

11

Probability and Statistics Introduction.

Random variable \rightarrow two types

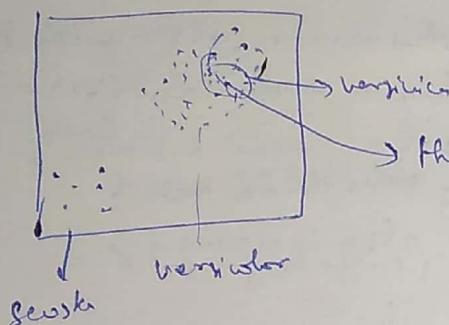
- discrete \rightarrow die is thrown: $\{1, 2, 3, 4, 5, 6\}$
- continuous \rightarrow height of student in cm.
 Suppose $(120 \text{ cm} - 190 \text{ cm})$
 Some say real
 numbers.

Outliers can also affect the probability.

height $\rightarrow \{120.1, 130.2, 146.5, 12.6, 184.8, 90.6\}$

brown
error.

short student

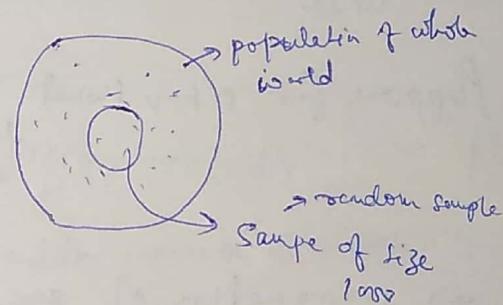


There is 80% of chance of being outlier
and 20% normal error.

Population and Sample.

Here we want to calculate population
avg height of people in the world, which is not
possible to calculate the mean with normal
process.

$$\bar{x} = \frac{1}{FB} \sum_{i=1}^{FB} h_i$$



So to overcome this we take random sample from the whole population
say 1000 random sample [which has 1 in 6 from India]
 $\bar{h} = \bar{x} = \frac{1}{1000} \sum_{i=1}^{1000} h_i$

This estimation is more easy and nearly
accurate

When size of sample increases

$$\bar{h} = \bar{x}$$

Estimation will be equal to population mean
of sample

Probability and Statistics introduction.

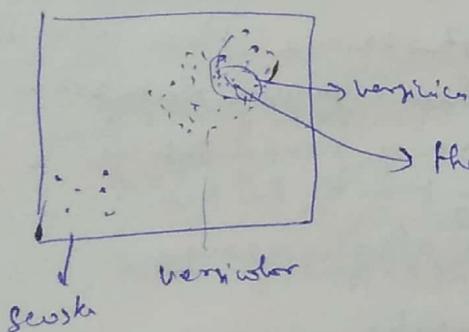
Random variable \rightarrow two types \rightarrow discrete \rightarrow die is Random: $\{1, 2, 3, 4, 5, 6\}$
 \rightarrow continuous \rightarrow height of student in cm.
 Suppose (120 cm - 190 cm)
 Some may real
 great waste.

Outliers can also affect the probability.

height $\rightarrow \{120.1, 130.2, 148.8, 12.6, 184.8, 90.6\}$

↳ human error.

↳ short student



there is 80% of chance of verifier
 and 20% verifier

Population and Sample.

Here we want to calculate population
 avg height of people in the world, which is not
 possible to calculate the area ~~at~~ with normal
 process.

$$\bar{x} = \frac{1}{n} * \sum_{i=1}^{FB} h_i$$

So to overcome this we take random sample from the whole population
 say 1000 random sample [which has 1 in 6 from India]

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

↳ This estimation is more easy and nearly
 accurate

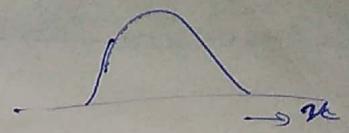
When size of sample increases

$$\bar{x} = \bar{h}$$

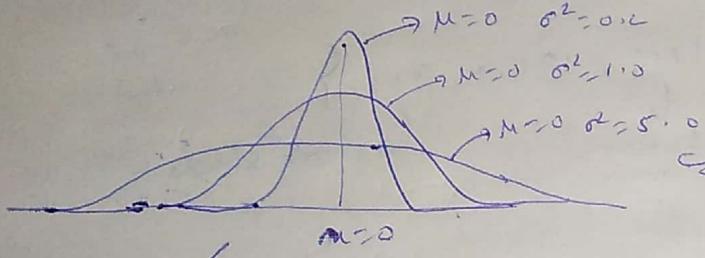
mean of estimation will be equal to population mean
 of sample

Gaussian and Normal Distribution Function.

If π is probability density function
then π is also gaussian distribution function.



We learn about distribution because this is a simple model and can be used, it can easily tell us how things are distributed randomly.



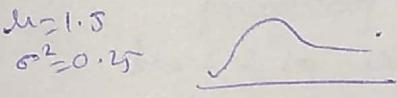
When we have idea of mean and variance, then with this we can draw gaussian distribution curve.

Given normal distribution is possible to find.

With increase in $\sigma^2 \rightarrow$ spread is also increases.

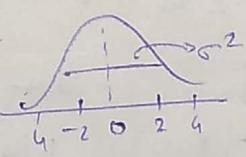
Generally peak is always at the mean value.

Suppose for set of flower. we have to make Gaussian distribution

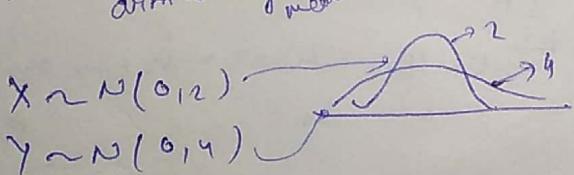


⇒ Parameters of gaussian distribution function (μ, σ^2)

$X \sim N(\mu, \sigma^2)$
 X follows normal distribution given ad mean variance



$X \sim N(0, 2)$



$X \sim N(0, 1)$

$Y \sim N(0, 4)$

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

Probabilistic random variable

Simplify this as -
(as $\mu=0$; $\sigma^2=1$;
 $\therefore \sigma=1$)

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

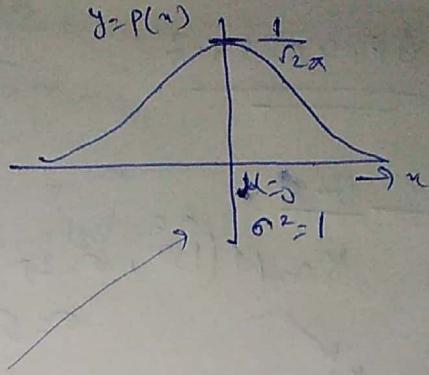
const const

Now, further simplifying it

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

$$\text{or } -x^2 + \exp(-x^2)$$

from this we can see



⇒ As x moves away from 0, $y \downarrow$

⇒ Symmetric

⇒ x moves away from 0, y reduces $\exp(-x^2)$

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

$$x=0 \quad y=1$$

$$x=1 \quad y = \exp(-1) = \frac{1}{e^1} = 0.3678 \quad \left\{ \text{from } 200n \right.$$

$$x=2 \quad y = \exp(-4) = \frac{1}{e^4} = 0.0183 \quad \left\{ \text{from } 100n \right.$$

$$x=3 \quad y = \exp(-9) = \frac{1}{e^9} = 0.000123 \quad \left\{ \text{from } 10n \right. \text{ drastic decrease in value.}$$

Q1 A ⇒ pdf is used to specify the probability of the random variable falling within particular range of value

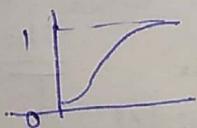
$$\text{if } \rightarrow P(1.1) = 0.59 \text{ from } \text{ex}(1-2)$$

$$\text{then } P(1.1) = 0.59$$

⇒ cdf will be continuous because there will not be sudden increase of probability at a single point.

