

Probability & statistics

- fundamental area
- histogram, PDF, CDF, mean, var, std-dev

↳ prob & stat

→ random variable :-

eg:- dice : six-sided $\rightarrow \{1, 2, 3, 4, 5, 6\}$
 experiment \rightarrow roll a ^{fair} dice ↳ Equally likely

r.v $\rightarrow X = \{ \quad \}$
 ↳ it can be anything $\{1, 2, \dots, 6\}$

→ tossing coin :-

r.v $\rightarrow Y = \{ H, T \}$
 ↳ Equally likely

→ always read eqn and code in English.

→ dice

* $P(X \text{ is even}) = \text{prob. that } X \text{ is even} = \frac{1}{2}$

Types of random variable :-

• Discrete random value :-

→ eg:- dice $\rightarrow X = \{1, 2, \dots, 6\}$

↳ a r.v which have set of outcomes

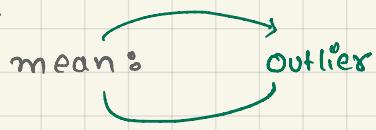
↳ a r.v can take one value in a finite set of values is called d.r.v.

• Continuous r.v :-

eg:- height of randomly picked student $\rightarrow 120\text{cm} \quad 190\text{cm}$

↳ here the set is not finite it's continuous so, C.r.v.

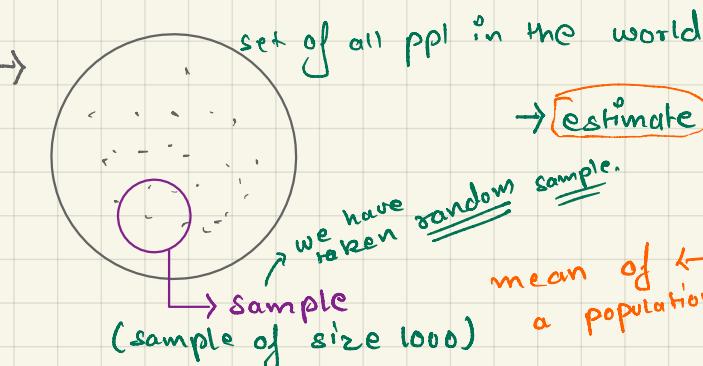
Outlier :-



Variance :-

MAD

Population & sample :-



→ → estimate the average/mean height of a human

mean of
a population

→ ←

mean of
a population

→

\bar{M}

$$\bar{M} = \frac{1}{TB} \sum_{i=1}^{TB} h_i \rightarrow \text{we can't find that.}$$

→ a "sample" is a subset of population.

mean of a sample → $\bar{h} = \text{Mest.} = \frac{1}{1000} \times \sum_{i=1}^{1000} h_i$

heights in my sample.

→ As sample size inc

$$\bar{x} = x$$

our sample mean will converge towards population mean.

Gaussian or Normal distribution and its PDF

PDF of a Gaussian dist r.v

→ dist. are very simple models for real world.

→ If you know the PDF you can plot CDF.

Let's assume :-

X : continuous r.v

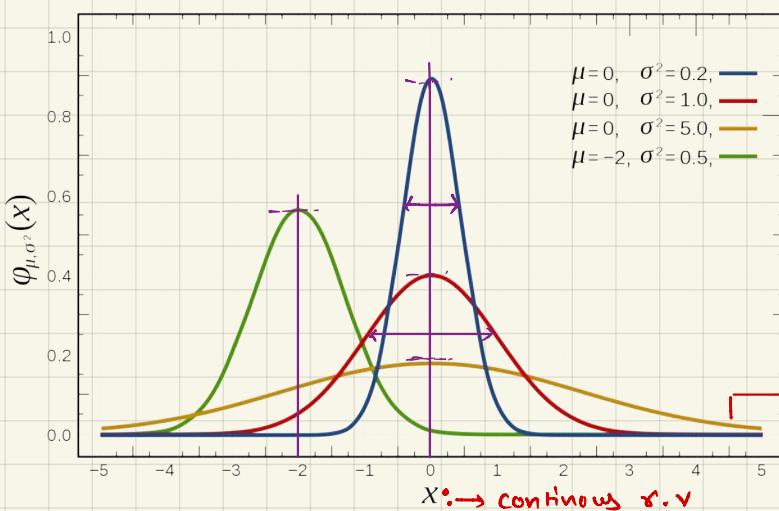
and

X PDF is look like "bell shaped"

then we will say 'X' has a distribution which is a gaussian dist.

→ why should we learn about this dist?

→ because lot of things in nature follow this dist. (not everything)
for eg:- P.L (petal-length) follow gaussian distribution.



→ If you tell me my 'X' (continuous r.v) follows Gaussian dist and if you give me it's mean (μ) and variance (σ^2) nothing else, just these two. Then I can tell you what it's PDF is. (without having one value/observation)

→ And if I know PDF, I can easily plot CDF.

