

Day - 76, Feb 14, 2025 (Falgun 2, 2081 B.S.)

## # Probability Distributions:

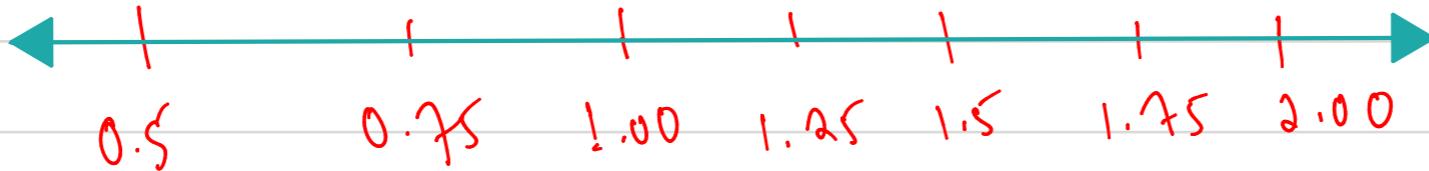
→ Describe the likelihood of the possible outcomes of a random event. Eg: probability of successful clinical trials -

Random Variables. Represents the values for the possible outcomes of a random event.

Discrete Random Variable → Has a countable possible values of P

Continuous Random Variable → Takes all the possible values in some range of numbers. It measure the outcome (continuous).

Continuous Values.



↳ Discrete Distributions represent discrete random variables.

↳ Continuous Distributions represent Continuous Random Variables.

Sample Space: The set of all possible values for a random variable.

• Sample Space for Single Coin Toss = {Heads, Tails}.

• Sample Space for Single Die Roll = {1, 2, 3, 4, 5, 6}

# Single Die Roll

• Sample Space = {1, 2, 3, 4, 5, 6}

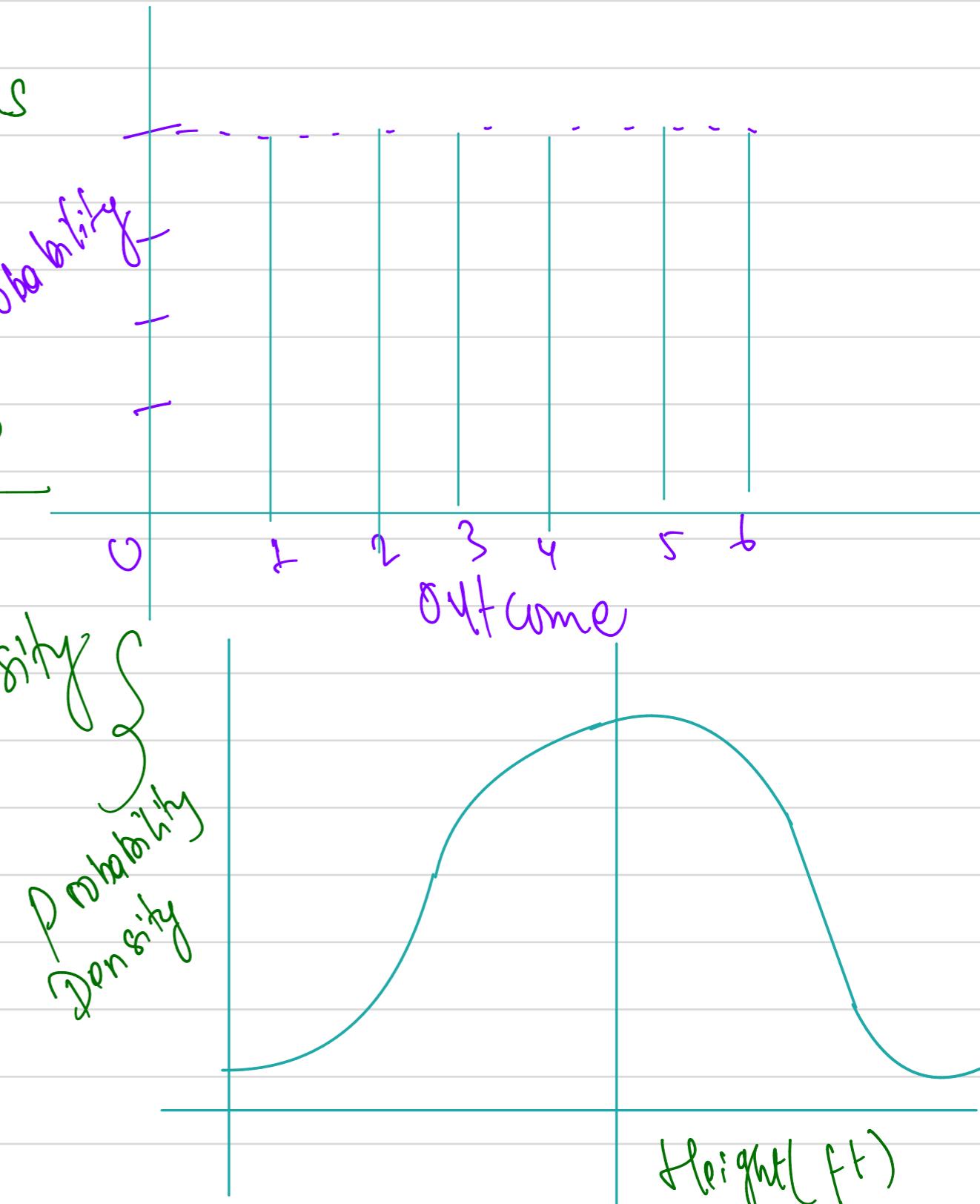
• probability of each outcome = 16.7%

# Infinite number of values = Continuous  
Random Variable

e.g.

$$5.67 + 5.66 + 5.15 + \dots + 20$$

So, we need distribution using continuous probability density function



## # Binomial Distribution

A discrete distribution that models the probability of events with only two possible outcomes, success or failure.