CDF of Gaussian Distribution.

↳ Cdf always lie b/w 0 and 1

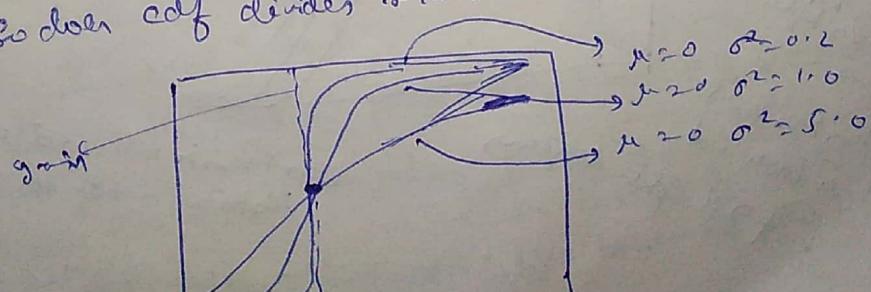


when variance is small → more close to y-axis
when variance is large → far away from y-axis.

↳ mean divides all data into two parts

↳ mean

so does cdf divides it into two equal percentage (probability)



$$6.8 - 9.5 - 9.9 \cdot 7$$

Standard Deviation



Normal distribution whose base is $\lambda = 3 \rightarrow$ non-negative (because σ)
lambda > 3 \rightarrow leptokurtic
lambda < 3 \rightarrow platykurtic

Support Curve

$$X \sim N(150, 5^2/25) \approx 5$$

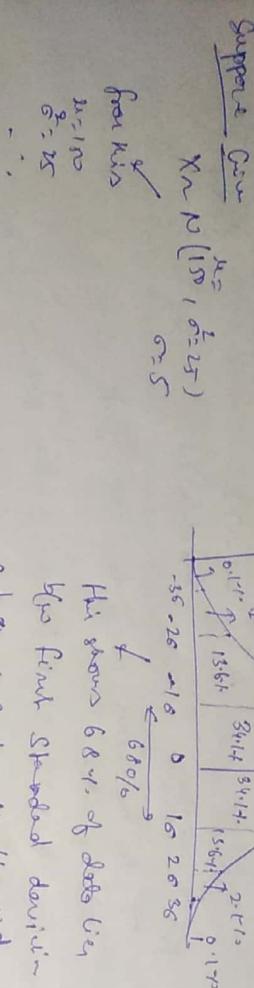
from his
 $\mu = 150$
 $\sigma = 25$

from 150-45 to 150+45

15 to 175

65% of people

live in world with height (150 cm to 175 cm)



Standardization and Standard Normal Variable (Z)

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

then X has random sample $\{x_1, x_2, x_3, \dots, x_n\}$
and variance s^2 .

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

then X has random sample $\{x_1, x_2, x_3, \dots, x_n\}$
then will be very useful while solving problems

$$\text{Now, with standardization: } z_i = \frac{x_i - \mu}{\sigma}, i = 1, 2, \dots, n,$$

$n \sim N(0, 1)$ \rightarrow the min value even value b/w 0 and 1 is

then will be very useful while solving problems

Symmetric distribution, Skewness and kurtosis

A histogram

\rightarrow deviation from symmetric normal distribution
of poly of x

Histogram is happening
on left side is slow
so, kurtosis on the
right side is slow

Excess kurtosis

When kurtosis is high
there is sharp peak
and when kurtosis is low
the peak is wide

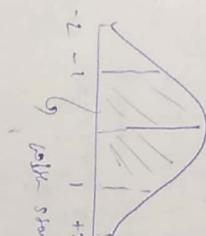
and kurtosis is 0.5
normal and kurtosis
is lower and kurtosis
is higher

when kurtosis is heavily tilted due to b/w
peak \rightarrow sharpener

Value is measured as tail length of
distribution with the freedom.

Kernel Density Estimation (KDE)

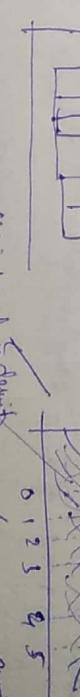
\rightarrow Using this we can know how a poly can be created using histogram



With standardization we would have $\mu = 0$ & $\sigma = 1$
and $z = \frac{x - \mu}{\sigma} = \frac{x - 0}{1} = x$



Density function



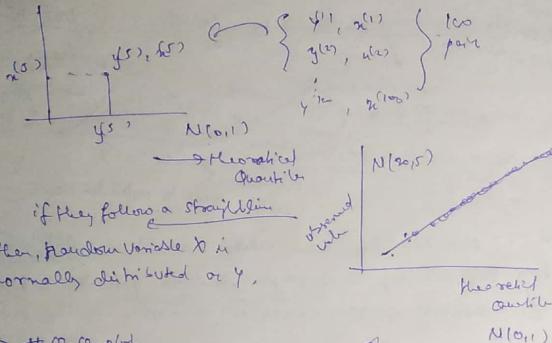
Density function

This is based on density and
function. Here, at one point (one x) our poly is created

Suppose there is a point x in between, density of point x is the sum of all
the kernels of poly below it.

At point P , mean + high kurtosis are present, so it has
higher peaks.

Step ③ plot Q-Q plot with $x^{(1)}, x^{(2)}, \dots, x^{(100)}$
 (percentile of $y^{(1)}, y^{(2)}, \dots, y^{(100)}$)



Code

```
# Q-Q-plot
import numpy as np
import pylab
import scipy.stats as stats
```

Adding standard distribution $N(0,1)$ mean variance 100 samples
 \Rightarrow std-normal = np.random.normal(loc=0, scale=1, size=100)

0 to 100 percentiles of std-normal

for i in range(0, 100):

print(i, np.percentile(std-normal, i))

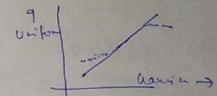
Q	P
0 - 3.29	
1 - -2.02	
2 - -0.82	
3 - -0.25	
4 - 0.07	
5 - 0.25	
6 - 0.44	
7 - 0.63	
8 - 0.82	
9 - 0.99	
10 - 1.18	
11 - 1.36	
12 - 1.54	
13 - 1.72	
14 - 1.90	
15 - 2.08	
16 - 2.26	
17 - 2.44	
18 - 2.62	
19 - 2.80	
20 - 2.98	
21 - 3.16	
22 - 3.34	
23 - 3.52	
24 - 3.70	
25 - 3.88	
26 - 4.06	
27 - 4.24	
28 - 4.42	
29 - 4.60	
30 - 4.78	
31 - 4.96	
32 - 5.14	
33 - 5.32	
34 - 5.50	
35 - 5.68	
36 - 5.86	
37 - 6.04	
38 - 6.22	
39 - 6.40	
40 - 6.58	
41 - 6.76	
42 - 6.94	
43 - 7.12	
44 - 7.30	
45 - 7.48	
46 - 7.66	
47 - 7.84	
48 - 8.02	
49 - 8.20	
50 - 8.38	
51 - 8.56	
52 - 8.74	
53 - 8.92	
54 - 9.10	
55 - 9.28	
56 - 9.46	
57 - 9.64	
58 - 9.82	
59 - 10.00	
60 - 10.18	
61 - 10.36	
62 - 10.54	
63 - 10.72	
64 - 10.90	
65 - 11.08	
66 - 11.26	
67 - 11.44	
68 - 11.62	
69 - 11.80	
70 - 11.98	
71 - 12.16	
72 - 12.34	
73 - 12.52	
74 - 12.70	
75 - 12.88	
76 - 13.06	
77 - 13.24	
78 - 13.42	
79 - 13.60	
80 - 13.78	
81 - 13.96	
82 - 14.14	
83 - 14.32	
84 - 14.50	
85 - 14.68	
86 - 14.86	
87 - 15.04	
88 - 15.22	
89 - 15.40	
90 - 15.58	
91 - 15.76	
92 - 15.94	
93 - 16.12	
94 - 16.30	
95 - 16.48	
96 - 16.66	
97 - 16.84	
98 - 17.02	
99 - 17.20	
100 - 17.38	

