

Probability & Statistics

① Random variable

i) Definition:

Let us understand with an example.

Eg:

Dice: Six sides $\rightarrow \{1, 2, 3, 4, 5, 6\}$



rolling dice can have random outcomes

Let's denote random variable by.

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



random
variable.

random outcomes

Eg: Tossing up a coin.

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



random
variable.

Dice roll :-

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

Where

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

↑

All are equally likely to occur.

$$\begin{aligned}\cdot P(X \text{ is even}) &= P(X=2) + P(X=4) + P(X=6) \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} \\ &= \frac{1}{2}\end{aligned}$$

General

$$P(X=x_1)$$

read as 'probability that random variable X is x_1 '.

Concise way: $P(x_1)$. [random variable implicitly]

1.2 Types of random variable

(i) Discrete random variable

can take a value from set of finite values.

Eg: dice roll

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

(ii) Continuous random variable:

can take any real value. Eg.

Height: 120cm - 190cm

$$Y = 162.45$$

$$Y = 132.65$$

Continuous random variable.

(2) Outlier

Eg:

y: weight of a student

{122.2, 146.4, 132.5, ... } 12.26 156.23, ... }

Outlier

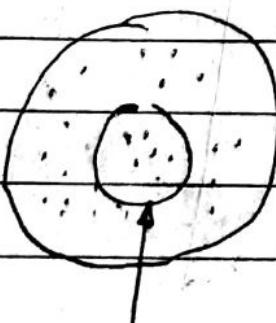
Outlier may be:

i) Genuine outlier

ii) Or caused by
human error

Observation error

(3) Population & Sample



& population [set of all population in the world]

sample: [subset of the population]

Eg: sample of 10 students

Mean of average height of population

FB

$$\mu = \frac{1}{FB} \sum_{i=1}^n h_i$$

[Since we cannot take population as whole,

Mean of average height of sample (1000)

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

[called random-sample]

sample is taken for analysis]

Typically, \bar{x} used to represent mean of the sample.

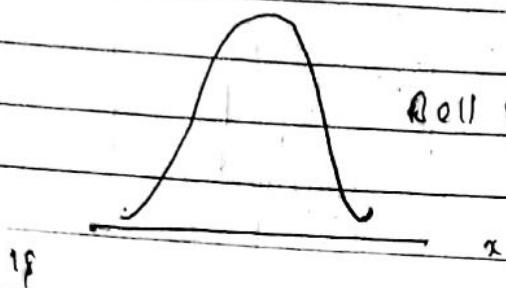
As sample size increases

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

Sample mean Population mean.

[sample is taken to estimate statistics [some property (like height) of population]]

(5) Gaussian Distribution / Normal Distribution



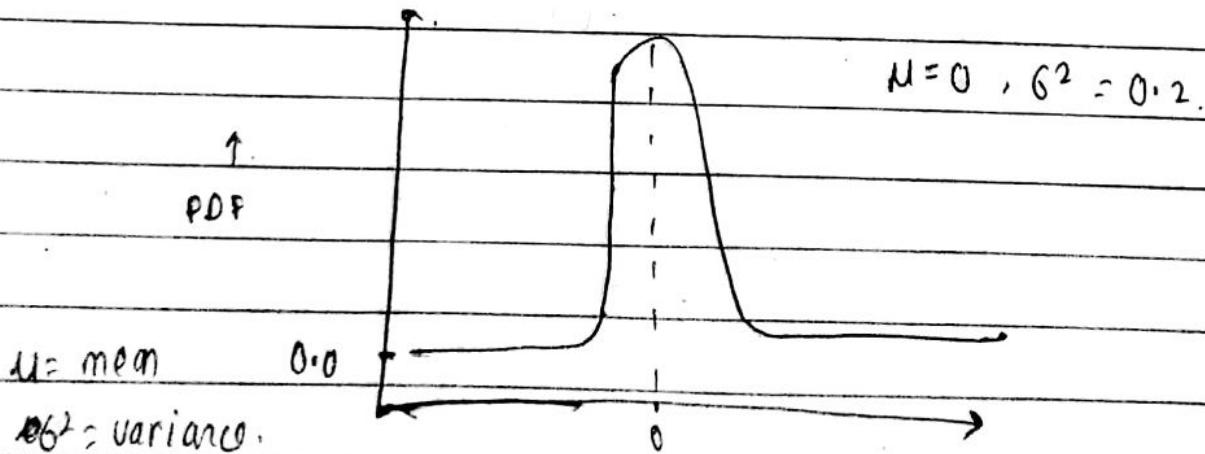
Bell shaped curve.

Probability density function (PDF) of a gaussian distributed random variable.

Why learn distributions?

↳ Helps in formulation of simple models giving details on various features and attributes

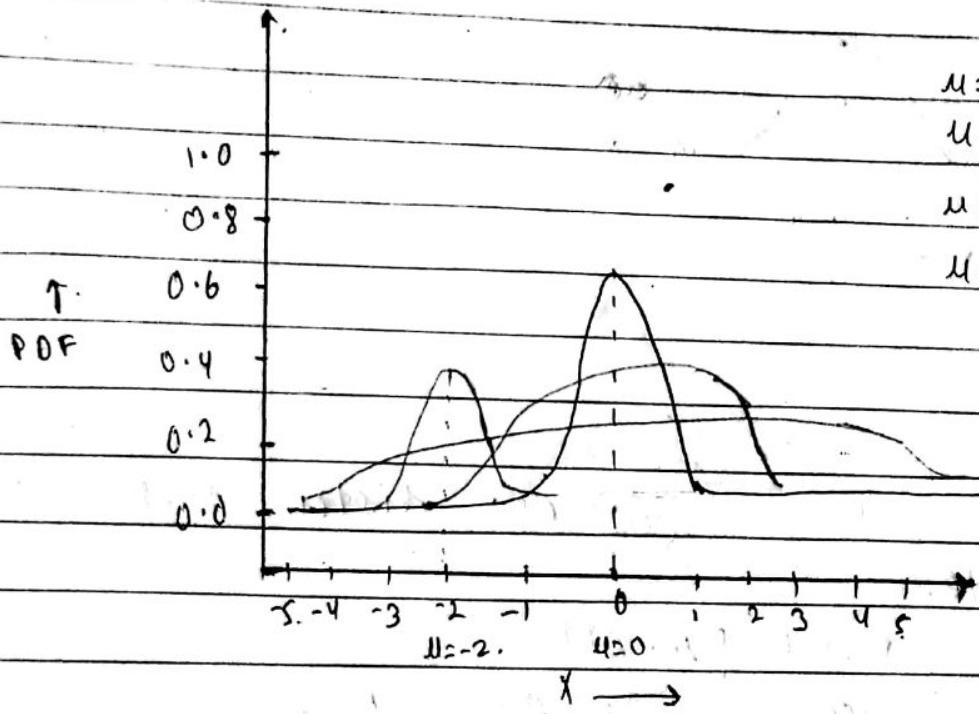
In general



* (Continuous Random Variable)

i.e. just by knowing σ , μ and σ^2 , we can find out PDF.

[μ and σ^2 are parameters of Gaussian distribution].



Legend:

$\mu = 0$	$\sigma^2 = 0.2$	-
$\mu = 0$	$\sigma^2 = 1.0$	-
$\mu = 0$	$\sigma^2 = 5.0$	-
$\mu = -2$	$\sigma^2 = 0.5$	-

Here, variance represents the spread.
Greater the spread, more pointed the PDF.

normal (Gaussian) distribution.

• $X \sim N(\mu, \sigma^2)$

↑ ↑ ↑ ↗ variance

random follows mean

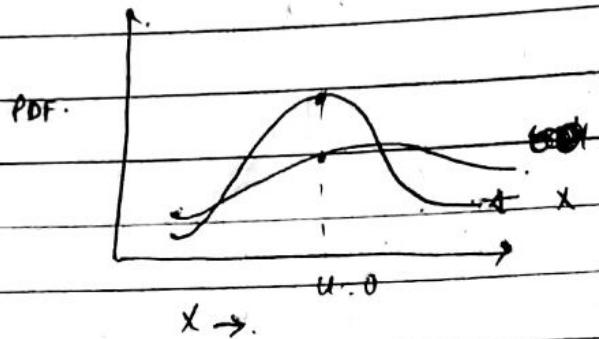
variable Parameters of Gaussian distribution

Eg. $X \sim N(0, 1)$

Scenario:

$$X \sim N(0, 2)$$

$$Y \sim N(0, 4)$$



Conclusion

If X is a Gaussian distributed random variable:

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

Just by knowing μ, σ^2 , we can calculate PDF.

And, it also lies at the peak of bell shaped curve.

Mathematically,

$$P(X=x) = P(x) = \frac{1}{\sqrt{2\pi} \cdot 6} \exp \left\{ -\frac{(x-4)^2}{2 \cdot 6^2} \right\}$$

↑ ↑
random value
variable

To simplify and understand,

Let's assume,

$$\mu=0, \sigma^2=1; \theta=1$$

$$\therefore P(x) = \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{x^2}{2} \right) = y \quad (\text{ignoring constants})$$

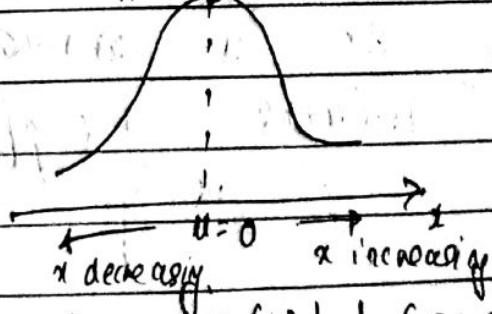
ignoring constants

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

hence, when we plot for

$y = \exp(-x^2)$, we get a bell shaped curve.

$$y = p(x)$$



$x \uparrow \text{ or } x \downarrow \rightarrow -x^2 \downarrow \text{ (always red)} \downarrow$

And hence $y \downarrow$ (also red)
i.e. $y = \exp(-x^2) \downarrow$

Conclusion from Gaussian distribution / normal distribution

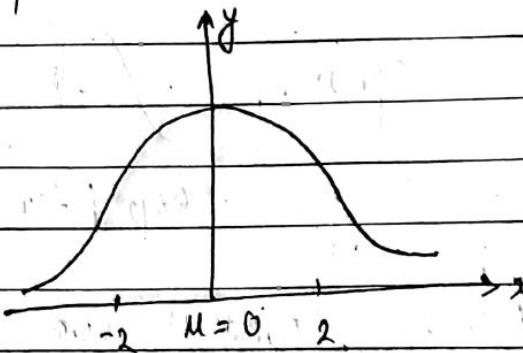
$$P(x) = y = \frac{1}{\sqrt{2\pi} \sigma^2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$



Simplified

$$P(u) = y = \exp(-u^2).$$

Graphical



#

(1)

Here, PDF is symmetric.

(2) As x moves away from μ , y reduces $\exp(-x^2)$
ie. y reduces very fast

i.e:

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

$$x=0$$

$$y=1$$

$$x=1 \quad 3^{1x}$$

$$y = \exp(-1) = 0.3678$$

$$x=2 \quad 3^{1.5x}$$

$$y = \exp(-4) = 0.0183$$

$$x=3 \quad 3^{1.5x}$$

$$y = \exp(-9) = 0.000123$$

It's clear here,

for x $1x$ increase, y $10x$, decrease.

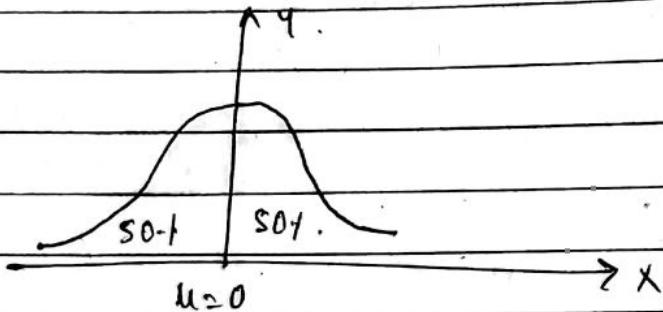
for x $1.5x$ increase, y $100x$ decrease.

(6) Quantitative distribution function (CDF) of Gaussian / Normal distribution :-

(7) Symmetric distribution, Skewness and Kurtosis :-

• Symmetric distribution :-

Ex: Gaussian distribution

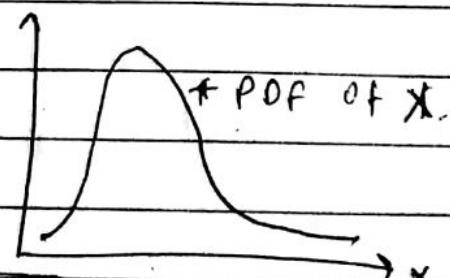
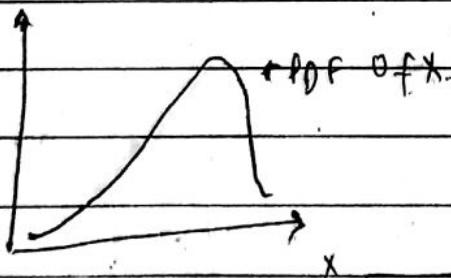


To quantify about non-symmetric distribution

(1) Skewness :-

• Negative skew

• Positive skew

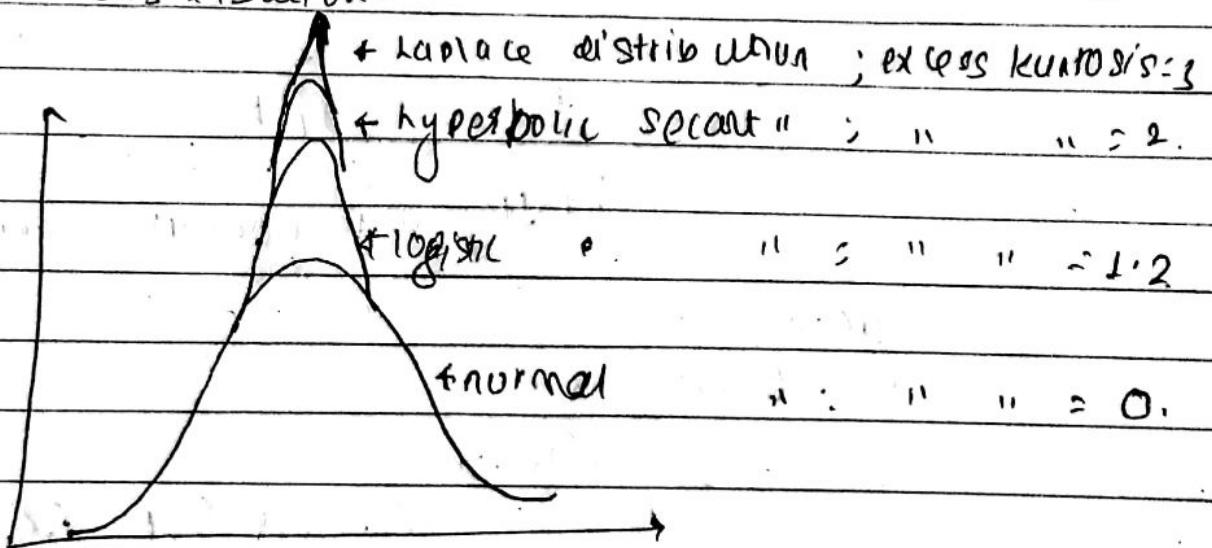


• left tail is longer | X : Random variable | • right tail is longer

Here, skewness tells us how different the distribution is from symmetric distribution.

