

Distributions

1 / 1

- A data distribution refers to the way in which data values are spread or distributed across a dataset.

Key characteristics of a Data Distributions :

- ① Cardinality (unique) and unique (for discrete columns)
- ② Count, Min and Max (For Numerical columns)
- ③ Central Tendency : This indicates the "center" of the data. Common measures of central tendency include the mean (average), median (middle value), and mode (most frequently occurring value).
- ④ Dispersion or Spread : This describes how spread out the data points are. Measures of dispersion include Range, variance, and Standard deviation.
- ⑤ Shape : This refers to the overall pattern or form of the distribution. Common shapes include normal (bell-shaped), skewed (lopsided), and uniform (Evenly distributed).

Data Distributions

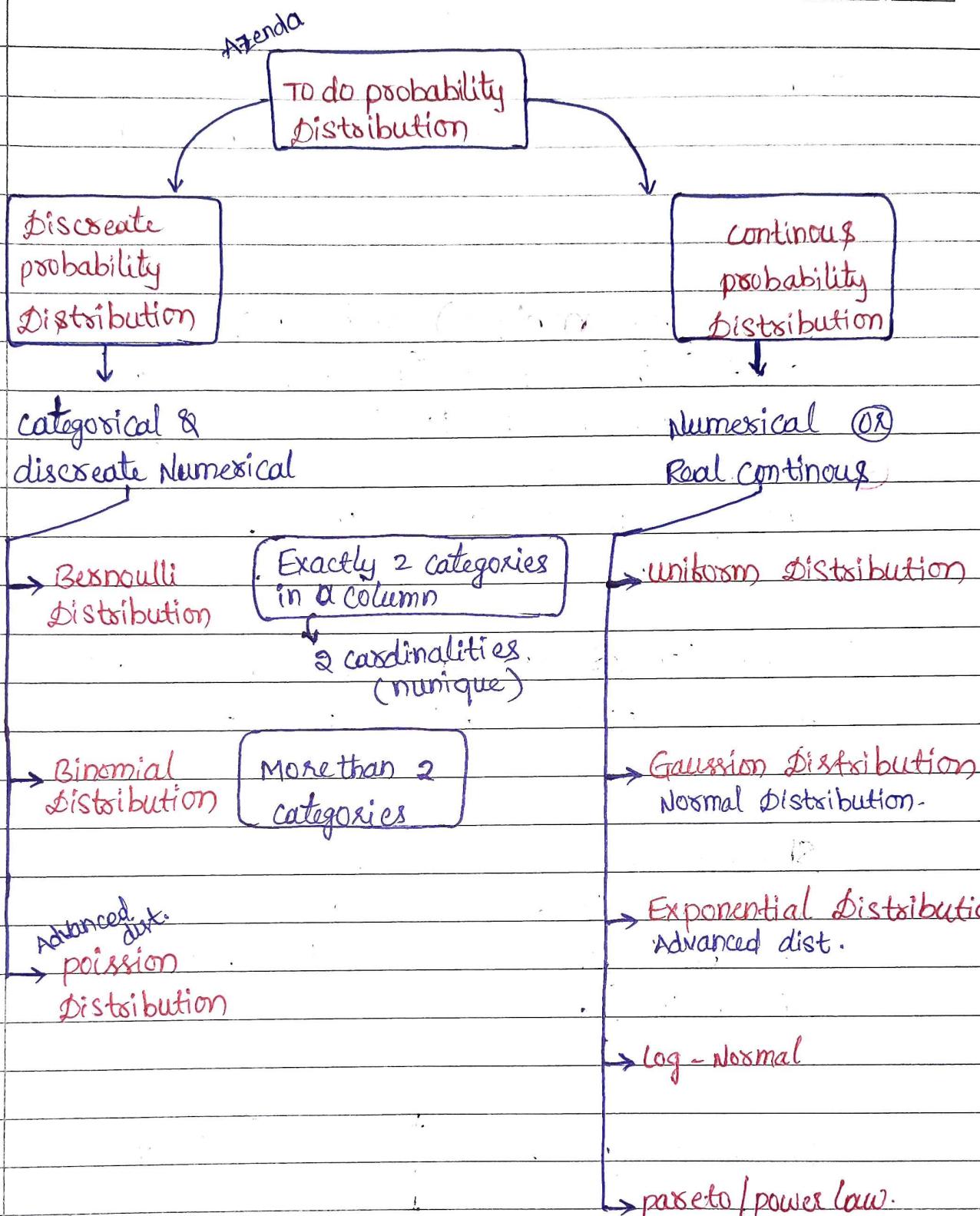
How to Analyse the data distribution ?