- Each event is independent
  - the probability of success is the same for each event.
  - Eg: tossing the coin for 10 times where each toss is 'f' or 's'.

## # Mutually Exclusive

Two outcomes are mutually exclusive if they cannot occur at the same time.

↳ Binomial Distributions Used in ML, Banking, medicine and investing.

### Random Experiment

- A process whose outcome cannot be predicted with certainty.
- More than one possible outcome and each outcome of the experiment depends on chance.

### # Binomial Experiment

- ↳ consists of a number of repeated trials
- ↳ Each trial has only two possible outcomes and the ' $p$ ' is success is the same for each trial.
- ↳ Each trial is independent.

## Binomial Experiment Examples:

- 10 repeated coin tosses
  - 2 possible outcomes: heads or tails
  - 'p' of success for each toss is the same: 50%.
  - the outcome of one coin toss does not affect the outcome of any other coin toss.
- 100 customer visits with two possible outcomes: return or not return
  - p of success for each customer visit is the same = 10%.
  - the outcome of one customer visit doesn't affect the outcome of any other customer visit.

## Binomial Distribution formula

$$P(X=k) = \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k}$$

•  $k$  refers to no. of Success

•  $n$  refers to the number of trials

•  $p$  refers to the probability of Success in a given trial.

$$P(X=k) = \frac{n!}{k!(n-k)!} * p^k * (1-p)^{n-k}$$

$\downarrow$

$$C_n^k ( n\text{-choose } k )$$

$C(n, k) \rightarrow$

$$P(X=0) = 0.729$$

$$P(X=2) = 0.027$$

$$P(X=3) = 0.001$$

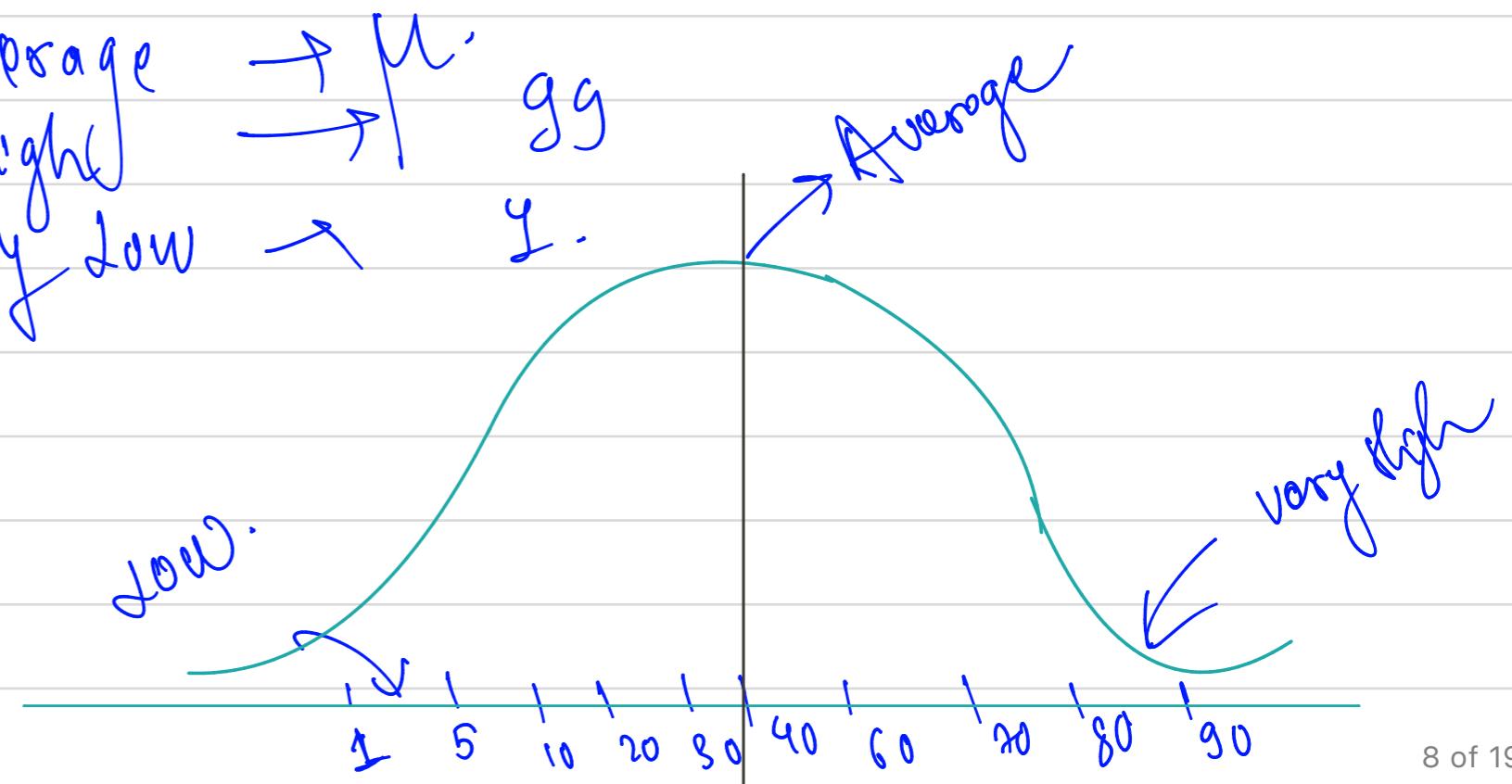
## # Normal Distribution:-

↳ While Discrete Probability Distributions like Binomial & Poisson Distribution can only take the discrete value into account.