KURTOSIS:

Measures peakness / sharpness of symmetric distribution.



⑧ Standard normal variate (z) and standardization:

$$z \text{ (SNV)}$$

$$z \sim N(0, 1)$$

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

$$\sigma^2 = 1$$

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

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

Hence
standardization:

$$x_i' = \frac{x_i - \mu}{\sigma}$$

$$x_i' \sim N(0, 1)$$

\hookrightarrow standard normal variate (z)

Hence,

Conclusion:

Given

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

Then, we can get
standard normal variate (z)
by subtracting every value of x from mean
and hence dividing by σ .

i.e.

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

With this
hence

$$z \sim N(0, 1)$$

The process is called standardization.

Why Standardization?

If we standardize, we get to know that:

if we have any normal distribution:

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

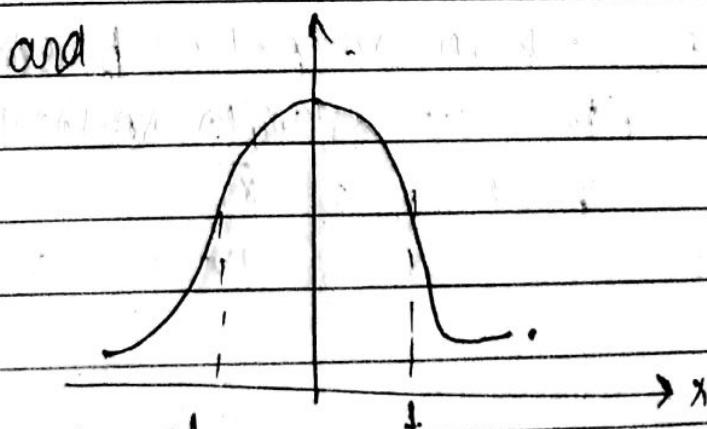
Then on standardization

via:

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

$$\therefore Z \sim N(0, 1)$$

So on standardization we get to have the normally distributed curve with $\mu = 0$ and $\sigma^2 = 1$.



Also 34.1% lies between -1 and 1.

68.2% lies between -1 and 1.

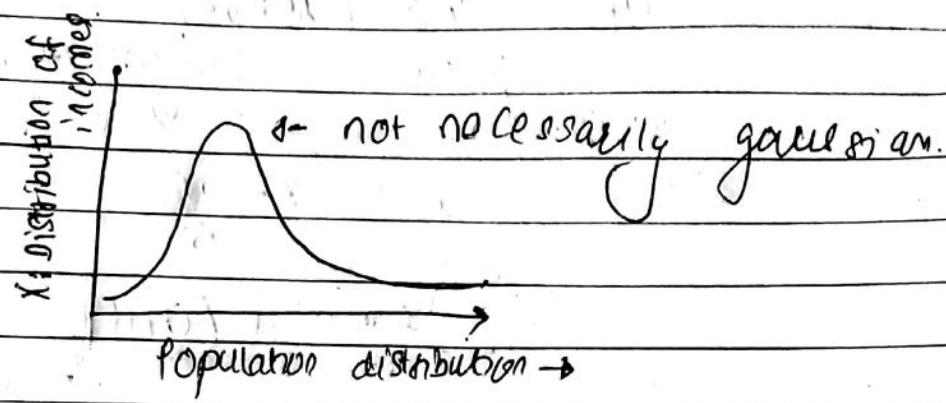
95.4% lies between -2 and 2.

⑨ Kernel density estimation

10 Sample distribution and central limit theorem

Sample distribution:-

Let's consider



Now let us take samples data from given population distribution.

let

(i) Random sample of size ($n=30$) $\rightarrow S_1$

$$\text{Let's find the sample mean of income}$$
$$\text{of } 81 = \bar{x}_1$$

$$H \rightarrow S_2$$

1120 11 11 11 11 11

$$\therefore \zeta = \bar{\zeta}_2$$

$$\therefore \zeta = \bar{\zeta}_2 \quad \text{and} \quad \eta = \bar{\eta}_2.$$

10. The following table gives the number of hours worked by each of the 1000 workers.

0 1 2 3 4 5 6 7 8 9

(a) $\text{H}_2 + \text{I}_2 \rightleftharpoons 2\text{HI}$ (exothermic) $\Rightarrow S$

P.S: All random sample: s_1, s_2, \dots, s_m taken are independent of each other.

Now,

$\bar{x}_1, \bar{x}_2, \dots, \bar{x}_m$: sample means



\bar{x}

All these \bar{x}_i where $i = 1 \text{ to } m$
will have distribution:

POI

Distribution of \bar{x}_i = Sampling distribution of sample means.



We have talked about it so far
to understand central limit theorem.

Central Limit Theorem

need $n \rightarrow \infty$ for $X = \text{popn. distribution}$
 to have finite μ and σ^2 .
 gaussian.

In case if
 has infinite μ and σ^2
 zipeto distribution

m samples of : s_1, s_2, \dots, s_m

size n_i

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

\downarrow

1 sample : $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_m$
 means

Here,

distribution of $\bar{x}_i = \text{Sampling distribution of sample mean}$

Q1

Q1

\bar{x}_i will be a
 gaussian distribution with μ
 and $\frac{\sigma^2}{n}$ as $n \rightarrow \infty$.

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

Conclusion :-

any distribution [need not to be gaussian]

$$x : \mu, \sigma^2$$

Let $n = 30$: s_1, s_2, \dots, s_m
(sample size)

Say $m = 1000$

so here we need to compute
 $30 \times 1000 = 30K$ people's income.

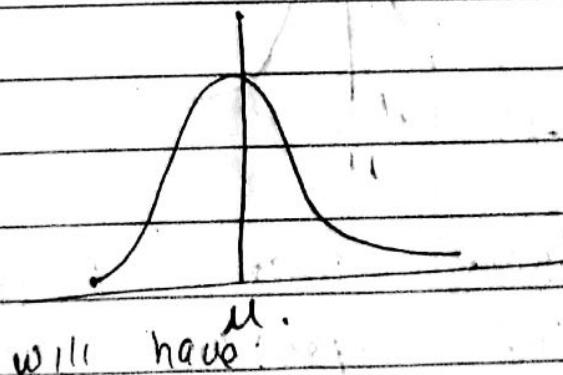
i.e:

$$: s_1, s_2, \dots, s_m$$

↓ ↓ ↓

Mean income $\rightarrow \bar{x}_1, \bar{x}_2, \dots, \bar{x}_m$
in each sample. $\hookrightarrow m = 1000$.

If I plot a distribution of it
will get gaussian distribution as



will have \bar{x}_i .

mean of $\bar{x}_i \approx \mu$ (popn. mean)

variation of $\bar{x}_i \approx \frac{\sigma^2}{n}$

such that $n \rightarrow \infty$.

But

Practically

According to rule of thumb:

We will get gaussian distribution for $n > 30$.

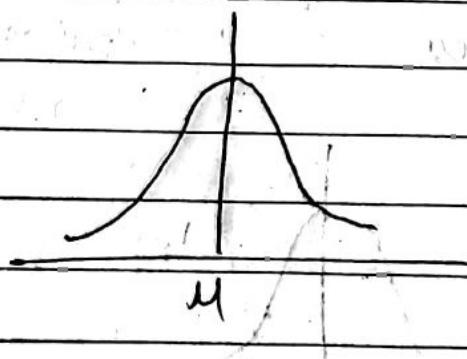
where

$n = \text{size of a sample}$.

Hence

for any distribution, we can estimate μ and σ^2

Just by doing simple sampling distribution of sample mean
And getting distribution as



And

CLT says that sampling distribution of sample mean represent gaussian/normal distribution with

$\mu = \text{popn mean}$

variance ~~σ^2~~ = $\frac{\sigma^2}{n}$ as $n \rightarrow \infty$.

OR $n \geq$ sample size.

(ii) Quantile - Quantile ($Q-Q$) Plot

Let us consider random variable ' x ' with 500 observations.

i.e.

$$x: x_1, x_2, x_3, \dots, x_{500}$$

Now, question is:

Is ' x ' Gaussian distribution?

To find out the answer we have.

' $Q-Q$ plot' [graphical method]

Steps for $Q-Q$ plot.

(1) Sort x_i 's in ascending order and compute percentile.

i.e.

$$x_1, x_2, \dots, x_{500}$$

↓ sorting

$$x'_1, x'_2, \dots, x'_{500}$$

↓ percentile computation

$$x'_{5}, x'_{10}, x'_{15}, \dots, x'_{500}$$

↓ ↓ ↓ ↓

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

↓ ↓ ↓

1st percentile 2nd percentile 100th percentile

value of x_i 's value of x_i 's " "

② Now let's create a random variable y .
such that

$$y \sim N(0, 1)$$

$y: y_1, y_2, \dots, y_{1000} \rightarrow 1000$ observations
from $N(0, 1)$.

\downarrow sort
 $y_1, y_2, \dots, y_{1000}$
1 percentile

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

③ Plotting of Q-Q plot using.

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

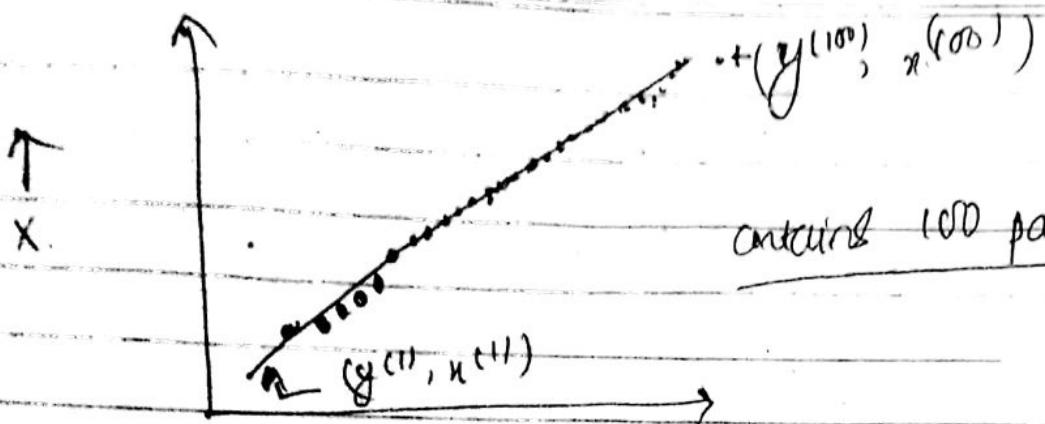
pairing like:

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

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

$$\vdots \quad \vdots \quad \vdots \\ (y^{(100)}, n^{(100)})$$

Q-Q plot



$$y = N(0, 1) \rightarrow (\text{Theoretical quantile})$$

Now, how (Q, Q) plot helps to find whether x is gaussian distribution or not:

Here each point corresponds to $(y^{(i)}, x^{(i)})$ & $i \rightarrow 1$ to 100.

If $(y^{(i)}, x^{(i)})$ & $i \rightarrow 1$ to 100 lie on a same straight line, then random variables x and y have similar distribution.

[Since 'y' is gaussian dist.
then 'x' is also gaussian distribution
since $(y^{(i)}, x^{(i)})$ lies on some straight line.]

Some code

```
① import numpy as np import random  
import scipy.stats as pylab  
std_normal = np.random.normal(loc=0, scale=1,  
size=1000)
```

It gives random obfuscation from
mean (0) and variance (\pm)

std_normal ~ $N(0, 1)$

calculating 0 to 100th percentile of
std_normal.

```
for i in range(0, 101):  
    print(i, np.percentile(std_normal, i))
```

② # Some import

for qq plot

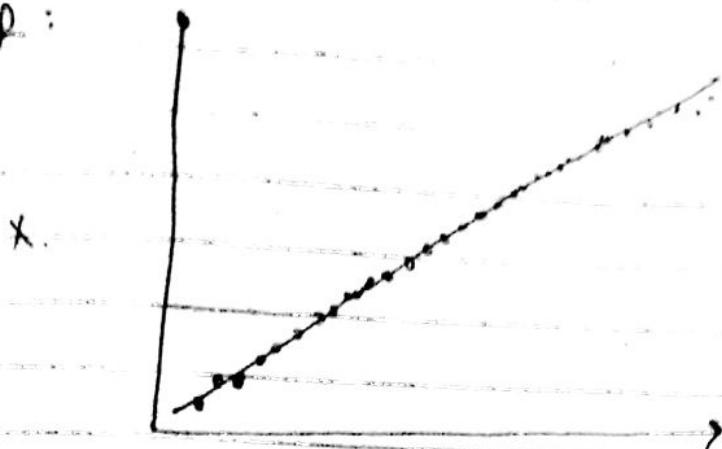
generate 100 samples from $N(20, 5)$

```
measurements = np.random.normal(loc=20, scale=5,  
size=100)
```

for qq plot

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

O/P :



$Y : N(0, 1)$ (\approx 1).

Here 'x' and 'y' are gaussian (previously said).

As my number of samples increases, Q-Q plot is more a condensed straight line.

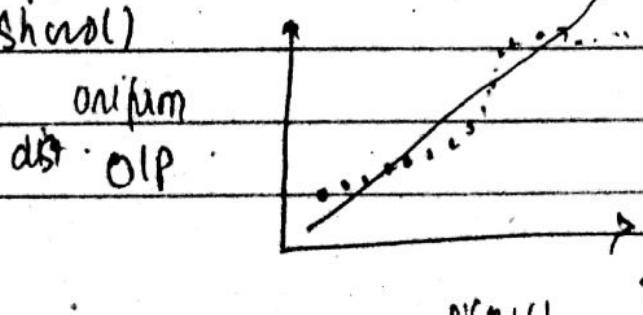
Limitation: If # Obs is small, it is hard to interpret Q-Q plot.