Non-visual

Understand the properties of distribution

visual

Shape, Spread, Outliers



why study Data Distributions ?

- * Understanding the distribution of data is crucial in a wide range of fields, including Statistics, Data Science, Economics, Finance, Health care, and more.
- * It informs decision-making, helps in selecting appropriate statistical tests $\&$ models, and provides insights into the underlying patterns and characteristics of a data set.
- * Following are the benefits of studying data distributions :

① Outliers and Missing values : Outliers are data points that are significantly different from the rest of the data. Knowledge of Data Distribution helps perform a better outlier and missing value treatment.

② Advance Data Analysis : Knowledge of data distribution and their PMFs / PDFs can help us with better insight. Generation from the given data and take data driven decisions.

③ Intersessional statistics : Knowledge of Distribution can help us select the appropriate Statistical Tests.

④ ML and DL : It also helps us identify appropriate Algorithms for Data Modelling in ML and DL .

Discrete Distributions

① Bernoulli Distribution: (cardinality = 2)

The Bernoulli distribution models a random experiment with two possible outcomes, often labeled as "Success" and "Failure". It is characterized by a single parameter, p , which represents the probability of success.

For eg :

Loan Approval - Suppose a bank is considering a loan application. The outcome can be either "Approved" (success)

or "Rejected" (failure) based on certain criteria.

Let's say probability of loan approval (p) is 0.8.

- What is the probability of Rejection?

$$\text{Ans: } 1-p = 1-0.8 = 0.2.$$

- Data follows Bernoulli distribution, we represent like:

$$X \sim B(1, p)$$

- probability Mass Function:

$$f(x) = \begin{cases} p & \text{if } x=1 \\ 1-p & \text{if } x=0 \end{cases}$$

- code:

```
def bernoulli_distribution_generator(p, size):
```

define the no. of trials (n) and probability of success (p)

```
return np.random.binomial(n=1, p=p, size=size)
```

```
bernoulli_dist = bernoulli_distribution_generator(p=0.3, size=10000)
```

```

from scipy import stats
# using scipy.stats to generate the Bernoulli distribution
prob_success = 0.30
var = stats.bernoulli.rvs(p=prob_success, size=10,000)

```

calculating PMF:

```

print("F(k=0):", stats.bernoulli.pmf(k=0, p=prob_success))
print("F(k=1):", stats.bernoulli.pmf(k=1, p=prob_success))

```

② Binomial Distribution:

- * The Binomial distribution models the number of successes (x) in a fixed number of independent Bernoulli trials (i.e. n), each with the same probability of success (p).

For Eg:

- Number of loan Approvals a day - Assuming the Bank receives 20 loan applications in a day (i.e; n), and each application has an independent probability of approval of $p = 0.8$.

- what is the probability of approving exactly 15 loans?

Ans:

- Data follows Binomial distribution:

$$X \sim B(n, p)$$

- Probability Mass Function:

$$f(x) = \binom{n}{x} (p)^x (1-p)^{n-x}$$

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

Code:

```
def binomial_distribution_generator(n, p, size):  
    # define the no. of trials (n) & the probability of success (p)  
    return np.random.binomial(n=n, p=p, size=size)
```

```
binomial_dist = binomial_distribution_generator(n=10, p=0.3, size=10000)
```

using Scipy.stats to generate the Binomial distribution

n_trials = 10

prob_success = 0.30

```
rvs = stats.binom.rvs(n=n_trials, p=prob_success, size=10000)
```

computing PMF & CDF:

"probability of 0 success out of 10 trials:" stats.binom.pmf(k=0, n=n_trials, p=prob_success)

"probability of 1 success out of 10 trials:"

```
stats.binom.pmf(k=1, n=n_trials, p=prob_success)
```

"probability of 2 success out of 10 trials:"

```
stats.binom.pmf(k=2, n=n_trials, p=prob_success)
```

"probability of <= 2 success out of 10 trials:"

```
stats.binom.cdf(k=2, n=n_trials, p=prob_success)
```

③ poisson Distribution :

- * The poisson distribution models the number of events that occurs in a fixed interval of time or space, assuming events occur at a constant average rate (λ) and are independent of the time since the last event.

For eg :

- Customers arrivals at a Bank Branch Every hour -
Suppose, on average, a bank branch receives 30 customers per hour. Here $\lambda = 30$ (i.e; Average no. of customers arriving per hour)
- What is the probability of 40 (i.e; x) customers arriving in the next hour?

