



ISDM (INDEPENDENT SKILL DEVELOPMENT MISSION)

BASIC & DESCRIPTIVE STATISTICS (MEAN, MEDIAN, MODE, VARIANCE)

CHAPTER 1: INTRODUCTION TO DESCRIPTIVE STATISTICS

1.1 WHAT IS DESCRIPTIVE STATISTICS?

DESCRIPTIVE STATISTICS IS A **BRANCH OF STATISTICS** THAT DEALS WITH SUMMARIZING AND ORGANIZING DATA IN A MEANINGFUL WAY. IT HELPS IN **UNDERSTANDING KEY CHARACTERISTICS** OF A DATASET WITHOUT MAKING PREDICTIONS OR GENERALIZATIONS.

DESCRIPTIVE STATISTICS INCLUDE:

- ✓ **MEASURES OF CENTRAL TENDENCY** – MEAN, MEDIAN, MODE (DESCRIBE WHERE THE DATA IS CENTERED).
- ✓ **MEASURES OF VARIABILITY** – VARIANCE, STANDARD DEVIATION, RANGE (DESCRIBE HOW SPREAD OUT THE DATA IS).
- ✓ **FREQUENCY DISTRIBUTIONS & GRAPHICAL REPRESENTATIONS** – HISTOGRAMS, BOX PLOTS, BAR CHARTS (REPRESENT DATA VISUALLY).

EXAMPLE:

A TEACHER CALCULATES THE **AVERAGE (MEAN)** TEST SCORE OF A

CLASS TO SUMMARIZE OVERALL PERFORMANCE INSTEAD OF ANALYZING EACH STUDENT'S SCORE INDIVIDUALLY.

CONCLUSION:

DESCRIPTIVE STATISTICS SIMPLIFY COMPLEX DATASETS, MAKING IT EASIER TO IDENTIFY PATTERNS, TRENDS, AND ANOMALIES.

CHAPTER 2: MEASURES OF CENTRAL TENDENCY

MEASURES OF CENTRAL TENDENCY HELP US IDENTIFY A SINGLE REPRESENTATIVE VALUE IN A DATASET.

2.1 MEAN (ARITHMETIC AVERAGE)

THE MEAN (OR AVERAGE) IS THE SUM OF ALL VALUES DIVIDED BY THE NUMBER OF VALUES IN A DATASET.

FORMULA:

$$\text{MEAN}(X^-) = \sum_{i=1}^n X_i / n$$

WHERE:

- X_i = EACH INDIVIDUAL VALUE
- n = TOTAL NUMBER OF VALUES

EXAMPLE:

A DATASET CONTAINS: **10, 20, 30, 40, 50**

$$\text{MEAN} = (10 + 20 + 30 + 40 + 50) / 5 = 30$$

USE CASES:

-  FINDING THE AVERAGE SALARY IN A COMPANY.
-  CALCULATING THE AVERAGE SALES PER MONTH.

CONCLUSION:

THE MEAN IS USEFUL WHEN ALL VALUES CONTRIBUTE **EQUALLY**, BUT IT IS SENSITIVE TO OUTLIERS (EXTREMELY HIGH OR LOW VALUES).

2.2 MEDIAN (MIDDLE VALUE OF ORDERED DATA)

THE MEDIAN IS THE MIDDLE VALUE WHEN DATA IS ARRANGED IN ASCENDING ORDER. IF THERE IS AN EVEN NUMBER OF VALUES, THE MEDIAN IS THE AVERAGE OF THE TWO MIDDLE VALUES.

EXAMPLE 1 (ODD NUMBER OF VALUES):

DATASET: 5, 10, 15, 20, 25

 **MEDIAN = 15 (MIDDLE NUMBER)**

EXAMPLE 2 (EVEN NUMBER OF VALUES):

DATASET: 4, 8, 12, 16

 **MEDIAN = $(8+12) / 2 = 10$**

USE CASES:

 FINDING THE MIDDLE SALARY IN A COMPANY.

 HOUSE PRICE ANALYSIS, WHERE EXTREME VALUES MAY SKEW THE MEAN.

CONCLUSION:

THE MEDIAN IS NOT AFFECTED BY OUTLIERS, MAKING IT MORE RELIABLE IN SKEWED DISTRIBUTIONS.

2.3 MODE (MOST FREQUENTLY OCCURRING VALUE)

THE MODE IS THE VALUE THAT APPEARS **MOST FREQUENTLY** IN A DATASET.

📌 EXAMPLE:

DATASET: 2, 4, 4, 6, 8, 8, 8, 10

✓ MODE = 8 (BECAUSE 8 APPEARS THE MOST TIMES).

📌 TYPES OF MODE:

- ◆ UNIMODAL – ONE MODE (E.G., {3, 4, 4, 5, 6} → MODE = 4)
- ◆ BIMODAL – TWO MODES (E.G., {2, 3, 4, 4, 6, 6, 7} → MODES = 4 & 6)
- ◆ MULTIMODAL – MORE THAN TWO MODES

📌 USE CASES:

- ✓ FASHION INDUSTRY USES MODE TO FIND MOST POPULAR CLOTHING SIZES.
- ✓ MARKETING TEAMS ANALYZE MODE TO FIND BEST-SELLING PRODUCTS.

💡 CONCLUSION:

MODE IS HELPFUL FOR CATEGORICAL DATA (E.G., MOST PREFERRED CAR BRAND) BUT MAY NOT ALWAYS EXIST OR BE UNIQUE IN NUMERICAL DATASETS.

📌 CHAPTER 3: MEASURES OF VARIABILITY (SPREAD)

MEASURES OF VARIABILITY HELP US UNDERSTAND HOW SPREAD OUT DATA POINTS ARE.

3.1 VARIANCE (MEASURE OF SPREAD FROM THE MEAN)

VARIANCE MEASURES HOW FAR EACH DATA POINT IS FROM THE MEAN. A HIGHER VARIANCE MEANS THE DATA IS MORE SPREAD OUT, WHILE A LOWER VARIANCE INDICATES DATA POINTS ARE CLOSER TO THE MEAN.

FORMULA:

$$\text{VARIANCE}(\Sigma_2) = \sum (x_i - \bar{x})^2 / N \quad \text{VARIANCE} (\sigma^2) = \frac{\sum (x_i - \bar{x})^2}{N}$$

WHERE:

- x_i = EACH VALUE
- \bar{x} = MEAN OF DATASET
- N = TOTAL NUMBER OF VALUES

EXAMPLE:

DATASET: **5, 10, 15** (MEAN = 10)

$$\begin{aligned} \text{VARIANCE} &= (5-10)^2 + (10-10)^2 + (15-10)^2 / 3 = (-5)^2 + (0)^2 + (5)^2 / 3 = 25 + 0 + 25 / 3 = 16.67 \\ \text{VARIANCE} &= \frac{(5-10)^2 + (10-10)^2 + (15-10)^2}{3} = \frac{(-5)^2 + (0)^2 + (5)^2}{3} = \frac{25 + 0 + 25}{3} = 16.67 \end{aligned}$$

USE CASES:

-  MEASURING INCOME INEQUALITY IN DIFFERENT CITIES.
-  ANALYZING STOCK PRICE FLUCTUATIONS OVER TIME.

CONCLUSION:

VARIANCE HELPS MEASURE DATA DISPERSION, BUT IT SQUARES THE DIFFERENCES, MAKING IT HARDER TO INTERPRET.

3.2 STANDARD DEVIATION (SQUARE ROOT OF VARIANCE)

STANDARD DEVIATION IS THE **SQUARE ROOT OF VARIANCE**, MAKING IT EASIER TO INTERPRET SINCE IT IS IN THE SAME UNIT AS THE DATA.

FORMULA:

STANDARD DEVIATION(Σ)=VARIANCE\TEXT{STANDARD DEVIATION}
(σ) = \SQRT{\TEXT{VARIANCE}}

📌 EXAMPLE:

USING THE PREVIOUS EXAMPLE,

STANDARD DEVIATION=16.67=4.08\TEXT{STANDARD DEVIATION} =
\SQRT{16.67} = 4.08

📌 USE CASES:

- ✓ RISK ASSESSMENT IN INVESTMENTS – HIGHER STANDARD DEVIATION MEANS HIGHER RISK.
- ✓ QUALITY CONTROL IN MANUFACTURING – DETECTS INCONSISTENCIES IN PRODUCT DIMENSIONS.

💡 CONCLUSION:

STANDARD DEVIATION PROVIDES A CLEAR MEASURE OF SPREAD, WITH REAL-WORLD APPLICATIONS IN FINANCE, SCIENCE, AND ENGINEERING.