③

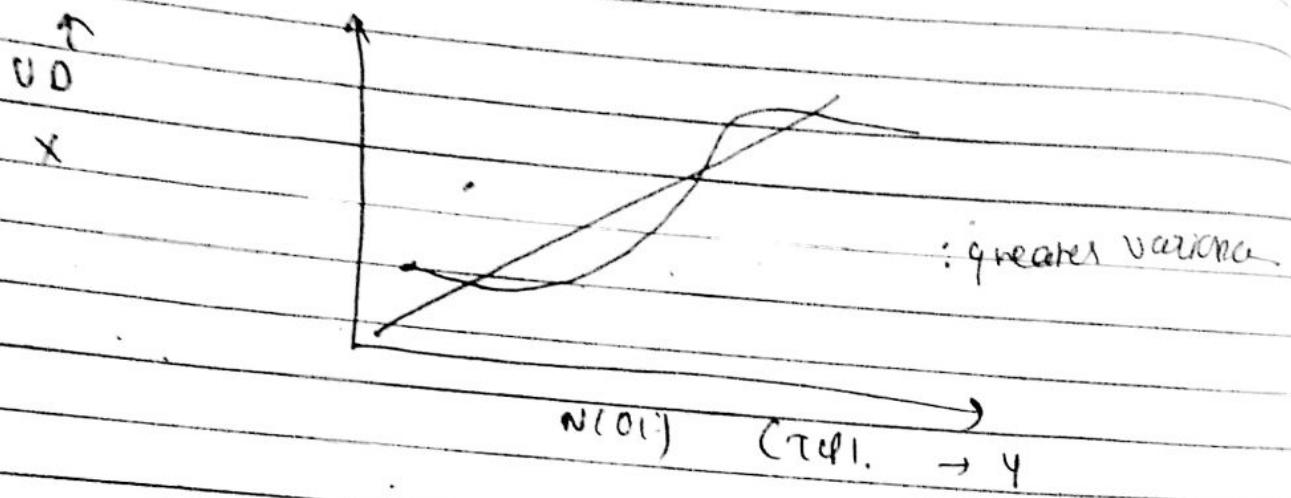
Let us take uniform distribution along y-axis

generate 100 sample from $N(0, 1)$
measurements = np.random.uniform(low=-1, high=1
size=100)

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



As we increase size of uniform dist. from 100 to 10000, we get dip in



hence X and Y are from different family.

Conclusion:-

i) Is $X \sim N(4, 6^2)$?

ii) Given any two random distribution.
Is $X \sim f_1$ & $Y \sim f_2$

broadest test. Does X and Y belong to same family?

(12) Now distribution are 0.9007

Till now we learnt a lot like
rv, pdf, gaussian, cdf, etc and soon.

Question is where we can use all these
mathematical concepts.

① # Data analysis: Answering questions about
data.

Q. Company \rightarrow Y42.

task 1: order T-shirts for all employees (100K)
(SML, XL)

Q. How many XL-T-shirts should you order?

X. Solution: Collect data for all 100K
employees. But cost more
time & money.

Solutions: from domain knowledge, we know
height $\geq 180\text{cm} \rightarrow$ wear XL t-shirt.

160s height $\leq 180\text{cm} \rightarrow$ " L t-shirt

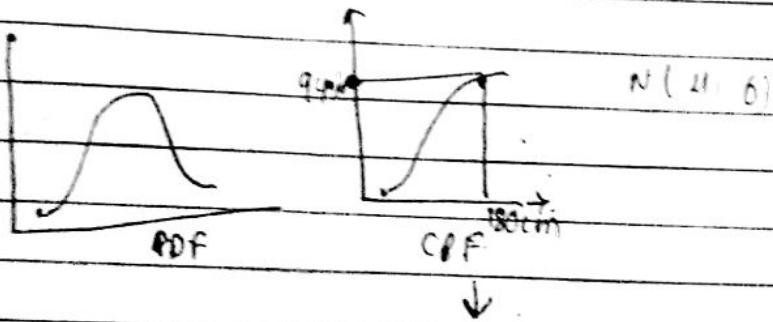
What we will do is:

1) Collect height of 500 random employees.

② Calculate μ and σ Ans
Example.

③ From doctor expert, we get
that height of people is
normally distributed.
i.e. heights $\sim N(\mu, \sigma)$.

Now we can easily compute
PDF and CDF.



$$\therefore P(h \geq 180\text{cm}) = ?$$

Hence $\pm 1\%$ employee will need
XL size t-shirt.

Normal gaussian dist. gives bell-shaped model
that is observed in many natural
phenomenon.

of salaries
of Sun (416)

? How we know salaries are gaussian distributed?

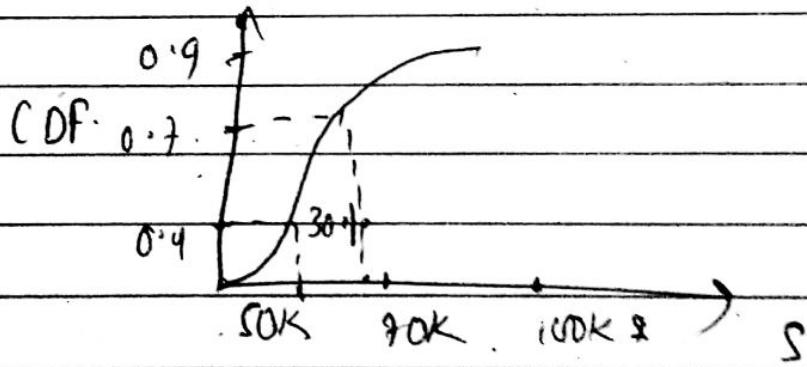
Ans : OG plot

Then, we can easily predict how many employees make a salary

(i) $> 100k \$$

(ii) $50k \$ \leq \text{salary} \leq 100k \$$

using CDF



It means 30% have salary
between 50k \\$ and 70k \\$

(13) Chebychev's inequality:

gaussian distribution:

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

→ 68% - 95% - 99.7% rule

Assume $\mu = 150 \text{ cm}$

$\sigma = 10 \text{ cm}$.

$$\begin{aligned} & \text{so } P(\text{height of student} = 95\%) \\ &= P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) \\ &= P(150 - 20 \leq X \leq 150 + 20) \\ &= P[130, 170] \end{aligned}$$

Hence, 95% student height lie in
[130, 170]

Similarly

$$68.27\% P(\mu - \sigma \leq X \leq \mu + \sigma)$$

$$99.71\% P(\mu - 3\sigma \leq X \leq \mu + 3\sigma)$$

Why Chebychev's inequality??

In case we don't know the type of distribution i.e. whether gaussian/uniform/...

But we know:

X : random variable
↳ salaries of individuals in a city
with

$$\mu = 90\text{K\$}$$

$$\sigma = 10\text{K\$}$$

- Q. What percentage of individual in range of
26. (a) $[20\text{K}, 60\text{K}]$?
36. (b) $[10\text{K}, 70\text{K}]$?

In case we knew it is gaussian distribution, ~~range will be~~

(a) 95%.

(b) 99.7%.

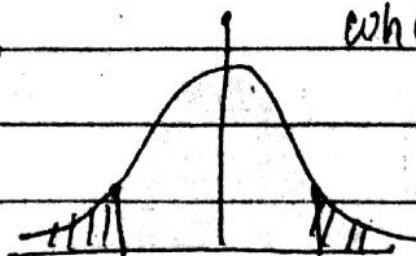
But,

here, we don't know the type of distribution incorporated.

Chebyshov's inequality?

Usage: when distribution is unknown

Diagram:



where: $\mu = \text{fixed}$

$\sigma = \text{non zero and finite}$

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

This is $k\sigma$ range around μ
 what % of points lies in \uparrow
 Now answering prev question

of salaries is

$$\mu = 40K$$

$$\sigma = 10K$$

(4) What % of individuals have salary in
 range of $[20K, 60K]$?
 $40 - 2 \times 10 \quad 40 + 2 \times 10$

$$\text{i.e. } \mu - 2\sigma, \mu + 2\sigma$$

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

$$\geq 1 - \frac{1}{4}$$

$$\geq \frac{3}{4} = 75\%$$

$\therefore \geq 75\% \text{ have salary } [20K, 60K]$

likewise.

(b) What if $P(10K, 2011)$
given: $P(\mu + 3\sigma \leq X \leq \mu + 6\sigma)$

$$P(1 - \frac{1}{9} \geq \frac{8}{9}) \geq \frac{800}{9}\%$$

Answer

More, chebsej inequality helps us to answer such question knowing not what the distribution is just by knowing its μ and σ .

14) Discrete and continuous uniform distribution

Uniform distribution types:

① of random variable (X) is discrete,
 \hookrightarrow discrete uniform distribution.

② " " " " " " " " continuous
 \hookrightarrow continuous "

① Discrete uniform distribution

for discrete random variable : PMF.
(Probability Mass fn.)

" continuous " " " " : PDF
(" density fn)

Here
let us consider dice throws
outcomes: 1, 2, 3, 4, 5, 6 $P = \frac{1}{6}$
parameter used in discrete uniform dist:

a : low

b : high

$$n = b - a + 1$$

In case of dice throw

$$a = 1$$

$$b = 6$$

$$n = b - a + 1 = 6 - 1 + 1 = 6$$

Basically $n = \text{no of outcomes}$

like

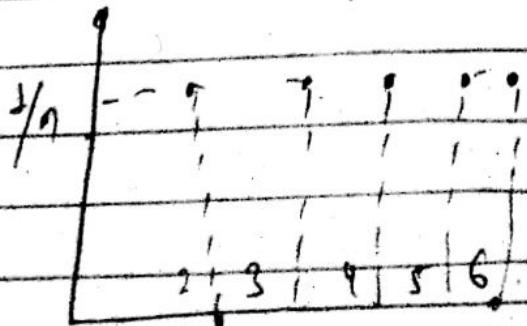
$N(4, 6)$	\rightarrow	$U(a, b)$
Gaussian dist		uniform dist

in $U(9, 15)$

All cases are equiprobable $(e = \frac{1}{n})$

feat. PMF

Eg:



$$\text{ie } U(2, 6) = \frac{1}{5}$$

$$\text{where } n = b - a + 1$$

(i) Discrete uniform

or discrete uniform

$U(a, b)$

$N = b - a + 1$

PMF: $\frac{1}{N}$

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

median

skewness: 0

(ii) Continuous uniform distribution

$U(a, b)$

Here, care taken in account of points

on the interval as well

i.e.



$$\text{PDF} = \frac{1}{b-a} \text{ for } x \in [a, b]$$

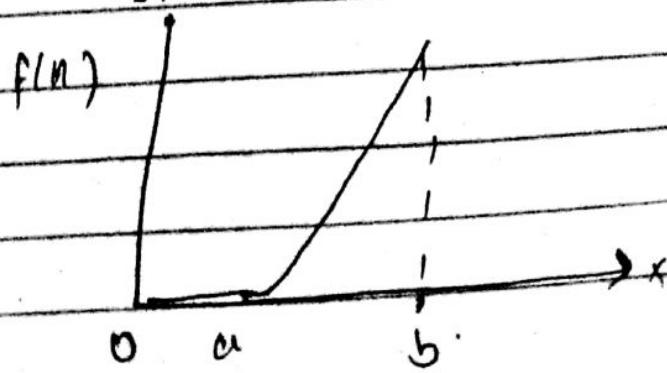
$$\text{mean} = \frac{a+b}{2}$$

$$\text{median} = \frac{a+b}{2}$$

$$\text{skewness} = 0$$

PDF: continuous random variable

CDF



(*) How to randomly sample data points
uniformly (deterministically)

We have dataset of 'n' points
 $d_{1:n} \text{ (points)}$

\downarrow We are trying to uniformly sample
'n' datapoints into 'm' datapoints
of $d_{1:m}$ where all 'n' datapoints selection
(sample) equal probability.

• random.random() generates random
values (any)

i.e.: import random.
print(random.random())

O/P: 0.79

any random value ranging from 0 to 1

• Loading Iris dataset with 150 points.

```
from sklearn import datasets  
iris = datasets.load_iris()  
d = iris.data  
d.shape
```

O/P: (150, 4)

Now to check internal working of random number generation (uniform distribution)

one at
time
greater
application

$$\begin{aligned} n &= 150 \\ m &= 30 \\ p &= m/n \\ \text{print}(p); \\ \text{sampled data} &= [] \end{aligned}$$

Here $p = 0.2$
Till random value ≤ 0.2 .
we are going to append dataset.

Finally print length.

for i in range(0, n):

$$a = \text{random.random}()$$

if $a \leq p$:

sampled_data.append(d[i, :])

len(sampled_data)

So, if result is 32

it means 20% of

equivalent

value for

O/P from 0

O/P: 0.2.

32 / 28 / 29 / 27 (always changing)

(1) Bernoulli and binomial distribution

(1) Bernoulli distribution:

- Discrete
- Parameters: $0 < p < 1$.
- used in event has 2 outcomes
- Pmf: $P(K=k) = p^k q^{1-k}$ for $K=0, 1$
 $q = 1-p$ for $K=0$

Ex: Mean = p

coin toss



(2) Binomial distribution:

Let us consider coin toss :- where

X ~ Bernoulli ($p = 0.5$)

when 'n' trials are conducted.

Ans.

Let y = no of times I get a head.

When I toss my fair coin 'n' times. (say $n = 10$)

I may get 10 heads

$\therefore y = \{0, 1, 2, \dots, 10\}$ on 10 trials

I may get '0' heads on 10 trials

hence,

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

\uparrow \uparrow
no. of trials prob of success in getting least
 performance

This is called binomial distribution.