If PL is G.D

then μ, σ^2

parameters of a Gaussian distb: μ, σ^2

Normal/gaussian

$$X \sim N(\mu, \sigma^2) \rightarrow X \text{ follows gaussian distb}$$

↓ ↓ ↓
follows mean variance

→ probability of finding x is this

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

Let assume,

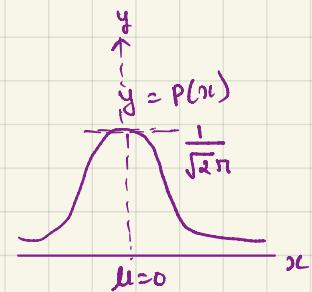
$$\mu = 0 \quad \& \quad \sigma^2 = 1$$

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

$$y = e^{-x^2}$$

if $x_c \uparrow$ then $y \downarrow \rightarrow$ the rate of reducing "y" is 'square' of 'exponent'

if $x_c \downarrow$, then $y \uparrow$



Conclusions

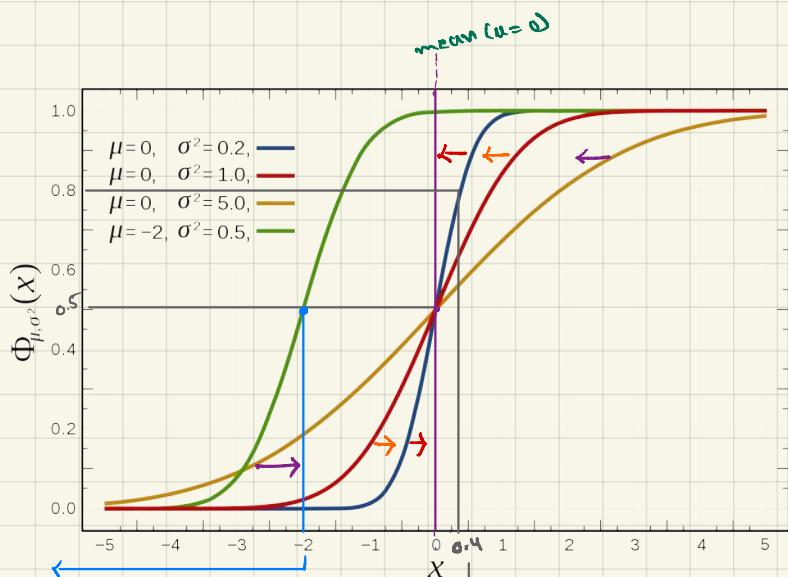
- As x_c moves away from μ ; $y \downarrow$
- symmetric along μ . \perp
- x_c moves away from μ ; 'y' reduces $\exp(-x^2)$

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

$$x_c = 0 \quad y = 1$$

$x_c = 1$	$y = \exp(-1)$	$= 0.3678$	$\rightarrow 20x \rightarrow \frac{0.018}{0.36}$
$\downarrow 2x$			
$x_c = 2$	$y = \exp(-4)$	$= 0.018$	$\downarrow 100x \downarrow$
$\downarrow 1.5x$			
$x_c = 3$	$y = \exp(-9)$	$= 0.000123$	

CDF



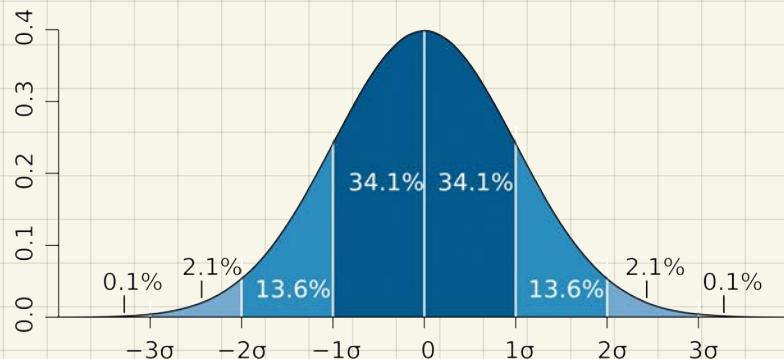
→ as σ increases our graph goes more away from mean.

→ Because our PDF (shown above) is symmetric our CDF is also symmetric. So, 50% our pt. lies one side of mean & 50% of pt. lies on other side.

(y-axis of CDF) → (means prob. of x having less than equal to 0.4 is 80%)

Std. deviation (σ)

→ If we know our $X \sim$ follows Gaussian dist' and the value of mean & std. deviation is known to us then what can we predict?



A plot of normal distribution (or bell-shaped curve) where each band has a width of 1 standard deviation –

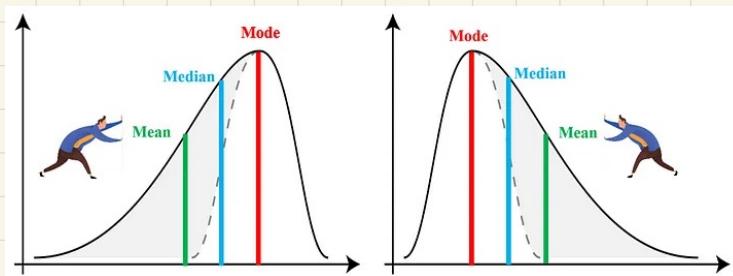
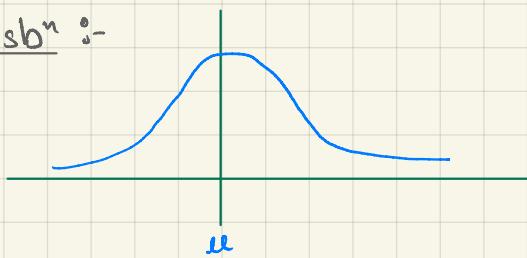
See also: 68–95–99.7 rule.