Random variable
 $X \hookrightarrow$ # Generate 100 samples from $N(20, 5)$
 measurements = np.random.uniform(loc=20, high=25, size=100)
 stats.probplot(measurements, dist="norm", plot=pylab)
 pylab.show()

cheatsheet normal
 if size is small
 Then we can see deviation
 but when we increase the size sample the the plot will be linear.

generate 100 samples for $N(20, 5)$

measurements = np.random.uniform(loc=20, high=25, size=100)
 stats.probplot(measurements, dist="norm", plot=pylab)
 pylab.show()



How/where to use distributions?

till now we have learned random variable, pdf, cdf, mean \rightarrow 68-95 rule

Distributions are generally use to do data analysis

Suppose

For an company are going to order t-shirts for look employees.

size of t-shirts can be \rightarrow XL, L, M, S

So, how we can't just go to every employee to know the t-shirt size.
 This would be time consuming.

So, to do this task in a simple way.

We will take the height of 500 people out of 6000 from our DB.

Mean, if height $> 180 \rightarrow$ XL there are general assumptions.
 heights 160 to 180 cm \rightarrow BM

So, Now when we have collected 500 heights

we will calculate mean and std.

$\hookrightarrow N(\mu, \sigma)$.

We will use Q-Q plot $N(0,1)$ to check if both are Gaussian distributed or not.

if yes,

-

suppose their cdf

says 99% of people are below 180cm.

\therefore P. of look = 1000 t-shirts are nearly equal.

This can be done for calculating salary range also.

Gaussian distribution are theoretical model

→ this happens in most of natural phenomena that's why we take it as reference.

Chebyshew's Inequality.

Till now, we studied Gaussian distribution $\rightarrow 68 - 95 - 99.7$ rule

Suppose $X \sim$ Height of Students

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

\downarrow Formed by mean and variance.

Then we can say.

↳ how much students lies b/w $\mu \pm 2\sigma$

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

↳ $[130, 170]$ width 40 lies b/w $[\mu \pm 2\sigma]$ distance $= 40$

But suppose, we don't know the distribution, all we know that we have finite mean and non-zero variance.

Can we calculate
how much

$\%$ of data lies within $\mu - 2\sigma$ & $\mu + 2\sigma$

$\%$ of data lies $\mu - 1\sigma$ & $\mu + 1\sigma$

Suppose we have to calculate $\%$ of individuals and we don't know the distribution, and we have

$$60k$$

$$40 - 10x2$$

$$40 + 2x10$$

Q1 What % of individual has salary in the range of $[20k, 60k]$?

Q2 [10k, 70k] ?

⇒ Chebyshew's inequality.

when we have $X \rightarrow \mathbb{R}$
don't know the distribution
finite mean $= \mu$
non-finite standard-dev = σ

probability that X lies b/w $(\mu - k\sigma)$ and $(\mu + k\sigma)$

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

$$\mu \pm k\sigma$$

$$X \leq \mu - k\sigma$$

$$X \geq \mu + k\sigma$$

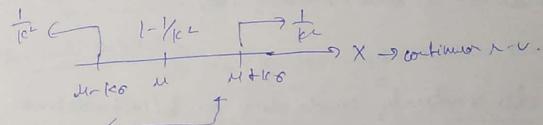
probability that X lies b/w $(\mu - k\sigma)$ and $(\mu + k\sigma)$ is

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

A1 $[20k, 60k]$

$$20k \text{ to } 60k \rightarrow P(20k \leq X \leq 60k) \geq 1 - \frac{1}{2^2} = \frac{3}{4} = 75\%$$

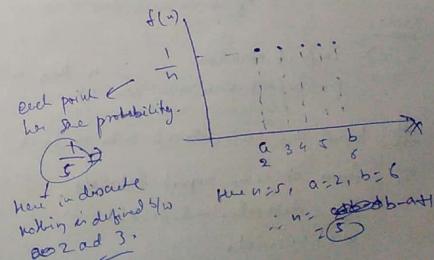
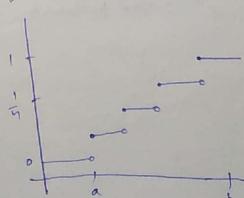
$$P(10k \leq X \leq 70k) \geq 1 - \frac{1}{3^2} = \frac{8}{9} = 90\%$$



Uniform Distribution

There are two types of uniform distribution

↳ probability
↳ function
↳ discrete → discrete random
variable
↳ uniform → uniform random variable



Each point has the probability $\frac{1}{n}$.
Here in discrete nothing is defined b/w a and b .
 $a=1, a=2, b=6$
 $n=5, n=6, b=4$

Notation used
parameters

$$\mu \in \{a, b\}$$

$$a \in \{-2, -1, 0, 1, 2\}$$

$$b \in \{-2, -1, 0, 1, 2\} \text{ s.t. } a < b$$

$$n = b - a + 1$$

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

$$CDF = \frac{[k] - a}{n}$$

mean = $\frac{a+b}{2}$
median
skewness = 0

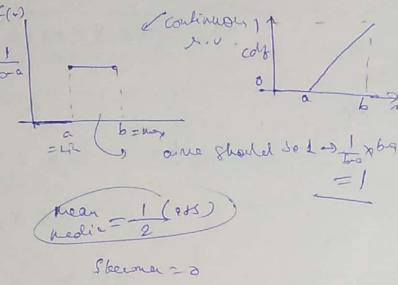
\hookrightarrow Uniform (continuous)

Notation = $U(a, b)$

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

cdf = $\begin{cases} 0 & \text{for } x < a \\ \frac{x-a}{b-a} & \text{for } a \leq x \leq b \\ 1 & \text{for } x > b \end{cases}$

mean = $\frac{a+b}{2}$
median = $\frac{a+b}{2}$
Skewness = 0



\approx The uniform distribution - all outcome is equally likely.

How do randomly sample data point (uniform distribution)

\Rightarrow Suppose we have 150 data points $\boxed{1 \ 1 \ 1 \ 1 \ 1 \dots 150}$
we have to select 30 data points out of 150

Set 1 \rightarrow $x_1, x_2, x_3, \dots, x_{150}$

Since we are doing this sampling in uniform probability
of each data should remain same, so fall in set 2.

Now,
 \Rightarrow import random
print(random.sample([1], 1)) → This will print value b/w [0, 1]
with equal probability.

