# Structure of the Presentation

1. What is statistics and why it is important for data scientists?
2. Important application of Mean and Geometric Mean
3. Introduction to probability concepts
4. Sampling method
5. Correlation and regression

# Statistics Defined

Friday, 28 June 2024     12:09 PM

**STATISTICS** is the science of collecting, organizing, presenting, analyzing, and interpreting data to assist in making more effective decisions.

⭐ Understanding statistical concepts is very important as these concepts are helpful in making informed decisions.
- Like business growth is associated with economic growth of a country
- Association of Fraud detection with corruption perception index
- Scientific productivity is related with literacy rate

For instance
- Difference between type of variable guide the selection of statistical technique such as Qualitative vs quantitative variable
- Measurement scale (nominal and ordinal vs interval and ratio)
- **Population vs sample data**
- **Type of** distribution (symmetric vs asymmetric)
- Relative vs absolute analytical measurements
- Sample size
- Dispersion in the data

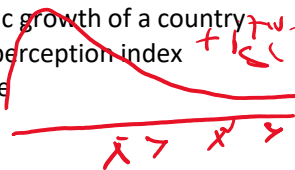Special Applications of Arithmetic Mean (Mean)

What is Geometric Mean and Why it is important?
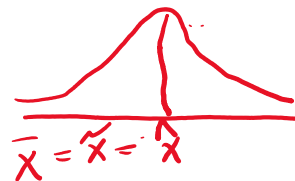Special Application of Geometric Mean

# AM Special Application Example

| Share price | Number of shares |
|-------------|------------------|
| 50          | 600              |
| 40          | 700              |
| 30          | 500              |

$$\frac{120}{3} = 40$$

$$40$$

1. You have an offer to sell all the shares at price 40/share. Will you sell the shares? Why or why not?

# GM Applications

|  | Earning Growth |
|---|---|
| Year 1 | 10% |
| year 2 | 15% |
| Year 3 | 20% |

Can you estimate the average annual growth in earning?

*(handwritten annotations in red):*

1000

$\frac{45}{3} = 15$

X

X 1000 × 10% = 1660 = 1600

comp

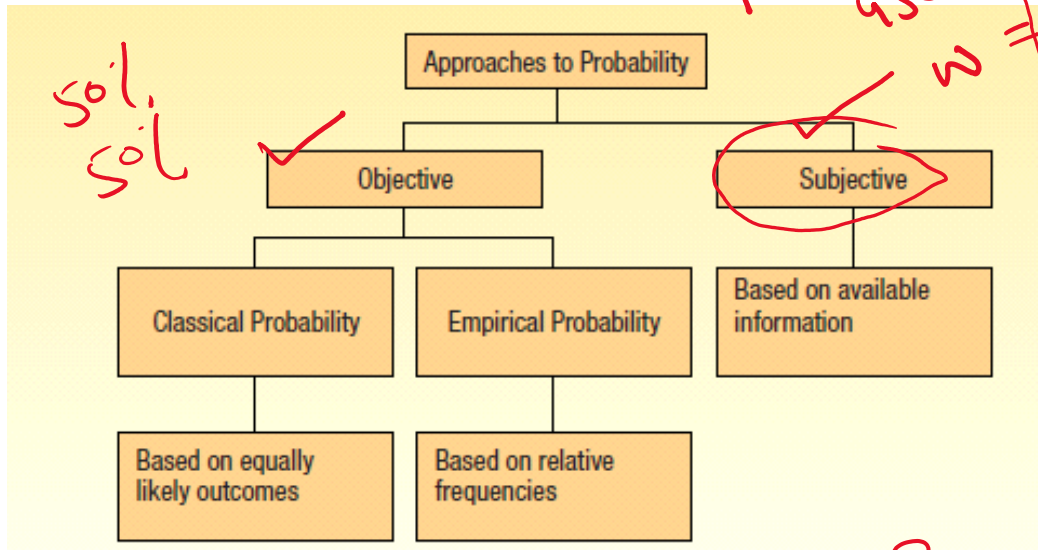10%

|  | bike | sugar |
|---|---|---|
| 2008 | Rs55,000 | Rs32 |
| 2023 | Rs160,000 | Rs140 |
| Change | 105000 | 108 |

Price of which product move faster?

# Survey of Probability Concepts (Beginning of Inferential Statistics)

Friday, 28 June 2024     12:09 PM

**PROBABILITY** A value between zero and one, inclusive, describing the relative possibility (chance or likelihood) an event will occur.