# Normal Distribution: → A continuous probability distribution that is symmetrical on both sides of mean and bell-shaped.

Eg: Maximum number score average  $\rightarrow \mu$ .  
Very less |||| High  $\rightarrow \mu$ . gg  
Very |||| Very low  $\rightarrow \mu$ . yy

↳ Used in Advanced Statistics  
like hypothesis testing.  
In terms of score.



## Features of Gaussian Distribution:

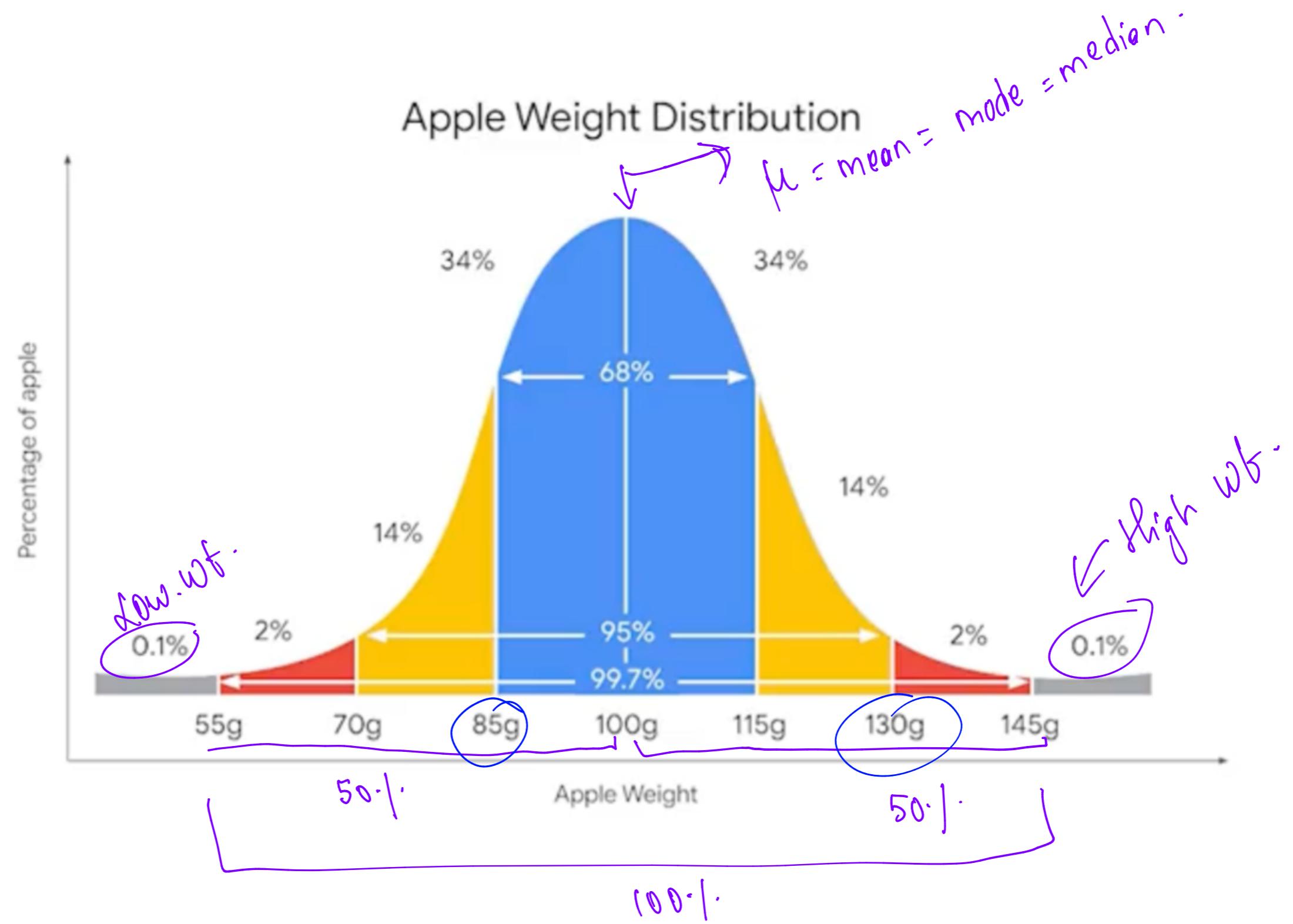
- ① Shape is a bell curve.
- ② Mean is located at the center of the curve
- ③ Curve is symmetrical on both sides of the center
- ④ The total area under the curve equals 1.

## # Standard Deviation:

Calculates the typical distance of a data point from the mean of your dataset.