Load the iris dataset with 150 point

\Rightarrow from sklearn import datasets
iris = datasets.load_iris()

d = iris.data
d.shape

\hookrightarrow $d/p \rightarrow (150, 4)$

Sample 30 points randomly from 150 point

$$\Rightarrow n = 150$$

$$m = 30 \quad \rightarrow \frac{30}{150} = 0.2$$

$$p = \frac{m}{n}$$

point(p)

sampled_data = [];

for i in range(0, n):

a = random.random() → This will print probabilities of occurring of

data where probability is $\leq p \rightarrow 0.2$.

if a $\leq p$:

sampled_data.append(d[i, :])

len(sampled_data) → 30

= 37 → it can vary around 30.

Bernoulli & Binomial distribution.

Bernoulli distribution

when we have only two outcome then bernoulli distribution is used.

Like tossing a coin (T, H)

~~if~~ or if occurrence of one outcome is P
then occurrence of another outcome is $(1-P) = q$.

Parameter $\rightarrow 0 \leq p \leq 1, p \neq 0$

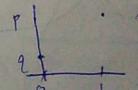
$$PMF = \begin{cases} 1-p & \text{for } k=0 \\ p & k=1 \end{cases}$$

$$CDF = \begin{cases} 0 & \rightarrow k=0 \\ 1-p & \rightarrow 0 \leq k < 1 \\ 1 & \rightarrow k \geq 1 \end{cases}$$

mean $\rightarrow P$

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

Variance $\rightarrow pq$,



Binomial distribution

↳ Number of times θ get a head when θ tosses a fair coin n times

$$Y \sim \text{Bin}(n, p)$$

Notation \rightarrow no. of heads $\sim \text{Bin}(n, p)$

$$\mathbb{P}[Y=k] = \binom{n}{k} p^k (1-p)^{n-k}$$

$$\text{Mean} = np$$

$$\text{Variance} = np(1-p)$$

Log Normal Distribution

Suppose we have to check if $Y \sim \text{log-normally distributed}$ or not.

So in a log normal distribution $Y \sim \mathcal{N}(\mu, \sigma^2)$

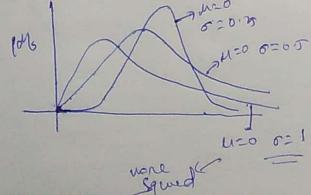
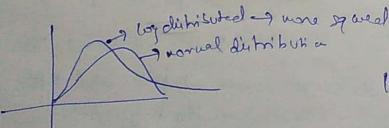
$$\text{false} \quad \log Y \sim \mathcal{N}(\mu, \sigma^2)$$

$$\text{Then } \ln(Y_1, Y_2, \dots, Y_n) \sim \mathcal{N}(0, 1)$$

will be normally distributed.

To check this - plot Q-Q plot b/w $\ln(Y_i)$ and $N(0, 1)$ \rightarrow if this follows linear plot then it is log normally distributed.

In most of the internet companies they work on log-normally distributed values, because very few times they get normally distributed values.

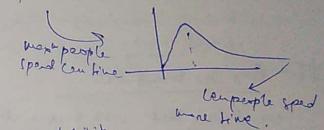


Notation $\rightarrow \text{lognormal}(\mu, \sigma^2)$

$$\text{pdf} \rightarrow \frac{1}{x\sigma\sqrt{2\pi}} e^{-\frac{(\ln x - \mu)^2}{2\sigma^2}}$$

↳ Comments on internet discussion forum follows a log normal distribution.

↳ Time spent by people in internet follows log normal distribution.



Poisson Loss distribution

It is a relationship b/w two question.

Pareto distribution

$x > 0 \rightarrow \text{scale}$

$\alpha > 0 \rightarrow \text{shape}$

y=probability

long tail
80% of data
in the region.
80-20 rule

$\alpha \downarrow \rightarrow \text{tail fattening}$

when $\alpha \rightarrow \infty$ then \rightarrow particle value probability is high.
long tail will reach to 0
at very high value.

direct delta function.

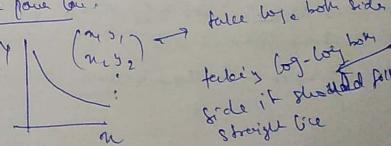
Loving.
pareto distribution is the distribution where, most population is ∞ value and small population having a very large value.

Application \rightarrow Hard disk drive error.

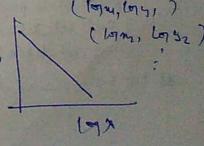
b) sizes of human settlements (few cities, many villages)

Log-log plot

To check power law,



false w.r.t both sides
fitter's log-log both sides if straight line



Or we can plot of (x, y)
 x having y as points distributed

↳ distribution of wealth in society

(C) pareto distribution.

↳ large portion of wealth is own by small portion of people.
 (80-20 rule)

Box-Cox transformation

We know that with gaussian distribution we can calculate a lot of things.
 So, the conversion of pareto distribution is necessary into gaussian distribution.
 because in machine learning gaussian distribution is most often used.

However, $\text{pareto} \sim X : [x_1, x_2, \dots, x_n]$ is not normal
 $\text{Gaussian} = Y : y_1, y_2, \dots, y_n$

So, to convert this box-cox is used

$$\begin{aligned} \textcircled{1} & \quad \text{box-cox}(x) = \text{lambda}(\lambda) \rightarrow \text{logarithm} \\ \textcircled{2} & \quad y_i = \begin{cases} \frac{x_i - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \log(x_i) & \text{if } \lambda = 0 \end{cases} \quad \text{if } \lambda = 0 \\ & \quad m_i = \log(x_i) \quad \text{else} \\ & \quad \text{The value we get after this conversion will be normally distributed.} \\ & \quad \forall i: 1-n \end{aligned}$$

[Code]

⇒ from scipy import stats
 import matplotlib.pyplot as plt

It we will generate some random variables from a non-normal distribution and make a probability plot for it, to show it is not normally distributed.

⇒ fig = plt.figure()
 $ax1 = fig.add_subplot(211)$
 $x = \text{stats.loggamma.rvs}(5, size=500) + 5$
 $prob = \text{stats.probplot}(x, dist=\text{stats.norm}, plot=ax1)$
 $ax1.set_kewel('') \rightarrow \text{qq-plot.}$
 $ax1.set_title("prob against normal distribution")$

too how change data using subplot so its close to normal.

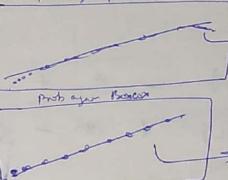
⇒ $ax2 = fig.add_subplot(212)$

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

prob = stats.probplot(xt, dist=stats.norm, plot=ax2)

ax2.set_title("prob. after Box-Cox transformation")

prob against normal



new data is denoted in front of the plot
 so not normally distributed

Correlation

⇒ Box-Cox won't work on every data, it's working or not can be checked using Q-Q plot.

⇒ Box-Cox transform power-law transforms into normal distribution.

Covariance

Suppose we have height
 & weight

	160	162
s1	150	154
s2	.	.
su	160	162

So here, we need to find
 what relation does x and y hold.

