



ISDM (INDEPENDENT SKILL DEVELOPMENT MISSION)

❖ DESCRIPTIVE STATISTICS: MEAN, MEDIAN, MODE, VARIANCE, STANDARD DEVIATION

❖ CHAPTER 1: INTRODUCTION TO DESCRIPTIVE STATISTICS

◆ 1.1 WHAT IS DESCRIPTIVE STATISTICS?

DESCRIPTIVE STATISTICS IS A BRANCH OF STATISTICS THAT SUMMARIZES AND DESCRIBES THE MAIN CHARACTERISTICS OF A DATASET. IT PROVIDES TOOLS TO ORGANIZE, VISUALIZE, AND INTERPRET DATA WITHOUT MAKING CONCLUSIONS BEYOND THE GIVEN DATASET.

THERE ARE TWO MAIN TYPES OF DESCRIPTIVE STATISTICS:

1. **MEASURES OF CENTRAL TENDENCY** – REPRESENT THE CENTER OF THE DATASET (MEAN, MEDIAN, MODE).
- 2 **MEASURES OF DISPERSION (SPREAD)** – SHOW HOW MUCH THE DATA VARIES (VARIANCE, STANDARD DEVIATION).

DESCRIPTIVE STATISTICS HELP IDENTIFY PATTERNS, DETECT OUTLIERS, AND UNDERSTAND DATA DISTRIBUTIONS. UNLIKE INFERENTIAL STATISTICS, WHICH MAKES PREDICTIONS ABOUT A

POPULATION, DESCRIPTIVE STATISTICS FOCUSES ONLY ON THE GIVEN DATASET.

 **EXAMPLE:**

A COMPANY COLLECTS THE MONTHLY SALARIES OF ITS EMPLOYEES AND USES DESCRIPTIVE STATISTICS TO CALCULATE THE AVERAGE SALARY, THE MOST COMMON SALARY, AND THE VARIATION IN SALARIES.

 **CONCLUSION:**

DESCRIPTIVE STATISTICS SIMPLIFIES COMPLEX DATA AND PROVIDES AN EASY-TO-UNDERSTAND SUMMARY OF CENTRAL VALUES AND DATA DISPERSION.

 **CHAPTER 2: MEASURES OF CENTRAL TENDENCY**

◆ **2.1 MEAN (ARITHMETIC AVERAGE)**

DEFINITION:

THE MEAN IS THE SUM OF ALL VALUES IN A DATASET DIVIDED BY THE TOTAL NUMBER OF OBSERVATIONS. IT REPRESENTS THE AVERAGE VALUE OF THE DATASET.

FORMULA FOR MEAN:

$$\text{Mean}(\mu) = \frac{\sum X}{N}$$

Where:

- ◆ X = Individual data points
- ◆ N = Total number of data points

EXAMPLE CALCULATION:

Consider the test scores of 5 students: 60, 70, 80, 90, 100.

$$\text{Mean} = \frac{60 + 70 + 80 + 90 + 100}{5} = \frac{400}{5} = 80$$

The average test score is 80.

ADVANTAGES OF MEAN:

- ✓ EASY TO CALCULATE AND UNDERSTAND.
- ✓ CONSIDERS ALL DATA POINTS.

DISADVANTAGES OF MEAN:

- ✗ SENSITIVE TO OUTLIERS (EXTREME VALUES).
- ✗ CAN BE MISLEADING IN SKEWED DISTRIBUTIONS.

📌 EXAMPLE:

INCOMES IN A COMPANY: \$30,000, \$35,000, \$40,000, \$50,000, AND \$1,000,000. THE MEAN SALARY IS \$231,000, WHICH IS MISLEADING BECAUSE ONE HIGH SALARY SKEWS THE AVERAGE.

💡 CONCLUSION:

THE MEAN IS USEFUL FOR SYMMETRICAL DATASETS BUT IS AFFECTED BY OUTLIERS, MAKING IT UNRELIABLE FOR SKEWED DATA.

◆ 2.2 MEDIAN (MIDDLE VALUE)

DEFINITION:

THE MEDIAN IS THE MIDDLE VALUE OF AN ORDERED DATASET. IF THE NUMBER OF OBSERVATIONS IS **ODD**, THE MEDIAN IS THE MIDDLE NUMBER. IF **EVEN**, THE MEDIAN IS THE AVERAGE OF THE TWO MIDDLE NUMBERS.

STEPS TO CALCULATE MEDIAN:

1. ARRANGE THE DATA IN ASCENDING ORDER.
2. IDENTIFY THE **MIDDLE VALUE**.
3. IF THERE ARE **TWO MIDDLE VALUES**, TAKE THEIR **AVERAGE**.

EXAMPLE CALCULATION (ODD DATASET):

DATASET: **10, 15, 20, 25, 30**

MEDIAN = **20** (MIDDLE VALUE).

EXAMPLE CALCULATION (EVEN DATASET):

Dataset: 5, 10, 15, 20, 25, 30

$$\text{Median} = \frac{15 + 20}{2} = 17.5$$

ADVANTAGES OF MEDIAN:

- ✓ NOT AFFECTED BY EXTREME VALUES (OUTLIERS).
- ✓ PROVIDES A BETTER MEASURE OF CENTER FOR **SKEWED DISTRIBUTIONS**.

DISADVANTAGES OF MEDIAN:

- ✗ DOES NOT CONSIDER ALL DATA VALUES.
- ✗ MORE COMPLEX TO CALCULATE THAN THE MEAN.

📌 EXAMPLE:

SALARIES: **\$30,000, \$35,000, \$40,000, \$50,000, \$1,000,000.**

THE MEDIAN SALARY IS \$40,000, WHICH IS A BETTER MEASURE THAN THE MEAN IN THIS CASE.

CONCLUSION:

THE MEDIAN IS PREFERRED FOR SKEWED DATA BECAUSE IT BETTER REPRESENTS THE CENTRAL VALUE WHEN OUTLIERS ARE PRESENT.

◆ **2.3 MODE (MOST FREQUENT VALUE)**

DEFINITION:

THE MODE IS THE MOST FREQUENTLY OCCURRING VALUE IN A DATASET. IT IS USEFUL FOR CATEGORICAL DATA (E.G., FAVORITE COLOR, MOST SOLD PRODUCT).

EXAMPLE CALCULATION:

DATASET: 3, 5, 7, 5, 8, 5, 9, 10

MODE = 5 (APPEARS MOST FREQUENTLY).