• Parameters

n = no. of trials

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

• PMF

$$P(K) = {}^n C_k p^k (1-p)^{n-k}$$

(17) Log normal distribution

If a random variable is log normally distributed:

$$\begin{cases} \ln(x) = \text{normal} \\ x = e^{\ln(x)} \end{cases}$$

$$x \sim \text{log-normal } (\mu, \sigma^2)$$

The,

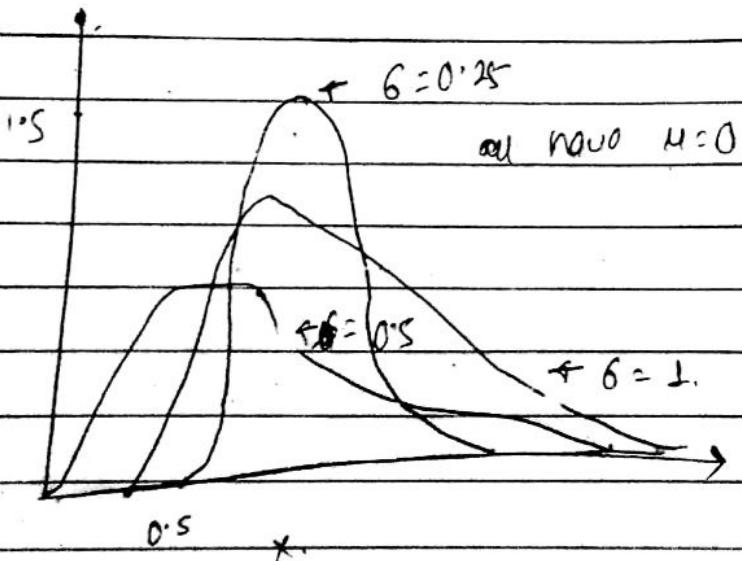
if then

natural
logarithm

$y = \ln(x)$ has a normal
distribution

Log normal

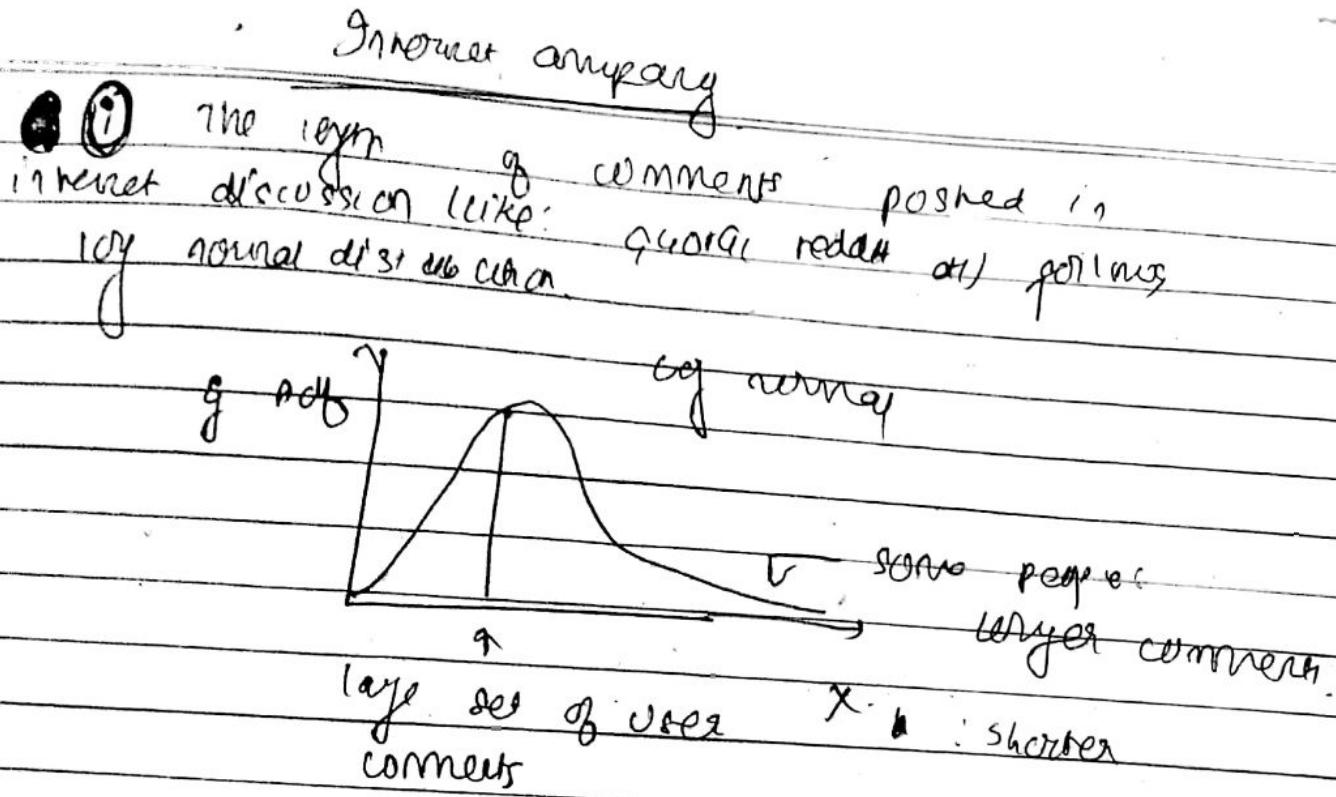
PDF.



All are skewed

and as σ increases spread has increased

* Occurrence and application



(ii) User dwell time on online articles follows a log normal distribution.

[Alternative: Pareto distribution]

Economics

(i) In economics, there is evidence that income of (97-99)% of population is distributed up normally. The distribution of high individuals follows a power law distribution.

(ii) In finance (stock market, ...) follows Black-Scholes model has a heavy tail.

e.g. normal dist.



$X \sim \text{log-normal}(\mu, \sigma)$

Q. Is X log normal?

To find out.

$X: n_1 \ n_2 \dots n_k$

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

$\ln(n_1) \ \ln(n_2) \ \dots \ \ln(n_k)$

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

$y_1 \ y_2 \ \dots \ y_n$

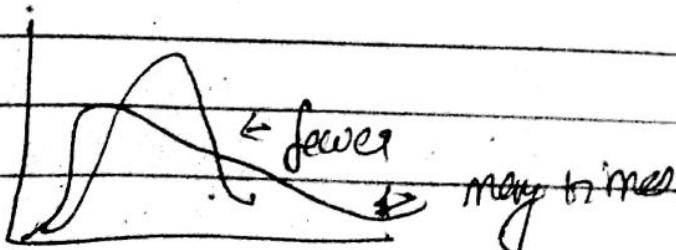
Check if y_i 's are gaussian using Q-Q plot

\downarrow

If yes

\underline{X} is log normal. (μ, σ)

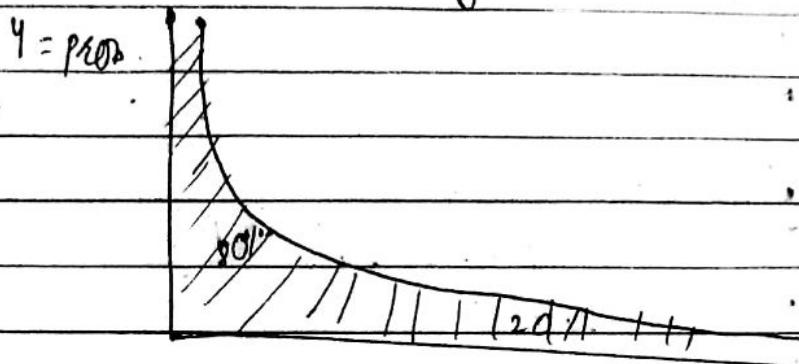
In most internet companies, we will get log normal distribution with skewed tail.



(18) Power law distribution, (Pareto distribution)

Before discussing on pareto distribution, let's focus on power law :-

Eg of power law graph



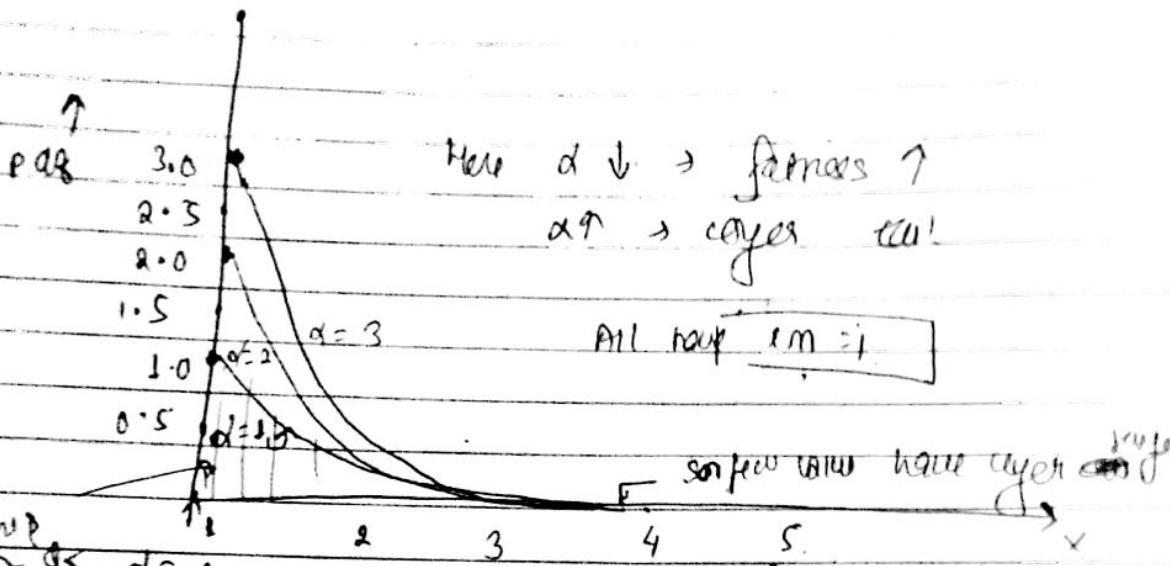
To the right: long tail. and rule
to the left are the few that dominate 80-20 rule.

↓
Pareto distribution:

Parameters:
1. $\alpha > 0$: Scale. = similar to $u \sim \gamma^{-1}$
2. $\gamma > 0$: shape = similar to b {gaussian distn}

Example

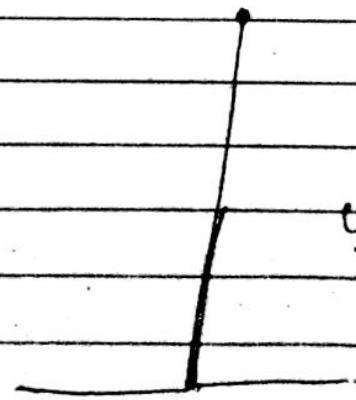
PDF



continuous random variable. As $d \rightarrow \infty$, distribution approaches $\delta(x - xm)$ where δ is a known δ function

| Dirac delta function |

i.e. $d \rightarrow \infty$

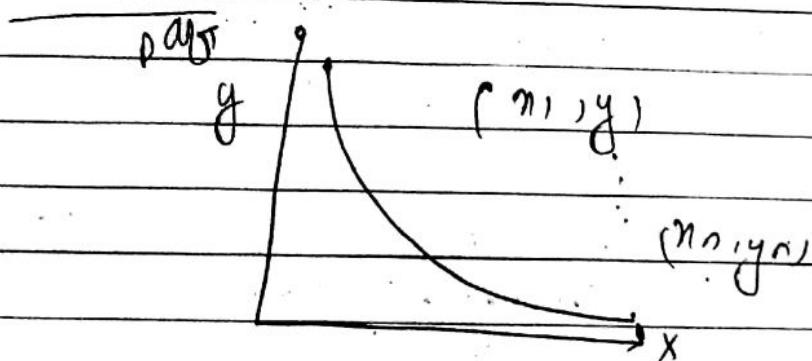


like δ , just appears called
at a given value.

App qualitatively

* To check for power law?

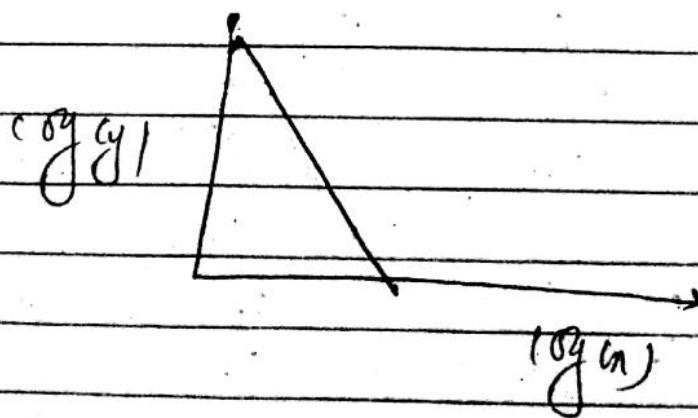
We have



for power law checking

$$n_i \rightarrow (\log n_i)$$
$$y_i \rightarrow (\log y_i).$$

make log-log plot (if we get a straight line like this, then this follows power law)



• To check for pareto distribution

• $X \rightarrow \text{obs}$

• $Y \sim \text{pareto distribution}$.

then, If X, Y give straight line (Q-Q Plot)
 X : pareto dist.

Power transform \rightarrow

(19) Box Cox transform.

Case I: $X \xrightarrow{\ln} Y$ we can ~~assess~~
 \rightarrow study new
(lognormal) (gaussian) since we know
a lot on gaussian dist.

K. $\ln X$ (assume that a feature is
gaussian dist):

Case II $X \rightarrow Y$?
power law ~~or pareto~~ How: Gaussian dist's
dist. we do it.

Here comes!

Power Transform

We use Box Cox transform under it.

Patho $\sim X$: $\{n_1, n_2, \dots, n_n\}$



We need to transform it to gaussian
to Gaussianity (g_1, g_2, \dots, g_n)

For we have to g_i 's from n_i .

Step 1: Apply box-cox (X) function.

$\{n_1, n_2, \dots, n_n\}$

Generate candidate (d).



How it generates? No need to understand.

Step 2: Now $y_i^d = \begin{cases} \frac{n_i^d - 1}{d} & \text{if } d \neq 0 \\ \ln(n_i) & \text{if } d = 0 \end{cases}$

i.e.

In case if $d = 0$

$\lambda \sim \text{log-normal}$

because $y_i^d = \ln(n_i)$ if $d = 0$

* look at also

scipy.stats.boxcox

boxcox (X , lmbda = None, alpha = None)

where:

X : n.d. array.

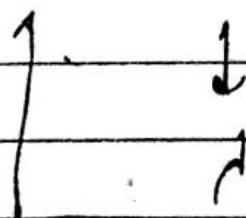
Returns

boxcox : nd array

i.e. box-cox power transformed array

i.e.