📌 CHAPTER 4: PRACTICAL APPLICATIONS & INTERPRETATION

4.1 COMPARING MEAN, MEDIAN, AND MODE

- ✓ SYMMETRIC DATA → MEAN ≈ MEDIAN ≈ MODE
- ✓ RIGHT-SKewed DATA (HIGH VALUES OUTLIERS) → MEAN > MEDIAN > MODE
- ✓ LEFT-SKewed DATA (LOW VALUES OUTLIERS) → MODE > MEDIAN > MEAN

📌 EXAMPLE:

SALARIES IN A COMPANY:

- IF MOST EMPLOYEES EARN \$50K, BUT THE CEO EARNS \$5M, THE MEAN SALARY WILL BE HIGH, BUT THE MEDIAN SALARY WILL BETTER REPRESENT EMPLOYEE EARNINGS.

CONCLUSION:

CHOOSING THE RIGHT MEASURE OF CENTRAL TENDENCY DEPENDS ON DATA DISTRIBUTION AND OUTLIERS.

CHAPTER 5: EXERCISES & ASSIGNMENTS

5.1 MULTIPLE CHOICE QUESTIONS

WHAT IS THE MEAN OF 5, 10, 15, 20?

(A) 10

(B) 15 

(C) 12

(D) 18

WHICH MEASURE IS BEST WHEN DATA HAS OUTLIERS?

(A) MEAN

(B) MEDIAN 

(C) MODE

(D) VARIANCE

5.2 PRACTICAL ASSIGNMENT

TASK:

1. CALCULATE MEAN, MEDIAN, MODE, AND VARIANCE FOR A DATASET.
2. VISUALIZE THE DATA USING HISTOGRAMS AND BOXPLOTS.

3. COMPARE RESULTS AND EXPLAIN WHICH MEASURE IS **MOST
USEFUL FOR THE DATASET.**

ISDM-NxT

PROBABILITY DISTRIBUTIONS (NORMAL, BINOMIAL, POISSON)

CHAPTER 1: INTRODUCTION TO PROBABILITY DISTRIBUTIONS

1.1 WHAT ARE PROBABILITY DISTRIBUTIONS?

A **PROBABILITY DISTRIBUTION** IS A MATHEMATICAL FUNCTION THAT DESCRIBES THE LIKELIHOOD OF DIFFERENT OUTCOMES IN A RANDOM EXPERIMENT. IT ASSIGNS PROBABILITIES TO **ALL POSSIBLE VALUES** THAT A RANDOM VARIABLE CAN TAKE.

IN REAL-WORLD SCENARIOS, PROBABILITY DISTRIBUTIONS HELP MODEL **UNCERTAINTIES** IN FIELDS LIKE **FINANCE, HEALTHCARE, ENGINEERING, AND ARTIFICIAL INTELLIGENCE**.

◆ KEY CONCEPTS IN PROBABILITY DISTRIBUTIONS:

✓ **RANDOM VARIABLE** – A VARIABLE WHOSE VALUES DEPEND ON THE OUTCOMES OF A RANDOM EVENT.

✓ **DISCRETE VS. CONTINUOUS DISTRIBUTIONS** – DISCRETE DISTRIBUTIONS DEAL WITH COUNTABLE VALUES, WHILE CONTINUOUS DISTRIBUTIONS WORK WITH UNCOUNTABLE VALUES.

✓ **MEAN & VARIANCE** – THE EXPECTED VALUE (MEAN) AND SPREAD (VARIANCE) DESCRIBE A DISTRIBUTION'S CHARACTERISTICS.

EXAMPLE:

A DICE ROLL FOLLOWS A **UNIFORM PROBABILITY DISTRIBUTION** WHERE EACH NUMBER (1 TO 6) HAS AN EQUAL PROBABILITY OF **1/6**.

CONCLUSION:

PROBABILITY DISTRIBUTIONS HELP ANALYZE UNCERTAINTY AND MAKE INFORMED PREDICTIONS IN VARIOUS APPLICATIONS.

CHAPTER 2: NORMAL DISTRIBUTION

2.1 WHAT IS NORMAL DISTRIBUTION?

THE **NORMAL DISTRIBUTION**, ALSO KNOWN AS THE **GAUSSIAN DISTRIBUTION**, IS ONE OF THE MOST COMMONLY USED CONTINUOUS PROBABILITY DISTRIBUTIONS. IT IS CHARACTERIZED BY A **BELL-SHAPED CURVE** WHERE MOST VALUES CLUSTER AROUND THE MEAN.

- ◆ **KEY PROPERTIES OF NORMAL DISTRIBUTION:**
- ✓ **SYMMETRIC** – THE LEFT AND RIGHT HALVES OF THE DISTRIBUTION ARE MIRROR IMAGES.
- ✓ **MEAN = MEDIAN = MODE** – THE HIGHEST PEAK IS AT THE CENTER (MEAN).
- ✓ **68-95-99.7 RULE** – DESCRIBES HOW DATA IS SPREAD ACROSS STANDARD DEVIATIONS.

EXAMPLE:

THE HEIGHTS OF ADULTS IN A POPULATION TYPICALLY FOLLOW A NORMAL DISTRIBUTION.

CONCLUSION:

THE NORMAL DISTRIBUTION IS WIDELY USED IN **STATISTICAL MODELING, FINANCE, MACHINE LEARNING, AND HYPOTHESIS TESTING**.

2.2 PROBABILITY DENSITY FUNCTION (PDF) OF NORMAL DISTRIBUTION

THE PROBABILITY DENSITY FUNCTION (PDF) OF A NORMAL DISTRIBUTION IS GIVEN BY:

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

WHERE:

- μ = MEAN
- σ^2 = VARIANCE
- σ = STANDARD DEVIATION

EXAMPLE:

A COMPANY WANTS TO DETERMINE THE PROBABILITY THAT AN EMPLOYEE'S SALARY FALLS WITHIN A CERTAIN RANGE IF SALARIES ARE NORMALLY DISTRIBUTED.

CONCLUSION:

UNDERSTANDING THE PDF HELPS IN COMPUTING PROBABILITIES FOR SPECIFIC VALUES IN NORMALLY DISTRIBUTED DATASETS.

2.3 STANDARD NORMAL DISTRIBUTION & Z-SCORES

THE STANDARD NORMAL DISTRIBUTION IS A NORMAL DISTRIBUTION WITH:

- MEAN (μ) = 0
- STANDARD DEVIATION (σ) = 1

A Z-SCORE MEASURES HOW FAR A VALUE IS FROM THE MEAN IN STANDARD DEVIATIONS:

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

📌 **EXAMPLE:**

IF A STUDENT SCORES **85** ON A TEST WHERE THE MEAN IS **75** AND THE STANDARD DEVIATION IS **5**, THE Z-SCORE IS:

$$Z = \frac{85 - 75}{5} = 2$$

THIS MEANS THE STUDENT'S SCORE IS **2 STANDARD DEVIATIONS ABOVE THE MEAN.**

💡 **CONCLUSION:**

Z-SCORES HELP COMPARE VALUES FROM DIFFERENT NORMAL DISTRIBUTIONS AND CALCULATE PROBABILITIES.

📌 **CHAPTER 3: BINOMIAL DISTRIBUTION**

3.1 WHAT IS BINOMIAL DISTRIBUTION?

THE BINOMIAL DISTRIBUTION IS A DISCRETE PROBABILITY DISTRIBUTION THAT MODELS THE NUMBER OF SUCCESSES IN N INDEPENDENT TRIALS, EACH WITH THE SAME PROBABILITY P.

◆ **KEY PROPERTIES OF BINOMIAL DISTRIBUTION:**

- ✓ **ONLY TWO OUTCOMES** – EACH TRIAL RESULTS IN EITHER SUCCESS OR FAILURE.
- ✓ **FIXED NUMBER OF TRIALS (N)** – THE NUMBER OF EXPERIMENTS IS PREDETERMINED.
- ✓ **CONSTANT PROBABILITY (P)** – THE PROBABILITY OF SUCCESS REMAINS THE SAME FOR EACH TRIAL.

💡 EXAMPLE:

FLIPPING A FAIR COIN **10 TIMES**, WHERE THE PROBABILITY OF HEADS (SUCCESS) IS **0.5**, FOLLOWS A BINOMIAL DISTRIBUTION.