TYPES OF MODE:

1. **UNIMODAL** – ONE MODE (E.G., 5, 6, 7, 7, 8 → MODE = 7)
2. **BIMODAL** – TWO MODES (E.G., 5, 5, 7, 8, 8, 9 → MODES = 5 & 8)
3. **MULTIMODAL** – MORE THAN TWO MODES.

ADVANTAGES OF MODE:

- ✓ WORKS WELL FOR CATEGORICAL DATA.
- ✓ NOT AFFECTED BY OUTLIERS.

DISADVANTAGES OF MODE:

- ✗ SOME DATASETS MAY NOT HAVE A MODE.
- ✗ NOT USEFUL FOR CONTINUOUS NUMERICAL DATA.

📌 EXAMPLE:

SURVEY OF FAVORITE ICE CREAM FLAVORS:

CHOCOLATE (10 VOTES), VANILLA (15 VOTES), STRAWBERRY (20 VOTES), MANGO (15 VOTES).

THE MODES ARE **STRAWBERRY & VANILLA** (BIMODAL DISTRIBUTION).

💡 CONCLUSION:

MODE IS USEFUL FOR CATEGORICAL DATA BUT LESS EFFECTIVE FOR CONTINUOUS DATASETS.

📌 CHAPTER 3: MEASURES OF DISPERSION (SPREAD)

◆ 3.1 VARIANCE (SPREAD OF DATA POINTS)

DEFINITION:

VARIANCE MEASURES THE SPREAD OF DATA POINTS FROM THE MEAN. HIGHER VARIANCE MEANS MORE SPREAD-OUT DATA, WHILE LOWER VARIANCE INDICATES CLOSER DATA POINTS.

FORMULA FOR VARIANCE:

Formula for Variance:

$$\sigma^2 = \frac{\sum(X_i - \mu)^2}{N}$$

Where:

- ◆ X_i = Each data value
- ◆ μ = Mean of the dataset
- ◆ N = Total number of observations

EXAMPLE CALCULATION:

DATASET: 2, 4, 6, 8, 10

MEAN = 6

VARIANCE = 8 (AFTER CALCULATIONS).

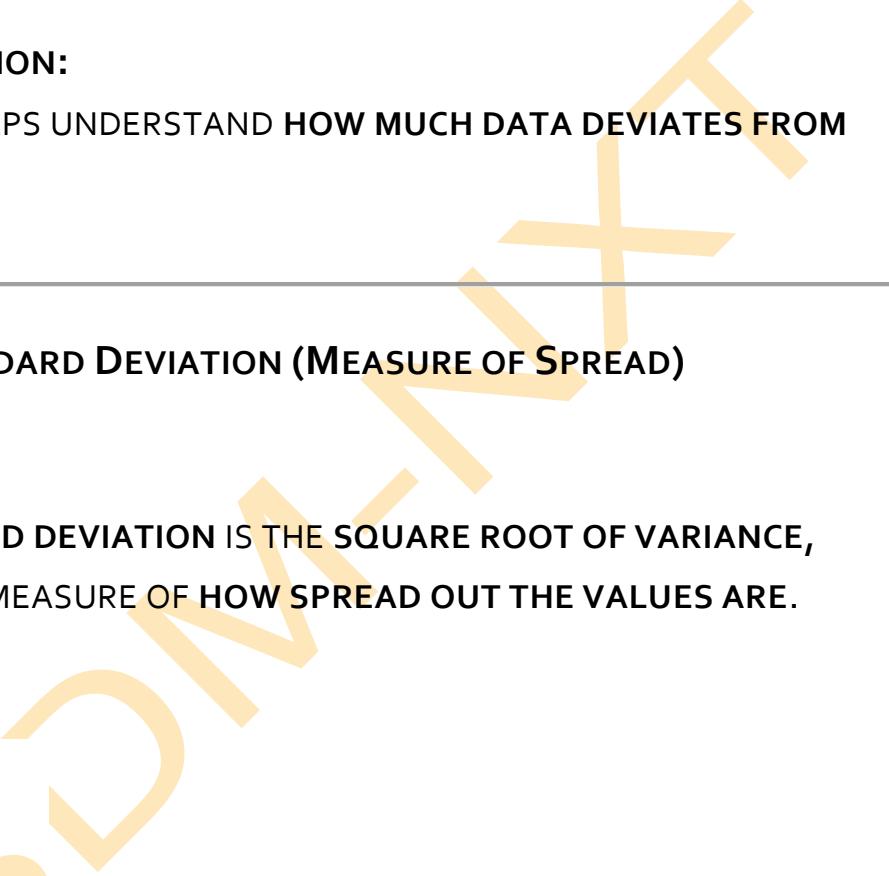
 **INTERPRETATION:**

✓ LOW VARIANCE – DATA POINTS ARE CLOSE TO THE MEAN.

✓ HIGH VARIANCE – DATA POINTS ARE WIDELY SPREAD.

 **CONCLUSION:**

VARIANCE HELPS UNDERSTAND HOW MUCH DATA DEVIATES FROM THE MEAN.

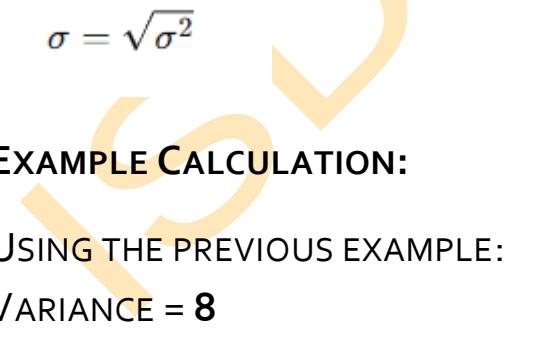
 **◆ 3.2 STANDARD DEVIATION (MEASURE OF SPREAD)**

DEFINITION:

THE STANDARD DEVIATION IS THE SQUARE ROOT OF VARIANCE, PROVIDING A MEASURE OF HOW SPREAD OUT THE VALUES ARE.

FORMULA:

$$\sigma = \sqrt{\sigma^2}$$

 **EXAMPLE CALCULATION:**

USING THE PREVIOUS EXAMPLE:

VARIANCE = 8

$$\text{Standard Deviation} = \sqrt{8} = 2.83$$

ADVANTAGES OF STANDARD DEVIATION:

- ✓ MORE INTUITIVE THAN VARIANCE.
- ✓ WIDELY USED IN STATISTICAL ANALYSIS AND PROBABILITY.