$Y = \text{boxcox}(X)$



automatically returning

Gaussian distribution

Let's look at the code

fig = plt.figure()

ax1 = fig.add_subplot(211)

from scipy import stats
import matplotlib.pyplot as plt

x = stats.loggamma.rvs(5, size=500) + 5

Scenario ↗

~~x~~ ~ loggamma (power law/parato
dist)

rarely used.

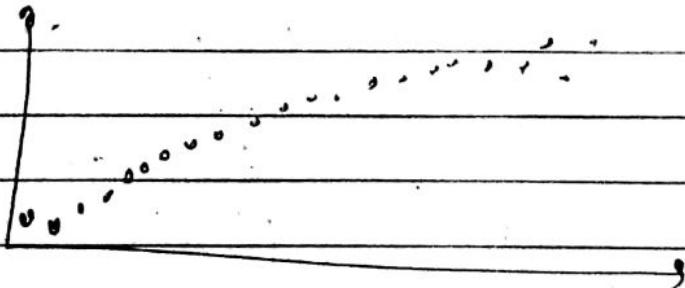
qqplot ↙ prob = stats.probplot(x, dist='stats.norm')

plot = ax1,

plt.show()

gives

(x) ~ loggamma



normal (y)

qqplot

* points deviates

X and Y don't have same distribution.

i.e. not gaussian.

b

Box Cox

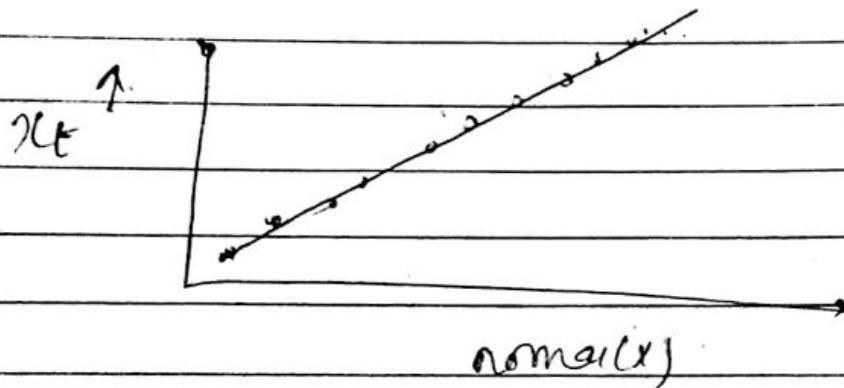
Box Cox.

• $\alpha_{x2} = \text{fig.add-subplot}(2, 2)$

• $x_t = \text{stats.boxcox}(X)$

$p105 = \text{stats.probplot}(x_t, dist='stats.norm')$
 $\text{plot} = \alpha_{x2}$.

plt.show()



Boxcox() best fm to
convert pureto dist to gaussian
distribution.

② Application of non gaussian distributions?

- Dist.

↓

Gaussian

↑

Non gaussian

- Uniform → used to generate random numbers

- Bernoulli, binomial

- Log-normal, pareto

There are 100's of distributions

? Why? ↓

if we have a random variable

(wellstudied gives us a

$X \sim \text{Distr}$) → theoretical model

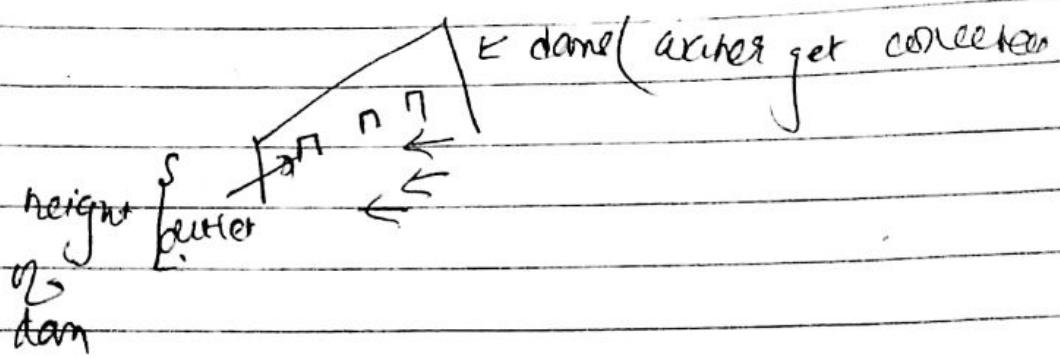
↓

Then with we can
~~not~~ have a lot of
inference / analysis /
understanding.

12:40

→ April 17.
558 [4.4 GB → 0.0]

* let us talk about
| 'weibull distribution'
| non gaussian
personal story



To mitigate accidents, CE have to determine
max. one day rainfall & for
dam height and no. of outlet suppression

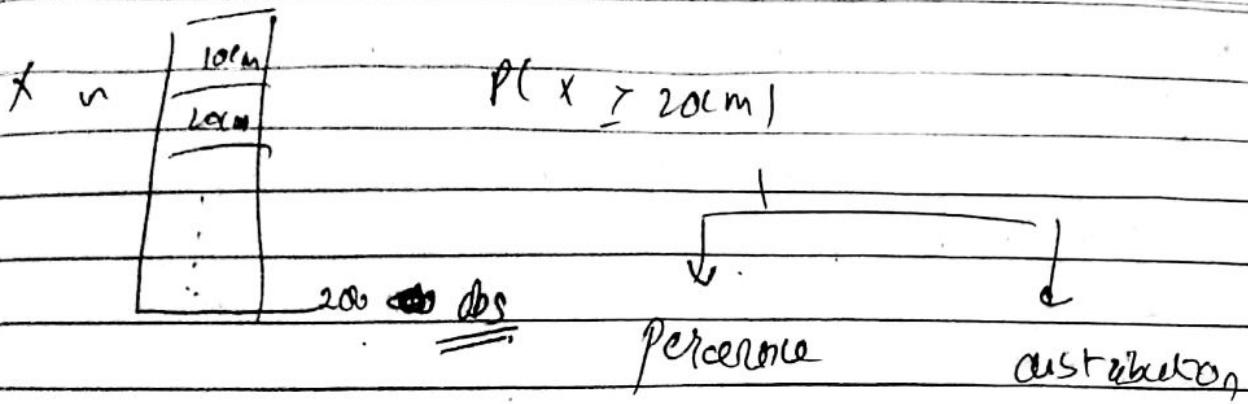
let $x = \text{max. one day rainfall}$

E

Now $P(x \geq 20\text{cm}) = ?$

In traditional time (pre computer era)

Let say, CE checks few 200. data points.



* percentage

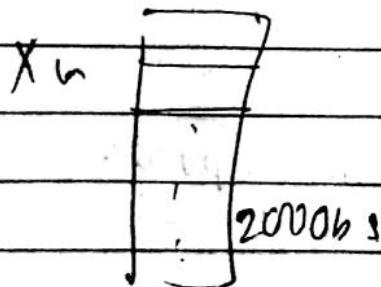
Let say only 1 obs meet $P(X \geq 20\text{cm})$

$$\therefore P(X \geq 20\text{cm}) = 1 - 0.5 = 0.5$$

$\frac{1}{200}$

(not trustworthy
can't trust)

* Theoretical model



→ fit a known and well studied dist/ theoretical model

Let say it fit Weibull distribution
(originally used for
particle physics)

Now X = max - daily rainfall

①

→ weibull (params)

③ pdf and cdf

②

$$④ P(X > 20\text{cm}) = 1 - P(X \leq 20\text{cm})$$

from cdf of X can be done

Why mathematical model is more mathematically?

Because in model all 200 data points fits the model.

Data points which are fitting only & fit.

So $P(x \geq 20m)$ is more trusted only here.

Today

Weibull \rightarrow Purely physical (originally)
 \rightarrow Civil Engg (now)

Trend: Distribution has greater application

[Non-gaussian dist ~~to~~ appear in terms of applications.]

(2)

D) Find the relationship between random variables :-

(a)

X: heights	$X = h$	$y = co.$
y: foreign	S ₁	160
	S ₂	150
	:	:
	S _n	140
		48

Q. relationship b/w X & Y

When X↑, does Y↑ ?

When X↓, does Y↓ ?

Measurement of relationships.

(A) Covariance

(B) Pearson correlation coeff.

(C) Spearman rank correlation coeff.

④ Covariance

$$\text{cov}(x_1y_1) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{var}(x_1) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\text{i.e. } \text{cov}(x_1x) = \text{var}(x)$$

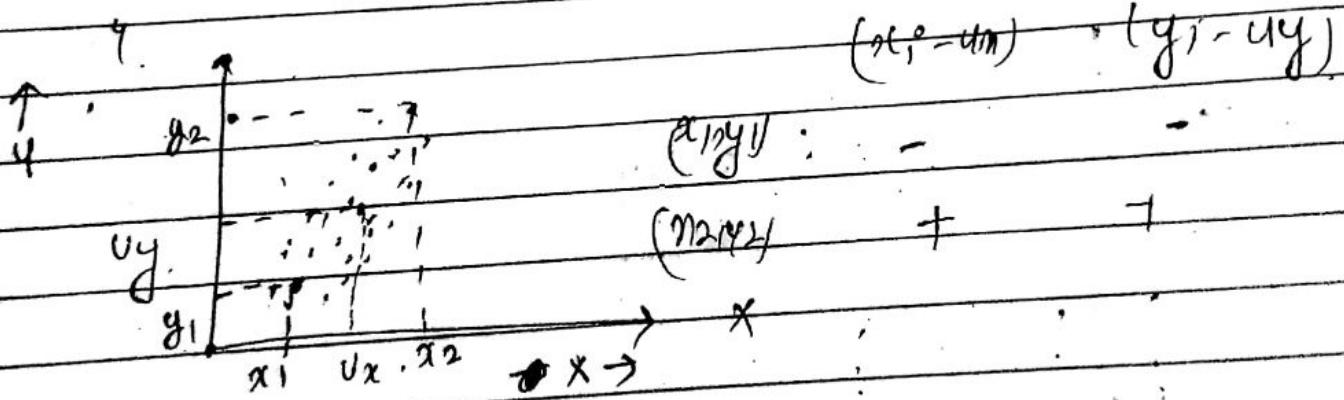
No. 20

$\text{cov}(x_1y_1) = \text{pos if } x_1, y_1$

$\text{cov}(x_1y_1) = \text{neg if } x_1, y_1$

1st scenario

$\text{cov}(x_1y_1) = \text{pos if } x_1, y_1$

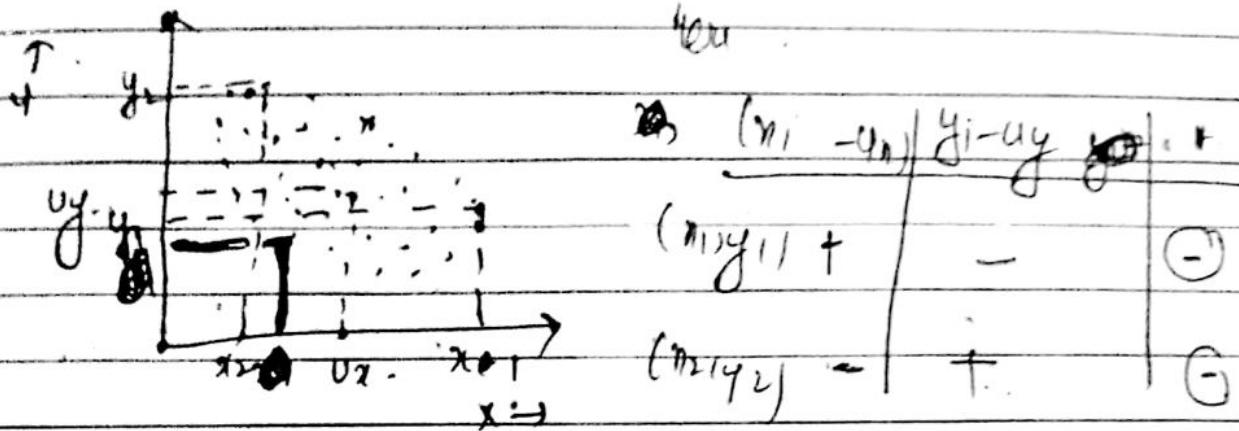


$$\text{cov}(x_1y_1) = \text{pos}$$

farmer the points form mean also (m_1, y_2) more positive the covariance

2nd scenario

$$\text{cov}(x_1y_1) = \text{ve} \text{ if } x \neq y$$



$$\text{cov}(x_1y_1) = \text{ve} \rightarrow \text{LT}, \text{ y↓}$$

Farmer the rain from mean (neg),
more negative the rainfall.

*

+ (less positive)

+

(more positive)

i.e.: Tells no relationship b/w x and y

Drawbacks

• S1 S7

+ height, weight (kg)

cov(x₁y₁)

→ weight (kg) ↑ cov(x₁, y₁)

height (cm)

↑ ↑ lbs.

i.e. Just by changing the metric of some set of measurements of some data, covariance may not match.

① Pearson correlation corr: 0 to 1

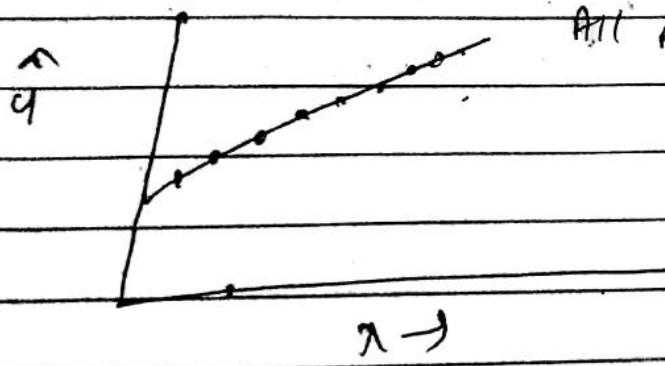
$$P_{x,y} = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y}$$

$$\text{where } \sigma_x = \sqrt{\text{var}(x)}$$

$\rightarrow \text{cov}(x,y)$

$x \uparrow, y \uparrow, \dots = \text{pos.} \quad \left. \begin{array}{l} \text{let covar} \\ \text{pos} \end{array} \right\}$
 $x \uparrow, y \downarrow, \dots = \text{neg.} \quad \left. \begin{array}{l} \text{let covar} \\ \text{neg} \end{array} \right\}$
covariance captures how much the +ve or
how much -ve.