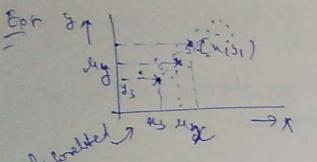
$\begin{array}{l} \text{Cov} \rightarrow \text{covariance} \\ \text{Pearson} \rightarrow \text{pearson correlation coefficient} \\ \text{Spearman} \rightarrow \text{spearman rank correlation coefficient} \end{array}$

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

$$\text{Var}(x) = \text{Cov}(x, x) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})$$

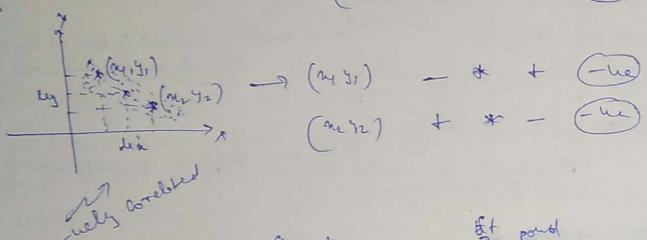
$$\text{if } \text{Cov}(x, y) = \text{true } x \uparrow, y \uparrow$$

$$\text{Cov}(x, y) = -\text{true } x \uparrow, y \downarrow$$



$$\text{Cov}(u_1, u_2) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_1)(y_i - \bar{y}_1)$$

$$\begin{aligned} (\mu_1 y_1) &= + * + (\text{true}) \\ (\mu_2 y_2) &= - * - (\text{true}) \end{aligned}$$



again suppose $\text{Cov}(\hat{u}_1, \hat{y}_2) \neq \text{Cov}(\hat{u}_1, y_1)$

if height in cm is height in ft and weight in pound not correlated with weight in kg then weight in ft and weight in pound will not be correlated.

Pearson Correlation Coefficient (PCC)

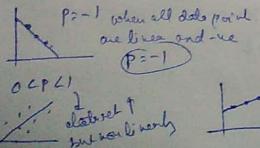
$$P_{xy} = \frac{\text{Cov}(u_1, u_2)}{\sigma_x \sigma_y}$$

In Covariance:

if $x \rightarrow y \uparrow$, $\text{Cov}(u_1, u_2) \rightarrow$ the but how much the if $x \rightarrow y \downarrow$, $\text{Cov}(u_1, u_2) \rightarrow$ we and we it will be it is never said.

So, to overcome this

$-1 \leq P \leq 1 \rightarrow$ PCC is always less than 1 and greater than -1 Pearson Correlation Coeff. is used.



$P=1$ when all data points are linear and we

need to check if it is linear or not

$-1 < P < 0 \rightarrow$ data are decreasing but not linear

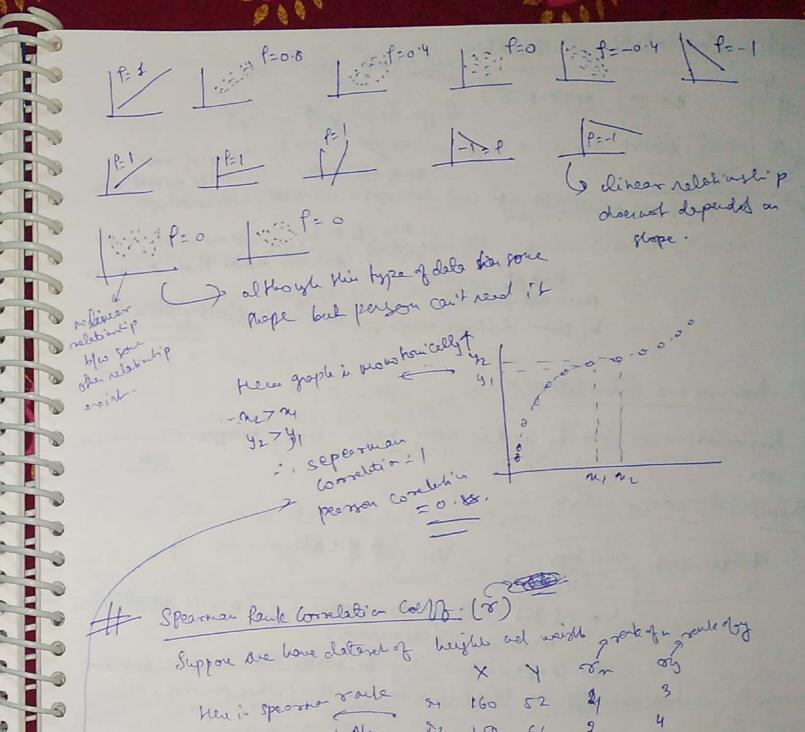
$P=0$ Pearson cannot

read any relation other than linear.

$P < 0$ Pearson correlation = 0.84

Pearson correlation = 0.67

Scanned by CamScanner



Spearman Rank Correlation Coeff. (S)

Suppose we have dataset of height and weight rank of a sample by

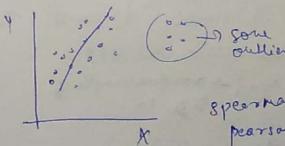
	X	Y	rank of x	rank of y
s1	160	52	4	3
s2	150	66	2	4
s3	170	68	5	5
s4	140	46	1	1
s5	188	57	3	2

Now in Spearman's rank we arrange the dataset A/T its ranks.



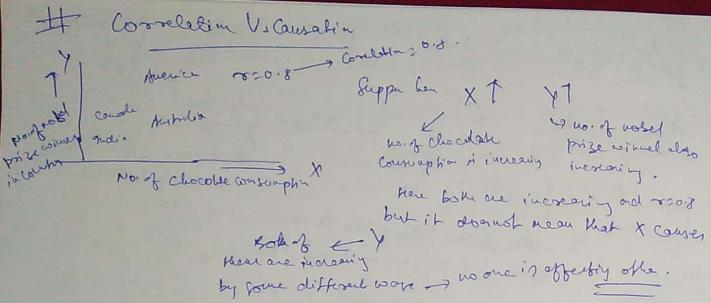
So as long as there is monotonicity in dataset. if it is linear or not

rank coefficient $r_s = \frac{1}{n(n-1)} \sum r_i^2 - \frac{n+1}{2}$



Spearman Correlation = 0.84

Pearson Correlation = 0.67



How to use Correlations?

Correlations are heavily used in many areas, like, ~~economics~~, etc.

Suppose in amazon website

if unique visitors in a day vs. \$ salary in a day
 ↳ if this increases then salary in a day will also increase. (correlated)
 So, company will try to ~~raise their website visibility~~ to get more no. of visitors per day.

Suppose in education

10x → Salary
 12x →
 eggs →
 Milk →
 phd →
 if we increase the study
 Salary might correlated increase.
 So education minister should focus on this.

Confidence Interval (CI)

Suppose ~~any distribution~~ X: height

↳ let random sample $\{x_1, x_2, \dots, x_{10}\}$ → random sample of size 10.

Estimated population mean of $X = \bar{x}$

for small sample set:
 $m \approx \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ → simple avg.
 This will be ~~more~~ equal to $\mu = \bar{x}$ when sample size increases.
 So here this $\bar{x} = \bar{x}$ is giving some exact (point) mean.
 But for some different data samples mean could change.
 So, we use Confidence Interval for this.

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

$\{150, 162, 151, 172, 168, 150, 171, 183, 165, 176\}$ → height of people in cm

$$\bar{x} = \bar{m} = \frac{1}{10} \sum_{i=1}^{10} m_i = 168.5 \text{ cm}$$

But CI $\Rightarrow \bar{x} \in [162, 174, 171]$ with 95% probability.