# mean + center

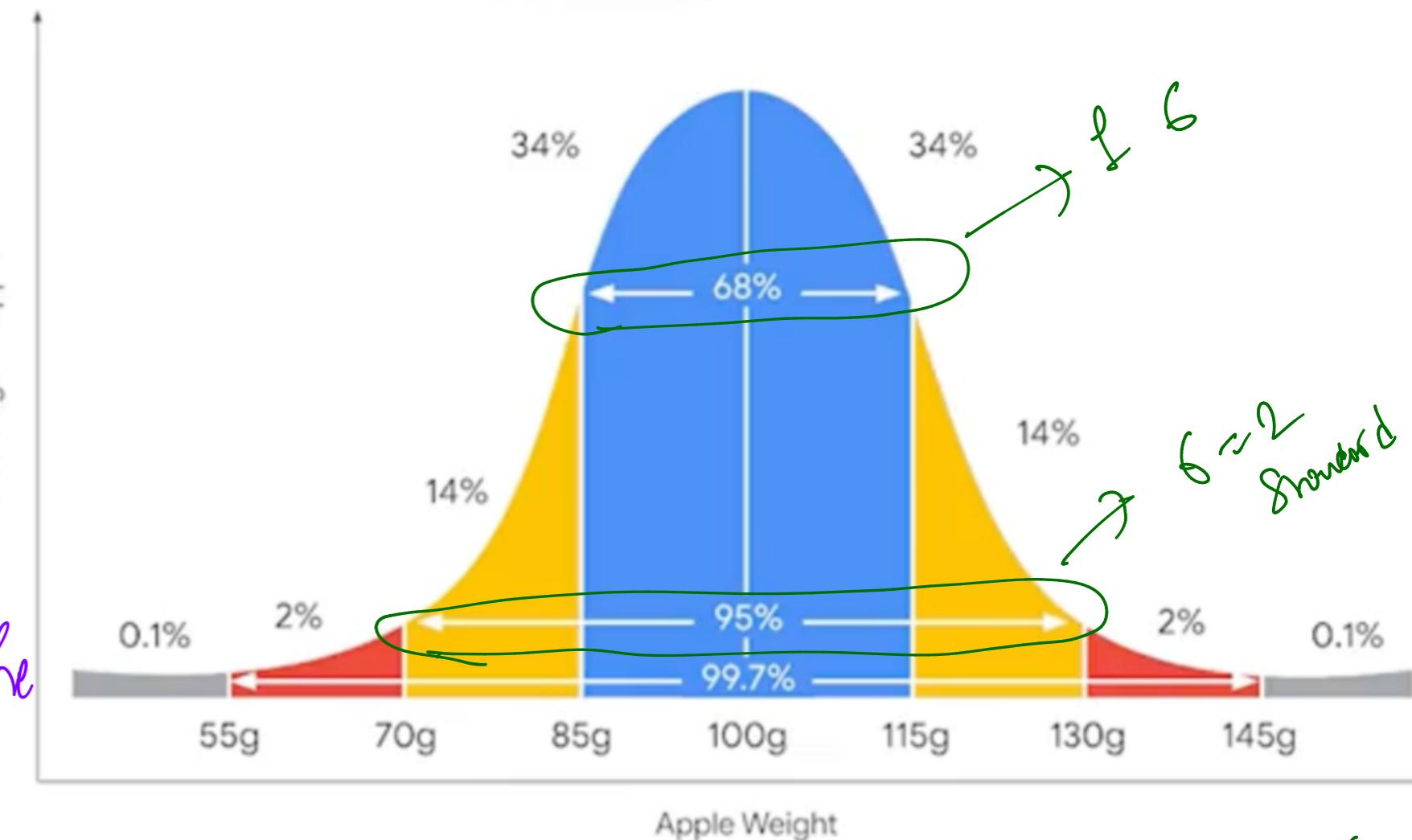
#  $\sigma \Rightarrow$  spread of data from the center or mean.



## Empirical Rule:

- 68% of values fall within 1 standard deviation of the  $\mu$ .
- 95% of values fall within 2 standard deviations of the  $\mu$ .
- 99.7% of values fall within 3 standard deviations of the  $\mu$ .