→ If we know $X \sim N(\mu, \sigma^2)$

$\mu = 0$
 $\sigma^2 = 4 \rightarrow \sigma = \pm 2$ $\xrightarrow{\text{obs.}}$ mean std. devi
 b/w 0 & 2 our 34.1% pt. lies.
 b/w -2 & 2 our 68.2% pt. lies.
 b/w -4 & 4 our 95% pt. lies

Symmetric, skew, and kurtosis distribution

Symmetric disbⁿ :-



$\mu \rightarrow$ mean \rightarrow location or middle

$\sigma^2 \rightarrow$ spreadness

Skew disbⁿ :- It tells you how dissimilar is the disbⁿ from symmetric disbⁿ.



There are four mathematical measures of relative skewness.

(a) Karl Pearson's Coefficient of Skewness

$$SK = \frac{\text{Mean} - \text{Mode}}{\text{S.D.}}$$

If mode is ill-defined, then we take

$$SK = \frac{3(\text{Mean} - \text{Median})}{\text{S.D.}}$$

(b) Bowley's Coefficient of Skewness

This is based on quartiles and median and is defined as

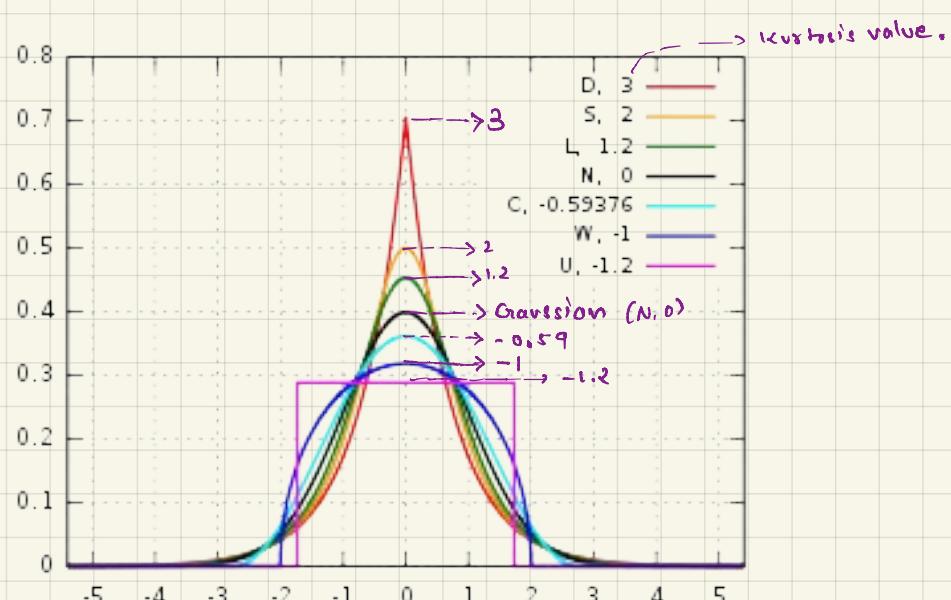
$$SK = \frac{Q_3 + Q_1 - 2 \text{Median}}{Q_3 - Q_1}$$

This formula is useful when the mode is ill-defined or the distribution has open end classes or unequal class-intervals.

If this value is between:

1. -0.5 and 0.5, the distribution of the value is almost symmetrical
2. -1 and -0.5, the data is negatively skewed, and if it is between 0.5 to 1, the data is positively skewed. The skewness is moderate.
3. If the skewness is lower than -1 (negatively skewed) or greater than 1 (positively skewed), the data is highly skewed.

Kurtosis :- How peaked your distribution is. (it measures that)



Standard normal Variate (z) and standardization :-

→ If my r.v (x) follows normal dist' and I want to standardize it because after standardization I can say that my 68% data lie b/w $\rightarrow -1 \pm 1$

having dim. 'n',
x (feature) follows ND & -----
or
95% data lie b/w $\rightarrow -2 \pm 2$

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

↳ How to read that $\rightarrow X$ (r.v) follows N.D & have mean (μ) & variance (σ^2)

$$X = [x_1, x_2, \dots, x_n] \quad (\text{assume})$$

all this are observation of this r.v.

Standardization \rightarrow method in which I can transform any mean & variance into a std. data.