📌 EXAMPLE:

TWO TEST SCORE DATASETS:

- **DATASET A:** 70, 72, 74, 76 (LOW STANDARD DEVIATION)
- **DATASET B:** 50, 60, 80, 90 (HIGH STANDARD DEVIATION)

DATASET B HAS GREATER VARIABILITY IN TEST SCORES.

💡 CONCLUSION:

STANDARD DEVIATION IS WIDELY USED IN DATA ANALYSIS TO MEASURE THE CONSISTENCY AND RELIABILITY OF DATA.

📌 SUMMARY & NEXT STEPS

✅ KEY TAKEAWAYS:

- ✓ MEAN, MEDIAN, AND MODE DESCRIBE THE CENTRAL TENDENCY OF DATA.
- ✓ VARIANCE AND STANDARD DEVIATION MEASURE DATA SPREAD.
- ✓ THE MEDIAN IS BETTER FOR SKEWED DATA, WHILE THE MEAN WORKS WELL FOR SYMMETRICAL DATA.

📌 NEXT STEPS:

- ◆ PRACTICE CALCULATING DESCRIPTIVE STATISTICS USING PYTHON (NUMPY, PANDAS).
- ◆ VISUALIZE DATA DISTRIBUTIONS USING HISTOGRAMS AND BOXPLOTS.

◊ INFERENTIAL STATISTICS: HYPOTHESIS TESTING, P-VALUES, CONFIDENCE INTERVALS

📌 CHAPTER 1: INTRODUCTION TO INFERENTIAL STATISTICS

◆ 1.1 WHAT IS INFERENTIAL STATISTICS?

INFERENTIAL STATISTICS IS A BRANCH OF STATISTICS THAT HELPS US DRAW CONCLUSIONS ABOUT A POPULATION BASED ON A SAMPLE. UNLIKE DESCRIPTIVE STATISTICS, WHICH SUMMARIZES DATA, INFERENTIAL STATISTICS USES PROBABILITY THEORY AND SAMPLING TECHNIQUES TO MAKE GENERALIZATIONS BEYOND THE COLLECTED DATA.

WHY IS INFERENTIAL STATISTICS IMPORTANT?

- ✓ IT HELPS PREDICT TRENDS AND MAKE DECISIONS BASED ON SAMPLE DATA.
- ✓ IT ENABLES RESEARCHERS TO TEST HYPOTHESES AND VALIDATE THEORIES.
- ✓ IT ALLOWS ESTIMATION OF POPULATION PARAMETERS USING CONFIDENCE INTERVALS.
- ✓ IT REDUCES DATA COLLECTION EFFORTS BY WORKING WITH A REPRESENTATIVE SAMPLE.

KEY COMPONENTS OF INFERENTIAL STATISTICS:

- ✓ HYPOTHESIS TESTING – DETERMINING IF AN ASSUMPTION ABOUT A POPULATION IS STATISTICALLY SIGNIFICANT.

✓ **P-VALUES – MEASURING THE STRENGTH OF EVIDENCE AGAINST A NULL HYPOTHESIS.**

✓ **CONFIDENCE INTERVALS – ESTIMATING THE RANGE IN WHICH A POPULATION PARAMETER LIKELY FALLS.**

📌 EXAMPLE:

A PHARMACEUTICAL COMPANY WANTS TO DETERMINE WHETHER A NEW DRUG LOWERS BLOOD PRESSURE. INSTEAD OF TESTING EVERY PATIENT IN THE WORLD, THEY SELECT A SAMPLE GROUP AND USE INFERENTIAL STATISTICS TO ESTIMATE THE DRUG'S EFFECT ON THE ENTIRE POPULATION.

💡 CONCLUSION:

INFERENTIAL STATISTICS IS ESSENTIAL IN SCIENTIFIC RESEARCH, BUSINESS ANALYTICS, MEDICINE, AND ECONOMICS FOR MAKING PREDICTIONS AND DATA-DRIVEN DECISIONS.

📌 CHAPTER 2: HYPOTHESIS TESTING

◆ 2.1 WHAT IS HYPOTHESIS TESTING?

HYPOTHESIS TESTING IS A STATISTICAL METHOD USED TO TEST ASSUMPTIONS OR CLAIMS ABOUT A POPULATION BASED ON A SAMPLE. IT HELPS DETERMINE WHETHER A HYPOTHESIS SHOULD BE ACCEPTED OR REJECTED.

KEY TERMS IN HYPOTHESIS TESTING:

✓ **NULL HYPOTHESIS (H_0) – ASSUMES NO EFFECT OR NO DIFFERENCE EXISTS IN THE POPULATION.**

✓ **ALTERNATIVE HYPOTHESIS (H_1 OR H_A) – SUGGESTS THERE IS A**

SIGNIFICANT EFFECT OR DIFFERENCE.

- ✓ **SIGNIFICANCE LEVEL (α)** – THE PROBABILITY OF REJECTING THE NULL HYPOTHESIS WHEN IT IS ACTUALLY TRUE (COMMONLY 0.05 OR 5%).
- ✓ **TEST STATISTIC** – A VALUE USED TO DECIDE WHETHER TO REJECT H_0 (E.G., T-STATISTIC, Z-SCORE).
- ✓ **DECISION RULE** – COMPARE THE P-VALUE WITH α TO DETERMINE IF WE REJECT OR FAIL TO REJECT H_0 .

❖ **EXAMPLE:**

A MANUFACTURER CLAIMS THEIR LIGHT BULBS LAST **10,000 HOURS**. A RESEARCHER TESTS A **RANDOM SAMPLE** OF BULBS TO CHECK IF THE AVERAGE LIFESPAN SIGNIFICANTLY DIFFERS FROM 10,000 HOURS.

◆ **2.2 STEPS IN HYPOTHESIS TESTING**

1 STATE THE HYPOTHESES (H_0 AND H_1)

- EXAMPLE:
 - H_0 : THE AVERAGE BATTERY LIFE OF A SMARTPHONE IS 12 HOURS.
 - H_1 : THE AVERAGE BATTERY LIFE IS NOT 12 HOURS.

2 CHOOSE A SIGNIFICANCE LEVEL (α)

- COMMON VALUES: 0.05 (5%) OR 0.01 (1%).
- LOWER α VALUES REDUCE FALSE POSITIVES BUT INCREASE FALSE NEGATIVES.