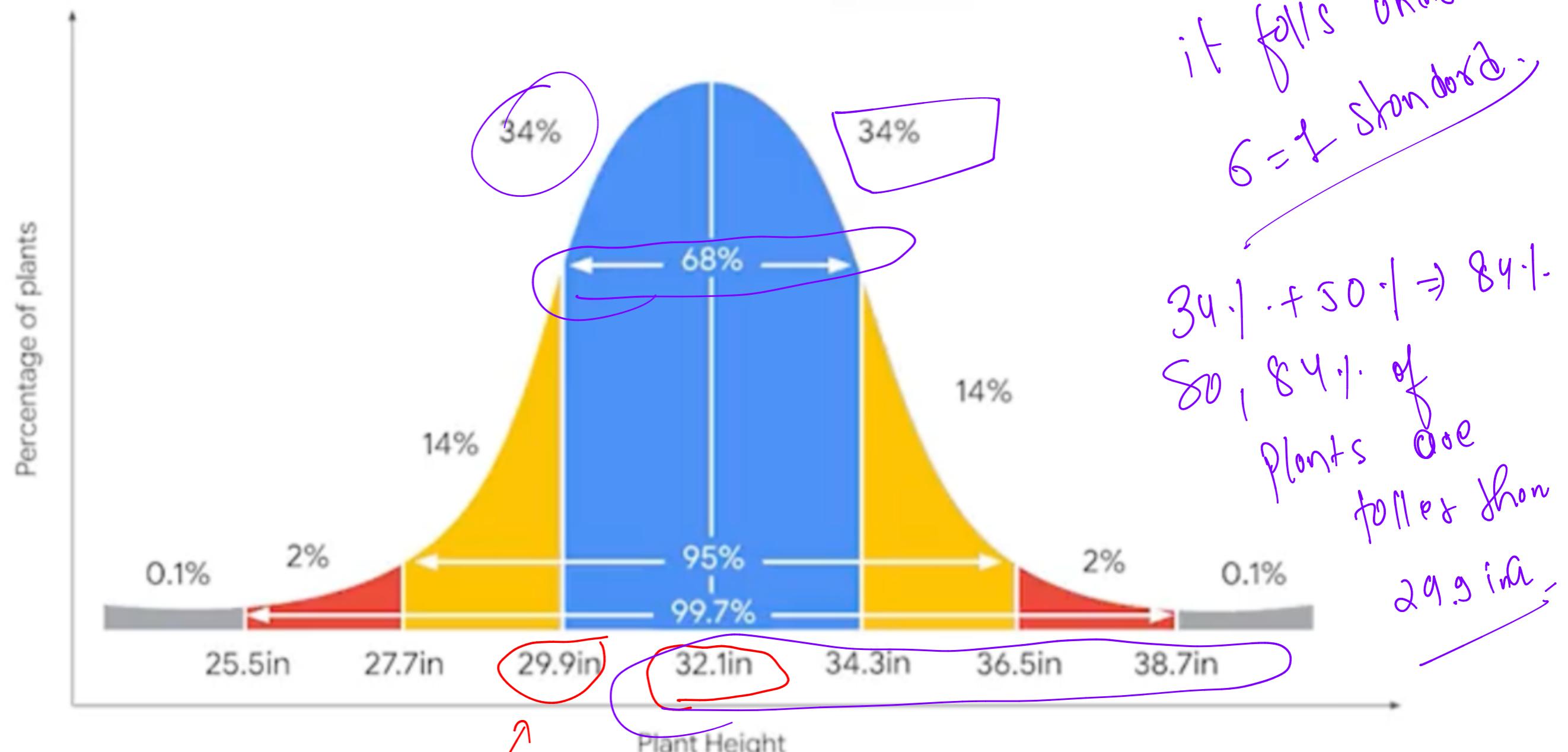
Apple Weight Distribution



$\rightarrow \delta = 2$  means of ranges between 70g and 130g.

$\rightarrow \delta = 1$  means 68% which means 85g and 115g between the ranges of apple weight ranges.

$S_1$ ,  $34.1\%$  from the center of mean for plant to be at least 29.9 inch tall because it falls under the  $6 = 1$  standard.



At least 29.9in tall?