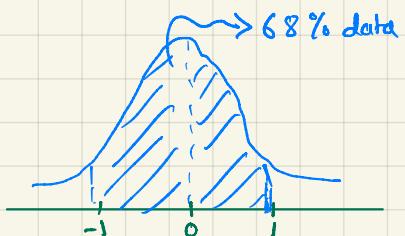
$$x_i' = \frac{x_i - \mu}{\sigma}$$

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

↳ standard Normal variate

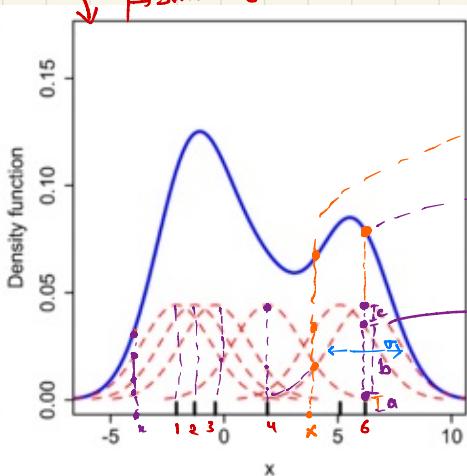
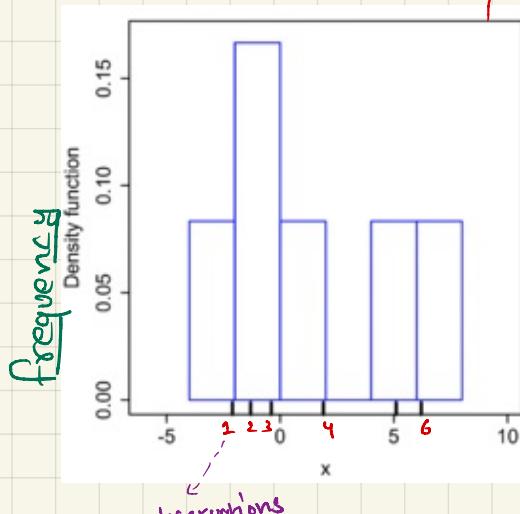
$$\Rightarrow X \sim N(\mu, \sigma^2) \xrightarrow{\text{standardization}} X' \sim N(0, 1)$$

$$\text{or} \quad Z \sim N(0, 1) = \frac{X - \mu}{\sigma}$$



Kernel density estimation :- a way of converting histogram into PDF.

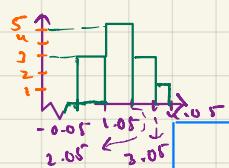
we can do smoothing using Kernel
↳ smooth form of histogram.



→ I want to calculate height of PDF at 'x'? \rightarrow Just by adding the kernels.
→ I calculated this value by adding $a + b + c$.
→ I have plotted a bell-shaped Kernel or gaussian Kernel here with center or mean at 6th point.

→ We have taken 'mean' in consideration what about Variance? \rightarrow The variance here is also called bandwidth of Kernel

Bandwidth



class limit $\rightarrow 0-1, 1-2, 2-3, 3-4$
 class boundaries $\rightarrow -0.05-1.05, 1.05-2.05, 2.05-3.05, 3.05-4.05$

e.g. assume

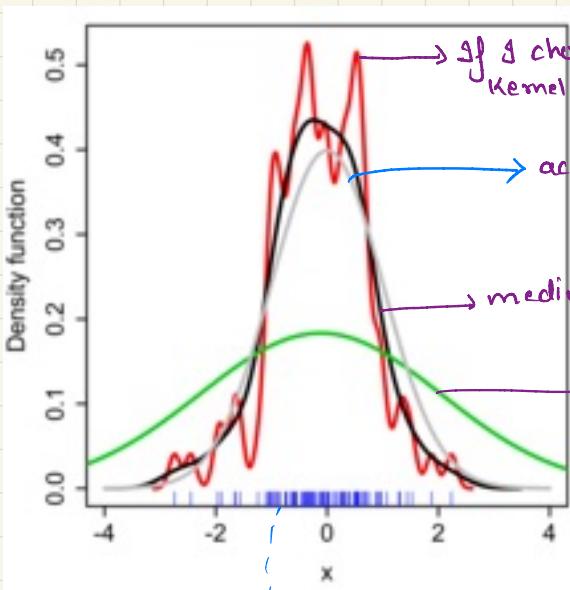
pr: {0, 0.5, 1, 1.2, 1.3, 1.4, 1.8, 2, 2.1, 2.8, 3.4}

→ histogram having bins ≈ 4

Since, our data points are integer, we can do to set the class boundaries subtract or add by 0.05

<u>Class boundaries</u>	<u>frequency</u>
-0.05 - 1.05	3
1.05 - 2.05	5
2.05 - 3.05	3
3.05 - 4.05	1

$\left\{ \begin{array}{l} 4 \text{ bins} \\ \text{observations} \end{array} \right.$



If I choose my b.w 'small' then my Kernel is very narrow or jaggeded (histogram with lots of bin).

→ actual distn.

→ medium b/w.

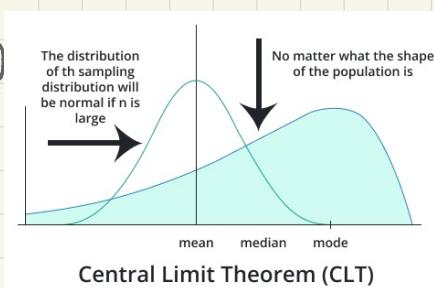
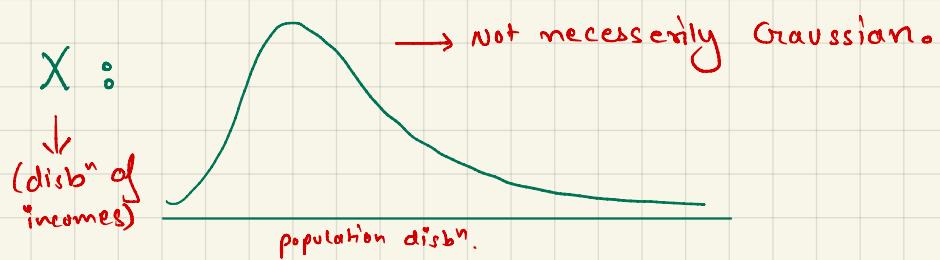
→ very large b/w.