3 SELECT THE APPROPRIATE STATISTICAL TEST

- **Z-TEST** (WHEN POPULATION VARIANCE IS KNOWN).
- **T-TEST** (WHEN POPULATION VARIANCE IS UNKNOWN).
- **CHI-SQUARE TEST** (FOR CATEGORICAL DATA).

4 CALCULATE THE TEST STATISTIC

- Example: If testing the mean, use the t-statistic formula:

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

where:

- \bar{x} = Sample Mean
- μ = Population Mean
- s = Sample Standard Deviation
- n = Sample Size

5 COMPARE THE P-VALUE WITH A

- IF P-VALUE < A, REJECT H_0 (THERE IS SIGNIFICANT EVIDENCE FOR H_1).
- IF P-VALUE > A, FAIL TO REJECT H_0 (INSUFFICIENT EVIDENCE TO SUPPORT H_1).

6 MAKE A CONCLUSION

- EXAMPLE: IF P-VALUE = 0.03 AND A = 0.05, WE REJECT H_0 AND CONCLUDE THERE IS SIGNIFICANT EVIDENCE THAT THE BATTERY LIFE IS DIFFERENT FROM 12 HOURS.

📌 EXAMPLE:

A UNIVERSITY WANTS TO KNOW WHETHER ONLINE STUDENTS SCORE HIGHER THAN IN-PERSON STUDENTS. THEY COMPARE TEST SCORES USING A T-TEST.

💡 CONCLUSION:

HYPOTHESIS TESTING HELPS VALIDATE CLAIMS IN SCIENTIFIC STUDIES, BUSINESS DECISIONS, AND QUALITY CONTROL.

📌 CHAPTER 3: UNDERSTANDING P-VALUES

◆ 3.1 WHAT IS A P-VALUE?

THE P-VALUE (PROBABILITY VALUE) QUANTIFIES THE STRENGTH OF EVIDENCE AGAINST THE NULL HYPOTHESIS (H_0). IT TELLS US HOW LIKELY WE WOULD OBTAIN THE OBSERVED DATA IF H_0 WERE TRUE.

INTERPRETING P-VALUES:

- ✓ P-VALUE < 0.05 – STRONG EVIDENCE AGAINST $H_0 \rightarrow$ REJECT H_0 .
- ✓ P-VALUE > 0.05 – WEAK EVIDENCE AGAINST $H_0 \rightarrow$ FAIL TO REJECT H_0 .
- ✓ P-VALUE < 0.01 – VERY STRONG EVIDENCE AGAINST H_0 (HIGHLY SIGNIFICANT RESULT).

📌 EXAMPLE:

A MEDICAL TRIAL TESTS WHETHER A NEW DRUG REDUCES BLOOD PRESSURE COMPARED TO A PLACEBO.

- IF THE P-VALUE = 0.02, WE REJECT H_0 (DRUG IS EFFECTIVE).

- IF THE P-VALUE = 0.10, WE FAIL TO REJECT H_0 (INSUFFICIENT EVIDENCE).

CONCLUSION:

THE P-VALUE IS CRUCIAL IN DECISION-MAKING, HELPING RESEARCHERS DETERMINE WHETHER RESULTS ARE STATISTICALLY SIGNIFICANT.

CHAPTER 4: CONFIDENCE INTERVALS

4.1 WHAT IS A CONFIDENCE INTERVAL?

A CONFIDENCE INTERVAL (CI) IS A RANGE OF VALUES THAT LIKELY CONTAINS THE TRUE POPULATION PARAMETER WITH A SPECIFIED CONFIDENCE LEVEL (E.G., 95% OR 99% CONFIDENCE).

WHY USE CONFIDENCE INTERVALS?

- ✓ THEY PROVIDE A RANGE OF ESTIMATES RATHER THAN A SINGLE VALUE.
- ✓ THEY ACCOUNT FOR SAMPLING VARIABILITY.
- ✓ THEY HELP IN DECISION-MAKING UNDER UNCERTAINTY.

◆ 4.2 CALCULATING A CONFIDENCE INTERVAL

FOR A POPULATION MEAN (μ), THE CONFIDENCE INTERVAL IS:

For a population mean (μ), the confidence interval is:

$$CI = \bar{x} \pm Z_{\alpha/2} \times \frac{s}{\sqrt{n}}$$

where:

- \bar{x} = Sample Mean
- s = Sample Standard Deviation
- n = Sample Size
- $Z_{\alpha/2}$ = Critical value from the Z-table (1.96 for 95% confidence, 2.58 for 99%)

EXAMPLE:

A SURVEY FINDS THE AVERAGE HEIGHT OF STUDENTS IN A UNIVERSITY IS 170 CM WITH A STANDARD DEVIATION OF 10 CM (SAMPLE SIZE = 100).

THE 95% CONFIDENCE INTERVAL IS:

The 95% confidence interval is:

$$170 \pm 1.96 \times \frac{10}{\sqrt{100}}$$

$$170 \pm 1.96 \times 1 = 170 \pm 1.96$$

$$CI = (168.04, 171.96)$$

INTERPRETATION:

"WE ARE 95% CONFIDENT THAT THE TRUE AVERAGE HEIGHT OF ALL STUDENTS IS BETWEEN 168.04 CM AND 171.96 CM."

CONCLUSION:

CONFIDENCE INTERVALS HELP ESTIMATE POPULATION PARAMETERS WITH A CERTAIN DEGREE OF CERTAINTY, REDUCING RELIANCE ON POINT ESTIMATES.

SUMMARY & NEXT STEPS

KEY TAKEAWAYS:

-  HYPOTHESIS TESTING HELPS DETERMINE WHETHER A CLAIM ABOUT A POPULATION IS STATISTICALLY SIGNIFICANT.
-  P-VALUES MEASURE THE STRENGTH OF EVIDENCE AGAINST THE NULL HYPOTHESIS.
-  CONFIDENCE INTERVALS PROVIDE A RANGE WITHIN WHICH A POPULATION PARAMETER IS LIKELY TO FALL.

NEXT STEPS:

- ◆ PRACTICE RUNNING HYPOTHESIS TESTS USING PYTHON'S SCIPY OR R.
- ◆ USE REAL-WORLD DATASETS TO CALCULATE P-VALUES AND CONFIDENCE INTERVALS.
- ◆ EXPLORE ADVANCED STATISTICAL METHODS LIKE ANOVA AND BAYESIAN INFERENCE. 