💡 CONCLUSION:

THE BINOMIAL DISTRIBUTION IS USEFUL IN SCENARIOS LIKE **QUALITY CONTROL, MEDICAL TESTING, AND SPORTS ANALYSIS.**

3.2 PROBABILITY MASS FUNCTION (PMF) OF BINOMIAL DISTRIBUTION

THE PROBABILITY MASS FUNCTION (PMF) OF A BINOMIAL DISTRIBUTION IS:

$$P(X=K) = (NK)P^k(1-P)^{N-K} = \text{BINOM}\{N\}\{K\} P^K (1 - P)^{N - K}$$

WHERE:

- **NN** = NUMBER OF TRIALS
- **KK** = NUMBER OF SUCCESSES
- **PP** = PROBABILITY OF SUCCESS IN A SINGLE TRIAL

💡 EXAMPLE:

A COMPANY TESTS **5 COMPUTER CHIPS**, AND EACH CHIP HAS A **90%** PROBABILITY OF FUNCTIONING CORRECTLY. THE PROBABILITY OF EXACTLY **4** FUNCTIONING CHIPS FOLLOWS:

$$P(X=4) = (54)(0.9)^4(0.1)^1 P(X = 4) = \text{BINOM}\{5\}\{4\} (0.9)^4 (0.1)^1$$

💡 CONCLUSION:

THE BINOMIAL PMF CALCULATES THE PROBABILITY OF ACHIEVING A SPECIFIC NUMBER OF SUCCESSES IN A SET OF TRIALS.

📌 CHAPTER 4: POISSON DISTRIBUTION

4.1 WHAT IS POISSON DISTRIBUTION?

THE POISSON DISTRIBUTION MODELS THE PROBABILITY OF A GIVEN NUMBER OF EVENTS OCCURRING IN A FIXED INTERVAL WHEN EVENTS HAPPEN INDEPENDENTLY AT A CONSTANT RATE.

◆ KEY PROPERTIES OF POISSON DISTRIBUTION:

- ✓ MODELS RARE EVENTS – USED WHEN EVENTS OCCUR RANDOMLY OVER TIME OR SPACE.
- ✓ ONLY ONE PARAMETER (λ) – REPRESENTS THE AVERAGE EVENT OCCURRENCE RATE.
- ✓ NO UPPER LIMIT – THEORETICALLY, ANY NUMBER OF EVENTS CAN OCCUR.

📌 EXAMPLE:

THE NUMBER OF CUSTOMER CALLS RECEIVED BY A SUPPORT CENTER IN AN HOUR FOLLOWS A POISSON DISTRIBUTION.

💡 CONCLUSION:

THE POISSON DISTRIBUTION IS WIDELY USED IN TRAFFIC ANALYSIS, NETWORK SECURITY, AND MEDICAL RESEARCH.

4.2 PROBABILITY MASS FUNCTION (PMF) OF POISSON DISTRIBUTION

THE POISSON PMF IS GIVEN BY:

$$P(X=k) = e^{-\lambda} \frac{\lambda^k}{k!}$$

WHERE:

- λ = AVERAGE NUMBER OF EVENTS PER INTERVAL

- $K = \text{NUMBER OF EVENTS OCCURRING}$

 **EXAMPLE:**

IF AN AIRPORT RECEIVES **2** EMERGENCY CALLS PER DAY ON AVERAGE, THE PROBABILITY OF RECEIVING EXACTLY **3** CALLS IN A DAY IS:

$$P(X=3) = e^{-2} (2)^3 / 3! P(X = 3) = \frac{e^{-2} (2)^3}{3!}$$

 **CONCLUSION:**

THE **Poisson PMF** HELPS PREDICT HOW OFTEN AN EVENT WILL OCCUR IN A GIVEN TIME FRAME.

 **CHAPTER 5: COMPARISON OF DISTRIBUTIONS**

FEATURE	NORMAL DISTRIBUTION	BINOMIAL DISTRIBUTION	POISSON DISTRIBUTION
TYPE	CONTINUOUS	DISCRETE	DISCRETE
USE CASE	MEASURING VALUES (E.G., HEIGHTS, TEST SCORES)	COUNTING SUCCESSES IN FIXED TRIALS	COUNTING RARE EVENTS OVER TIME
SHAPE	BELL-SHAPED	VARIABLES (SKEWED FOR SMALL P, SYMMETRIC FOR P ≈ 0.5)	RIGHT-SKewed

 **EXAMPLE:**

- **NORMAL DISTRIBUTION – HEIGHTS OF STUDENTS IN A UNIVERSITY.**

- **BINOMIAL DISTRIBUTION** – NUMBER OF DEFECTIVE PRODUCTS IN A BATCH.
- **POISSON DISTRIBUTION** – NUMBER OF EMAILS RECEIVED PER HOUR.

 **CONCLUSION:**

EACH DISTRIBUTION IS USED FOR DIFFERENT TYPES OF **REAL-WORLD PROBLEMS AND DATA ANALYSIS TASKS.**

 **CHAPTER 6: EXERCISES & ASSIGNMENTS**

6.1 MULTIPLE CHOICE QUESTIONS

 **WHICH DISTRIBUTION MODELS CONTINUOUS DATA?**

- (A) BINOMIAL
- (B) POISSON
- (C) NORMAL 

 **WHICH PARAMETER REPRESENTS SUCCESS PROBABILITY IN A BINOMIAL DISTRIBUTION?**

- (A) N
- (B) P 
- (C) λ

6.2 PRACTICAL ASSIGNMENT

 **TASK 1: SIMULATE A BINOMIAL, POISSON, AND NORMAL DISTRIBUTION USING PYTHON.**

 **TASK 2: CALCULATE THE PROBABILITY OF SPECIFIC OUTCOMES USING GIVEN DISTRIBUTIONS.**

🌟 CONCLUSION: THE IMPORTANCE OF PROBABILITY DISTRIBUTIONS

PROBABILITY DISTRIBUTIONS HELP MODEL UNCERTAINTY AND MAKE PREDICTIONS IN FINANCE, SCIENCE, ENGINEERING, AND AI. UNDERSTANDING THEM ENABLES BETTER DECISION-MAKING, STATISTICAL INFERENCE, AND MACHINE LEARNING APPLICATIONS.





HYPOTHESIS TESTING & CONFIDENCE INTERVALS

📌 CHAPTER 1: INTRODUCTION TO HYPOTHESIS TESTING

1.1 WHAT IS HYPOTHESIS TESTING?

HYPOTHESIS TESTING IS A **STATISTICAL METHOD** USED TO DETERMINE WHETHER THERE IS ENOUGH EVIDENCE IN A SAMPLE TO SUPPORT OR REJECT A CLAIM ABOUT A POPULATION. IT HELPS ANALYSTS AND RESEARCHERS MAKE **DATA-DRIVEN DECISIONS** BY EVALUATING ASSUMPTIONS USING SAMPLE DATA.

KEY USES OF HYPOTHESIS TESTING:

- ✓ COMPARING TWO OR MORE GROUPS (E.G., A/B TESTING IN MARKETING).
- ✓ DETERMINING IF A NEW DRUG IS MORE EFFECTIVE THAN AN EXISTING ONE.
- ✓ ASSESSING BUSINESS PERFORMANCE BEFORE AND AFTER A STRATEGIC CHANGE.

📌 EXAMPLE:

A COMPANY LAUNCHES A NEW MARKETING CAMPAIGN AND WANTS TO TEST IF IT SIGNIFICANTLY INCREASES SALES. HYPOTHESIS TESTING CAN HELP DETERMINE IF THE INCREASE IS DUE TO THE CAMPAIGN OR JUST RANDOM CHANCE.

💡 CONCLUSION:

HYPOTHESIS TESTING ELIMINATES GUESSWORK AND ENSURES DECISIONS ARE BACKED BY STATISTICAL EVIDENCE RATHER THAN INTUITION.

📌 CHAPTER 2: KEY CONCEPTS IN HYPOTHESIS TESTING

2.1 NULL HYPOTHESIS (H_0) AND ALTERNATIVE HYPOTHESIS (H_1)

◆ NULL HYPOTHESIS (H_0):

- A DEFAULT ASSUMPTION THAT THERE IS NO SIGNIFICANT DIFFERENCE OR NO EFFECT.
- IT REPRESENTS THE STATUS QUO OR EXISTING BELIEF.

◆ ALTERNATIVE HYPOTHESIS (H_1 OR H_A):

- A CLAIM THAT CONTRADICTS H_0 , SUGGESTING THERE IS A SIGNIFICANT EFFECT OR DIFFERENCE.

📌 EXAMPLE:

A SCHOOL TESTS WHETHER A NEW TEACHING METHOD IMPROVES STUDENT PERFORMANCE.

- H_0 : THE NEW METHOD DOES NOT IMPROVE SCORES.
- H_1 : THE NEW METHOD DOES IMPROVE SCORES.

💡 CONCLUSION:

A HYPOTHESIS TEST EVALUATES WHETHER TO REJECT OR FAIL TO REJECT H_0 BASED ON SAMPLE DATA.

2.2 TYPES OF ERRORS IN HYPOTHESIS TESTING

HYPOTHESIS TESTING IS NOT FOOLPROOF, AND ERRORS CAN OCCUR.

◆ TYPE I ERROR (FALSE POSITIVE):

- REJECTING H_0 WHEN IT IS ACTUALLY TRUE.