→ observations.

→ To find the write b.w what does 'seaborn' do is?

They choose a very small 'b.w' initially & see how much jaggeded is our PDF is and then they inc' the 'b.w' till the time the line become smooth enough.

Sampling distn and central limit theorem (CLT)



I took independently.
 It doesn't depend on one another

→ I take 'random sample' of size $n=30$ $\rightarrow S_1 \rightarrow \bar{x}_1$

again I take 'random sample' of size $(n=30) \rightarrow S_2 \rightarrow \bar{x}_2$

⋮

Sample mean
 \uparrow
 $\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_m \rightarrow \bar{x}_m$

→ $\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_m \rightarrow m$ sample means

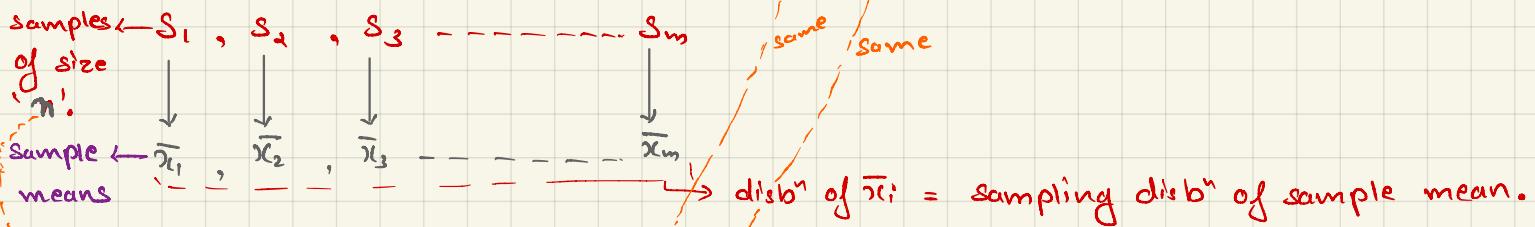
$\rightarrow \bar{x}_i \sim \text{distb}$

\hookrightarrow distb of sample means

\rightarrow distb of $\bar{x}_i = \text{(Sampling distb)} \text{ of } \text{(sample-mean)}$

CLT (only for finite mean & variance)

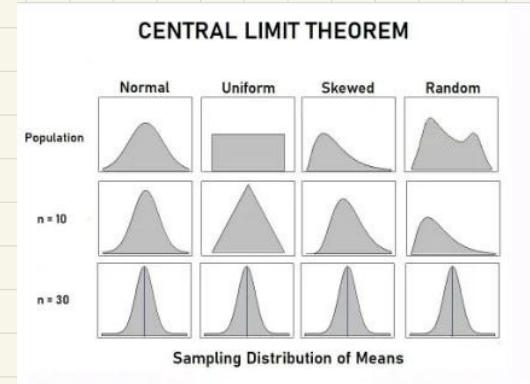
If X : (population distb): have finite μ & σ^2



CLT^o-

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

\uparrow gaussian distb



\Rightarrow e.g:- $X: \mu, \sigma^2$ (income)

any distb
need not be Gaussian.

Let $n=30$

mean income in each sample

$$\mu \approx \text{mean of } \bar{x}_i \text{'s} = \bar{x}_u$$

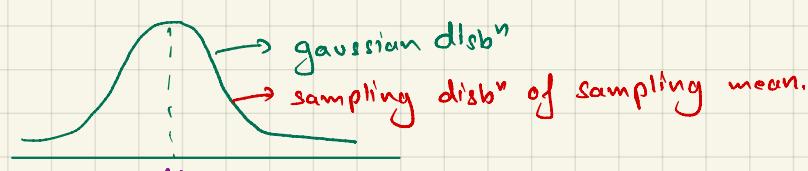
$$\frac{\sigma^2}{n} \approx \text{variance of } \bar{x}_i \text{'s} \approx \sum_{i=1}^m \frac{(\bar{x}_i - \bar{x}_u)^2}{m}$$

\hookrightarrow rule of thumb $\rightarrow n \geq 30$

Let,
 $m = 1000$

$$\frac{1000 \times 30}{m} = 30K \text{ (people income)}$$

$$\rightarrow \bar{x}_u = \frac{\bar{x}_1 + \bar{x}_2 + \dots + \bar{x}_m}{m}$$



Note:- By looking at '30K' data points / observations. I am able to estimate the whole population $\xrightarrow{\text{mean}}$ $\xrightarrow{\text{variance}}$

\rightarrow To understand the population mean or variance (μ, σ^2) of any distb (not necessarily gaussian) I just need to know they are finite or well defined.