◊ PROBABILITY DISTRIBUTIONS: NORMAL, BINOMIAL, POISSON

📌 CHAPTER 1: INTRODUCTION TO PROBABILITY DISTRIBUTIONS

◆ 1.1 WHAT IS A PROBABILITY DISTRIBUTION?

A PROBABILITY DISTRIBUTION IS A FUNCTION THAT DESCRIBES HOW THE VALUES OF A RANDOM VARIABLE ARE DISTRIBUTED. IT PROVIDES THE PROBABILITIES OF OCCURRENCE OF DIFFERENT POSSIBLE OUTCOMES IN AN EXPERIMENT. PROBABILITY DISTRIBUTIONS PLAY A FUNDAMENTAL ROLE IN STATISTICS, MACHINE LEARNING, AND DATA SCIENCE, ENABLING US TO MODEL REAL-WORLD SCENARIOS, MAKE PREDICTIONS, AND ANALYZE UNCERTAINTY.

TYPES OF PROBABILITY DISTRIBUTIONS

PROBABILITY DISTRIBUTIONS ARE BROADLY CLASSIFIED INTO:

1. DISCRETE PROBABILITY DISTRIBUTIONS – DEALS WITH DISCRETE VARIABLES (E.G., NUMBER OF HEADS IN COIN TOSSES).
 - BINOMIAL DISTRIBUTION – MODELS THE NUMBER OF SUCCESSES IN A FIXED NUMBER OF TRIALS.
 - POISSON DISTRIBUTION – MODELS THE NUMBER OF EVENTS OCCURRING IN A FIXED INTERVAL OF TIME OR SPACE.

2. CONTINUOUS PROBABILITY DISTRIBUTIONS – DEALS WITH CONTINUOUS VARIABLES (E.G., HEIGHT OF INDIVIDUALS).

- **NORMAL DISTRIBUTION (GAUSSIAN DISTRIBUTION)** – USED TO MODEL NATURAL PHENOMENA SUCH AS TEST SCORES, STOCK PRICES, AND HEIGHTS.

📌 EXAMPLE:

- THE NORMAL DISTRIBUTION CAN MODEL IQ SCORES, AS THEY ARE SYMMETRICALLY DISTRIBUTED AROUND THE MEAN.
- THE BINOMIAL DISTRIBUTION CAN MODEL THE PROBABILITY OF GETTING HEADS IN REPEATED COIN TOSSES.
- THE POISSON DISTRIBUTION CAN MODEL THE NUMBER OF CUSTOMER ARRIVALS AT A SERVICE DESK IN AN HOUR.

💡 CONCLUSION:

PROBABILITY DISTRIBUTIONS ARE ESSENTIAL IN STATISTICAL MODELING, PROVIDING INSIGHTS INTO THE BEHAVIOR OF RANDOM VARIABLES.

📌 CHAPTER 2: NORMAL DISTRIBUTION

◆ 2.1 WHAT IS THE NORMAL DISTRIBUTION?

THE NORMAL DISTRIBUTION (GAUSSIAN DISTRIBUTION) IS A BELL-SHAPED PROBABILITY DISTRIBUTION THAT IS SYMMETRIC AROUND THE MEAN. IT IS WIDELY USED IN STATISTICS BECAUSE MANY NATURAL AND HUMAN-MADE PHENOMENA FOLLOW THIS DISTRIBUTION.

PROPERTIES OF NORMAL DISTRIBUTION:

- ✓ **BELL-SHAPED & SYMMETRIC** – THE LEFT AND RIGHT SIDES ARE MIRROR IMAGES.
- ✓ **MEAN, MEDIAN, AND MODE ARE EQUAL** – THE HIGHEST POINT OCCURS AT THE MEAN.
- ✓ **DEFINED BY TWO PARAMETERS** – MEAN (μ) AND STANDARD DEVIATION (σ).
- ✓ **EMPIRICAL RULE (68-95-99.7 RULE):**

- 68% OF DATA FALLS WITHIN **1 STANDARD DEVIATION** ($\mu \pm \sigma$).
- 95% OF DATA FALLS WITHIN **2 STANDARD DEVIATIONS** ($\mu \pm 2\sigma$).
- 99.7% OF DATA FALLS WITHIN **3 STANDARD DEVIATIONS** ($\mu \pm 3\sigma$).

📌 EXAMPLE:

- HEIGHTS OF PEOPLE IN A POPULATION FOLLOW A **NORMAL DISTRIBUTION**.
- IQ SCORES ARE NORMALLY DISTRIBUTED WITH **MEAN = 100** AND **STANDARD DEVIATION = 15**.

💡 CONCLUSION:

THE NORMAL DISTRIBUTION IS FUNDAMENTAL IN STATISTICS, APPEARING IN REAL-WORLD DATASETS, HYPOTHESIS TESTING, AND MACHINE LEARNING.

◆ 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 = Standard deviation
- e = Euler's number (approximately 2.718)

📌 EXAMPLE CALCULATION IN PYTHON:

```
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIB.PYTHON AS PLT  
FROM SCIPY.STATS IMPORT NORM  
  
# DEFINE MEAN AND STANDARD DEVIATION  
MU, SIGMA = 100, 15  
  
# GENERATE DATA POINTS  
X = NP.LINSPACE(MU - 4*SIGMA, MU + 4*SIGMA, 100)  
  
# COMPUTE NORMAL DISTRIBUTION
```

```
PDF = NORM.PDF(X, MU, SIGMA)
```

```
# PLOT THE NORMAL DISTRIBUTION
```

```
PLT.PLOT(X, PDF, LABEL="NORMAL DISTRIBUTION", COLOR='BLUE')
```

```
PLT.XLABEL("X VALUES")
```

```
PLT.YLABEL("PROBABILITY DENSITY")
```

```
PLT.TITLE("NORMAL DISTRIBUTION (M=100, Σ=15)")
```

```
PLT.LEGEND()
```

```
PLT.SHOW()
```

CONCLUSION:

THE PDF OF A NORMAL DISTRIBUTION HELPS IN CALCULATING THE LIKELIHOOD OF DIFFERENT VALUES OCCURRING WITHIN THE DATASET.

◆ 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 many standard deviations a data point is from the mean:

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

📌 EXAMPLE:

📌 Example:

If a student scores 130 in an IQ test where $\mu = 100$ and $\sigma = 15$, the Z-score is:

$$Z = \frac{130 - 100}{15} = 2$$

This means the score is 2 standard deviations above the mean.