Ans :

- Data follows poisson distribution,
 $x \sim p(\lambda)$

+ probability Mass function :

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

- Typically, poisson distribution has following characteristics :

- ① It should exhibit a right-skewed, unimodal shape.
- ② Variance is approximately equal to the mean.

code:

```
def poisson_distribution_generator(lam, size):
```

define expected no. of events occurring in a fixed-time interval (lam).

```
return np.random.poisson(lam=lam, size=size)
```

poisson_dist = poisson_distribution_generator(lam=30, size=10,000)

using Scipy.stats to generate the poission Distribution

lambda_ = 30

var = stats.poisson.rvs(mu=lambda_, size=10,000)

computing PMF and CDF

"probability of 20 events happening, given that on an average 30 events happen:"

stats.poisson.pmf(k=20, mu=lambda_)

"probability of 25 events happening, given that on an average 30 events happen:"

stats.poisson.pmf(k=25, mu=lambda_)

"probability of 35 events happening, given that on an average 30 events happen:"

stats.poisson.pmf(k=35, mu=lambda_)

"probability of <= 35 events happening, given that on an average 30 events happen:"

stats.poisson.cdf(k=35, mu=lambda_)

Multinomial Distribution :

- * If you have a discrete column in your data set with more than two distinct values (i.e; a cardinality greater than 2), the distribution of this column would be a specific type of discrete probability distribution known as a Multinomial Distribution.
- * In a multinomial distribution, each observation falls into one of several categories, with each category having a specified probability.
- * In your case, with three distinct values, you have 3 categories.

+ Conclusion :

These examples illustrate how each of these distributions can be applied in different scenarios to model random processes with specific characteristics.

- * The Bernoulli distribution can model binary outcomes like pass/fail.
- * The Bernoulli distribution can model the number of successes in a fixed number of trials (E.g; success in a tutoring program)
- * The Poisson distribution can model the occurrences of events (E.g; student arrives at a library) over a specific time frame.

Continuous Distribution

① Uniform Distribution:

* In a uniform distribution, all outcomes are equally likely within a specified range. For, E.g;

- Rolling a Fair Six-Sided Die - When rolling a fair six-sided die, each face has an equal probability of $1/6$. This is an example of a discrete uniform distribution.

- Data follows uniform distribution;

$$X \sim U(\alpha, \beta)$$

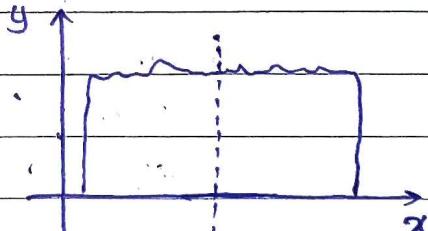
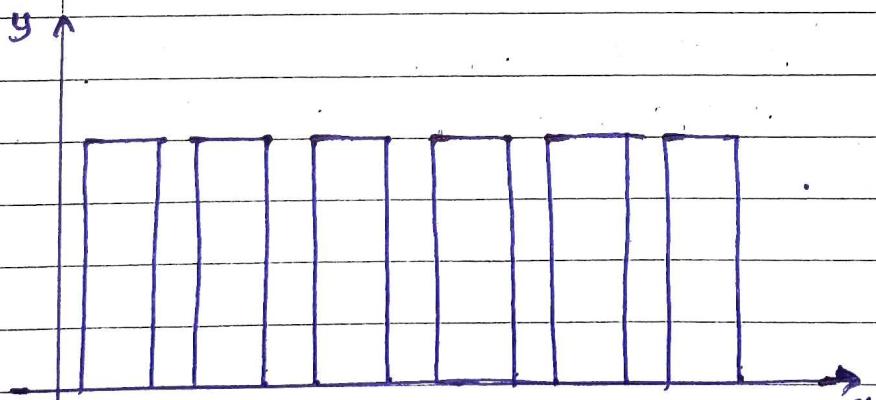
- probability density function:

$$f(x) = \begin{cases} \frac{1}{\beta-\alpha} & \text{if } \alpha \leq x \leq \beta \\ 0 & \text{otherwise} \end{cases}$$

Code:

```
def uniform_distribution_generator(min, max, size):
    return np.random.uniform(low=min, high=max, size=size)
```

uniform_dist = uniform_distribution_generator(min=10, max=40, size=10000)



{
 symmetric
 mean = median
 uniform shape (flat)
 properties

Using `scipy.stats` to generate the uniform distribution:

`min = 10`

`max = 40`

`var = stats.uniform.rvs(loc=min, scale=max, size=10000)`

computing PDF and CDF:

"`F(x=0):`", `stats.uniform.pdf(x=0, loc=min, scale=max)`

"`F(x=20):`", `stats.uniform.pdf(x=20, loc=min, scale=max)`

"`F(x=30):`", `stats.uniform.pdf(x=30, loc=min, scale=max)`

"`F(x<=30):`", `stats.uniform.cdf(x=30, loc=min, scale=max)`

percent point Function (inverse of cdf):