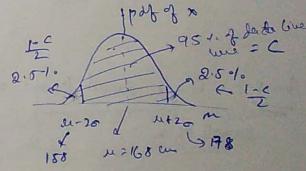
↳ gives a range of value with confidence of 95% probability in which data lies.

Computing Confidence Interval given the underlying distribution

Suppose we have heights that is gaussian distributed

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

$$\text{left } (\bar{x} = 168 \text{ cm})$$



⇒ Value of \bar{x} , range of height in the range of 95% probability.

→ Since it is Gaussian distributed.

↳ 95% data will lie $\bar{x} \pm (\mu - 2\sigma, \mu + 2\sigma)$
 $[158, 178]$

∴ $C = 95\%$. So, Value of 90%, 80%, 70% probability can be known by chart of normal distribution.

CI for mean of a normal random variable

Suppose $X \sim N$ with population mean of μ and std-deviation of σ .
 follows
 $\{x_1, x_2, x_3, \dots, x_n\}$ -> having sample of size $n=10$
 $\{180, 162, 158, 172, 168, 150, 170, 183, 165, 176\}$

What is the 95% CI of μ .

Case 1 Suppose when $\sigma = 5$ is known {we know pop. std-dev}.

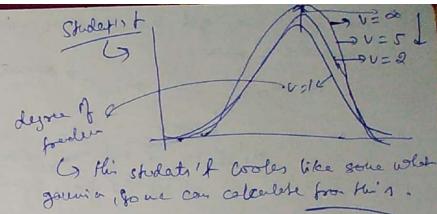
Here, in this case, we will use CLT, since CLT is used in a lot of cases, but since we have one sample CLT is used.

$$\text{CLT: } \bar{x} = \text{sampled mean} = \frac{1}{10} \sum_{i=1}^{10} x_i$$

$\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$
 Suppose follows normal dist. with population std-deviation $(\frac{\sigma}{\sqrt{n}})$
 $\therefore \bar{x} \in [\bar{x} - 2\frac{\sigma}{\sqrt{n}}, \bar{x} + 2\frac{\sigma}{\sqrt{n}}]$ with 95% Confidence
 near best $\bar{x} = 168.5$ $\sigma = 5$ $n = 10$ $\frac{\sigma}{\sqrt{n}} = 1.6$

Case 2: if we don't know σ (pop std-dev).

Sample of n in this case t -distribution is used.
 $\bar{x} \sim t(n-1)$
 Sample mean follows t -dist with degree of freedom $(n-1)$
 σ is unknown



Confidence Interval using Bootstrapping

Using this we can calculate CI for median, mean, std etc.

Suppose we only have sample size and don't know the distribution or anything.

Let sample of size n : $S = \{x_1, x_2, \dots, x_n\}$ for $n=10$

Once we have the sample using discrete uniform r.v.

We create different sample size of length $m \leq n$

$$\therefore u(1, n) = \frac{1}{10} \sqrt{\frac{1}{2} \left(\frac{n(n+1)}{12} \right)}$$

$$s_1 = (x_1^{(1)}, x_2^{(1)}, x_3^{(1)}, \dots, x_m^{(1)}) \quad m \leq n$$

This is bootstrap sample. Sampling with replacement from S .

$$s_2 = \{x_1^{(2)}, x_2^{(2)}, \dots, x_m^{(2)}\}$$

$$s_3 = \{x_1^{(3)}, x_2^{(3)}, \dots, x_m^{(3)}\}$$

$$\vdots$$

$$s_{10} = \{x_1^{(10)}, x_2^{(10)}, x_3^{(10)}, \dots, x_m^{(10)}\}$$

Here we generate k different samples from sample S with repetition. Then calculate median for all samples k for CI of median.

Suppose median of sample $s_1 \rightarrow m_1$,

$$s_2 \rightarrow m_2$$

$$s_3 \rightarrow m_3$$

$$\vdots$$

$$s_{10} \rightarrow m_{10}$$

Now once we got the median, we arrange the median in ascending order.

$$\begin{aligned} & \text{Suppose } k=100 \\ & m_1, m_2, \dots, m_{100} \rightarrow \text{Median} \\ & m'_1 \leq m'_2 \leq m'_3 \dots \leq m'_{100} \\ & \text{1st CI } (m'_1, m'_{10}) \\ & 25 \quad m'_1 \quad m'_{95} \quad 75 \end{aligned}$$

To get 95% CI. out of 100 sample medians we take median from m'_{12} to m'_{95} .

Similarly we can calculate for various others.

Code \Rightarrow import numpy

```
from pandas import read_csv
from sklearn.utils import resample
from sklearn.metrics import accuracy_score
from matplotlib import pyplot
```

load dataset

$\Rightarrow x = \text{numpy.array}([130, 162, 158, 172, 168, 150, 171, 183, 165, 176])$

Configure bootstrap

$\Rightarrow n_iterations = 1000$

$n_size = \text{int}(\text{len}(x))$

Run bootstrap

$\text{medians} = \text{list}()$

for i in range(n_iterations):

prepare train & test sets

s = resample(x, n_samples=n_size);

m = numpy.median(s);

print(m)

medians.append(m)

Plot scores

pyplot.hist(medians)

pyplot.show()

Confidence Intervals

$$\alpha = 0.95$$

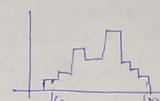
$$p = ((1.0 - \alpha)/2.0) * 100 \rightarrow 25$$

$$\text{lower} = \text{numpy.percentile}(\text{medians}, p)$$

$$p = (\alpha + ((1.0 - \alpha)/2.0)) * 100$$

$$\text{upper} = \text{numpy.percentile}(\text{medians}, p)$$

print('95.0 confidence interval', lower, upper)



95.0 confidence interval 161.5 and 176.0

Hypothesis testing methodology, Null-hypothesis, p-value

Suppose we have 500 students ~~each~~ in two classes each.

So, out of 500 students in each class, we take 50-50 samples from two class. And noted

down the heights of each student.

	cl1	cl2
1	150	162
2	152	156
:	:	:
50	160	162

So, the question here is

Is there any difference

b/w the height of students

b/w two classes.

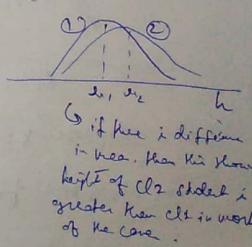
So, to know this we do a better work.

Suppose the mean of cl1 is μ_1
cl1 — μ_1

So to do hypothesis testing we do the following steps.

① Choosing the test-statistic

(Now we have chosen mean. $(\mu_2 - \mu_1)$)



If μ_2 is different in mean than μ_1 , then the mean height of cl2 student is greater than μ_1 in most of the case.

② Null Hypothesis (H_0)

(Here we will follow proof by contradiction)

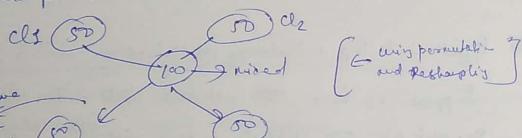
$H_0 \rightarrow$ no difference in μ_1 and μ_2 .

Alternative Hypothesis (H_1) \rightarrow difference in μ_1 and μ_2 .

3) P-value: Check if prob. of observation ($\mu_2 - \mu_1$) \leq null hypothesis is true.

Now, from the observation we know that there is a difference b/w mean of two class is greater than or equal to 10 cm, from all the 1000 students of two class.