💡 CONCLUSION:

Z-SCORES ALLOW COMPARISON ACROSS DIFFERENT NORMAL DISTRIBUTIONS AND ARE USED IN HYPOTHESIS TESTING AND STANDARDIZATION.

📌 CHAPTER 3: BINOMIAL DISTRIBUTION

◆ 3.1 WHAT IS THE BINOMIAL DISTRIBUTION?

THE BINOMIAL DISTRIBUTION MODELS THE NUMBER OF SUCCESSES IN A FIXED NUMBER OF INDEPENDENT TRIALS, WHERE EACH TRIAL HAS TWO POSSIBLE OUTCOMES (SUCCESS OR FAILURE).

PROPERTIES OF BINOMIAL DISTRIBUTION:

- ✓ **FIXED NUMBER OF TRIALS (N)** – THE NUMBER OF EXPERIMENTS IS PREDETERMINED.
- ✓ **TWO POSSIBLE OUTCOMES** – SUCCESS (1) OR FAILURE (0).
- ✓ **CONSTANT PROBABILITY OF SUCCESS (P)** – THE PROBABILITY OF SUCCESS REMAINS THE SAME ACROSS TRIALS.
- ✓ **INDEPENDENT TRIALS** – THE OUTCOME OF ONE TRIAL DOES NOT AFFECT ANOTHER.

📌 EXAMPLE:

- TOSSING A COIN 10 TIMES AND COUNTING THE NUMBER OF HEADS.
- PREDICTING THE NUMBER OF DEFECTIVE ITEMS IN A BATCH OF 50 MANUFACTURED PRODUCTS.

 CONCLUSION:

THE BINOMIAL DISTRIBUTION IS USED TO MODEL DISCRETE PROBABILITY EVENTS IN FIELDS LIKE QUALITY CONTROL, FINANCE, AND GENETICS.

◆ 3.2 BINOMIAL PROBABILITY FORMULA

The probability of exactly k successes in n trials is given by the Binomial Formula:

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

where:

- $\binom{n}{k} = \frac{n!}{k!(n-k)!}$ (Binomial coefficient)
- p = Probability of success
- $1 - p$ = Probability of failure

 EXAMPLE CALCULATION IN PYTHON:

```
FROM SCIPY.STATS IMPORT BINOM
```

```
# PARAMETERS
```

```
N = 10 # NUMBER OF TRIALS
```

```
P = 0.5 # PROBABILITY OF SUCCESS
```

COMPUTE PROBABILITY OF GETTING EXACTLY 5 HEADS IN 10 COIN TOSSES

```
PROB_5_HEADS = BINOM.PMF(5, N, P)
```

```
PRINT(F"PROBABILITY OF 5 HEADS IN 10 TOSSES:  
{PROB_5_HEADS:.4F}")
```

 CONCLUSION:

THE BINOMIAL FORMULA HELPS COMPUTE THE PROBABILITY OF ACHIEVING A SPECIFIC NUMBER OF SUCCESSES IN REPEATED INDEPENDENT TRIALS.

 CHAPTER 4: POISSON DISTRIBUTION

◆ 4.1 WHAT IS THE POISSON DISTRIBUTION?

THE POISSON DISTRIBUTION MODELS THE NUMBER OF EVENTS OCCURRING IN A FIXED INTERVAL OF TIME OR SPACE, ASSUMING THE EVENTS HAPPEN AT A CONSTANT RATE INDEPENDENTLY OF EACH OTHER.

PROPERTIES OF POISSON DISTRIBUTION:

- ✓ DESCRIBES RARE EVENTS – USED WHEN THE PROBABILITY OF AN EVENT OCCURRING IS LOW.
- ✓ SINGLE PARAMETER (λ) – REPRESENTS THE AVERAGE NUMBER OF OCCURRENCES.
- ✓ NON-NEGATIVE INTEGER VALUES – THE VARIABLE CAN TAKE VALUES 0, 1, 2, 3, ...

 **EXAMPLE:**

- THE NUMBER OF CALLS RECEIVED AT A CALL CENTER PER HOUR.
- THE NUMBER OF CUSTOMER ARRIVALS AT A RESTAURANT PER MINUTE.

 **CONCLUSION:**

POISSON DISTRIBUTION IS WIDELY USED IN QUEUE THEORY, RISK ANALYSIS, AND NETWORK TRAFFIC MODELING.

ISDM-NXT

◊ CORRELATION & REGRESSION ANALYSIS

📌 CHAPTER 1: INTRODUCTION TO CORRELATION & REGRESSION ANALYSIS

◆ 1.1 WHAT IS CORRELATION & REGRESSION ANALYSIS?

CORRELATION AND REGRESSION ANALYSIS ARE TWO FUNDAMENTAL STATISTICAL TECHNIQUES USED IN DATA SCIENCE, MACHINE LEARNING, AND BUSINESS ANALYTICS TO UNDERSTAND RELATIONSHIPS BETWEEN VARIABLES.

- CORRELATION ANALYSIS MEASURES THE STRENGTH AND DIRECTION OF THE RELATIONSHIP BETWEEN TWO VARIABLES.
- REGRESSION ANALYSIS HELPS IN PREDICTING ONE VARIABLE BASED ON ANOTHER BY MODELING THE RELATIONSHIP.

THESE TECHNIQUES ARE WIDELY USED IN VARIOUS DOMAINS SUCH AS FINANCE, HEALTHCARE, MARKETING, ECONOMICS, AND ENGINEERING TO ANALYZE TRENDS AND MAKE DATA-DRIVEN DECISIONS.

WHY ARE CORRELATION & REGRESSION IMPORTANT?

- ✓ UNDERSTAND DATA RELATIONSHIPS – IDENTIFY HOW VARIABLES ARE ASSOCIATED.
- ✓ PREDICT FUTURE OUTCOMES – REGRESSION MODELS PREDICT VALUES BASED ON PAST DATA.
- ✓ FEATURE SELECTION – HELPS IN DETERMINING WHICH

VARIABLES IMPACT AN OUTCOME.

✓ OPTIMIZING BUSINESS STRATEGIES – USED IN SALES FORECASTING, RISK ASSESSMENT, AND MARKETING OPTIMIZATION.

 **EXAMPLE:**

A RETAIL COMPANY WANTS TO ANALYZE IF THERE IS A CORRELATION BETWEEN ADVERTISING SPEND AND SALES REVENUE. IF A STRONG RELATIONSHIP IS FOUND, REGRESSION ANALYSIS CAN BE USED TO PREDICT FUTURE SALES BASED ON ADVERTISING BUDGETS.