- EXAMPLE: A DOCTOR DIAGNOSES A HEALTHY PERSON WITH A DISEASE.
- ◆ TYPE II ERROR (FALSE NEGATIVE):
 - FAILING TO REJECT H_0 WHEN IT IS ACTUALLY FALSE.
 - EXAMPLE: A FAULTY FIRE ALARM FAILS TO DETECT A FIRE.

📌 EXAMPLE:

IN DRUG TESTING, A TYPE I ERROR WOULD MEAN APPROVING AN INEFFECTIVE DRUG, WHILE A TYPE II ERROR WOULD MEAN REJECTING A BENEFICIAL DRUG.

💡 CONCLUSION:

MINIMIZING THESE ERRORS IS CRUCIAL FOR MAKING ACCURATE AND RELIABLE DECISIONS.

📌 CHAPTER 3: HYPOTHESIS TESTING PROCESS

3.1 STEPS IN HYPOTHESIS TESTING

- ✓ STEP 1: DEFINE THE HYPOTHESES (H_0 AND H_1).
- ✓ STEP 2: CHOOSE A SIGNIFICANCE LEVEL (α).
- ✓ STEP 3: SELECT A TEST STATISTIC (E.G., T-TEST, CHI-SQUARE TEST).
- ✓ STEP 4: COMPUTE THE P-VALUE AND COMPARE WITH α .
- ✓ STEP 5: MAKE A DECISION (REJECT OR FAIL TO REJECT H_0).

📌 EXAMPLE:

A WEBSITE RUNS AN A/B TEST TO COMPARE TWO HOMEPAGE DESIGNS.

- H_0 : THE NEW DESIGN DOES NOT INCREASE USER ENGAGEMENT.

- H_1 : THE NEW DESIGN DOES INCREASE USER ENGAGEMENT.

IF THE P-VALUE < 0.05, THE COMPANY REJECTS H_0 AND CONCLUDES THAT THE NEW DESIGN IS BETTER.

CONCLUSION:

A STRUCTURED HYPOTHESIS TESTING PROCESS ENSURES ACCURATE AND UNBIASED DECISION-MAKING.

CHAPTER 4: TYPES OF HYPOTHESIS TESTS

4.1 PARAMETRIC VS. NON-PARAMETRIC TESTS

◆ PARAMETRIC TESTS:

- ASSUMES DATA FOLLOWS A NORMAL DISTRIBUTION.
- EXAMPLES: T-TEST, ANOVA, Z-TEST.

◆ NON-PARAMETRIC TESTS:

- DOES NOT REQUIRE NORMALITY ASSUMPTIONS.
- EXAMPLES: MANN-WHITNEY U TEST, KRUSKAL-WALLIS TEST, CHI-SQUARE TEST.

CONCLUSION:

CHOOSING THE RIGHT TEST DEPENDS ON THE DATA TYPE AND DISTRIBUTION.

4.2 COMMON HYPOTHESIS TESTS

- ◆ **Z-TEST:** USED FOR COMPARING POPULATION MEANS WHEN THE SAMPLE SIZE IS LARGE (>30).

◆ **T-TEST:** USED WHEN THE SAMPLE SIZE IS SMALL (<30). TYPES INCLUDE:

✓ **ONE-SAMPLE T-TEST** – COMPARES A SAMPLE MEAN TO A KNOWN POPULATION MEAN.

✓ **INDEPENDENT T-TEST** – COMPARES MEANS BETWEEN TWO INDEPENDENT GROUPS.

✓ **PAIRED T-TEST** – COMPARES MEANS OF THE SAME GROUP BEFORE AND AFTER A CHANGE.

◆ **ANOVA (ANALYSIS OF VARIANCE):**

- COMPARES MORE THAN TWO GROUP MEANS (E.G., TESTING THREE MARKETING CAMPAIGNS).

◆ **CHI-SQUARE TEST:**

- USED FOR CATEGORICAL DATA (E.G., TESTING GENDER PREFERENCE IN PRODUCT CHOICES).

📌 **EXAMPLE:**

A UNIVERSITY WANTS TO KNOW IF STUDENTS FROM DIFFERENT MAJORS PERFORM DIFFERENTLY ON AN ENTRANCE EXAM. THEY USE **ANOVA** TO TEST FOR SIGNIFICANT DIFFERENCES.

💡 **CONCLUSION:**

USING THE APPROPRIATE TEST ENSURES ACCURATE STATISTICAL INFERENCES.

📌 **CHAPTER 5: INTRODUCTION TO CONFIDENCE INTERVALS**

5.1 WHAT IS A CONFIDENCE INTERVAL?

A CONFIDENCE INTERVAL (CI) IS A RANGE OF VALUES THAT LIKELY CONTAINS THE TRUE POPULATION PARAMETER. IT PROVIDES AN ESTIMATE OF WHERE THE TRUE MEAN OR PROPORTION LIES.

✓ A 95% CONFIDENCE INTERVAL MEANS THAT IF WE REPEATED THE TEST MULTIPLE TIMES, 95% OF THE INTERVALS WOULD CONTAIN THE TRUE VALUE.

 **EXAMPLE:**

A RESEARCHER ESTIMATES THAT THE AVERAGE WEIGHT OF ADULTS IN A CITY IS **70 KG ± 2 KG** WITH 95% CONFIDENCE. THIS MEANS THE TRUE AVERAGE IS LIKELY BETWEEN **68 KG AND 72 KG**.

 **CONCLUSION:**

CONFIDENCE INTERVALS HELP QUANTIFY UNCERTAINTY IN STATISTICAL ESTIMATES.

5.2 FORMULA FOR CONFIDENCE INTERVALS

FOR A POPULATION MEAN (μ):

$$CI = \bar{X} \pm (Z_{\alpha/2} \times \frac{\sigma}{\sqrt{n}})$$

WHERE:

- \bar{X} = SAMPLE MEAN
- $Z_{\alpha/2}$ = CRITICAL VALUE FROM Z-TABLE
- σ = STANDARD DEVIATION
- n = SAMPLE SIZE

 **EXAMPLE:**

A COMPANY SURVEYS 100 EMPLOYEES ABOUT JOB SATISFACTION.

THE AVERAGE RATING IS **4.2/5** WITH A STANDARD DEVIATION OF **0.5**.
THE **95% CONFIDENCE INTERVAL** FOR JOB SATISFACTION IS
CALCULATED AS:

$$\text{CI} = 4.2 \pm (1.96 \times 0.5 / \sqrt{100}) \\ \text{CI} = 4.2 \pm 0.098 = 4.2 \pm 0.098 = (4.1, 4.3) = (4.1, 4.3)$$

THIS MEANS THE TRUE JOB SATISFACTION LEVEL IS LIKELY BETWEEN **4.1 AND 4.3**.

💡 CONCLUSION:

CONFIDENCE INTERVALS PROVIDE AN ESTIMATE WITH A MARGIN OF ERROR, MAKING DATA-DRIVEN DECISIONS MORE RELIABLE.

📌 CHAPTER 6: EXERCISES & ASSIGNMENTS

6.1 MULTIPLE CHOICE QUESTIONS

WHICH OF THE FOLLOWING DESCRIBES THE NULL HYPOTHESIS (H_0)?

- (A) THE CLAIM WE ARE TESTING
- (B) THE ASSUMPTION OF NO EFFECT
- (C) A RANDOM GUESS

WHICH TEST COMPARES TWO INDEPENDENT SAMPLE MEANS?

- (A) Z-TEST
- (B) T-TEST
- (C) CHI-SQUARE TEST

WHAT DOES A 95% CONFIDENCE INTERVAL MEAN?

- (A) THE PROBABILITY THAT THE SAMPLE MEAN IS 95% CORRECT
- (B) THE RANGE IN WHICH THE TRUE POPULATION MEAN LIES 95%

OF THE TIME

(C) A PREDICTION OF FUTURE VALUES

6.2 PRACTICAL ASSIGNMENT

📌 **TASK 1:** CONDUCT A HYPOTHESIS TEST ON A DATASET

COMPARING TWO GROUPS (E.G., MALE VS. FEMALE AVERAGE SALARIES).

📌 **TASK 2:** CALCULATE A **95%** CONFIDENCE INTERVAL FOR AN AVERAGE TEST SCORE IN A DATASET.

🌟 CONCLUSION: THE IMPORTANCE OF STATISTICAL TESTING

HYPOTHESIS TESTING AND CONFIDENCE INTERVALS ARE FUNDAMENTAL TOOLS IN DATA-DRIVEN DECISION-MAKING. THEY HELP BUSINESSES, RESEARCHERS, AND POLICYMAKERS MAKE RELIABLE AND STATISTICALLY BACKED CONCLUSIONS.