Now to check for 100 students \Rightarrow from one class set 50 for one.



After mixing we again resample. If 50-50 students \Rightarrow then calculate the mean.

Suppose $\mu_1 - \mu_2 = S_1 \Rightarrow$ from first resample.
We repeat against resample again suppose $\mu_1 - \mu_2 = S_2 \Rightarrow$ 2nd resample
we continue. $\mu_1 - \mu_2 = S_{1000} \Rightarrow$ 1000th resample

We get $S_1, S_2, S_3, \dots, S_{1000}$
↓
so this are order.

$S_1^{(1)}, S_2^{(1)}, S_3^{(1)}, \dots, S_{1000}^{(1)}$

Suppose $X = \mu_1 - \mu_2 \geq 10$ cm lies at 850th place

So, ~~prob.~~ prob. that $(\mu_1 - \mu_2 = x)$ greater than 850th $= \frac{850}{1000} \times 100$
 $= 20 \%$

So, null hypothesis is right.

Let just say. if 10cm lies at 200 place.

$$\text{then prob that } X \geq 10 \text{cm} = \frac{200}{1000} \times 100$$

2.1. 154-

\therefore Null hypothesis is wrong \therefore check for alternative hypothesis.

So, if p-value ($X \geq 10 \text{cm} | H_0$) ≥ 0.9
 $\therefore H_0 \rightarrow$ accepted

if p-value ($X \geq 10 \text{cm} | H_0$) ≤ 0.01
↓
How & from observation false $\therefore H_0 \rightarrow$ rejected.
↓
true

Hypothesis testing: Coin intuition.

Ex: Here, we have given a coin and we have to check if the coin is biased towards heads or not.

If it is biased towards heads then $P_{head} \geq 0.5$
If it is not biased towards heads then $P_{head} \leq 0.5$

For this we flip the coin 5 times and count heads $\Rightarrow X \rightarrow$ test-statistic
after performing the expt suppose we get 5 heads. \rightarrow
flip, flip, flip, flip, flip
n n n n n $\therefore X=5$
(This is observational expt.)

$$P(X=5 | \text{coin is not biased}) = P(\text{obs} | H_0)$$

Here, we will check, but the hypothesis is coin is not biased towards heads.

not biased $\Rightarrow H_0 \rightarrow$ we flip 5 time
 $\therefore P(\text{head}) = \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{32} = 0.03 = 3\%$

But this shows that when a coin is flipped 5 times, there are 3% probability that all five head occur.
 $\text{here } 3 \approx 2.5\%$ So, our hypothesis is false, and as by observation we got all five head, coin is biased towards head.

Action we need to choose hypothesis and testing very carefully, because this changes the whole scenario.

Suppose we take the no. of times coin was tossed to 3.
 Now, observation = $\frac{1}{3}x_1 + \frac{2}{3}x_2 = \frac{1}{3} \times 600 = 200 \geq 5\%$
 Tossed with it is not biased = $\frac{1}{2}x_1 + \frac{1}{2}x_2 = \frac{1}{2} \times 600 = 300 \geq 5\%$
 Which is true.
 Now, H_0 is true. \therefore Coin is not biased.

(Q1) This observation is done only once, because, don't it n times would cost you a lot of money.

P-value is defined for single experiment.

KS test for similarity of two distribution.

Here, D_n stands for Kolmogorov-Smirnov Test.

This test is used where we have to check if two distributions are similar or not.

Suppose, we have, two sets of random variables here,

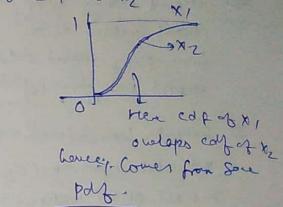
$X_1 = [x_{11}, x_{12}, \dots, x_{1n}]$ -> variable.

$X_2 = [x_{21}, x_{22}, \dots, x_{2n}]$ -> variable

Please we have to check if X_1 and X_2 come from same distribution or not.

So, for this to happen, for two distribution to be same.

CDF of X_1 should lie on each other to CDF of X_2 .
 if n_1, n_2 are large, then distance between two CDF must be 0.



If we have small datasets, or large then there should be a minimum Distance(D). After which our hypothesis will fail.

Now, we have two R.V.s x_1 and x_2

$H_0 \rightarrow$ Null Hypothesis is \rightarrow the ~~two~~ R.V. x_1 and x_2 come from same distribution.

Now, $D_{n_1, n_2} = \sup |F_{1, n_1}(x) - F_{2, n_2}(x)|$

place $n_1 = 1000$
 $n_2 = 500$
 $x = 0.05$

$D_{1000, 500} > 0.049$ for rejection of H_0 .

for rejection of H_0 , $D_{n_1, n_2} > C(\alpha)$

supremum function.

empirical distribution function.

now go to

$$D_{n_1, n_2} > C(\alpha) = \sqrt{\frac{-1}{2} \ln \alpha}$$

Where

$$C(\alpha) = \sqrt{-1 \ln \alpha}$$

To reject H_0 $D_{n_1, n_2} > C(\alpha) = \sqrt{\frac{-1}{2} \ln \alpha}$

Let $D_{n_1, n_2} = D$

$$D = \sqrt{\frac{n_1}{n_1 + n_2}} > \sqrt{\frac{-1}{2} \ln \alpha}$$

$$D^2 \left(\frac{n_1}{n_1 + n_2} \right) > \frac{1}{2} \ln \alpha$$

$$D^2 \left(\frac{n_1}{n_1 + n_2} \right) < \exp \left[-2 \sqrt{\frac{n_1}{n_1 + n_2}} \right] \times \alpha$$

p-value.

α	0.10	0.05	0.01
$C(\alpha)$	1.073	1.224	1.328 / 1.52

Code snippet of k-s test

```

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

# Generate Gaussian(Normal) r.v. of x.
x = stats.norm.rvs(size=1000); This will generate r.v. of size 1000.
sns.set_style('whitegrid'); Scaling factor for std. of data.
sns.kdeplot(np.array(x), bw=0.5); It takes data as 1D array.
plt.show()

```

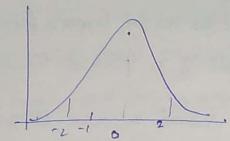
\Rightarrow # Applying k-s test table 1-Domino or Stray

```

stats.kstest(x, 'norm')

```

Out:
 (statistic = 0.0230, pvalue = 0.7554)
 $D_{\text{stat}} \approx 0.75 \rightarrow \therefore \text{Hypothesis accepted.}$



Y - Continuous Uniform Distribution ($0, 1$)

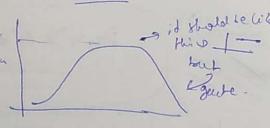
```

y = np.random.uniform(0.1, 1.0000); generates 1000 uniform random
sns.kdeplot(np.array(y), bw=0.1); Variable.
plt.show()

```

\Rightarrow stats.kstest(y, 'norm') it can be any distribution as we want.

Out:
 (statistic = 0.50159, pvalue = 0.0)
 $D_{\text{stat}} \approx 0.5 \rightarrow \text{Hypothesis fails here. Most}$
 $y \text{ is triangularly distributed.}$



Hypothesis Testing: Another Example

Part 1: Determine if population means of heights of people in these two cities are equal or not.

Suppose City 1 has population of 1 million (C_1)
 City 2 has population of 2 million (C_2)

Since average height to 1 million people is a tough task, so, we take 50 samples from two cities.