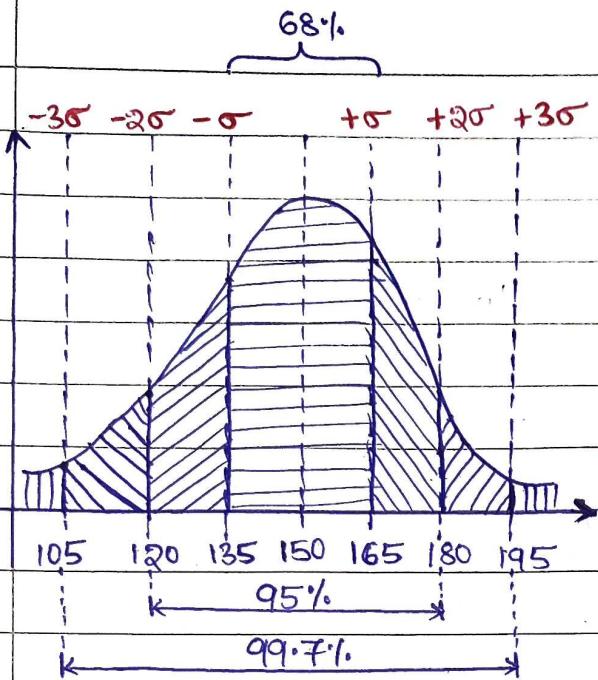
`stats.uniform.ppf(q=0.25, loc=min, scale=max)`

② Gaussian/Normal Distribution: 68-95-99.7% Rule

* The Normal distribution is characterized by its Bell-shaped curve and is widely used in various fields to model a wide range of natural phenomena.

* This distribution can be described by its mean and standard deviation, providing insights into the central tendency and spread of the data.

* Important - If a data follow normal distribution, it will follow 68-95-99.7% Rule. i.e; 68% of total data will lie in 1 std range, 95% of the total data will lie in 2 std range and 99.7% of the total data will lie in 3 std range.



$$\sigma = 15$$

Properties:

- Symmetric
- mean = median
- Bell shape

Note: How do we identify the type of distribution?
→ shape of the distribution

- Data follows Normal Distribution :

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

- Probability Density Function :

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

Code:

```
def normal_distribution_generator(mean, std, size):
    return np.random.normal(loc=mean, scale=std, size=size)
```

normal_dist = normal_distribution_generator(mean=10., std=40, size=100)

using `scipy.stats` to generate the normal distribution :

mean = 10

std = 40

vax = `stats.norm.rvs(loc=mean, scale=std, size=10000)`

computing PDF and CDF :

"F(x=0)", `stats.norm.pdf(x=0, loc=mean, scale=std)`

"F(x=20)", `stats.norm.pdf(x=20, loc=mean, scale=std)`

"F(x=30)", `stats.norm.pdf(x=30, loc=mean, scale=std)`

"F(x<=30)", `stats.norm.cdf(x=10, loc=mean, scale=std)`

percent point function (Inverse of cdf) :

`stats.norm.ppf(q=0.5, loc=min, scale=max)`

③ Exponential Distribution :

The exponential distribution models the time until an event occurs in a poisson process, where events occurs continuously and independently at a constant average rate (λ).

For eg.:

- * Time between customer arrivals.
- * Modeling waiting times, such as the time between phone calls at a call center.
- * Reliability engineering, eg; time until failure of a component.

- Data follows exponential distribution :

$$X \sim E(\lambda)$$

- Probability Density Function :

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

code :

```
def exponential_distribution_generator(lam , size ):  
    # scale (beta) is inverse of lambda parameter  
    return np.random.exponential(scale = 1/lam , size = size )
```

```
exponential_dist = exponential_distribution_generator(lam=5, size=1000)
```

using scipy.stats to generate the normal distribution :

lambda_ = 5

```
var = stats.expon.rvs(scale = 1/lambda_ , size = 1000 )
```

computing PDF and CDF :

"F(x=0.5):", stats.expon.pdf(x = 0.5, scale = 1/lambda)

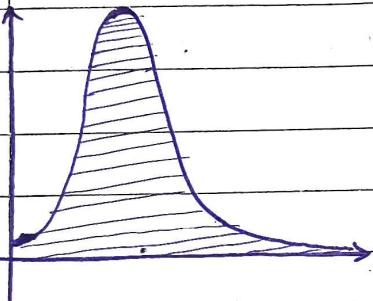
"F(x=0.75):", stats.expon.pdf(x = 0.75, scale = 1/lambda)

"F(x=0.90):", stats.expon.pdf(x = 0.90, scale = 1/lambda)

"F(x<=1):", stats.expon.cdf(x = 1, scale = 1/lambda)

percent point function (inverse of cdf) :

stats.expon.ppf(q = 0.994, scale = 1/lambda)



④ Log-Normal Distribution :

- * It is right-skewed and typically used for variables that are positive and have a long right tail.
- * The log-normal distribution describes a random variable whose logarithm is normally distributed.