DATA VISUALIZATION WITH MATPLOTLIB, SEABORN, POWER BI, AND TABLEAU

CHAPTER 1: INTRODUCTION TO DATA VISUALIZATION

1.1 WHAT IS DATA VISUALIZATION?

DATA VISUALIZATION IS THE PROCESS OF **REPRESENTING COMPLEX DATA IN GRAPHICAL FORMATS SUCH AS CHARTS, GRAPHS, AND INTERACTIVE DASHBOARDS.** IT HELPS ORGANIZATIONS **INTERPRET LARGE DATASETS EFFICIENTLY, UNCOVER HIDDEN TRENDS, AND SUPPORT DATA-DRIVEN DECISION-MAKING.**

- ◆ **KEY BENEFITS OF DATA VISUALIZATION:**
- ✓ ENHANCES DATA INTERPRETATION BY MAKING PATTERNS AND TRENDS VISIBLE.
- ✓ IMPROVES DECISION-MAKING THROUGH CLEAR INSIGHTS.
- ✓ ENABLES QUICK IDENTIFICATION OF OUTLIERS, CORRELATIONS, AND TRENDS.
- ✓ FACILITATES BETTER STORYTELLING IN BUSINESS REPORTS AND PRESENTATIONS.

EXAMPLE:

A COMPANY VISUALIZING **MONTHLY SALES PERFORMANCE** USING A BAR CHART CAN QUICKLY IDENTIFY PEAK SALES PERIODS AND LOW-PERFORMING MONTHS.

CONCLUSION:

DATA VISUALIZATION BRIDGES THE GAP BETWEEN **RAW DATA AND**

MEANINGFUL INSIGHTS, MAKING DATA ACTIONABLE AND ACCESSIBLE FOR BUSINESSES.

📌 CHAPTER 2: DATA VISUALIZATION TECHNIQUES

2.1 COMMON TYPES OF DATA VISUALIZATIONS

- ◆ **BAR CHARTS** – IDEAL FOR CATEGORICAL COMPARISONS (E.G., SALES BY REGION).
- ◆ **LINE GRAPHS** – BEST FOR TIME-SERIES DATA (E.G., STOCK PRICE TRENDS).
- ◆ **HISTOGRAMS** – USED TO DISPLAY FREQUENCY DISTRIBUTIONS OF NUMERICAL DATA.
- ◆ **SCATTER PLOTS** – HELPS IN IDENTIFYING RELATIONSHIPS BETWEEN TWO VARIABLES.
- ◆ **PIE CHARTS** – DISPLAYS PROPORTIONS (THOUGH NOT IDEAL FOR DETAILED ANALYSIS).
- ◆ **HEATMAPS** – SHOWS CORRELATIONS BETWEEN VARIABLES USING COLOR CODING.

📌 EXAMPLE:

A LINE GRAPH CAN HELP VISUALIZE WEBSITE TRAFFIC TRENDS OVER TIME, REVEALING SEASONAL FLUCTUATIONS.

💡 CONCLUSION:

CHOOSING THE RIGHT VISUALIZATION DEPENDS ON THE TYPE OF DATA AND INSIGHTS NEEDED.

📌 CHAPTER 3: DATA VISUALIZATION USING MATPLOTLIB

3.1 WHAT IS MATPLOTLIB?

MATPLOTLIB IS A PYTHON LIBRARY USED FOR CREATING STATIC, ANIMATED, AND INTERACTIVE VISUALIZATIONS. IT IS HIGHLY CUSTOMIZABLE AND INTEGRATES WELL WITH NUMPY AND PANDAS.

3.2 KEY FEATURES OF MATPLOTLIB

- ✓ PROVIDES LINE, BAR, SCATTER, AND PIE CHARTS.
- ✓ ALLOWS CUSTOMIZATION OF LABELS, COLORS, AND ANNOTATIONS.
- ✓ SUPPORTS SUBPLOTTING FOR MULTIPLE VISUALIZATIONS IN ONE FIGURE.

3.3 COMMON MATPLOTLIB FUNCTIONS

📌 LINE CHART EXAMPLE:

```
IMPORT MATPLOTLIB.PYPLOT AS PLT
```

```
# DATA
```

```
MONTHS = ['JAN', 'FEB', 'MAR', 'APR', 'MAY']
```

```
SALES = [25000, 27000, 30000, 28000, 32000]
```

```
# PLOT
```

```
PLT.PLOT(MONTHS, SALES, MARKER='O', LINESTYLE='-', COLOR='B')
```

```
PLT.XLABEL("MONTHS")
```

```
PLT.YLABEL("SALES ($)")
```

```
PLT.TITLE("MONTHLY SALES TRENDS")
```

```
PLT.GRID(TRUE)
```

PLT.SHOW()

📌 **EXAMPLE:**

A SCATTER PLOT IN MATPLOTLIB HELPS VISUALIZE THE RELATIONSHIP BETWEEN MARKETING SPEND AND REVENUE.

💡 **CONCLUSION:**

MATPLOTLIB IS A POWERFUL AND FLEXIBLE VISUALIZATION TOOL, COMMONLY USED FOR STATIC VISUALIZATIONS IN PYTHON.

📌 **CHAPTER 4: DATA VISUALIZATION USING SEABORN**

4.1 WHAT IS SEABORN?

SEABORN IS BUILT ON TOP OF MATPLOTLIB AND PROVIDES STATISTICAL VISUALIZATIONS WITH IMPROVED AESTHETICS.

4.2 KEY FEATURES OF SEABORN

- ✓ BUILT-IN THEMES FOR VISUALLY APPEALING PLOTS.
- ✓ STATISTICAL FUNCTIONS FOR CORRELATIONS AND DISTRIBUTIONS.
- ✓ HEATMAPS, PAIR PLOTS, AND VIOLIN PLOTS FOR ADVANCED ANALYSIS.

4.3 COMMON SEABORN FUNCTIONS

📌 **HISTOGRAM (DISTRIBUTION PLOT) EXAMPLE:**

IMPORT SEABORN AS SNS

IMPORT MATPLOTLIB.PYTHON AS PLT

SAMPLE DATA

```
DATA = SNS.LOAD_DATASET("TIPS")  
  
# HISTOGRAM  
  
SNS.HISTPLOT(DATA['TOTAL_BILL'], KDE=TRUE, BINS=20,  
COLOR='BLUE')  
  
PLT.TITLE("DISTRIBUTION OF TOTAL BILL AMOUNT")  
  
PLT.XLABEL("TOTAL BILL ($)")  
  
PLT.YLABEL("FREQUENCY")  
  
PLT.SHOW()
```

 **EXAMPLE:**

A HEATMAP IN SEABORN HELPS VISUALIZE CORRELATIONS BETWEEN DIFFERENT NUMERICAL FEATURES IN A DATASET.

 **CONCLUSION:**

SEABORN ENHANCES DATA VISUALIZATION IN PYTHON WITH RICH, EASY-TO-INTERPRET PLOTS THAT MAKE STATISTICAL ANALYSIS MORE EFFECTIVE.

 **CHAPTER 5: DATA VISUALIZATION USING POWER BI**

5.1 WHAT IS POWER BI?

POWER BI IS A BUSINESS INTELLIGENCE (BI) TOOL BY MICROSOFT THAT ALLOWS USERS TO CREATE INTERACTIVE DASHBOARDS AND REPORTS.

5.2 KEY FEATURES OF POWER BI

- ✓ DRAG-AND-DROP FUNCTIONALITY FOR CREATING REPORTS.
- ✓ CONNECTIVITY WITH DATABASES, EXCEL, AND CLOUD PLATFORMS.
- ✓ INTERACTIVE DASHBOARDS WITH REAL-TIME UPDATES.
- ✓ ADVANCED AI-POWERED INSIGHTS AND PREDICTIVE ANALYTICS.

5.3 CREATING A POWER BI DASHBOARD

- ◆ STEP 1: IMPORT DATA FROM EXCEL, SQL SERVER, OR CLOUD SOURCES.
- ◆ STEP 2: USE BAR CHARTS, LINE GRAPHS, AND KPIs TO VISUALIZE KEY METRICS.
- ◆ STEP 3: ADD FILTERS AND SLICERS FOR USER INTERACTIVITY.
- ◆ STEP 4: PUBLISH AND SHARE DASHBOARDS VIA POWER BI SERVICE.

📌 EXAMPLE:

A RETAIL COMPANY USES POWER BI TO TRACK SALES PERFORMANCE ACROSS DIFFERENT STORES IN REAL-TIME.

💡 CONCLUSION:

POWER BI IS A POWERFUL BI TOOL THAT TRANSFORMS DATA INTO INTERACTIVE, SHAREABLE DASHBOARDS.

📌 CHAPTER 6: DATA VISUALIZATION USING TABLEAU

6.1 WHAT IS TABLEAU?

TABLEAU IS A DATA VISUALIZATION SOFTWARE THAT HELPS USERS ANALYZE AND PRESENT DATA THROUGH INTERACTIVE CHARTS AND REPORTS.