Approaches to Probability

- Objective
  - Classical Probability
    - Based on equally likely outcomes
  - Empirical Probability
    - Based on relative frequencies
- Subjective
  - Based on available information

Important concepts in probability:
- Events
- Independent Events
- Dependents Event
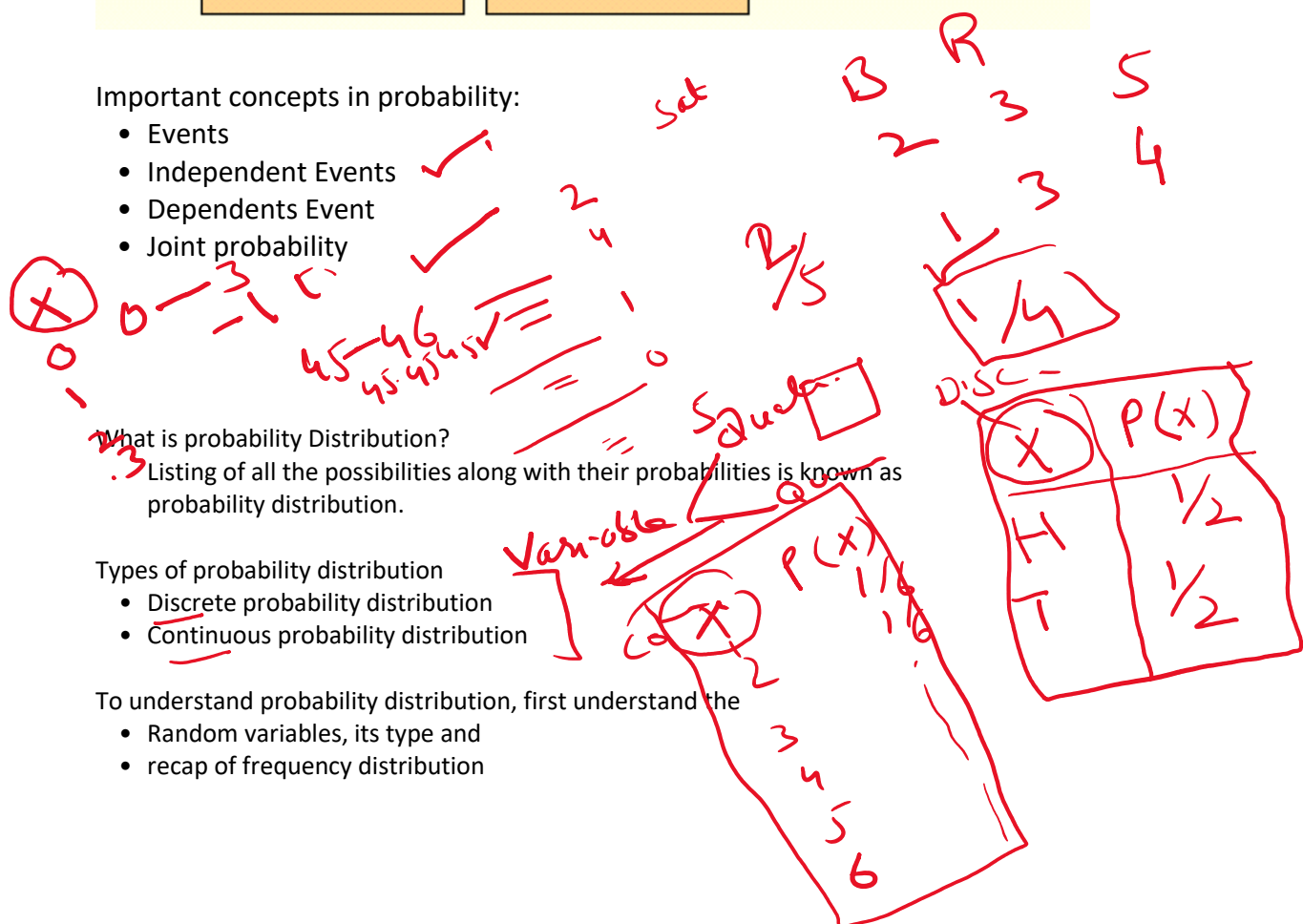- Joint probability

What is probability Distribution?
   Listing of all the possibilities along with their probabilities is known as probability distribution.

Types of probability distribution
- Discrete probability distribution
- Continuous probability distribution

To understand probability distribution, first understand the
- Random variables, its type and
- recap of frequency distribution

# Sampling Methods

Probability sampling vs Non-probability sampling

Concept of generalizability and its relationship with sampling methods

Probability sampling techniques
- Simple random sampling

- Systematic random sampling

- Stratified random sampling

- Cluster sampling

# Correlation and Regression (The Heart of Machine Learning)

Friday, 28 June 2024   8:00 PM

What is correlation analysis?
A group of techniques to measure the relationship between two variables (Independent and dependent variables).