## Explanation:

If you want to determine the percentage of plants taller than 29.9 inches for your backyard landscape design, you can use the Empirical Rule of normal distribution.

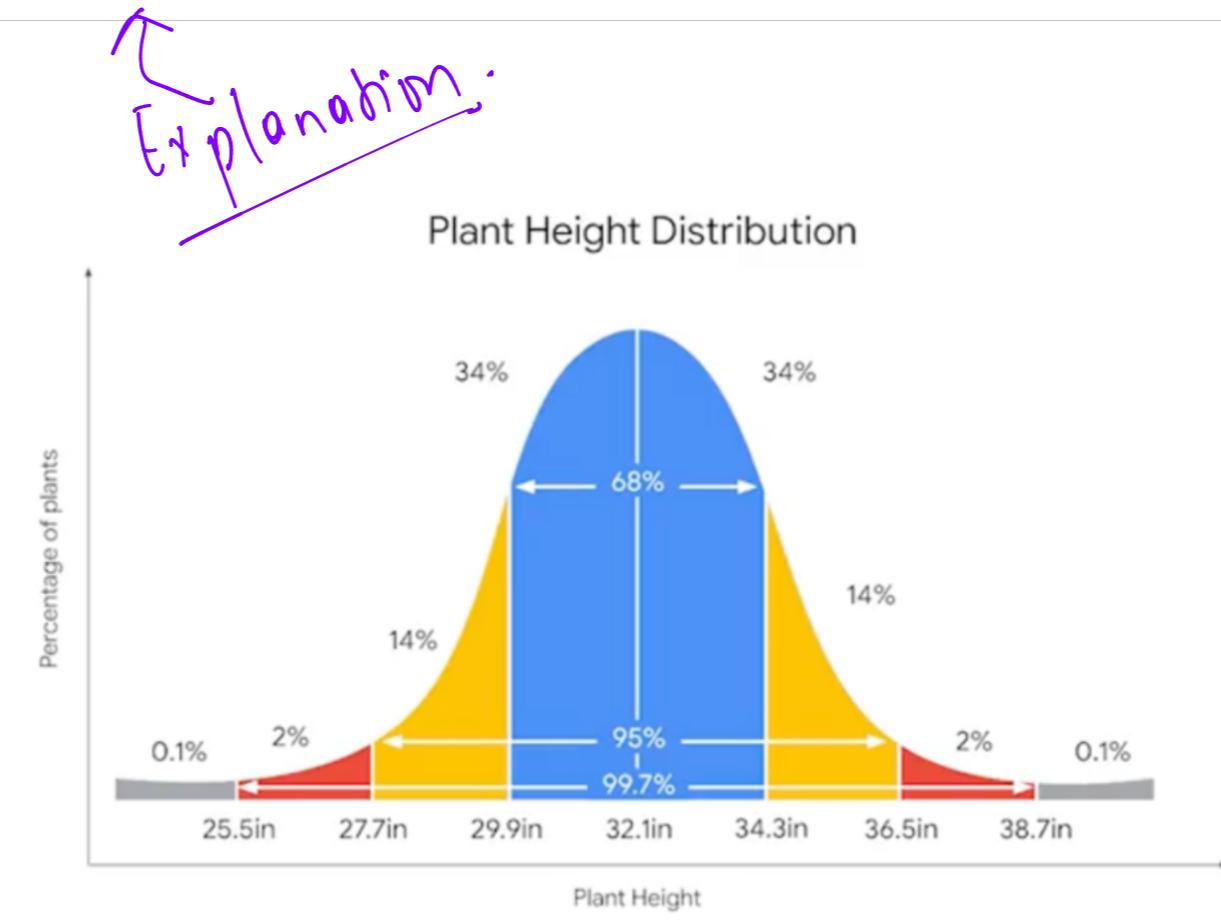
Since 29.9 inches is one standard deviation below the mean, the Empirical Rule states that 68% of values fall within one standard deviation of the mean. This means 34% of values lie between 29.9 and the mean. Additionally, in a normal distribution, 50% of values are above the mean.