6.2 KEY FEATURES OF TABLEAU

- ✓ DRAG-AND-DROP INTERFACE FOR CREATING CHARTS.
- ✓ INTEGRATION WITH MULTIPLE DATA SOURCES (EXCEL, DATABASES, APIs).
- ✓ AI-POWERED INSIGHTS FOR PREDICTIVE ANALYSIS.
- ✓ STORYTELLING AND DASHBOARD FUNCTIONALITIES.

6.3 CREATING A TABLEAU DASHBOARD

- ◆ STEP 1: CONNECT TO A DATA SOURCE (EXCEL, SQL, GOOGLE SHEETS).
- ◆ STEP 2: CREATE VISUALIZATIONS USING BAR CHARTS, LINE GRAPHS, AND MAPS.
- ◆ STEP 3: BUILD A DASHBOARD WITH MULTIPLE CHARTS AND KPIs.
- ◆ STEP 4: PUBLISH THE DASHBOARD ON TABLEAU SERVER OR TABLEAU PUBLIC.

 **EXAMPLE:**

A FINANCIAL ANALYST USES TABLEAU TO VISUALIZE MARKET TRENDS AND FORECAST STOCK PRICES.

 **CONCLUSION:**

TABLEAU IS A LEADING BI TOOL FOR CREATING INTERACTIVE DATA VISUALIZATIONS AND ANALYTICS DASHBOARDS.

 **CHAPTER 7: EXERCISES & ASSIGNMENTS**

7.1 MULTIPLE CHOICE QUESTIONS

WHICH PYTHON LIBRARY IS USED FOR CREATING STATISTICAL VISUALIZATIONS?

(A) MATPLOTLIB

(B) SEABORN 

(C) POWER BI

 WHICH VISUALIZATION TOOL ALLOWS INTERACTIVE DASHBOARDS?

(A) POWER BI 

(B) MATPLOTLIB

(C) NUMPY

 WHAT TYPE OF CHART IS BEST FOR VISUALIZING TIME-SERIES DATA?

(A) PIE CHART

(B) LINE CHART 

(C) HISTOGRAM

7.2 PRACTICAL ASSIGNMENT

 **TASK 1:** CREATE A LINE CHART IN MATPLOTLIB TO VISUALIZE MONTHLY SALES TRENDS.

 **TASK 2:** GENERATE A HEATMAP IN SEABORN TO SHOW CORRELATION BETWEEN NUMERICAL FEATURES.

 **TASK 3:** BUILD AN INTERACTIVE POWER BI DASHBOARD TO ANALYZE FINANCIAL DATA.

 **TASK 4:** CREATE A TABLEAU DASHBOARD SHOWCASING E-COMMERCE SALES PERFORMANCE.

CONCLUSION

DATA VISUALIZATION IS KEY TO MAKING SENSE OF COMPLEX DATA.
USING MATPLOTLIB AND SEABORN FOR PYTHON VISUALIZATIONS,

AND POWER BI AND TABLEAU FOR INTERACTIVE DASHBOARDS,
PROFESSIONALS CAN GAIN DEEP INSIGHTS AND MAKE DATA-DRIVEN
DECISIONS EFFICIENTLY.





CORRELATION & REGRESSION ANALYSIS

📌 CHAPTER 1: UNDERSTANDING CORRELATION & REGRESSION ANALYSIS

1.1 WHAT IS CORRELATION & REGRESSION ANALYSIS?

CORRELATION AND REGRESSION ANALYSIS ARE STATISTICAL TECHNIQUES USED TO EXAMINE THE RELATIONSHIP BETWEEN TWO OR MORE VARIABLES. THESE METHODS HELP IN UNDERSTANDING HOW VARIABLES INFLUENCE EACH OTHER AND PREDICTING FUTURE OUTCOMES BASED ON EXISTING DATA.

- ◆ **CORRELATION** MEASURES THE STRENGTH AND DIRECTION OF A RELATIONSHIP BETWEEN TWO VARIABLES.
- ◆ **REGRESSION** HELPS MODEL THE RELATIONSHIP AND PREDICT VALUES BASED ON INDEPENDENT VARIABLES.

📌 EXAMPLE:

A COMPANY MIGHT ANALYZE THE CORRELATION BETWEEN ADVERTISING SPEND AND SALES REVENUE TO DETERMINE HOW EFFECTIVELY THEIR MARKETING EFFORTS INFLUENCE CUSTOMER PURCHASES.

💡 CONCLUSION:

BOTH **CORRELATION & REGRESSION** PLAY A CRUCIAL ROLE IN **DATA SCIENCE, BUSINESS ANALYTICS, AND DECISION-MAKING**, HELPING ORGANIZATIONS MAKE DATA-DRIVEN DECISIONS.

📌 CHAPTER 2: INTRODUCTION TO CORRELATION ANALYSIS

2.1 WHAT IS CORRELATION?

CORRELATION MEASURES HOW STRONGLY TWO VARIABLES ARE RELATED AND WHETHER THEY MOVE TOGETHER. THE CORRELATION COEFFICIENT (R) IS A NUMERICAL VALUE BETWEEN **-1 AND +1**, REPRESENTING THE RELATIONSHIP STRENGTH AND DIRECTION.

- ◆ **POSITIVE CORRELATION ($R > 0$):** BOTH VARIABLES INCREASE TOGETHER.
- ◆ **NEGATIVE CORRELATION ($R < 0$):** ONE VARIABLE INCREASES WHILE THE OTHER DECREASES.
- ◆ **NO CORRELATION ($R \approx 0$):** NO RELATIONSHIP EXISTS BETWEEN THE VARIABLES.

📌 EXAMPLE:

- **STRONG POSITIVE CORRELATION:** MORE STUDY HOURS ↑ → HIGHER TEST SCORES ↑
- **STRONG NEGATIVE CORRELATION:** INCREASED SPEED ↑ → LESS TRAVEL TIME ↓
- **NO CORRELATION:** A PERSON'S HEIGHT DOES NOT AFFECT THEIR PHONE NUMBER.

💡 CONCLUSION:

CORRELATION HELPS IN UNDERSTANDING VARIABLE RELATIONSHIPS, BUT IT DOES NOT IMPLY CAUSATION—JUST BECAUSE TWO VARIABLES ARE CORRELATED DOESN'T MEAN ONE CAUSES THE OTHER.

📌 CHAPTER 3: TYPES OF CORRELATION

3.1 PEARSON'S CORRELATION COEFFICIENT (R)

PEARSON'S R IS THE MOST WIDELY USED CORRELATION MEASURE AND IS BASED ON LINEAR RELATIONSHIPS BETWEEN VARIABLES.

✓ FORMULA:

$$R = \frac{\sum(X - \bar{X})(Y - \bar{Y})}{\sqrt{\sum(X - \bar{X})^2} \sqrt{\sum(Y - \bar{Y})^2}}$$

✓ VALUES INTERPRETATION:

- $R = +1 \rightarrow$ PERFECT POSITIVE CORRELATION
- $R = -1 \rightarrow$ PERFECT NEGATIVE CORRELATION
- $R = 0 \rightarrow$ NO CORRELATION

📌 EXAMPLE:

IF $R = 0.85$ BETWEEN ADVERTISING SPEND AND REVENUE, IT INDICATES A STRONG POSITIVE CORRELATION—INCREASING ADS LEADS TO HIGHER REVENUE.

3.2 SPEARMAN'S RANK CORRELATION

USED FOR ORDINAL (RANKED) DATA, INSTEAD OF ASSUMING A LINEAR RELATIONSHIP LIKE PEARSON'S R.

📌 EXAMPLE:

IF STUDENTS ARE RANKED BASED ON MATH AND SCIENCE SCORES, SPEARMAN'S CORRELATION MEASURES HOW WELL THESE RANKINGS MATCH.

CONCLUSION:

CHOOSING THE RIGHT CORRELATION METHOD DEPENDS ON WHETHER DATA IS LINEAR, RANKED, OR NONLINEAR.

CHAPTER 4: INTRODUCTION TO REGRESSION ANALYSIS

4.1 WHAT IS REGRESSION?

REGRESSION IS A PREDICTIVE MODELING TECHNIQUE THAT ESTIMATES THE RELATIONSHIP BETWEEN A DEPENDENT VARIABLE (Y) AND ONE OR MORE INDEPENDENT VARIABLES (X).

- ◆ **SIMPLE LINEAR REGRESSION** – RELATIONSHIP BETWEEN ONE INDEPENDENT VARIABLE (X) AND A DEPENDENT VARIABLE (Y).
- ◆ **MULTIPLE REGRESSION** – RELATIONSHIP BETWEEN MULTIPLE INDEPENDENT VARIABLES ($X_1, X_2, X_3\dots$) AND A DEPENDENT VARIABLE (Y).