 **CONCLUSION:**

CORRELATION AND REGRESSION ANALYSIS ARE ESSENTIAL TOOLS FOR EXPLORING DATA PATTERNS, RELATIONSHIPS, AND PREDICTING OUTCOMES.

◆ **1.2 DIFFERENCES BETWEEN CORRELATION & REGRESSION**

FEATURE	CORRELATION ANALYSIS	REGRESSION ANALYSIS
PURPOSE	MEASURES THE STRENGTH & DIRECTION OF A RELATIONSHIP BETWEEN TWO VARIABLES	PREDICTS DEPENDENT VARIABLE BASED ON INDEPENDENT VARIABLE(S)
CAUSALITY	DOES NOT IMPLY CAUSATION	HELPS IN ESTABLISHING CAUSE-AND-EFFECT RELATIONSHIPS

TYPES	PEARSON, SPEARMAN, KENDALL	LINEAR, LOGISTIC, POLYNOMIAL, MULTIPLE REGRESSION
OUTPUT	CORRELATION COEFFICIENT (R)	REGRESSION EQUATION ($Y = A + BX$)
EXAMPLE	FINDING RELATIONSHIP BETWEEN STUDY TIME & EXAM SCORES	PREDICTING EXAM SCORES BASED ON STUDY HOURS

 **EXAMPLE:**

- **CORRELATION ANALYSIS:** DETERMINING IF THERE IS A RELATIONSHIP BETWEEN EXERCISE TIME AND CALORIE BURN.
- **REGRESSION ANALYSIS:** BUILDING A MODEL TO PREDICT CALORIES BURNED BASED ON EXERCISE DURATION.

 **CONCLUSION:**

CORRELATION QUANTIFIES RELATIONSHIPS, WHILE REGRESSION GOES ONE STEP FURTHER BY MAKING PREDICTIONS BASED ON DATA.

 **CHAPTER 2: UNDERSTANDING CORRELATION ANALYSIS**

◆ **2.1 WHAT IS CORRELATION?**

CORRELATION IS A STATISTICAL MEASURE THAT EXPRESSES THE STRENGTH AND DIRECTION OF THE RELATIONSHIP BETWEEN TWO VARIABLES. IT TELLS US HOW CHANGES IN ONE VARIABLE ARE ASSOCIATED WITH CHANGES IN ANOTHER.

TYPES OF CORRELATION:

1 POSITIVE CORRELATION (+) – WHEN ONE VARIABLE INCREASES, THE OTHER ALSO INCREASES.

- EXAMPLE: HIGHER STUDY HOURS LEAD TO HIGHER GRADES.

2 NEGATIVE CORRELATION (-) – WHEN ONE VARIABLE INCREASES, THE OTHER DECREASES.

- EXAMPLE: INCREASE IN EXERCISE LEADS TO WEIGHT LOSS.

3 NO CORRELATION (0) – NO SIGNIFICANT RELATIONSHIP BETWEEN TWO VARIABLES.

- EXAMPLE: THE NUMBER OF PETS OWNED VS. MONTHLY SALARY.

MEASURING CORRELATION: PEARSON'S CORRELATION COEFFICIENT (R)

THE MOST COMMON CORRELATION MEASURE IS PEARSON'S CORRELATION COEFFICIENT (R), WHICH RANGES FROM -1 TO +1.

CORRELATION COEFFICIENT (R)	RELATIONSHIP STRENGTH
+1.0	PERFECT POSITIVE CORRELATION
+0.5 TO +0.9	STRONG POSITIVE CORRELATION
+0.1 TO +0.5	WEAK POSITIVE CORRELATION
0	NO CORRELATION
-0.1 TO -0.5	WEAK NEGATIVE CORRELATION
-0.5 TO -0.9	STRONG NEGATIVE CORRELATION

-1.0	PERFECT NEGATIVE CORRELATION
------	------------------------------

 **EXAMPLE:**

A MARKETING TEAM ANALYZES THE CORRELATION BETWEEN SOCIAL MEDIA SPENDING AND CUSTOMER ENGAGEMENT TO DETERMINE IF INCREASING AD BUDGETS WILL BOOST INTERACTIONS.

 **CONCLUSION:**

CORRELATION HELPS UNDERSTAND ASSOCIATIONS BETWEEN VARIABLES, BUT IT DOES NOT ESTABLISH CAUSALITY.

◆ **2.2 TYPES OF CORRELATION ANALYSIS**

- ✓ **PEARSON CORRELATION** – MEASURES LINEAR RELATIONSHIPS BETWEEN CONTINUOUS VARIABLES.
- ✓ **SPEARMAN RANK CORRELATION** – MEASURES MONOTONIC RELATIONSHIPS (WORKS WELL FOR RANKED DATA).
- ✓ **KENDALL TAU CORRELATION** – USED FOR ORDINAL (RANKED) DATA WITH SMALL SAMPLE SIZES.

 **EXAMPLE:**

A FINANCIAL ANALYST USES PEARSON CORRELATION TO STUDY THE RELATIONSHIP BETWEEN STOCK PRICES AND INTEREST RATES.

 **CONCLUSION:**

CHOOSING THE RIGHT CORRELATION METHOD DEPENDS ON THE DATA TYPE AND DISTRIBUTION.

CHAPTER 3: REGRESSION ANALYSIS

3.1 WHAT IS REGRESSION ANALYSIS?

REGRESSION ANALYSIS IS A PREDICTIVE MODELING TECHNIQUE THAT ESTABLISHES A RELATIONSHIP BETWEEN A DEPENDENT VARIABLE (OUTCOME) AND ONE OR MORE INDEPENDENT VARIABLES (PREDICTORS).

TYPES OF REGRESSION ANALYSIS:

1 SIMPLE LINEAR REGRESSION – MODELS A RELATIONSHIP BETWEEN ONE INDEPENDENT VARIABLE AND ONE DEPENDENT VARIABLE ($Y = A + BX$).

- EXAMPLE: PREDICTING HOUSE PRICES BASED ON SQUARE FOOTAGE.

2 MULTIPLE LINEAR REGRESSION – MODELS A RELATIONSHIP BETWEEN ONE DEPENDENT VARIABLE AND MULTIPLE INDEPENDENT VARIABLES ($Y = A + BX_1 + CX_2 + DX_3$).

- EXAMPLE: PREDICTING EMPLOYEE SALARY BASED ON EXPERIENCE, EDUCATION, AND CERTIFICATIONS.