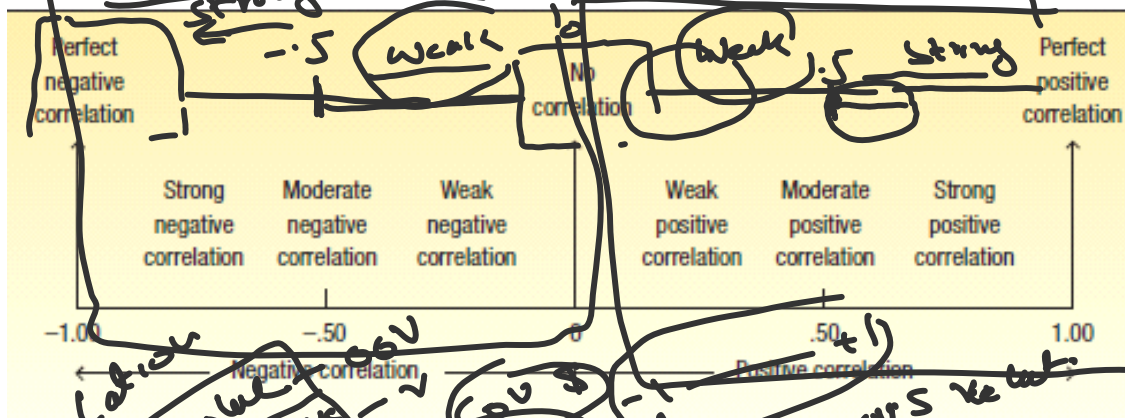Two approaches to access the relationship between two variables
   Graphical approach (scatter plot) using MS Excel and Python
   Mathematical approach using MS Excel and Python

CORRELATION COEFFICIENT A measure of the strength of the linear relationship between two variables.

CHARACTERISTICS OF THE CORRELATION COEFFICIENT
1. The sample correlation coefficient is identified by the lowercase letter r.
2. It shows the direction and strength of the linear relationship between two interval- or ratio-scale variables.
3. It ranges from -1 up to and including +1.
4. A value near 0 indicates there is little relationship between the variables.
5. A value near 1 indicates a direct or positive relationship between the variables.
6. A value near -1 indicates inverse or negative relationship between the variables.

| Perfect negative correlation | | | No correlation | | | Perfect positive correlation |
|---|---|---|---|---|---|---|
| Strong negative correlation | Moderate negative correlation | Weak negative correlation | | Weak positive correlation | Moderate positive correlation | Strong positive correlation |
| -1.00 | | -.50 | 0 | | .50 | 1.00 |
| | | Negative correlation | | | Positive correlation | |

Let's take an example of correlation analysis using:
   • Graphical approach
   • Mathematical approach

What is R-square and how do we interpret it?
Does r represent absolute or relative measure of association?

What is Regression Analysis?
What is simple linear regression?
What is multiple linear regression?
Do we always confront with linear relationships?

