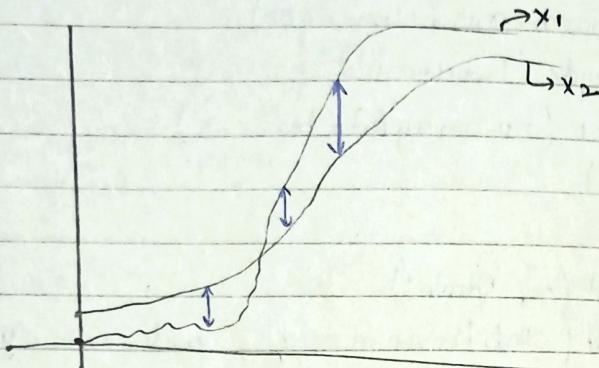
A K-S test for similarity of 2 distribution

2 R.V $X_1: [x_1^1, x_2^1, \dots, x_n^1]$

$X_2: [x_1^2, x_2^2, \dots, x_m^2]$

is distribution X_1 is same as X_2

first we plot CDF of both X_1 & X_2



H_0 : X_1 & X_2 have same distribution

$$\text{Test statistic} = D_{n,m} = \sup_x \left| \underbrace{F_{1,n}(x)}_{\text{CDF}} - \underbrace{F_{2,m}(x)}_{\text{CDF}} \right|$$

maximum difference b/w CDF of X_1 & X_2

If $n \times m$ are large then $D_{n,m} = 0$ and hence H_0 is true.

The null hypothesis is rejected at level α if

$$D_{n,m} > c(\alpha) \sqrt{\frac{n+m}{nm}}$$

when $n \times m$ are small we get higher cutoff to reject H_0

where $n \times m$ as sizes of X_1 & X_2 . get a small cutoff

α	0.10	0.05	0.025	0.01	0.005	0.001
$c(\alpha)$	1.22	1.36	1.48	1.63	1.73	1.95

to reject H_0 .

$$c(\alpha) = \sqrt{\frac{-1}{2} \ln \left(\frac{\alpha}{2} \right)}$$

code

```
import numpy as np  
import seaborn as sns  
from scipy import stats  
import matplotlib.pyplot as plt
```

```
x = stats.norm.rvs(size=1000)  
sns.set_style('whitegrid')  
sns.kdeplot(np.array(x), bw=0.5)  
plt.show()
```

stats.ks_2samp(x, 'norm')

O/P : Ks_2sampResult(statistic=0.021308, pvalue=0.75424)

```
y = np.random.uniform(0, 1, 10000);
```

stats.ks_2samp(y, 'norm')

Ks_2sampResult(statistic=0.501596, pvalue=0.0).

* Another example :- Hypothesis testing

Task :- two cities C₁ & C₂ determine if population mean of heights of people living in these cities is same or not.

C₁ C₂

problem is there are lot of people we can't calculate by going to each people in city and calculate height.

① Design experiment :- Randomly pick 50 people from each city

$\begin{matrix} C_1 \\ \text{So } \left[\begin{matrix} h_1 \\ h_2 \\ \vdots \\ h_{50} \end{matrix} \right] \end{matrix}$ $\begin{matrix} C_2 \\ \text{So } \left[\begin{matrix} h_1 \\ h_2 \\ \vdots \\ h_{50} \end{matrix} \right] \end{matrix}$ → Samples of size 50

Sample mean $M_1 = \frac{h_1 + h_2 + \dots + h_{50}}{50}$

Sample mean $M_2 = h_1 + h_2 + \dots + h_{50}$

(i) Let statistics : $\mu_1 - \mu_2 = x$

$$\text{Let } \mu_1 = 162 \quad \mu_2 = 167$$

$$|\mu_1 - \mu_2| = 5$$

(ii) Null Hypothesis :- There is no difference in population means of height of people of 2 cities.

(iii) p-value : $p(x = 5 \text{ cm} | H_0)$

diff in sample mean
of sample size of 50

How to calculate $p(x = 5 \text{ cm} | H_0)$:- using Resampling & Permutations

* proportional sampling :-

Let we have g. matrix d with 5 elements

$$d = \begin{array}{|c|c|c|c|c|} \hline & 2.0 & 6.0 & 1.2 & 5.8 & 20.0 \\ \hline 1 & 2 & 3 & 4 & 5 = n \\ \hline \end{array}$$

Task :- I want to pick an element amongst these n-elements s.t. the probability of picking an element is proportional to d_i 's

If we randomly pick a value, each of the value is equally probable. But we don't want to pick randomly. But we pick element s.t. picking the 5th element is very large because 5th element is maximum.

Similar, the prob. of picking the 3rd element must be extremely small because 3rd element is least.

So, how do we do it.

Step 1a : $S = \sum_{i=0}^5 d_i$ { sum of all values }

$$S = 2 + 6 + 1.2 + 5.8 + 20 = 35$$

Step 1b : $\rightarrow d'_i = d_i / S$ { Normalizing the values }

$$d'_1 = 0.0571, d'_2 = 0.17118, d'_3 = 0.034, d'_4 = 0.165, d'_5 = 0.571$$

Step 0 i.e.: Cumulative Normalized sum

$$\begin{array}{ll} d_1 = 0.0571 & \tilde{d}_1 = d_1 = 0.0571 \\ d_2 = 0.171428 & \tilde{d}_2 = \tilde{d}_1 + d_2 = 0.228528 \\ d_3 = 0.0343 & \tilde{d}_3 = \tilde{d}_2 + d_3 = 0.262828 \\ d_4 = 0.1657 & \tilde{d}_4 = \tilde{d}_3 + d_4 = 0.428528 \\ d_5 = 0.5714 & \tilde{d}_5 = \tilde{d}_4 + d_5 = 1.00 \end{array}$$

Step 2: Sample one value from uniform Random value b/w 0

$r = \text{numpy.random.uniform}(0.0, 1)$

let $r=0.6$

Step 3: If $r \leq \tilde{d}_1$

return 1

else if $r \leq \tilde{d}_2$

return 2

else if $r \leq \tilde{d}_3$

return 3

Proof:- prob of picking up d_i

= prob. of r lying b/w \tilde{d}_{i-1} & \tilde{d}_i

$$d'_i \propto d_i \quad \int d'_i = \frac{d_i}{S}$$