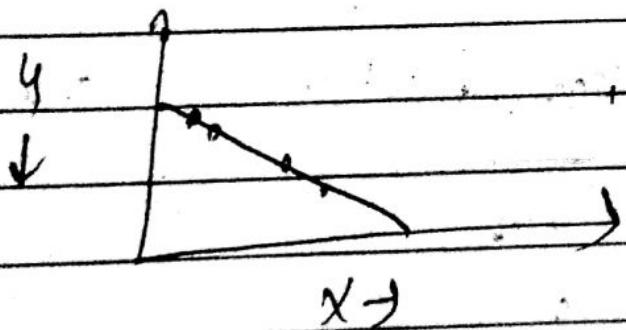
Graphs



All points lies on straight line.

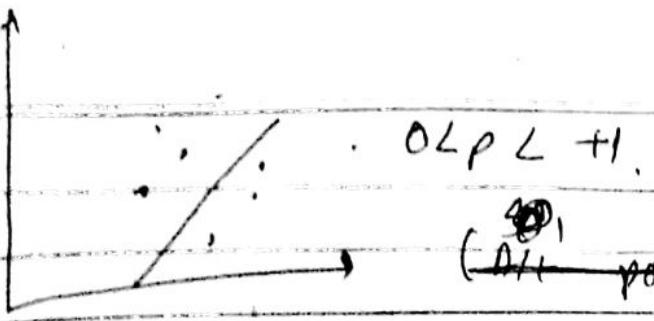
$$P = +1.$$

(Since x_1, y_1)



$$P = -1.$$

$\sin P(-x_1, y_1)$



~~(All)~~ point \to no linear relationship.

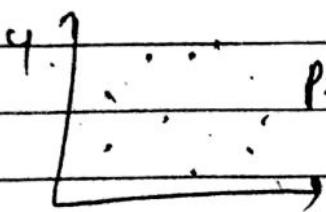
between x and y) so,
And

$X \uparrow, Y \uparrow.$



$\uparrow \downarrow$

$X \uparrow, Y \downarrow.$

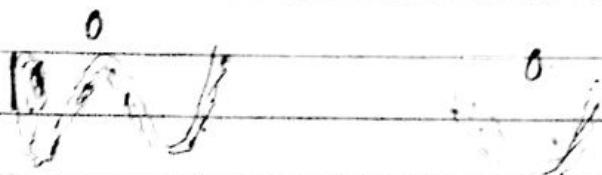
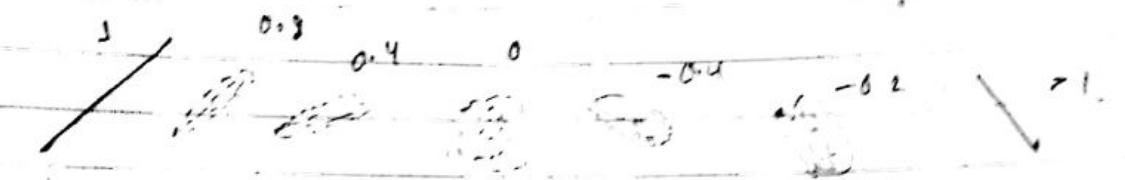


NO relationship at all

$\left\{ \begin{array}{l} \text{when } p \neq 0 \text{ there is only a linear relationship between } x \text{ and } y. \\ \text{when } p = 0 \text{ there is no relationship between } x \text{ and } y. \end{array} \right.$

corr = 100

Some example values of P.



$$(Y = \sin(X) + 0.7)$$

i.e. PCC only useful in case of linear relationships
between X and Y.

So we introduce another
concept.

(C) Spearman Rank correlation coeff.
(r)



mer, we calculate rank

Obs	Sg	X	Y	r _x	r _y
1	160	52	4	4	2
2	150	78	2	2	4
3	170	68	5	5	5
4	140	46	1	1	1
5	158	51	3		2

~~R_{xy}~~ is calculated as

$$r_x = 1 - \frac{(\text{sum of } r_{xy} + x)}{n}$$

$$r_y = 1 - \frac{(\text{sum of } r_{xy} - y)}{n}$$

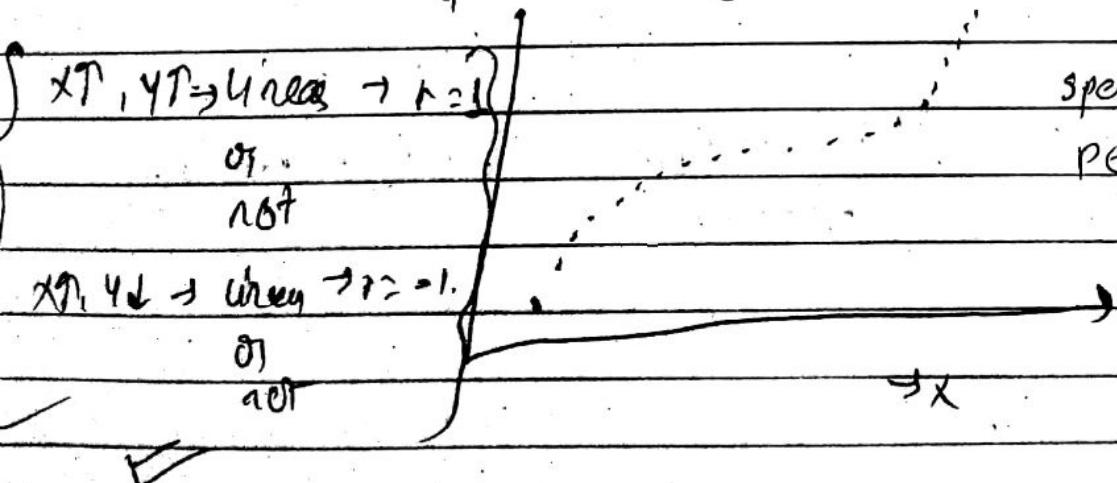
and so on.

Summary r_y [assumed $r_y = 1 - \dots$] on
the basis of accuracy (order of Y).

Now,

$$r = \text{Pcc.}(r_x, r_y)$$

Let see what happens when we apply &
consider strictly non decreasing graph



$$\text{Spearman's } r = 1$$

$$\text{Pearson's } P_{cc} = 0.88$$

io: regardless of linear or not

$$r = 1 \quad \left[\sum x_i^2, \sum y_i^2 \right]$$

$$r = -1 \quad \left[\sum x_i y_i \right]$$

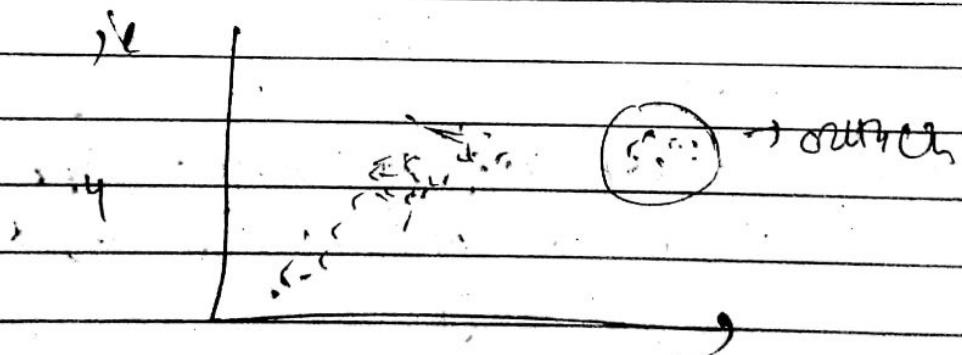
↓

But

bec won't be so exact..

* In case of outliers

r' is more robust than
'Pcc'

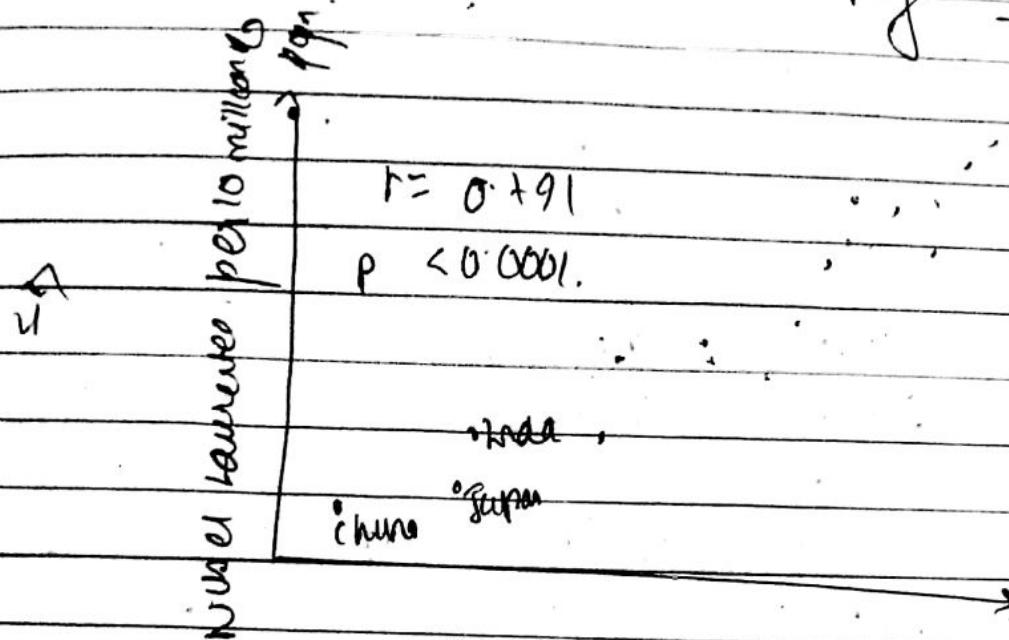


$r = 0.84$ → more robust.

$PCC = 0.67$

(22)

"Correlation" doesn't imply "causation"



Chocolate consumption (kg/capita) \rightarrow X

Not: $[x \rightarrow y]$ \rightarrow good correlation

$r = 0.791$
(Spearman rank)

But x causes y / y causes x
i.e. more chocolate consumption produces more
nobel laureates. That's absurd

i.e. if x and y are not necessary. that
correlated,

x causes y or
 y causes x .

23

Applications of correlation

pos neg.

[just correlation;
not causation]

(1) Is salary correlated with sq. footage of your home?

i.e. if salary ↑ → sq. footage ?

so, brokerage will show house according to salary.

That doesn't mean salary causes
sq. footage of home or vice versa.

(2) Is # years of education correlated with income?

x y
if $x \uparrow \rightarrow y \uparrow$

if correlated in this way,
education ministry will focus on
good education.

③ climber angon.com

④ time spent in 24 hrs : \rightarrow money spent in
last 24 hrs

*

4.

if $x \uparrow$

$\rightarrow 4 \uparrow$

⑤ # unique visitors in q : \$ sales in Q3

2y

x

4

\$ $x \uparrow \rightarrow 4 \uparrow$

100k \rightarrow f. 1M

120k \rightarrow f. 1.6M

good for
company.

⑥

No drama:

usage of a grey indicator is reduction in blood sugar

$x \uparrow \rightarrow 4 \uparrow$

if $x \uparrow \rightarrow 4 \uparrow \rightarrow$ good by

medical practitioner

(2) Confidence Interval Intro: -
C.C.I)

μ = popn. mean
 \bar{x} = sample mean

Let us consider a sample ' X ' of size '10' such that:

$$X = \{x_1, x_2, \dots, x_{10}\} \rightarrow 820 \text{ people}$$

$$X: 180, 162, 158, 172, 168, 150, 171, \\ 183, 165, 176$$

\therefore POINT ESTIMATE of $\mu = \bar{x} = \frac{1}{10} \sum_{i=1}^{10} x_i$

$$\therefore \bar{x} = 168.5 \text{ cm}$$

↑

Point estimation of μ as \bar{x} is less precise.

Despite saying

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

↑

↑

~~Confidence interval
C.C.I)~~

INTERVAL

CONFIDENCE

[is more precise]

6: 32

② Compute CI given underlying distribution.

Let

height $\sim \sigma$

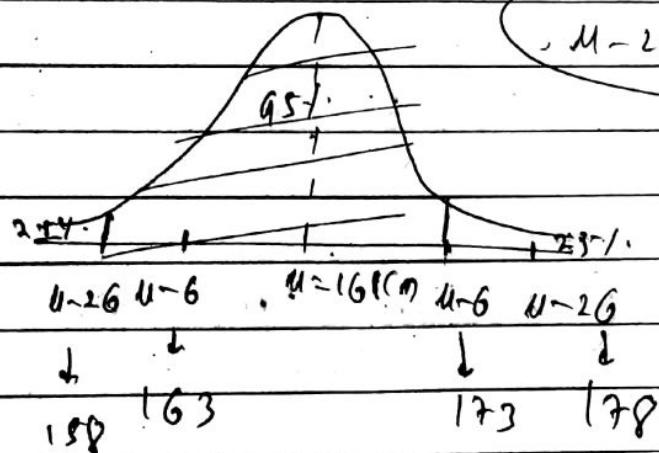
$X \sim N(4, 6)$

Let

$\mu = 16 \text{ cm}$

$\sigma = 5 \text{ cm}$

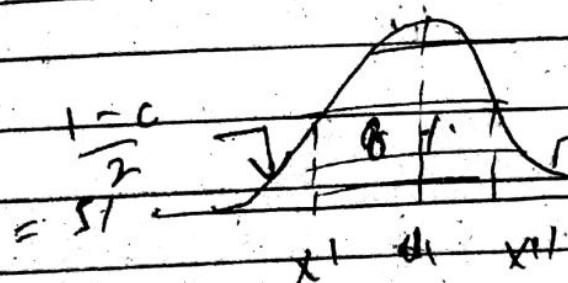
$M=26, m=26.95\%$



$\therefore 95\% \text{ of my obs. } [158, 178] \text{ with } 95\% \text{ prob.}$

95%, 68.27%, 68.1% are okay.

What about $C = 90\%$



where n^1, n^{11} are
find out using
normal
dist. table.

$\therefore \text{lie in } [x^1, x^{11}] \text{ with } 90\% \text{ confidence}$
lower bound \rightarrow upper bound

③ C.I. for mean of a R.V

Let X in F. with popn mean μ and std dev of σ .

$\{x_1, x_2, \dots, x_{10}\} \rightarrow$ sample of size 10 + n = 10

$\{180/62, 181/171, 162 \text{ Cm} \text{ in cm}\}$ thus sample
 $\{150/171, 183/161, 171\}$

⑨ What is 95% C.I. of μ ?

Given: $\sigma = 8\text{cm}$ (we know popn. std. dev.)