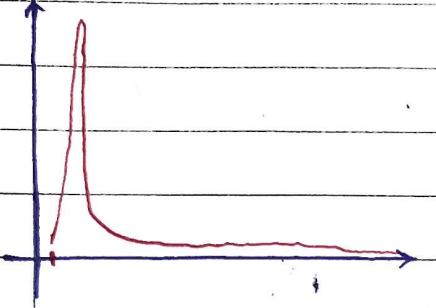
* Applications :

- ① Modeling variables like income, asset prices, and biological measurements that cannot be negative.
- ② In Finance, modeling stock prices.

code :

```
def lognormal_distribution_generator(mean, std, size):
    return np.random.lognormal(mean=mean, sigma=std, size=size)
```

`lognormal_dist = lognormal_distribution_generator(mean = 0, std = 2, size = 10000)`



- * The way location, scale, and shape parameters work in Scipy for the Log-Normal distribution is confusing.
- * Suppose a normally distributed random variable X has mean μ and standard deviation σ .
- * Then $Y = \exp(X)$ is lognormally distributed with $s = \sigma$ and $scale = \exp(\mu)$.

using Scipy.stats to generate the normal distribution :

mean = 0

std = 2

`var = stats.lognorm.rvs(s=std, scale=np.exp(mean), size=10000)`

computing PDF & CDF :

`F(x=0.5): stats.lognorm.pdf(x=0.5, s=std, scale=np.exp(mean))`

`F(x=0.75): stats.lognorm.pdf(x=0.75, s=std, scale=np.exp(mean))`

`F(x=2): stats.lognorm.pdf(x=2, s=std, scale=np.exp(mean))`

`F(x<=2): stats.lognorm.cdf(x=2, s=std, scale=np.exp(mean))`

`stats.lognorm.ppf(q=0.636, s=std, scale=np.exp(mean))`

⑤ pareto distribution : [power-law] 80:20 distribution

- * The power-law distribution, also known as a pareto distribution, is characterized by a long tail of rare events.
- * It's a heavy-tailed distribution. It describes phenomena where large events are rare but have a substantial impact.
- * It is also known as the "80-20 rule".

Applications:

- * Modeling city sizes, where a few cities have a very large population while most have relatively few inhabitants.
- * In social sciences, representing the distribution of income

⑥ Wealth.

code:

```
def pareto_distribution_generator(alpha, size):
    return np.random.pareto(a=alpha, size=size)
```

```
pareto_dist = pareto_distribution_generator(alpha=2, size=10000)
```

using scipy.stats to generate the pareto distribution.

alpha = 2

```
var = stats.pareto.rvs(b=alpha, size=10000)
```

computing PDF and CDF :

$F(x=0.5)$: stats.pareto.pdf($x=0.5, b=\text{alpha}$)

$F(x=0.75)$: stats.pareto.pdf($x=0.75, b=\text{alpha}$)

$F(x=2)$: stats.pareto.pdf($x=2, b=\text{alpha}$)

$F(x \leq 2)$: stats.pareto.cdf($x=2, b=\text{alpha}$)

Verify the Type of Distribution - QQ plot :

Normal distribution

{ stats.probplot (normal_dist, dist = stats.norm, plot = plt)

The above code will help you check observed distribution with Theoretical distribution.