To find the percentage of plants taller than 29.9 inches, add these two values:

$$34\% \text{ (between 29.9 and the mean)} + 50\% \text{ (above the mean)} = 84\%.$$

Thus, 84% of your plants meet the height requirement. This quick estimation method helps analyze data patterns efficiently. As a future data professional, understanding the normal distribution will be key to interpreting and making data-driven decisions.

Source: The power of Statistics  
Courseware Course Provided  
by Google.



## # Standardize Data using Z-scores:

- Z-score is a measure of how many  $\sigma$  below or above the population mean a data point is.
- Z-score gives idea on how far the data is from the  $\mu$ .
- Z-score is 0 if the value is equal to the mean.
- Z-score is +ve if the value is greater than the  $\mu$ .
- Z-score is -ve if the value is less than the  $\mu$ .

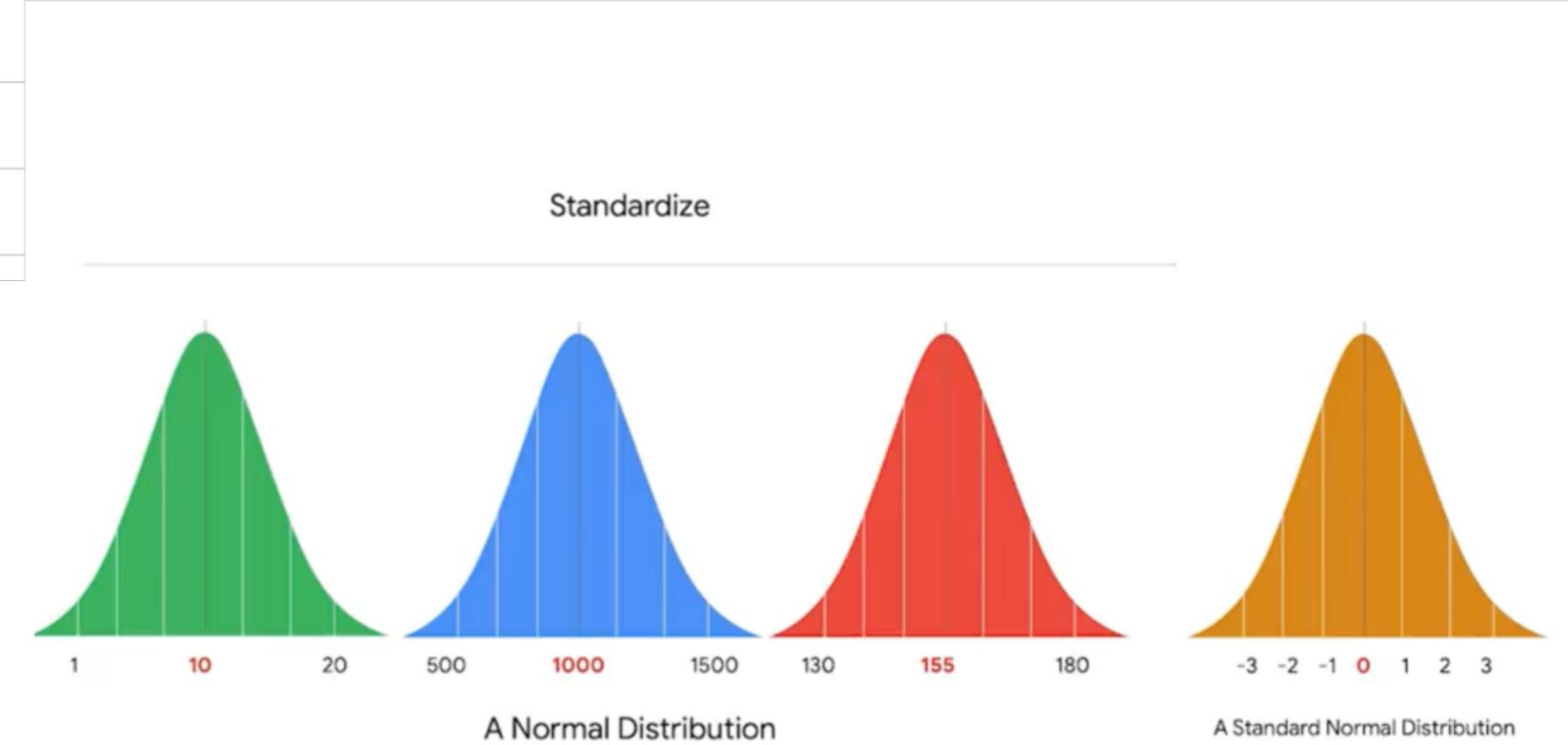
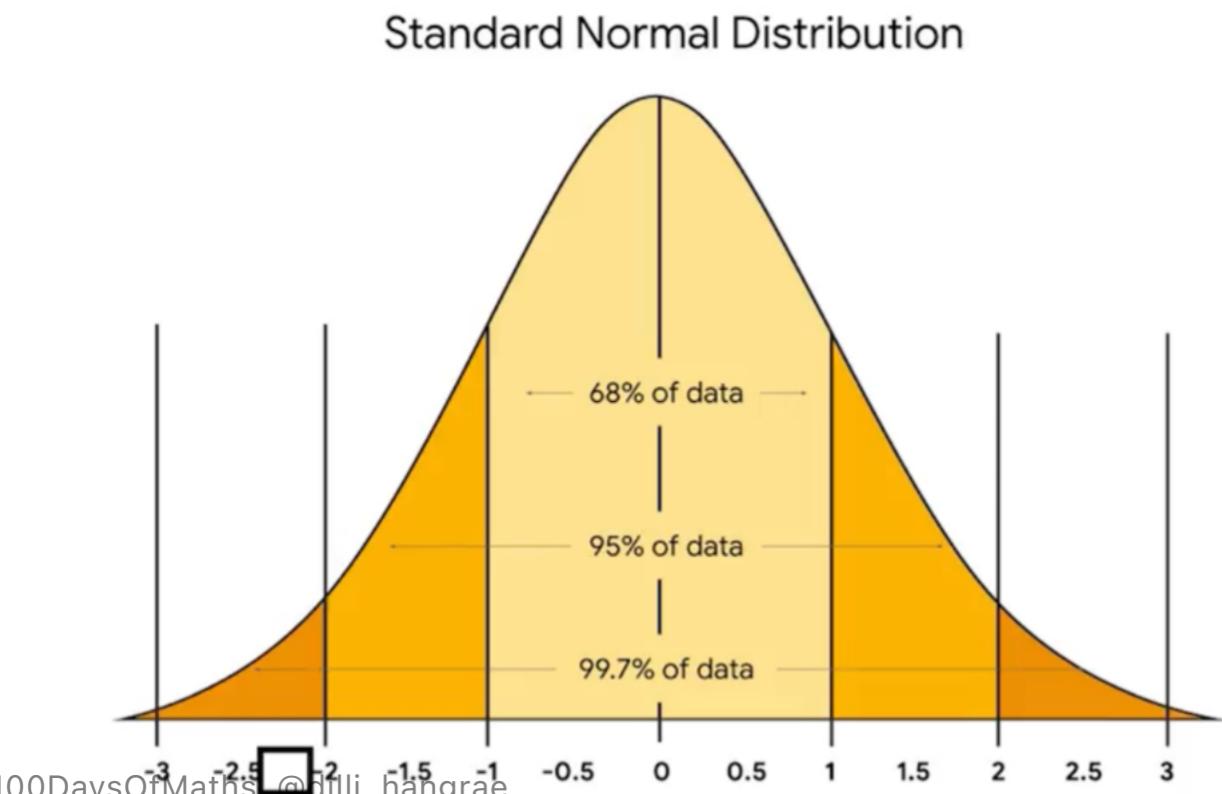
# Standardization, the process of putting different variables in the same scale.