EXAMPLE:

- **SIMPLE REGRESSION:** PREDICTING HOUSE PRICE (Y) BASED ON SQUARE FOOTAGE (X).
- **MULTIPLE REGRESSION:** PREDICTING HOUSE PRICE (Y) USING SQUARE FOOTAGE, LOCATION, AND NUMBER OF BEDROOMS (X_1, X_2, X_3).

CONCLUSION:

REGRESSION IS ESSENTIAL FOR FORECASTING, DECISION-MAKING, AND RISK ANALYSIS IN INDUSTRIES LIKE FINANCE, HEALTHCARE, AND MARKETING.

📌 CHAPTER 5: LINEAR REGRESSION ANALYSIS

5.1 SIMPLE LINEAR REGRESSION

A STATISTICAL TECHNIQUE TO MODEL THE RELATIONSHIP BETWEEN TWO CONTINUOUS VARIABLES (X AND Y).

✓ EQUATION OF SIMPLE LINEAR REGRESSION:

$$Y = MX + B$$

WHERE:

- **Y** = DEPENDENT VARIABLE (WHAT WE PREDICT)
- **X** = INDEPENDENT VARIABLE (PREDICTOR)
- **M** = SLOPE (CHANGE IN Y PER UNIT CHANGE IN X)
- **B** = INTERCEPT (Y-VALUE WHEN X=0)

📌 EXAMPLE:

A COFFEE SHOP MODELS SALES REVENUE (Y) BASED ON NUMBER OF CUSTOMERS (X).

✓ INTERPRETATION:

- IF **M = 2**, IT MEANS THAT FOR EVERY ADDITIONAL CUSTOMER, REVENUE INCREASES BY \$2.

💡 CONCLUSION:

SIMPLE LINEAR REGRESSION IS USEFUL FOR TREND ANALYSIS AND FORECASTING WHEN ONE PREDICTOR INFLUENCES AN OUTCOME.

5.2 MULTIPLE LINEAR REGRESSION

EXTENDS SIMPLE REGRESSION BY INCLUDING **MULTIPLE INDEPENDENT VARIABLES.**

✓ **EQUATION OF MULTIPLE REGRESSION:**

$$Y = B_0 + B_1X_1 + B_2X_2 + \dots + B_NX_N \quad Y = B_0 + B_1X_1 + B_2X_2 + \dots + B_NX_N$$

📌 **EXAMPLE:**

PREDICTING HOUSE PRICES BASED ON SQUARE FOOTAGE (X_1), NUMBER OF BEDROOMS (X_2), AND LOCATION SCORE (X_3).

✓ **INTERPRETATION:**

EACH B COEFFICIENT MEASURES THE IMPACT OF EACH INDEPENDENT VARIABLE.

💡 **CONCLUSION:**

MULTIPLE REGRESSION ALLOWS MORE ACCURATE PREDICTIONS BY CONSIDERING MULTIPLE FACTORS.

📌 **CHAPTER 6: REGRESSION MODEL EVALUATION**

6.1 ASSESSING REGRESSION MODEL PERFORMANCE

✓ **R-SQUARED (R^2):** MEASURES HOW WELL THE MODEL EXPLAINS VARIANCE IN Y .

- $R^2 = 0.90 \rightarrow 90\%$ OF Y 'S VARIATION IS EXPLAINED BY X VARIABLES.

✓ **MEAN ABSOLUTE ERROR (MAE):** MEASURES AVERAGE ERROR IN PREDICTIONS.

✓ **ROOT MEAN SQUARED ERROR (RMSE):** PENALIZES LARGE PREDICTION ERRORS MORE HEAVILY.

 **EXAMPLE:**

A SALES FORECASTING MODEL WITH $R^2 = 0.85$ IS CONSIDERED HIGHLY ACCURATE IN PREDICTING REVENUE.

 **CONCLUSION:**

EVALUATING REGRESSION MODELS ENSURES ACCURACY AND RELIABILITY BEFORE APPLYING THEM IN REAL-WORLD DECISION-MAKING.

 **CHAPTER 7: EXERCISES & ASSIGNMENTS****7.1 MULTIPLE CHOICE QUESTIONS**

WHICH CORRELATION COEFFICIENT REPRESENTS A PERFECT POSITIVE CORRELATION?

- (A) -1
- (B) 0
- (C) +1
- (D) 0.5

WHICH REGRESSION METHOD PREDICTS Y USING MULTIPLE INDEPENDENT VARIABLES?

- (A) SIMPLE LINEAR REGRESSION
- (B) MULTIPLE REGRESSION
- (C) LOGISTIC REGRESSION
- (D) CORRELATION

WHAT DOES A CORRELATION OF R = -0.85 INDICATE?

- (A) STRONG NEGATIVE RELATIONSHIP
- (B) NO RELATIONSHIP

-
- (C) WEAK POSITIVE RELATIONSHIP
 - (D) PERFECT POSITIVE CORRELATION
-

7.2 PRACTICAL ASSIGNMENT

📌 **TASK 1:**

USE PYTHON OR EXCEL TO CALCULATE THE CORRELATION BETWEEN ADVERTISING SPEND AND SALES REVENUE IN A DATASET.

📌 **TASK 2:**

BUILD A SIMPLE LINEAR REGRESSION MODEL TO PREDICT HOUSE PRICES BASED ON SQUARE FOOTAGE.

📌 **TASK 3:**

ANALYZE A DATASET USING MULTIPLE REGRESSION AND INTERPRET HOW DIFFERENT FACTORS INFLUENCE THE OUTCOME.

🌟 CONCLUSION: THE ROLE OF CORRELATION & REGRESSION IN DATA ANALYTICS

CORRELATION & REGRESSION ANALYSIS ARE ESSENTIAL TOOLS IN DATA SCIENCE, BUSINESS INTELLIGENCE, AND RESEARCH. UNDERSTANDING THEIR PRINCIPLES ENABLES ANALYSTS TO DETECT PATTERNS, FORECAST TRENDS, AND MAKE DATA-DRIVEN DECISIONS.



ASSIGNMENT:

CONDUCT EDA (EXPLORATORY DATA ANALYSIS) ON A REAL-WORLD DATASET & CREATE INSIGHTS USING DATA VISUALIZATION TOOLS

ISDM-Nxt

SOLUTION: CONDUCTING EXPLORATORY DATA ANALYSIS (EDA) & CREATING INSIGHTS USING DATA VISUALIZATION TOOLS

OBJECTIVE:

THE GOAL OF THIS ASSIGNMENT IS TO **PERFORM EXPLORATORY DATA ANALYSIS (EDA)** ON A REAL-WORLD DATASET, CLEAN AND PREPROCESS THE DATA, GENERATE SUMMARY STATISTICS, VISUALIZE PATTERNS, AND EXTRACT INSIGHTS USING **PYTHON, PANDAS, MATPLOTLIB, AND SEABORN.**

◆ **STEP 1: INSTALL & IMPORT REQUIRED LIBRARIES**

BEFORE WE BEGIN, ENSURE THAT THE REQUIRED PYTHON LIBRARIES ARE INSTALLED. IF NOT, INSTALL THEM USING:

PIP INSTALL PANDAS NUMPY MATPLOTLIB SEABORN

IMPORT NECESSARY LIBRARIES

IMPORT PANDAS AS PD

IMPORT NUMPY AS NP

IMPORT MATPLOTLIB.PYPLOT AS PLT

IMPORT SEABORN AS SNS

◆ **STEP 2: LOAD THE DATASET**

FOR THIS EXAMPLE, WE WILL USE THE **IRIS DATASET** (A WELL-KNOWN DATASET IN DATA SCIENCE). YOU CAN USE ANY DATASET RELATED TO SALES, CUSTOMER BEHAVIOR, OR HEALTHCARE AS PER YOUR INTEREST.

📌 LOAD DATASET USING PANDAS

```
# LOAD THE DATASET
```

```
DF = PD.READ_CSV("IRIS.CSV") # REPLACE WITH YOUR DATASET FILE
```

```
# DISPLAY THE FIRST 5 ROWS
```

```
DF.HEAD()
```

🔍 OBSERVATIONS:

- WE CHECK THE FIRST FEW ROWS TO **UNDERSTAND COLUMN NAMES AND DATA STRUCTURE.**
- IF THE DATASET IS NOT AVAILABLE IN CSV FORMAT, WE CAN LOAD IT FROM DATABASES (SQL) OR APIs.

◆ STEP 3: UNDERSTAND THE DATASET

📌 CHECK DATASET INFORMATION

```
# DISPLAY DATASET INFORMATION
```

```
DF.INFO()
```

🔍 KEY INSIGHTS:

- THIS STEP TELLS US **DATA TYPES, MISSING VALUES, AND MEMORY USAGE OF THE DATASET.**

- IF THERE ARE MISSING VALUES OR INCORRECT DATA TYPES, WE WILL HANDLE THEM IN THE NEXT STEPS.