uniform distribution

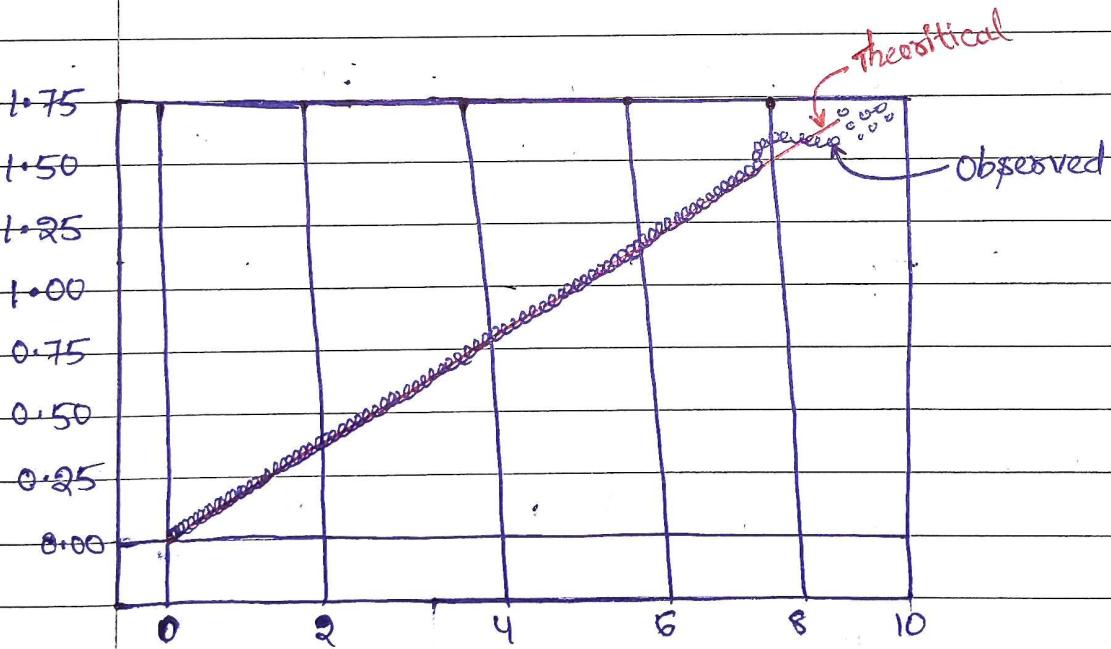
{ stats.probplot (uniform_dist, dist = stats.uniform, plot = plt)

observed distribution

theoretical distribution

Exponential distribution

{ stats.probplot (exponential_dist, dist = stats.expon, plot = plt)



Feature Engineering :

(*) Feature Scaling (①) Rescaling ↴

- * Feature scaling is all about removing the scale difference in Numerical columns.
- * Feature Scaling, is a process of transforming data into a common format ② Scale.
- * This ensures that data from different sources ③ with different units of measurement can be compared ④ combined more effectively.

① Min-Max Scaling (Normalization) :

After Rescaling $\rightarrow \min = 0$
 $\max = 1$

$$\Rightarrow \frac{x - \min}{\max - \min}$$

code:

```
def min_max_scalar(data):
    min_value = np.min(data)
    max_value = np.max(data)
    return (data - min_value) / (max_value - min_value)
```

② Z-transformation (Standardization) :

After Rescaling $\rightarrow \mu = 0$
 $\sigma = 1$

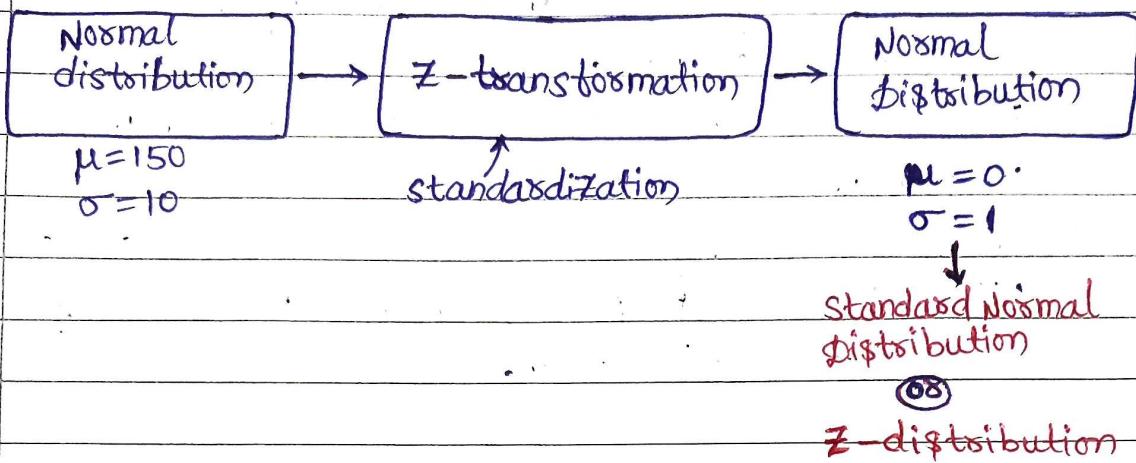
$$\Rightarrow \frac{x - \mu}{\sigma}$$

code: def z_transformer(data):

```
mean = np.mean(data)
std = np.std(data)
```

```
return (data - mean) / (std)
```