Z-score helps to understand the relationship between different datasets.

## Z-Score Application in Anomaly Detection Application

- ① Fraud in financial Transactions
- ② Flaws in manufacturing products
- ③ Intrusions in Computer Networks.

# Standardize



↑ Standardization helps to perform and analyze the distribution efficiently & smoothly

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

score                      mean  
 ↓                          ↓  
 $\sigma$                        $\mu$

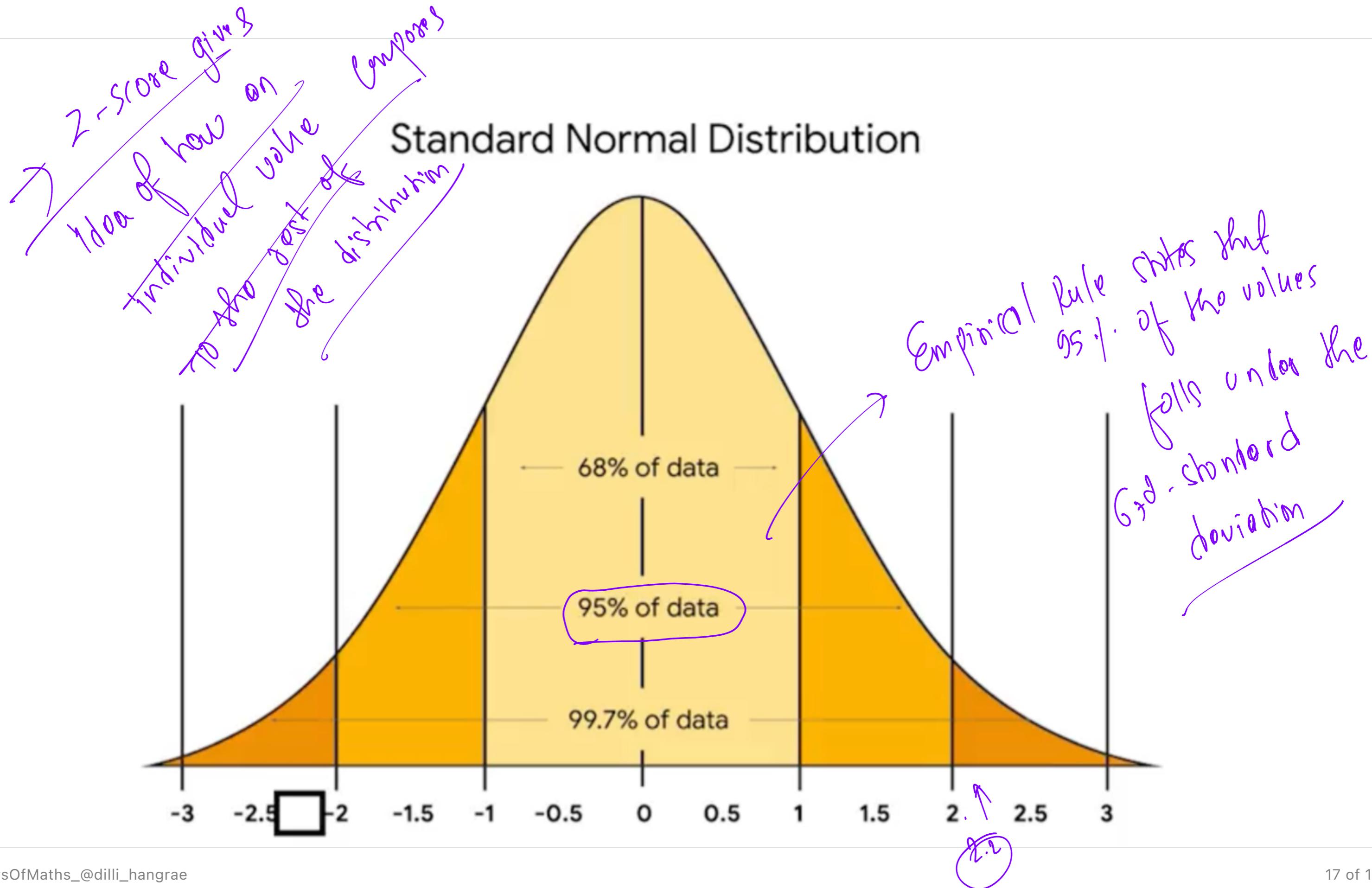
$$z = \frac{(133 - 100)}{15}$$

$$z = 2.2$$

$z = 2.2$  tells us that your test score is 2.2 standard deviations above the mean or average score.

So, the z-score of data value equals to the  $\mu$ . When  $z$ -score =  $\mu$ . A z-score is a measure of how many standard deviations below or above the population mean or data point is -

## Standard Normal Distribution



## Example:

Say you score in 85, you want to find out if that's a good score relative to the rest of the class.

Whether or not it's a good score depends on the mean and standard deviation of all exam scores.

Suppose the exam scores are normally distributed with a mean score of 90 and a standard deviation of 4, you can use the formula to calculate the z-score of a raw score of 85.

Your z-score is yours raw score, 85 minus the mean score 90, divided by the standard deviation 4.

This is 85 minus 90 divided by 4 equals -5, divided by 4 equals 1.25. Your z-score of -1.25 tells you that your exam score of 85 is 1.25 standard deviations below the mean or average exam score. Z-scores give you an idea of how individual values compared to the mean.

As a data professional, you'll use z-scores to help you better understand the relationship between specific values in your data set.

## Introduction to Sampling:

Types of Statistics

→ Descriptive statistics summarize the main features of a dataset.

→ Inferential statistics use sample data to draw conclusions about a larger population.

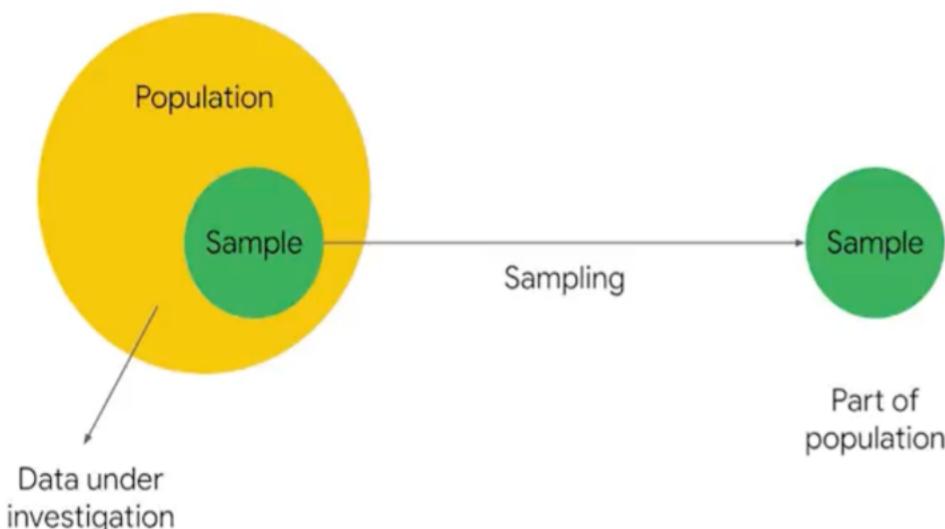
#Sampling - Process of drawing a subset of data from a population.

## Questions Answered by Sampling:

A representative sample accurately reflects the characteristics of a population. If a sample doesn't accurately reflect the characteristics of a population, then the inferences will be likely to be unreliable and inaccurate.

## Questions answered by sampling

- How many products in an app store do we need to test to feel confident that all the products are secure from malware?
- How do we select a sample of users to run an effective A/B test for an online retail store?
- How do we select a sample of customers of a video streaming service to get reliable feedback on the shows they watch?



This subset is your sample, then you can

Sampling should be drawn randomly or unbiased.

Q finding the Average height of people in Nepal we use Sampling?