Randomly taken \rightarrow

c_1	c_2
h_1'	h_1'
h_2'	h_2'
h_3'	h_3'
\vdots	\vdots
h_{50}'	h_{50}'

Sample height of 50 people \rightarrow

$$\bar{h}_1 = \frac{h_1' + h_2' + \dots + h_{50}'}{50}$$

sample height of n people \rightarrow

$$\bar{h}_2 = \frac{h_1' + h_2' + \dots + h_{50}'}{50}$$

$\therefore \bar{h}_2 = 162 \text{ cm}$ $\bar{h}_1 = 167 \text{ cm}$

\therefore Then these are the observed values.

Now,

$$\text{Test statistic} = \bar{h}_2 - \bar{h}_1 = n = 162 + 167 = 5 \text{ cm}$$

\therefore difference in ~~population~~ mean for cities

Null-Hypothesis (H_0): There is no difference in population mean.

Compute for. $P(n=5 \text{ cm} | H_0)$

probability of observing a diff. of 5 cm in sample mean heights of sample size 50 b/w C_1 and C_2 , if there is no population difference in mean heights.

Case 1 $P(n=5/H_0) = 0.2 = 20\%$ (Suppose)

There is 20% chance of observing a difference of 5m in sample mean heights of G_1 & G_2 with sample size of 50, if there is no population mean difference.

Here 20% is significant.

\therefore assumption must be true.

\therefore accept H_0 .

Case 2 $P(n=5/H_0) = 0.03 = 3\%$

$P(\text{obs}/\text{assumption}) = 3 \rightarrow$ small

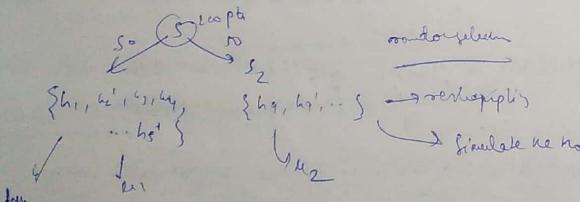
\Rightarrow assumption must be incorrect

\Rightarrow reject $H_0 \rightarrow$ accept H_1

Resampling and permutation tests

Step 1 $\rightarrow S = \{h_1, h_2, h_3, \dots, h_{50}, h'_1, h'_2, \dots, h'_{50}\}$

4 Use \rightarrow concatenating the two sets.



- after first resampling:
- (1) $h_1 - h_2 \rightarrow d_1 \leftarrow \text{diff}_1$
 - repeat (2) $h_2 - h_3 \rightarrow d_2 \rightarrow d_3$
 - + (3) $h_3 - h_4 \rightarrow d_4 \rightarrow d_5$
 - :
 - (n) $h_n - h_1 \rightarrow d_n \rightarrow d_{n-1}$

(last n-1 = 49)

Step 3 sort d_i

$$d_1' \leq d_2' \leq d_3' \dots \leq d_{50}'$$

Simulated diff.

Case 2 give observed diff = 5m.

Suppose H_0 is true

$$d_1' \leq d_2' \leq d_3' \dots \leq d_{50}' \leq d_m' \dots \leq d_{50}'$$

80% of sim diff. $\leq 5m$

20% of simulated diff.

$\geq 5m$ \rightarrow observed diff.

: P-value $P(\text{obs-diff}/\text{assumption}) = 20\% > 5\%$

\therefore assumption = true
accept H_0

if $P(\text{obs-diff}/\text{assumption}) = 3\% < 5\%$

assumption false
reject H_0

Hypothesis testing (hard to use)

6 we have seen t-test that tells us about if the variable follows same distribution or not

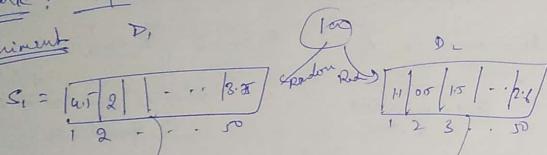
\Rightarrow Suppose we have two drugs here D_1 and D_2 in which D_1 claims that it can reduce fever in an about 6 hours and D_2 claims that it can reduce fever before D_1 .

(in the market) $\rightarrow D_1$

\rightarrow new drug.

So, here b'out of Indra has to decide whether D_2 is claiming right or not. Should it be giving ~~the same~~ license to D_2 company to sell it or not. So, it can be tested using hypothesis testing.

Here,
Task : exp to determine if claim is true or not.
Experiment: D_1



So, they decided to experiment this n=50, no sampler of people, and the finding is also noted.
 Suppose mean time ≈ 4 hrs.
 mean time between $d_{12} = 2$ hrs.
 & this is the observed value.

Since we can't rule the life of any human.
 \therefore we will use d or p-value $\leq 1\%$.

Hypothesis testing : (1) H_0 : D_1 and D_2 take the same time to reduce fever.
 (2) Test statistic is $K = d_{12} - d_{11} = 4 - 2 = 2$ (observation).

Notes
 $P(X \geq 2 | H_0) = \text{Very small}$ (Suppose 1/10)

P. value
 If there is no difference in D_1 and D_2 , the prob. of obs. that $X \geq 2$ is very small ($1/10$)

\therefore No ad obs do not agree with each other.
 $\therefore H_0$ is rejected (incorrect)

for e-commerce company like amazon, we can use $L \leq 5\%$.
 because for a little loss of money can be tolerate.

Proportional Sampling (probability sampler)

Suppose we have taken 5 random elements

$$d = \begin{bmatrix} 1.0 \\ 0.6 \\ 1.2 \\ 0.8 \\ 2.0 \end{bmatrix} \quad i=1, 2, 3, 4, 5 = n$$

Goal is to pick an element amongst the n elements such that probability of picking an element is proportional to $d_i^{1/2}$.
 sum of all the elements = $\sum d_i = 3.5$.

Step 1: (a) Compute the sum $\Rightarrow S = \sum_{i=1}^n d_i = 3.5$.

(b) Normalizing using the sum.

$$d'_i = d_i / S$$

$$\begin{aligned} d'_1 &= \frac{1.0}{3.5} = 0.2857 \\ d'_2 &= \frac{0.6}{3.5} = 0.171428 \\ d'_3 &= \frac{1.2}{3.5} = 0.342857 \\ d'_4 &= \frac{0.8}{3.5} = 0.228571 \\ d'_5 &= \frac{2.0}{3.5} = 0.571429 \end{aligned}$$

0 to 1
sum to L

(c) Cumulative Normalized sum.

$$\tilde{d}_i = \text{Cumulative Normalized Value}$$

Step 2: Sample one random uniform value (x) $\Rightarrow \text{unif}(0, 1)$
 $x = \text{numpy.random.uniform}(0, 1)^{\text{new}}$

Let $\sigma = 0.6$

Step 3

prop sampling

if $\sigma \leq d_1$
reborn 1
else if $\sigma \leq d_2$
reborn 2
else if $\sigma \leq d_3$
reborn 3

$$\mu = \frac{\sigma - 0.6}{\sigma} \cdot d_5$$

\therefore pick d_5

But, prob. of picking d_5 ,

$$= \text{prob of } \sigma \text{ being } \leq d_5$$
$$= \frac{d_5 - 0.6}{d_5} \cdot d_5^{-1}$$
$$\propto d_5^{-1}$$
$$\therefore d_5' = \frac{d_5}{5}$$

that \Rightarrow 3-Candidate selection example.