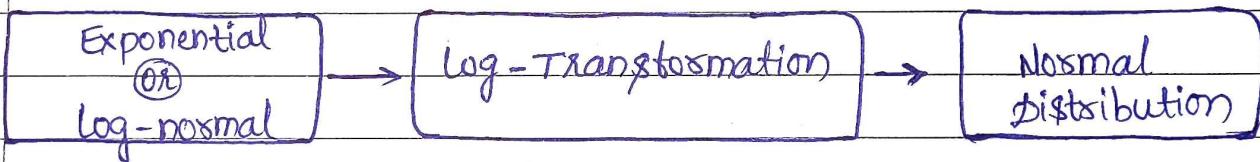
Note:



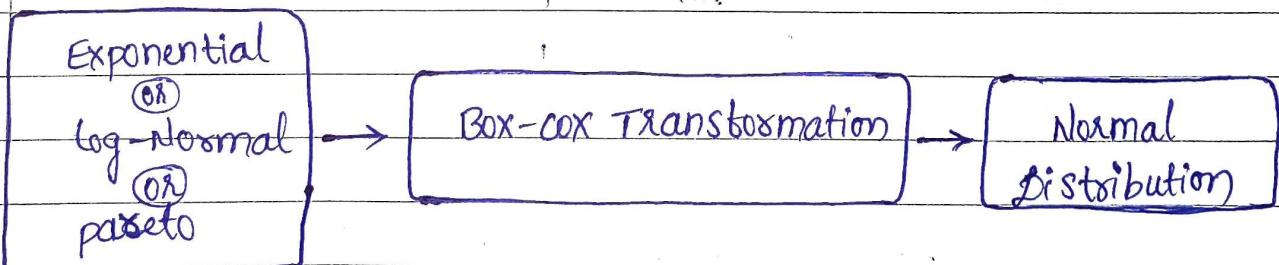
Data Transformation (Reshaping) :

- It alters the shape of the distribution.

① Log - Transformation :



② Box - cox Transformation :



Case Study - Budget Allocation :

Scenario: Setting up a Garment Store

your friend is embarking on the journey of setting up a garment store. Launching a new business with limited funds demands meticulous planning, resourcefulness, and a sharp understanding of the market.

For the initial phase, your friend has decided to focus solely on the product - Jeans, and is investing 1,00,000 INR in inventory.

Question: How should the available budget be allocated to inventory?

Answer: Determining the optimal budget allocation can be a bit challenging without a thorough market analysis. A primary suggestion might be to allocate 50,000 INR for both male and female customers.

However, what if there's a scarcity of male customers? Fortunately, your friend has you, a trusted best friend who ~~also~~ happens to be a successful Data Scientist. Budget Allocation can rely on many factors like:

① Gender split: Since your friend is offering jeans, it's essential to consider the Gender Split in potential customers. This will help in determining how to allocate the Budget effectively.

② sizing considerations: Different body types and sizes should be taken into account. It might be beneficial to allocate a portion of the budget for a diverse range of sizes.

Task - Data driven Budget Allocation :

Let's assume you with the help of an expert got the following information about jeans sizes and Height ranges for both the genders. Consider the following chart for better Budget Allocation.

Group	Small	Medium	Large
Female Heights	[120, 136.61]	[136.61 to 153.40]	[154.41 to 170]
Male Heights	[99.87, 128.87]	[129.97 to 160.07]	[160.07 to 190.17]

Given data :

- ① Import libraries : Numpy, pandas, matplotlib.pyplot, seaborn
- ② Load data : df = pd.read_csv('data/height-gender-data.csv')
- ③ Description about data :
The data set containing 2 columns (Height & Gender) and 10,000 rows.
df.columns - Height, Gender
df.shape - (10000, 2).

⇒ Gender Split - Budget Allocation for Male & Female :

Total Budget - ₹1,00,000/-

Non-visual Analysis :

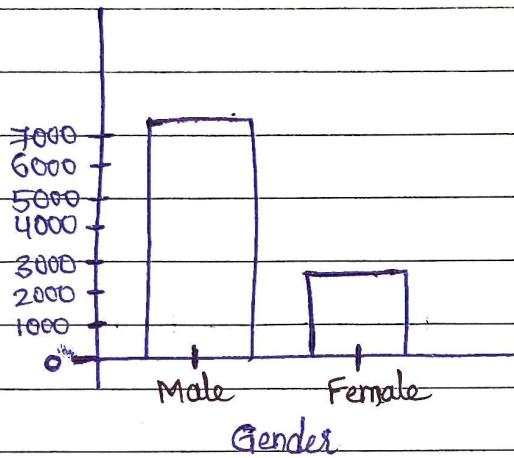
```
df['gender'].value_counts()
```

male 7114

Female 2886

Visual Analysis :

fig, ax = plt.subplots(figsize=(3,3), facecolor='light yellow', constrained_layout=True)
 sns.countplot(data=df, x='gender', ax=ax)



- Based on visual & non-visual analysis Budget Allocation for Male & Female Accordingly.