$$r = \frac{n\Sigma xy - \Sigma x \Sigma y}{\sqrt{[n\Sigma x^2 - (\Sigma x)^2][n\Sigma y^2 - (\Sigma y)^2]}}$$
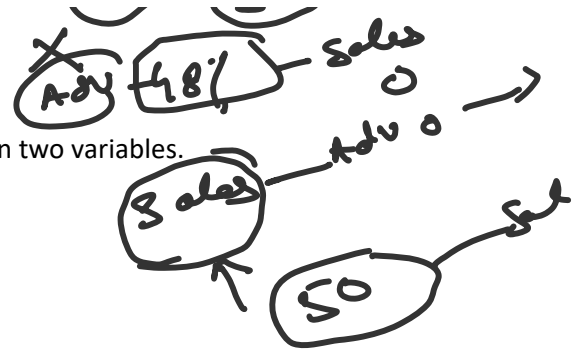
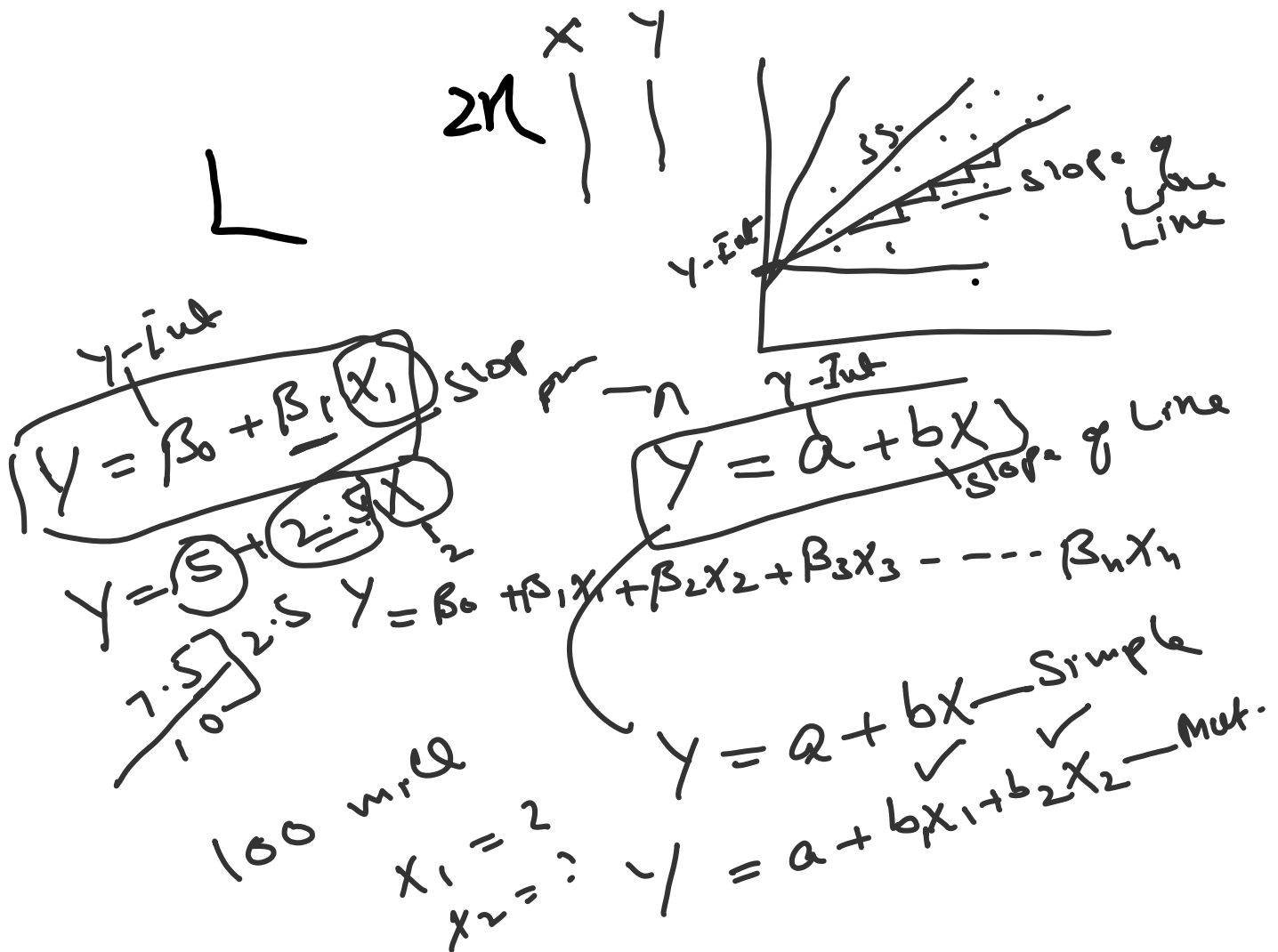Do we always confront with linear relationships?

Regression Analysis
An equation that expresses the linear relationship between two variables.
Role of Intercept and Slope coefficient

- Regression analysis is used for prediction

- What are the various applications of regression analysis specifically in data science?

- What is the role of regression analysis in building machine learning models?

2n

L

X   Y



Y-Int

$Y = \beta_0 + \beta_1 X_1$  Slope  $-n$   Y-Int  $Y = a + bX$  Slope q Line

$Y = \textcircled{5} + \textcircled{2.5}X$

7.5  2.5   $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 ---- \beta_n X_n$

10

100 mrll

$X_1 = 2$
$X_2 = ?$

$Y = a + bX$ — Simple

$Y = a + bX_1 + b_2 X_2$ — Mult.

Null Hyp $H_0$: ——— ———
Alternate $H_1$: ———

| Hyp | Beta Coff | t-value z-valu 1.967 | P-Value |
|-----|-----------|----------------------|---------|
| H1  | 0.75      |                      |         |
| H2  | 0.02      | 2.5                  | 0.001   |
| H3  |           | 1.58                 | 0.09 N Sig |

$H_0: \bar{X} = 0$  X

$H_1: \bar{X} \neq 0$  ✓

$H_0: \bar{X} = 0$  ✓

$H_1: \bar{X} \neq 0$  X

$\alpha = 5\%$