Quantile-Quantile (Q-Q) plot:-

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

Q) Is 'X' Gaussian dist? → \rightarrow QQ Plot (graphical method)
 ↗ statistical testing (KS, AD)

Steps for Q-Q plot

→ sort x_i 's & compute percentiles

x_1, x_2, \dots, x_{500}

↓ sort (asc.)

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

↓ Percentiles

① ② ③ ⑩ → 100th percentile
 $x'_5, x'_{10}, x'_{15}, \dots, x'_{500}$

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

1st- percentile
 value of x_i 's

→ std. Normal / gaussian distn

→ $Y \sim N(0, 1)$ → (create this)

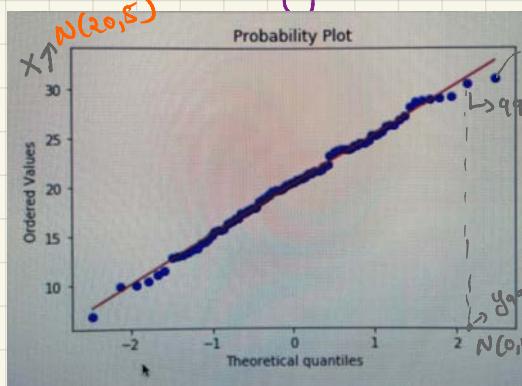
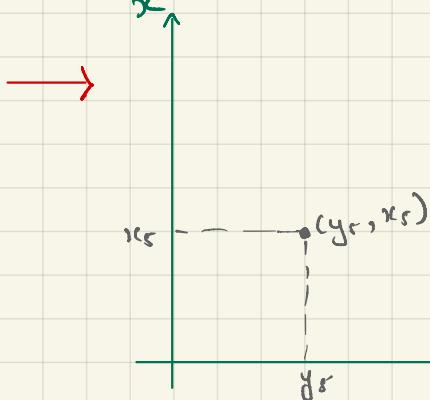
$y_1, y_2, \dots, y_{1000} \rightarrow 1000$ obs. from $N(0, 1)$

↓ sort (asc.)

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

↓ Percentiles

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



$(y^{(i)}, x^{(i)}) \forall i: 1 \rightarrow 100$

$y^{(i)}, x^{(i)} \forall i: 1 \rightarrow 100$
 lie on a st. line

then X & Y have
 a similar distn.

⇒ Q-Q plot

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

→ Theoretical quantiles

Q-Q Plot

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

# N(0,1)
std_normal = np.random.normal(loc = 0, scale = 1, size=1000)

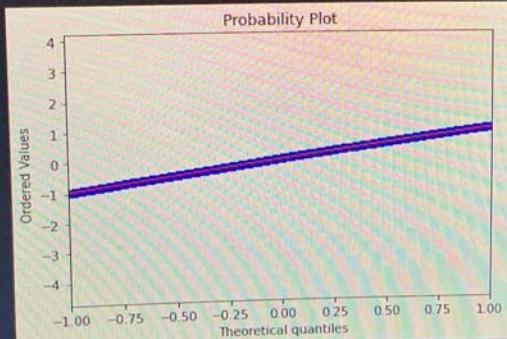
# 0 to 100th percentiles of std-normal
for i in range(0,101):
    print(i, np.percentile(std_normal,i))
```

```
0 -2.9941560608546496
1 -2.331155489997273
2 -2.102462010297869
3 -1.9030070976001443
4 -1.7878574029822945
5 -1.6383831947251253
6 -1.5935117657666127
7 -1.5315554277910335
8 -1.4362998197560297
9 -1.394737614240833
10 -1.3132161092517054
94 1.4626263448353378
95 1.5608572685665627
96 1.6966233988361676
97 1.9147750543381148
98 2.0157626321951363
99 2.246560915109006
100 2.7918505373969498
```

```
[ ] import matplotlib.pyplot as plt

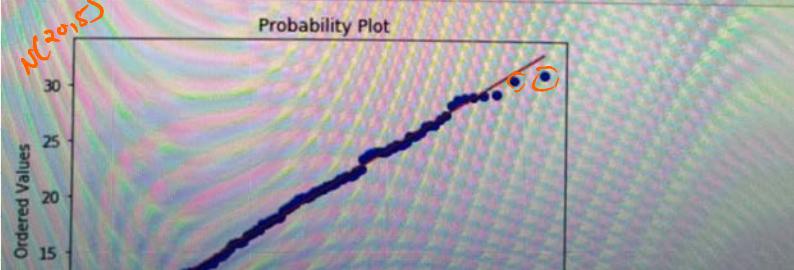
# generate 100 samples from N(20,5)
measurements = np.random.normal(loc = 0, scale = 1, size=10000)
#try size=1000

plt.xlim(-1,1) Q-Q plot
stats.probplot(measurements, dist="norm", plot=pylab)
pylab.show()
```



```
# generate 100 samples from N(20,5)
measurements = np.random.normal(loc = 20, scale = 5, size=100)
#try size=1000

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



→ generate random no. from a normal/gaussian dist^b having mean → 0 & std. dev. → 1 & generate '1000' observation/sample. and put all value in std-normal.