- Male - ₹114 out of 10,000 → ₹ 71,140/- ~ ₹1000/-
- Female - 2886 out of 10,000 → ₹ 28860/- ~ ₹9000/-

⇒ Budget Allocation for Small, Medium & Large sizes :

Groupby Gender :

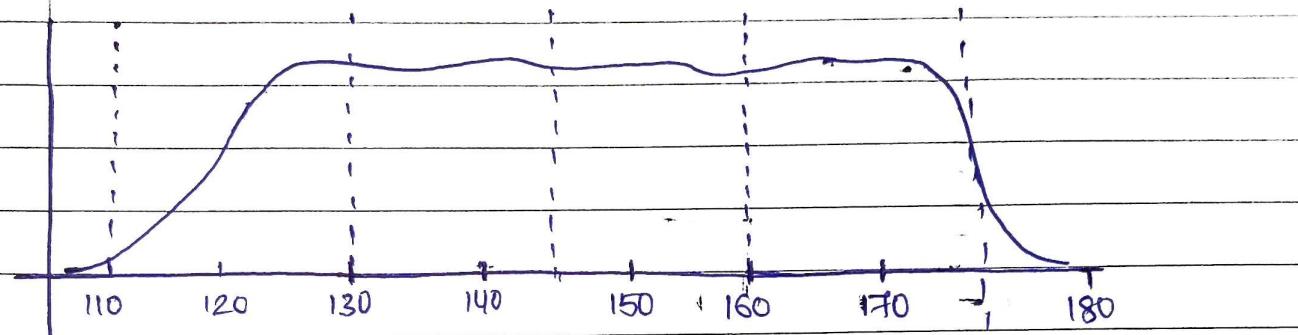
grouped_df = df.groupby('gender')

female_height = grouped_df.get_group('Female')['Height']

male_height = grouped_df.get_group('Male')['Height']

- Female :

```
fig, axs = plt.subplots(figsize=(8,3), facecolor='lightyellow',  
                      constrained_layout=True)  
axs.set_title('Histogram plot')  
sns.kdeplot(data=female_height, ax=axs, fill=True)
```



→ Female Height follows the uniform distribution -

Budget = ₹ 29,000/-

- Based on this it's better to allocate Equal sum of all the sizes.
 - small - $29000/3 = 9666/-$
 - medium - $29000/3 = 9666/-$
 - large - $29000/3 = 9666/-$

- Male :

```
fig, axs = plt.subplots(figsize=(8,3), facecolor='lightyellow',  
                      constrained_layout=True)
```

axs.set_title('kde plot')

```
sns.kdeplot(data=male_height, ax=axs, fill=True)
```

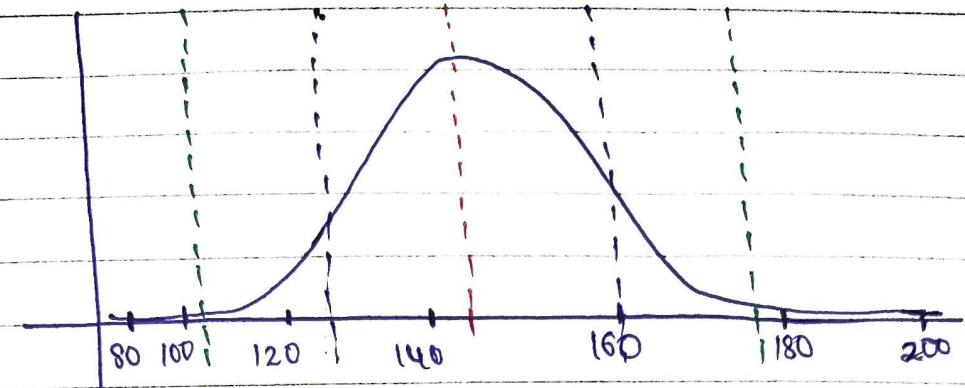
axs.axvline(110, ls='--', color='green')

axs.axvline(130, ls='--', color='orange')

axs.axvline(145, ls='--', color='red')

axs.axvline(160, ls='--', color='orange')

axs.axvline(175, ls='--', color='green')



⇒ Male Height follows Gaussian distribution - Budget = ₹ 71,000/-

- Based on this, it can follow 68-95-99.7% Rule
 - Small - 15.85% of ₹ 71000 = ₹ 11,253.5/-
 - Medium - 68% of ₹ 71000 = ₹ 48,280/-
 - Large - 15.85% of ₹ 71000 = ₹ 11,253.5/-