As per C.I.

$$\bar{x} = \text{sample mean} = \frac{1}{10} \sum_{i=1}^{10} x_i = 168\text{ cm}$$

So, for C.I.

$$\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$$

Sample mean

known popn. mean
digit

~~std popn. std dev~~

~~std sample std dev~~

$$\therefore \mu \in \left[\bar{x} - \frac{2\sigma}{\sqrt{n}}, \bar{x} + \frac{2\sigma}{\sqrt{n}} \right] \quad \begin{cases} n=26, \\ \sigma=8 \end{cases}$$

as (n - 95%) confidence

i. Using sample of size 10

\therefore popn mean (μ) $\in [85.30, 171.60]$ with
95% confidence

| Case II: if we don't know σ (popn. std. dev.)

to find ch of μ

↓

Consider sample of size 'n'

In this case, we use t-distribution
also called Student's t-distr

$\therefore \bar{x} \sim t(n-1)$

↑ ↑ ↗
sample mean t-distr degrees of freedom

* ~~Student's~~ Student's t-distr was created to
study ch of mean of RV when ' σ '
is not given.

Conclusion

Estimate CR

Estimate CR of μ if σ known

Case I: σ known, \rightarrow CR $\rightarrow N(\mu, \frac{\sigma^2}{n})$

Case II: σ unknown \rightarrow t dist($n-1$)



Estimate CR of μ if σ unknown

" " " median "

" " " 90th percentile "



using

bootstrap CR \rightarrow modern computer

(9:3)

(11:18)

(3:18)

④ Confidence Interval (CI)
using bootstrapping

Let

 $X \sim F$

↑

any distribution

Let $S^1 = \text{sample of size } (n) = 10$

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

- ↑ discrete

↓
processed by uniform
random numbers generator,
 $U(1/n)$.

random sample w \leftarrow $S_1 : \{x_1^{(1)}, x_2^{(1)}, x_3^{(1)}, \dots, x_m^{(1)}\}$

Size 'm' generated
from 'S' using

 $U(1, n)$

may be sampling with repetition

Hence

Conclusion

Estimate CR of μ by a RV

estimate CR of μ by a RV

Case I: σ known, \rightarrow CR $\rightarrow N(\mu, \frac{\sigma^2}{n})$

Case II: σ unknown \rightarrow t dist($n-1$)

Q

Estimate CR of σ by a RV

" " " median "

" " " 90th percentile "



using

bootstrap CR \rightarrow modern computer

(6)

(9:13)

(11:18)

(11:18)

(4)

Confidence Interval (CI)

using bootstrapping

Let

$X \sim F$

T

any distribution

Let $S^1 = \text{sample of size } n_1 = 10$

$$\{x_1, \dots, x_{n_1}\}$$

discrete.

↓
Process by uniform
random numbers generator,
 $U(0,1)$.

random sample w

$$S_1 : \{x_1^{(1)}, x_2^{(1)}, x_3^{(1)}, \dots, x_m^{(1)}\}$$

Size 'm' generated

from 'S' using

$\sigma(t, n)$

such that $m \leq n$.

may be sampling with repetition

Hence

$S = \{n_1, n_2, \dots, n_k\}$

Step

↓ using sampling replacement

$\{n_1^{(1)}, n_2^{(1)}, \dots, n_m^{(1)}\} \rightarrow m_1$

Step

$\{n_1^{(2)}, n_2^{(2)}, \dots, n_m^{(2)}\} \rightarrow m_2$

bootstrap

samples

Size
 $m \leq n$

size : $n_1, n_2, \dots, n_m \rightarrow m_k$

while

$m_1 = \text{median of samples}$

$m_2 = \dots \text{ " } 2$

$m_k = \dots \text{ " } k$

arrow

Let we obtain 1000 median using
bootstrap samples

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

sort

$m_1' \leq m_2' \leq m_3' \dots \leq m_{1000}'$ (increasing)

$m_{125}' \leq \bar{m}_{950}' \rightarrow 750 \text{ order}$,
 m_{125}'

15:47

$\therefore 95\% \text{ CI of mean of } X_i$

$$\{n_{LR}, n_{RR}\}$$

for S and variance, just calculate them
from bootstrap sample and find the

CI from sorted ones.

Hence, Bootstrap technique is non parametric
technique where we don't make any assumption
about dist. of data.

Code

```
import numpy
from pandas import read_csv
from sklearn import resample
from numpylib import eyelet
```

Load dataset

```
x = numpy.array([100, 162, ..., 176])
```

Configure resample

```
n_iteration = 1000
```

(K)

```
n_size = int(len(x))
```

\rightarrow size of sample (n)

run bootstrap

median = np.median

for i in range(n_iterations):

prepone train and test sets

s = resample(x, n_samples=n_size);

m = numpy.median(y);

medians.append(m)

plot scores

pyplot.hist(medians)

pyplot.show()

confidence interval

alpha = 0.95

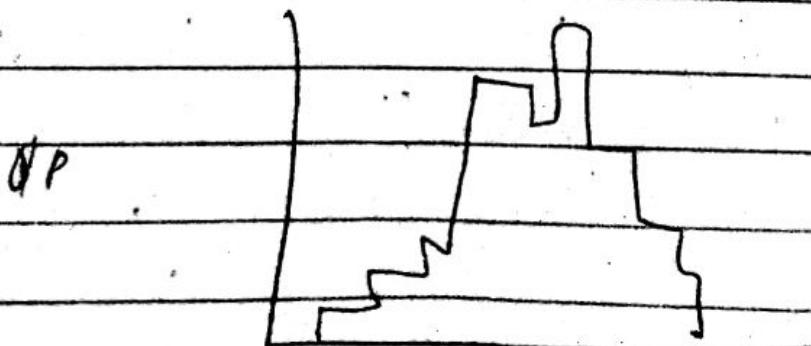
p = ((1.0 - alpha) / 2.0) * 100

lower = numpy.percentile(medians, p)

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

upper = numpy.percentile(medians, p)

print('%.1f confidence interval %.1f and
%.1f' % (alpha * 100, lower, upper))



95.0 confidence interval 152.0 and 162.0

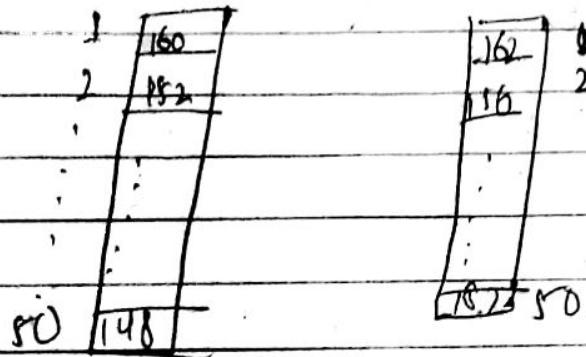
22

Hypothesis testing

i)

M180

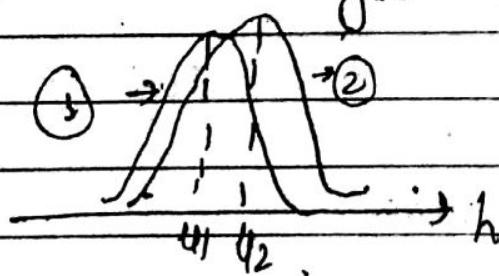
Let us assume heights of 50 students in
C1 : \bar{x}_1 C2 : \bar{x}_2 classroom L 42.



Q9 Is there a diff in heights of
students in C1 and C2?



Let us draw histograms



$$H_2 > H_1 \Rightarrow \text{height}_{C2} > \text{height}_{C1}$$

But \bar{x}_2 is slightly greater than \bar{x}_1

How can we be so sure that $\text{height}_{C2} > \text{height}_{C1}$???

To correctly estimate such a confidence.

Hypothesis testing.

Step 1 Choosing a test - Statistics

$\mu_2 = \text{mean height of C12}$
 $\mu_1 = \text{... of C11 students}$

$\therefore \text{Test } \bar{x} - x = (\mu_2 - \mu_1)$
statistic

Step 2 Null hypothesis (H_0). [Prove by contradiction]

$H_0 = \text{no difference in}$
 $\mu_1 \text{ and } \mu_2$.

Alternative hypothesis (H_1):

$H_1 = \text{difference in } \mu_1 \text{ and } \mu_2$

Step 3 P-value:

Prob of obs ($\mu_2 - \mu_1$) provided if
null hyp is true.

Assume H_0 is true

$$\text{i.e. } \sqrt{u_2 - u_1} = 0$$

~~Case-1~~ If $p\text{-value} = 0.9$ $H_0 - u_1$
Prob of 10cm is 0.9
provided that H_0 is true.

It means that assuming $\sqrt{u_2 - u_1} = 0$,
there is 90% probability that $u_2 - u_1 = 10\text{cm}$.

Hence, H_0 is accepted.

~~Case-2~~ If $p\text{-value} = 0.05$ It means
Prob of 10cm is 0.5%.
if H_0 is true.

In this case, if $H_0 = \text{true}$. (assumed)
But since prob of $H_0 - u_1 = 5\%$. (small)

H_0 is rejected.

Therefore

H_1 is accepted (taken into account)

(11)

intuition

using

Hypothesis testing coin toss & expt.^(a)

Ques: Given a coin, determine if coin is biased towards heads or not.

~~Basic probability~~ Biased towards head, $P(H) \geq 0.5$
 not " " " " " $P(H) = 0.5$

i) design exp:

flip a coin 5 times and
count #heads = X . [Test
(r.v) statistics]

ii) perform exp

f f f f f
 ↓ ↓ ↓ ↓ ↓
 H H H W W

$\boxed{X = 5 \rightarrow \text{Observation by exp}}$

so,

$P(X=5 \mid \text{coin is not biased towards}) = P(\text{Obs} \mid \text{H})$

↑
idea

Obs

↑
assumption

Null hypothesis (H_0)

H_0 : coin is not biased towards head.

$$P(X=5 \mid H_0) = \frac{1}{\binom{10}{5}} = \frac{1}{32} \approx 0.03 \approx 3\%$$

five heads in
five tosses coin is
not biased towards
head.

$$\left\{ \begin{array}{l} S f f f f \\ \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \end{array} \right.$$

$$P(H) = 1/2 = 0.5$$

$$\therefore P(X=5 \mid H_0) = 3\%$$

It means, there is 3% chance of getting
5 heads in 10 flips if the coin is not
biased towards head.

$$p\text{ value} \rightarrow P(\text{obs} \mid \text{assumption}) = 3\%$$

assumed small

According to rule of thumb if $< 5\%: H_0$
not considered

If $P(\text{obs} / \text{no}) < s/$

then

no may be incorrect

↓

alternation is not true

↓

reject $H_0 \rightarrow$ reject coin is not biased
towards head

↓

accept coin is biased towards
head.

i.e we accept "alternative hypothesis"

expt 2: flip coin 8 times.

count # heads = x & rest - stem & g.

Perform

expt

f f +

H H H

$X = 3$

Observation

$$P(\text{obs} / \text{assumption}) = \frac{1}{2^8} = \frac{1}{256} = 12.51\%.$$

no

2^3

8

→ coin is not biased towards head.

P-value
here, we accept H_0

~~(*)~~ Careful conclusion!

We must be careful in

- (i) design of expt
- (ii) $H_0 \rightarrow$ null hypothesis.
- (iii) design X.



so that
 $p(\text{obs} | H_0)$ is easy and
feasible.

~~H_0~~ ~~(*)~~ Resampling and permutation testing

~~(*)~~ (iii)

G_1

C_{12}

: [] : []
so [] so []

$$d_1 - d_2 = \Delta$$

ie ~~(*)~~ all data merged into JUNIOR ✓

Now we randomly sample

X Y

$$\textcircled{1} : \begin{bmatrix} \cdot \\ \cdot \\ \vdots \\ \cdot \\ \cdot \end{bmatrix} \quad \begin{bmatrix} \cdot \\ \cdot \\ \vdots \\ \cdot \\ \cdot \end{bmatrix}_{\text{so}}$$

$$u_1 - u_2 \rightarrow d_1$$

$$\textcircled{2} : \begin{bmatrix} \cdot \\ \cdot \\ \vdots \\ \cdot \\ \cdot \end{bmatrix} \quad \begin{bmatrix} \cdot \\ \cdot \\ \vdots \\ \cdot \\ \cdot \end{bmatrix}_{\text{so}}$$

$$u_1 - u_2 \rightarrow d_2$$

$$\textcircled{10k} : u_1 - u_3 \rightarrow d_{10k}$$

Sum d, d' 's in increasing order.
we get.

$$d'_1, d'_2, d'_3, \dots, \Delta, \dots, d'_{10k}$$

$\overbrace{\quad \quad \quad}^{n \cdot l}$

How much n.l.

By how much n.l. d's greater
than A.

Final may $n \neq$ gives p-value.

* most agrees

$$n \neq = 51$$

$$P = 0.05$$

Or

$$P = 0.01$$

: modified user

23

K8 Test for similarity of two distribution

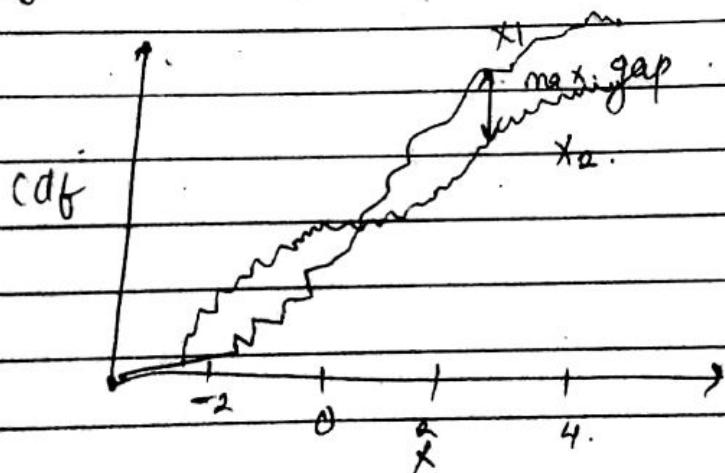
~~Defn:~~ To find out whether two random variables have same distribution or not.

Let X_1 and X_2 are two random variables with sample size 'n' and 'm' respectively.

[Hypothesis testing] :-

(i) Let the null hypothesis be
 $H_0: X_1$ and X_2 have same distribution.

Plotting CDF of X_1 and X_2



So by KS testing

Test statistic: $D_{n,m} = \sup |F_{1,n}(x) - F_{2,m}(x)|$

Where, $F_{1,n}(x) = \text{CDF of random variable } X_1$
 with 'n' sample.