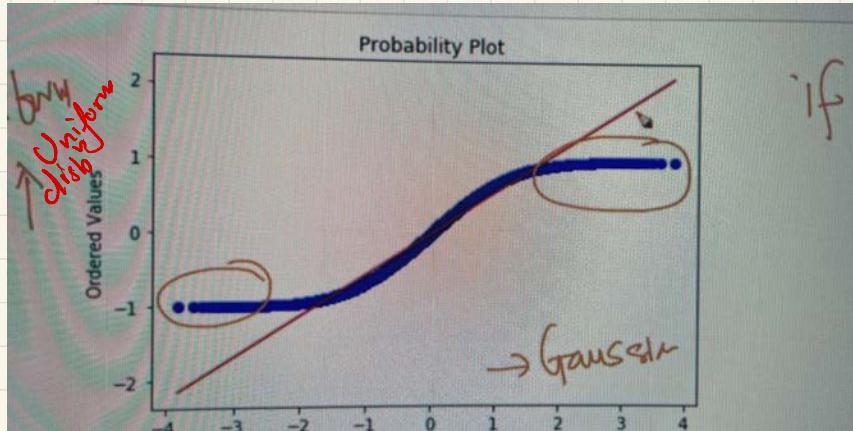
my theoretical
value goes from
(-3 to 3) approximately

use this library for plotting.
→ it says I want to compare it with 'norm' or std-normal variate.

→ if we incⁿ size → 1000
↳ we will see more & more pt. lie on std line.

→ if we decⁿ size → 50
↳ we will see more & more deviation from line.

∴ If no. of observation/sample is small,
it's hard to interpret Q-Q plot.



- suppose if we generate data from "uniform" disb
- if we incⁿ size here from 100 → ¹⁰⁰⁰
↳ we will see more deviation
- if we decⁿ size here 100 → 80
↳ we will see our deviation is less. [Just opp. the above concept]

Q-Q plot

- is $x \sim N(\mu, \sigma^2) \rightarrow$ is x coming from Normal disbⁿ
 - is $x, y: r.v \left\{ \begin{array}{l} \text{does } x \& y \text{ have} \\ \text{the same disb?} \end{array} \right.$
- } we can answer both by using Q-Q plot.

How/where to use distributions?

→ r.v, pdf, cdf, Gaussian → 68-95-99.7 rule

probability → data analysis → answering questions about data

Q) Company → xyz

Task & order t-shirts for all employees → 100K
size S, M, L, XL

a) How many 'XL' t-shirts should you order?

- ① collect data for all 100K employees,
- ② typically ppl. height $> 180\text{cm} \rightarrow$ XL t-shirt

domain knowledge.

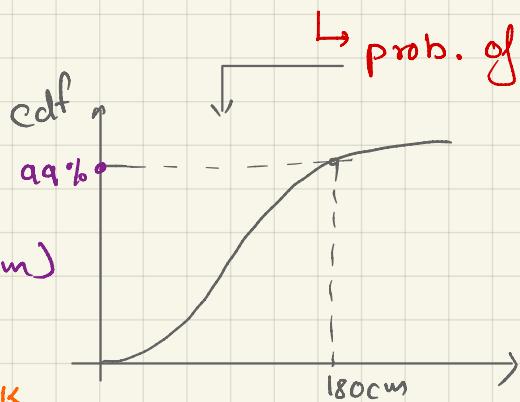
$[160\text{cm}, 180\text{cm}] \rightarrow$ L t-shirt (100K)

here I assume my sample (μ & σ) are approx \approx population (μ, σ)

→ collect heights of 500 random employees

$\frac{\text{mean, std. dev.}}{\mu}$

→ somebody told me heights $\sim N(\mu, \sigma)$ (doctor told me)



↳ prob. of (height $> 180\text{cm}$) that I can calculate

↳ how? → I know, $\mu + \sigma$ and I know it follows normal distn then I can plot CDF & PDF.

so, prob. (height $\geq 180\text{cm}$)

is 1% $\rightarrow 80$,

$$100K \times 1\% = 1K$$

↳ I will order '1K' of 'XL' t-shirt for employee.

assumptions

- i) you can predict size of t-shirt from heights.
- ii) height is a Gaussian distn & I am finding μ, σ for sample.

→ Gaussian distn → Theoretical model

↓
observed in many natural phenomena.

Q) Salaries $\xrightarrow{(a)}$ of all ppl. (100K)

↳ if $s \sim N(\mu, \sigma)$

using CDF I can answer that.

↳ how many employees make a salary $> 100K$?

↳ " " " " " " " " $\$50K \leq s \leq \$70K$

→ How do you know your s (salary) → Gaussian?

↳ Q-Q test (or other technique)

⇒ If salaries are not gaussian distn we can't use CDF.

→ Gaussian distn → give us theoretical model.

↳ to check that we can use 'Q-Q' plot.

Chebyshov's Inequality → helps us to find the outliers for a feature.

X : Gaussian - dist \rightarrow 68 - 95 - 99% rule

$X \leftarrow$ heights of students

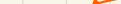
→ If we know ①, ② & ③

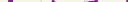
$$X \sim N(\mu, \sigma^2)$$

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

assume, $\mu = 180$, $T = 10$

→ 95% of students heights lie in the [130, 170]

Q) What if I don't know the dist'n.
 we can al

But I know (μ, σ)  We can also calculate this by using "CLT".