📌 CHECK SUMMARY STATISTICS

GENERATE DESCRIPTIVE STATISTICS

DF.DESCRIBE()

🔍 KEY INSIGHTS:

- THIS PROVIDES MEAN, MEDIAN, MIN, MAX, AND PERCENTILES FOR NUMERICAL COLUMNS.
- HELPS IDENTIFY OUTLIERS AND CHECK DATA DISTRIBUTION.

◆ STEP 4: HANDLING MISSING VALUES & DATA CLEANING

MISSING VALUES CAN DISTORT ANALYSIS, SO WE NEED TO DETECT AND HANDLE THEM PROPERLY.

📌 CHECK FOR MISSING VALUES

COUNT MISSING VALUES PER COLUMN

DF.ISNULL().SUM()

📌 HANDLE MISSING VALUES (IF ANY)

FILL MISSING VALUES WITH MEAN FOR NUMERICAL COLUMNS

DF.FILLNA(DF.MEAN(), INPLACE=TRUE)

FILL MISSING VALUES IN CATEGORICAL COLUMNS WITH MODE

DF.FILLNA(DF.MODE().ILOC[0], INPLACE=TRUE)

📌 REMOVE DUPLICATE ROWS (IF ANY)

DROP DUPLICATE ROWS

```
DF.DROP_DUPLICATES(INPLACE=TRUE)
```

🔍 KEY INSIGHTS:

- HANDLING MISSING VALUES ENSURES ACCURATE ANALYSIS.
- REMOVING DUPLICATES PREVENTS BIAS IN PATTERNS.

◆ STEP 5: DETECT & HANDLE OUTLIERS

📌 VISUALIZE OUTLIERS USING BOXPLOT

PLOT BOXPLOT TO DETECT OUTLIERS

```
PLT.FIGURE(figsize=(8,5))
```

```
SNS.BOXPLOT(data=DF)
```

```
PLT.SHOW()
```

📌 REMOVE OUTLIERS USING THE IQR METHOD

```
Q1 = DF.quantile(0.25)
```

```
Q3 = DF.quantile(0.75)
```

```
IQR = Q3 - Q1
```

DEFINE LOWER AND UPPER BOUNDS

```
LOWER_BOUND = Q1 - 1.5 * IQR
```

```
UPPER_BOUND = Q3 + 1.5 * IQR
```

```
# REMOVE OUTLIERS
```

```
DF_CLEANED = DF[(DF >= LOWER_BOUND) & (DF <= UPPER_BOUND)]
```

🔍 KEY INSIGHTS:

- BOXPLOTS HELP IDENTIFY EXTREME VALUES (OUTLIERS) IN THE DATASET.
- REMOVING OUTLIERS IMPROVES ACCURACY OF PREDICTIVE MODELS.

◆ STEP 6: UNIVARIATE ANALYSIS (INDIVIDUAL FEATURES)

6.1 DISTRIBUTION OF NUMERICAL DATA

📌 HISTOGRAM OF NUMERICAL FEATURES

```
# PLOT HISTOGRAM FOR NUMERICAL COLUMNS
```

```
DF.HIST(FIGSIZE=(10, 6), BINS=20, EDGECOLOR="BLACK")
```

```
PLT.SHOW()
```

📌 KERNEL DENSITY ESTIMATE (KDE) PLOT

```
# VISUALIZING DISTRIBUTION USING KDE PLOT
```

```
SNS.KDEPLOT(DF['SEPAL_LENGTH'], SHADE=True, COLOR="BLUE")
```

```
PLT.SHOW()
```

🔍 KEY INSIGHTS:

- HISTOGRAMS HELP US UNDERSTAND THE SPREAD OF DATA.

- KDE PLOTS SHOW IF DATA IS **NORMALLY DISTRIBUTED OR SKEWED.**

- ◆ **STEP 7: BIVARIATE ANALYSIS (RELATIONSHIP BETWEEN TWO FEATURES)**

7.1 CORRELATION HEATMAP

📌 CHECK CORRELATION BETWEEN NUMERICAL VARIABLES

```
# COMPUTE CORRELATION MATRIX
```

```
CORR_MATRIX = DF.CORR()
```

```
# VISUALIZE WITH HEATMAP
```

```
PLT.FIGURE(figsize=(8,6))
```

```
SNS.HEATMAP(CORR_MATRIX, annot=True, cmap="COOLWARM")
```

```
PLT.SHOW()
```

🔍 KEY INSIGHTS:

- **STRONG POSITIVE CORRELATION** INDICATES A DIRECT RELATIONSHIP BETWEEN FEATURES.
- **STRONG NEGATIVE CORRELATION** SUGGESTS INVERSE RELATIONSHIPS.

7.2 SCATTER PLOTS

📌 PLOT RELATIONSHIPS BETWEEN TWO NUMERICAL VARIABLES

```
# SCATTER PLOT TO CHECK RELATIONSHIPS
```

```
SNS.SCATTERPLOT(X=DF['SEPAL_LENGTH'], Y=DF['PETAL_LENGTH'],  
HUE=DF['SPECIES'])
```

```
PLT.SHOW()
```

🔍 KEY INSIGHTS:

- HELPS DETERMINE LINEAR OR NON-LINEAR RELATIONSHIPS.
- CAN REVEAL CLUSTERS OR PATTERNS IN DATA.

◆ STEP 8: MULTIVARIATE ANALYSIS (INTERACTIONS BETWEEN MULTIPLE FEATURES)

8.1 PAIRPLOT (VISUALIZING ALL RELATIONSHIPS AT ONCE)

```
# PAIRPLOT OF NUMERICAL FEATURES
```

```
SNS.PAIRPLOT(DF, HUE="SPECIES", DIAG_KIND="KDE")
```

```
PLT.SHOW()
```

🔍 KEY INSIGHTS:

- PAIRPLOTS HELP IN IDENTIFYING CLUSTERS AND FEATURE INTERACTIONS.
- CATEGORICAL FEATURES (LIKE SPECIES) CAN BE USED TO DIFFERENTIATE CLASSES.

8.2 BOXPLOT GROUPED BY CATEGORIES

```
# BOXPLOT OF SEPAL LENGTH BY SPECIES
```

```
SNS.BOXPLOT(X=DF["SPECIES"], Y=DF["SEPAL_LENGTH"])
```

```
PLT.SHOW()
```

🔍 KEY INSIGHTS:

- HELPS COMPARE DISTRIBUTIONS ACROSS DIFFERENT CATEGORIES.

◆ STEP 9: CREATING SUMMARY INSIGHTS

📌 SUMMARIZE KEY FINDINGS FROM EDA

✓ DATA QUALITY CHECK

- NO MISSING VALUES AFTER PREPROCESSING.
- NO DUPLICATE RECORDS FOUND.
- HANDLED OUTLIERS USING IQR METHOD.

✓ KEY TRENDS & PATTERNS

- STRONG CORRELATION BETWEEN PETAL LENGTH & SEPAL LENGTH.
- SOME SPECIES HAVE DISTINCT FEATURE DISTRIBUTIONS, USEFUL FOR CLASSIFICATION.

✓ DATA DISTRIBUTION OBSERVATIONS

- MOST FEATURES FOLLOW A NORMAL DISTRIBUTION.
- OUTLIERS WERE PRESENT IN SOME COLUMNS AND WERE HANDLED.

📌 RECOMMENDATIONS BASED ON INSIGHTS

- IF USING MACHINE LEARNING, CONSIDER FEATURE SCALING BEFORE MODEL TRAINING.

- CONSIDER DIMENSIONALITY REDUCTION IF DATASET IS LARGE.
-

📌 FINAL THOUGHTS & NEXT STEPS

🚀 SUMMARY OF STEPS PERFORMED:

- 1 LOADED DATASET & PERFORMED BASIC EXPLORATION.
- 2 CLEANED MISSING VALUES & REMOVED DUPLICATES.
- 3 HANDLED OUTLIERS USING IQR METHOD.
- 4 CONDUCTED UNIVARIATE ANALYSIS USING HISTOGRAMS AND KDE PLOTS.
- 5 EXPLORED RELATIONSHIPS USING CORRELATION HEATMAPS, SCATTER PLOTS, AND BOXPLOTS.
- 6 SUMMARIZED KEY INSIGHTS AND BUSINESS IMPLICATIONS.

📌 NEXT STEPS:

- ◆ USE THIS ANALYSIS TO TRAIN A MACHINE LEARNING MODEL.
- ◆ CREATE AN INTERACTIVE DASHBOARD USING POWER BI OR TABLEAU.
- ◆ PERFORM FEATURE ENGINEERING TO IMPROVE PREDICTIVE MODELS.