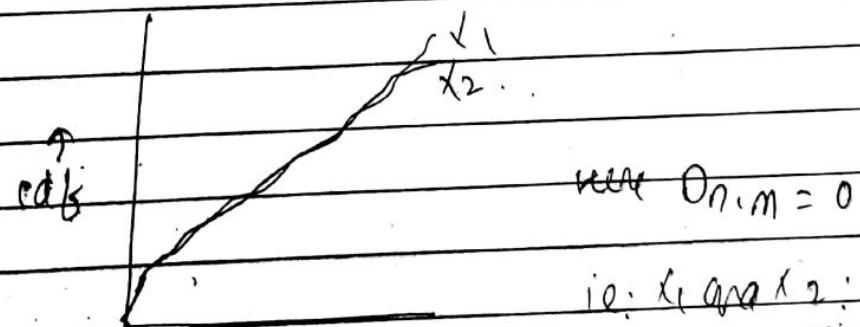
$f_{2|m}(x) = \text{cdf of random variable } x_2$
 when in sample

$D_{n,m} = \sup | \quad |$
 = supremum function representing
 max. gap b/w cdf of x_1 and x_2 .

in case

If $n \ggg$
 $m \ggg$,

then plot will be



i.e. x_1 and x_2 :
 same distribution

SG. w/o 11 true.

iii) Consider $\alpha!$ w/p representing significance level.

α	0.10	0.05	0.025	0.01
$c(\alpha)$	1.073	1.224	1.358	1.714

w/w $c(\alpha) \sqrt{\frac{1}{2} \ln \alpha}$.

(1)
(2) Now (3)

$$D_{\text{obs}} > C(\text{Cal}) \sqrt{\frac{n+m}{nm}} \quad (\text{number of } S)$$

(calculated
from graph)
Then H_0 is rejected.

else accepted.

g.

$$n = 1000$$

$$m = 8000$$

$$\text{for } \alpha = 0.05$$

$$C(x) = 1.224$$

$$\text{then } D_{\text{obs}} > 0.047 \quad (\text{after Lehmann})$$

True

Lehmann is rejected at 0.05 significance level.

i.e. X_1 and X_2 don't have same distribution.

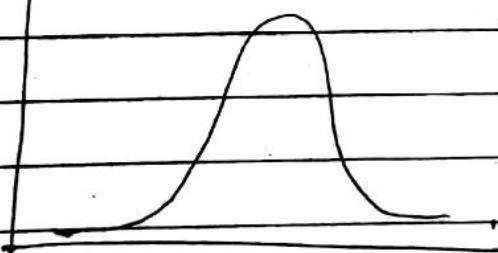
Code for KS testing

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

Code :- # generate gaussian at r, u, n

```
x = stats.norm.rvs (size=1000),  
# generates a normal dist with 1000 values  
sns.set_style ('whitegrid')  
sns.kdeplot (np.array(x), bw=0.5)  
plt.show()
```

O/P :



• stats.kstest (x, 'norm')

O/P :

KSTEST RESULT (statistic = 0.021..., p-value = 0.75...)

↑
norm

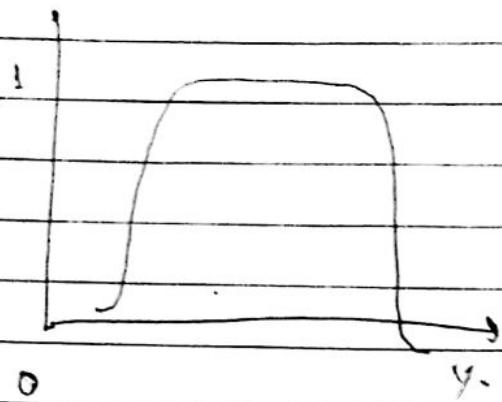
↑
 $c(\alpha) \sqrt{\frac{n+m}{n}}$

Since. $D(n,m) \subset C(\alpha)$ //

∴ X is a normal distribution

Task 2: Continuous uniform dist. (0,1)

$y \sim \text{ap-random-uniform}(0,1, 100)$
size: n=100 (np.array(y)) bw: 0.5
ct: snow()



stats::kstest(y, 'norm')

exp:
KStestResult(statistic=0.5..., pvalue=0.2)

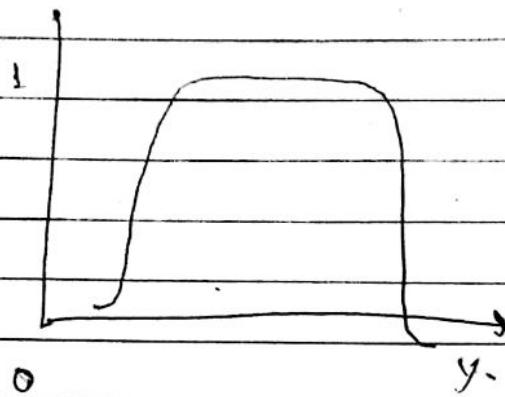
$$H_0: D_{n,m} > c(\alpha) \sqrt{\frac{n+m}{m}}$$

if H_0 rejected

if y is not a normal distribution

Task 2: continuous uniform dist (0,1)

$y = np.random.uniform(0, 1, 1000)$
 $y_{\sin} \sim \text{indep. CA}(\text{np.array}(y), bw=0.3)$,
plt.show()



stats.kstest(y, 'norm')

out:
KSTESTRESULT (statistic = 0.5..., pvalue = 0.0)

$$D_{n,m} \geq c(\alpha) \sqrt{\frac{\log m}{n}}$$

i.e. H_0 rejected

i.e. y is not a normal distribution.

(3:41)

(6:59)

(9:38)

(11:58)

(15:00)

7:28

(24) Hypothesis Testing continued.

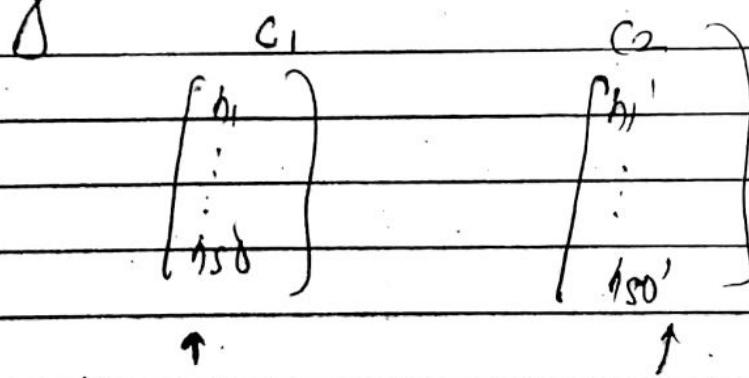
Another example :-

Q. ask:

determining if the pop. mean of heights
of people in city C_1 and C_2 is
same or not.

Step 1: Designing experiment :

not us randomly 1000 $\underline{50}$ sample from
each city. or.



sample height of 50 people in C_1 and C_2

$\therefore u_1 = \text{sample height in } C_1 = 162\text{cm} (\text{let say})$
 observed value $\{ \therefore u_2 = \dots \dots \text{, } C_2 = 167\text{cm} (\text{if})$

Step 2 Test statistic $112 - 41 = 2 = 167 - 162 = 5\text{cm}$

$$\begin{array}{c} \downarrow \\ c_2 \end{array} \quad \begin{array}{c} \downarrow \\ c_1 \end{array}$$

10

1 P

Step 3: Null hypothesis

H_0 : There is no difference in population means.

Step 4: Compute criterion

$$P(X = 5\text{cm} | H_0)$$

↓

We have to compute

Mean prob. of observing a difference in sample mean height of c_1 and c_2 if there is no popn. diff. in mean height

We will compute next

Let's assume first to make conclusion,

Step 5: Assumption & conclusion.

$$\text{Ques } P(X = 5 | H_0) = 0.2 = 20\% \text{ > S.I.}$$

i.e. There is 20% chance of observing dfn of 5cm in sample mean height of c_1 and c_2 (of size 50) if there is

no pop^{mean} diff

: $P(\text{obs} | \text{assumption}) = 20\% = \text{significant}$

Mean: assumption must be true
ie accept H_0

Qn 1. $P(X \geq S | H_0) = 0.03 = 3\% < 5\%$

plausible $P(\text{obs} | \text{assumption}) = 3\% (\text{very small})$

Now,

- assumption must be incorrect
- reject H_0
- Accept H_1

Now, for actual computation of p-value

We use something known as nonparametric

f permutation testing.

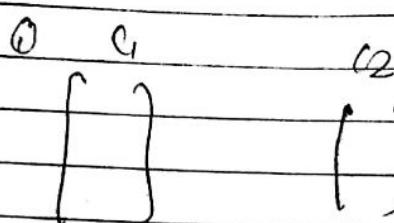
0:29 2:10 5:47 8:21 11:05 14:36 16:50

av-1. Running and Permutation Testing

$P(X \geq 5/n_0)$ → compute it.
i.e.

$P(\text{obs} / \text{assumption})$

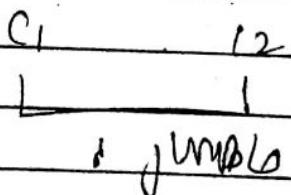
What we have



for sample so sample.

(2) $X = \mu_2 - \mu_1 = 5\text{cm}$

Step 1 :



$$S: \{h_1, h_2, \dots, h_n, b_1, b_2, \dots, b_s\}$$

Step 2 ~~randomly~~ $\rightarrow S_1 \{h_1, h_2, h_3, h_4, \dots\} = u_1$

~~randomly~~ ① $S \rightarrow S_2 \{h_3, h_4, h_2, \dots\} = u_2$

to satisfy H₀

$$\mu_2 - \mu_1 = 3\text{cm} > d_1$$

?

Resampling Repeated K times

(1) $u_1 - u_1 = 3\text{cm} = d_1$

Repeat (2) $u_2 - u_1 = -2\text{cm} = d_2$

repeat (3) $u_2 - u_1 = 1\text{cm} = d_3$

(K) $u_2 - u_1 = 6\text{cm} = d_K$

Step 3 : Sort d's.

$$d_1 \leq d_2 \leq d_3 \leq \dots \leq d_K$$

simulated differences \rightarrow increasing order

Case 2: Obs-dif $= x = 5\text{cm}$

$$P(d_{\text{diff}} \geq 5\text{cm}/60) = ?$$

Cut K=1000

Let say

$$d_1 \leq d_2 \leq d_3 \leq \dots \leq 5\text{cm} \leq d_{801} \leq \dots \leq d_{1000}$$

801 of sim. differences $\leq 5\text{cm}$ (obs. diff)

201 of sim. diff $\geq 5\text{cm}$ (obs. d.)

$P(\text{Obs. dip} \geq 5^\circ \text{ cm} | \text{assumption}) = 20\%$ \rightarrow Significant

∴ Assumption must be true
Accept H₀

Ans G

$$d_1' \leq d_{21} \dots d_{10}' \leq \text{sum of } g \leq d_{10}$$

9.1% 3.1%

$P(\text{Obs. dip} \geq 5\text{cm} | H_0) = 3\%$
may $\leq 5\%$.

∴ Assumption Chd must be incorrect

∴ Reject H₀

(11) 311 ~~200~~ 7:28 19:30 15:48 15:32 22:15

(25) How to use hypothesis testing ??

(1) KS test 1980 hypothesis testing \rightarrow to check if two random variables have some distribution or not.

Application.

(1) Assigning drug / medicine

Copy c₁



Cold) drug D₁

Copy c₂

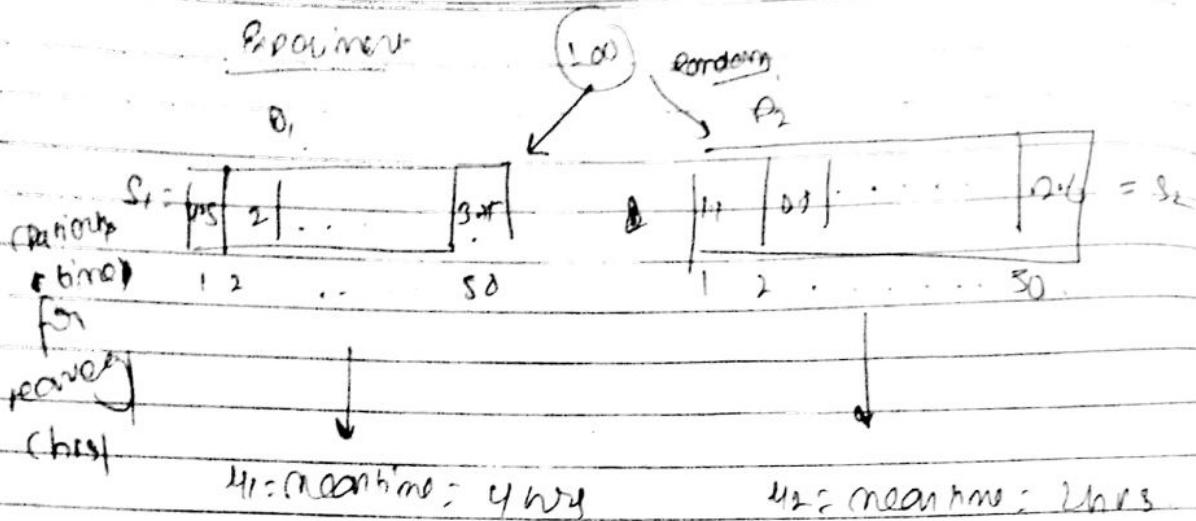


Drug D₂ (new).

Reduces fever in
4 hrs

Claim: Reduces fever faster
than D₁.

Task: experiment \rightarrow to determine if claim is true or not



To be sue about observation, let apply

Hypothesis testing

(i) $H_0: D_2$ and D_1 take same time to reduce fever

(ii) test statistic: $X = \bar{D}_2 - \bar{D}_1 = 2 \text{ hrs (obs)}$

(iii) $P(X \geq 2 | H_0) = 0.5 \text{ max} \approx 1.0 \text{ (assume)}$

↓
no rejected

mean D_2 has faster recovery than

↓
for actual computation

↓
faster resampling and permutation testing

In case of medicine

we must take in account the significance level (α)

Since it is matter of life and death
if $p\text{-val.} < \alpha$

reject H₀
else accept H₀

for medicine

$\alpha = 1\%$ or even 0.1% .

Normaly

$\alpha = 5\%$.

(ii) Economic α

similar w^g application

where $\alpha = 3\%$ or 5% .