\downarrow ↳ 'non-zero' and 'finite'
finite

→ Can I say: $x\%$ of data lies within $\mu - 2\sigma$ & $\mu + 2\sigma$

y% of $\mu - 1.5\sigma$ to $\mu + 1.5\sigma$

Can I find $x\%$ & $y\%$?

why?

Salaries of individuals → don't know the distⁿ

(μ , σ) \longrightarrow (I know)

→ \$10K
→ \$40K

no -
P

Q) What %age of individual have salary in the range of [20K, 60K]?

→ Here, comes the concept → Chebyshev's inequality:

$\rightarrow X$: r.v finite mean = μ

non-zero & finite std-dev = +

don't know
the dish".

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

$$\begin{aligned} X &\geq \mu + k\sigma \\ X &\leq \mu - k\sigma \end{aligned}$$

→ Let, assume our distn look like that

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

OR

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

Let's go to Ques → b/w salary $\rightarrow [20K, 60K]$

$$20K = \mu - 2\sigma$$

$$40K = \mu$$

$$60K = \mu + 2\sigma$$

from chebyshev's inequality

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

here $k=2$

$$P(20K < X < 60K) \geq 1 - \frac{1}{2^2}$$

$$P(20K < X < 60K) \geq 0.75$$

→ it means, atleast 75% ppl. salaries lie in this region.

in gaussian distn

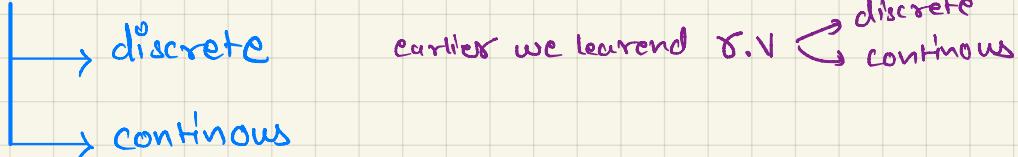
→ In this we can't say accurately like 95% ppl salaries is lie in that region. S can only give "range".

a) $\rightarrow [10K \leq X \leq 70K] \geq 1 - \frac{1}{9} \approx 90\% \rightarrow$ greater than 90% ppl. salaries lie in that region.

<https://www.learndatasci.com/glossary/chebyshevs-inequality/>

link to understand more.

Uniform distribution \rightarrow here uniform means, the value of 'x' w.r.t 'y' is const.

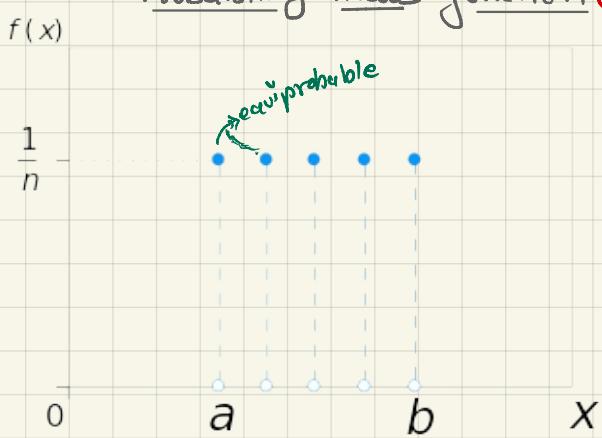


Discrete uniform distb

PDF \rightarrow continuous r.v

Pmf \rightarrow discrete r.v

Probability mass function (Pmf)



notation :- $U(a, b)$

parameter :- $a \in \text{integer}$
 $b \in \text{integer}$
 $b \geq a$

of outcomes $\leftarrow n = b - a + 1$
 we have.

PDF :- $\frac{1}{n}$

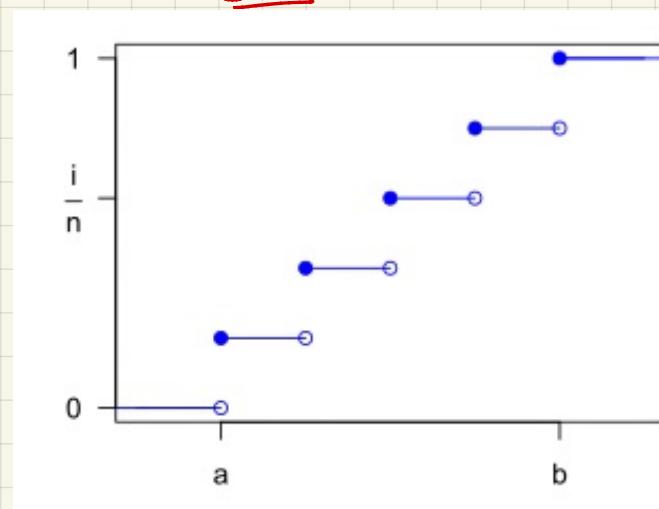
e.g.: In a dice $\rightarrow a = 1$

$\downarrow b = 6$

$\downarrow n = 6 - 1 + 1 = 6$ \rightarrow probability of every outcome $= \frac{1}{n} = \frac{1}{6}$
 ↑
 total outcomes

\Rightarrow if you know the parameters you know all the information about the distribution :- $N(\mu, \sigma^2) \rightarrow$ gaussian

CDF $U(a, b) \rightarrow$ uniform



\rightarrow non-smooth funt