3 POLYNOMIAL REGRESSION – MODELS CURVED RELATIONSHIPS BETWEEN VARIABLES.

- EXAMPLE: PREDICTING CUSTOMER SATISFACTION OVER TIME (NON-LINEAR TRENDS).

4 LOGISTIC REGRESSION – USED FOR CATEGORICAL PREDICTIONS (BINARY CLASSIFICATION: YES/NO, PASS/FAIL).

- EXAMPLE: PREDICTING WHETHER A LOAN APPLICANT WILL DEFAULT (YES/NO).

 **EXAMPLE:**

A HEALTHCARE COMPANY USES MULTIPLE REGRESSION TO PREDICT DIABETES RISK BASED ON BMI, AGE, AND DIET.

 **CONCLUSION:**

REGRESSION IS WIDELY USED IN FORECASTING, RISK ANALYSIS, AND DECISION-MAKING.

- ◆ 3.2 INTERPRETING REGRESSION RESULTS
- ✓ REGRESSION COEFFICIENT (B_1, B_2, \dots) – INDICATES THE STRENGTH OF PREDICTORS.
- ✓ INTERCEPT (A) – THE EXPECTED OUTCOME WHEN ALL PREDICTORS ARE ZERO.
- ✓ R-SQUARED (R^2) – MEASURES HOW WELL THE MODEL FITS THE DATA (RANGES FROM 0 TO 1).
- ✓ P-VALUE – DETERMINES STATISTICAL SIGNIFICANCE ($P < 0.05$ IS CONSIDERED SIGNIFICANT).

 **EXAMPLE:**

A SALES TEAM RUNS A REGRESSION MODEL TO PREDICT REVENUE BASED ON ADVERTISING BUDGET AND STORE LOCATION. A HIGH R^2 INDICATES A STRONG MODEL FIT.

 **CONCLUSION:**

REGRESSION RESULTS MUST BE INTERPRETED CAREFULLY TO ENSURE ACCURATE PREDICTIONS AND DECISION-MAKING.

📌 CHAPTER 4: PRACTICAL APPLICATIONS OF CORRELATION & REGRESSION

◆ 4.1 REAL-WORLD APPLICATIONS

- ✓ FINANCE: PREDICTING STOCK PRICES, RISK ASSESSMENT.
- ✓ HEALTHCARE: DISEASE PREDICTION, PATIENT SURVIVAL ANALYSIS.
- ✓ MARKETING: CUSTOMER SEGMENTATION, AD CAMPAIGN EFFECTIVENESS.
- ✓ SPORTS ANALYTICS: PREDICTING TEAM PERFORMANCE BASED ON PLAYER STATISTICS.
- ✓ REAL ESTATE: ESTIMATING HOUSE PRICES BASED ON LOCATION, SIZE, AND AMENITIES.

📌 EXAMPLE:

A BANKING INSTITUTION USES REGRESSION TO PREDICT LOAN APPROVAL CHANCES BASED ON CREDIT SCORE, INCOME, AND DEBT.

💡 CONCLUSION:

CORRELATION AND REGRESSION DRIVE DATA-DRIVEN DECISION-MAKING IN EVERY INDUSTRY.

📌 SUMMARY & NEXT STEPS

✓ KEY TAKEAWAYS:

- ✓ CORRELATION QUANTIFIES RELATIONSHIPS, BUT DOES NOT IMPLY CAUSATION.
- ✓ REGRESSION MODELS RELATIONSHIPS AND MAKES

PREDICTIONS.

- ✓ DIFFERENT TYPES OF CORRELATION AND REGRESSION ARE USED BASED ON DATA TYPE AND PURPOSE.
- ✓ REAL-WORLD APPLICATIONS SPAN FINANCE, HEALTHCARE, MARKETING, AND RISK ASSESSMENT.

📌 NEXT STEPS:

- ◆ APPLY CORRELATION AND REGRESSION TECHNIQUES USING PYTHON (PANDAS, SCIKIT-LEARN).
- ◆ WORK ON REAL-WORLD DATASETS FROM KAGGLE OR OPEN-SOURCE REPOSITORIES.
- ◆ MASTER ADVANCED REGRESSION TECHNIQUES LIKE RIDGE, LASSO, AND BAYESIAN REGRESSION.



◊ DATA STORYTELLING & VISUALIZATION WITH POWER BI & TABLEAU

📌 CHAPTER 1: INTRODUCTION TO DATA STORYTELLING & VISUALIZATION

◆ 1.1 WHAT IS DATA STORYTELLING?

DATA STORYTELLING IS THE PROCESS OF COMMUNICATING INSIGHTS DERIVED FROM DATA USING VISUAL AND NARRATIVE TECHNIQUES. IT COMBINES DATA ANALYTICS, VISUALIZATION, AND STORYTELLING TO HELP DECISION-MAKERS UNDERSTAND AND ACT UPON THE INSIGHTS EFFECTIVELY.

DATA STORYTELLING GOES BEYOND SIMPLE CHARTS AND DASHBOARDS—IT PROVIDES CONTEXT, INTERPRETATION, AND ACTIONABLE INSIGHTS IN A STRUCTURED AND ENGAGING WAY.

KEY COMPONENTS OF DATA STORYTELLING:

- ✓ DATA – THE RAW INFORMATION COLLECTED AND ANALYZED.
- ✓ VISUALS – GRAPHS, CHARTS, AND DASHBOARDS THAT ILLUSTRATE TRENDS AND PATTERNS.
- ✓ NARRATIVE – THE EXPLANATORY CONTEXT THAT MAKES THE DATA MEANINGFUL.

◆ 1.2 IMPORTANCE OF DATA STORYTELLING IN BUSINESS INTELLIGENCE

IN TODAY'S DATA-DRIVEN WORLD, ORGANIZATIONS COLLECT VAST AMOUNTS OF DATA, BUT RAW DATA ALONE IS MEANINGLESS. DATA

STORYTELLING BRIDGES THE GAP BETWEEN COMPLEX ANALYTICS AND BUSINESS DECISION-MAKING.

- ✓ ENHANCES DATA-DRIVEN DECISION-MAKING.
- ✓ SIMPLIFIES COMPLEX INSIGHTS INTO ACTIONABLE INFORMATION.
- ✓ ENGAGES STAKEHOLDERS BY MAKING DATA COMPELLING AND RELATABLE.
- ✓ IMPROVES DATA LITERACY ACROSS TEAMS.

 **EXAMPLE:**

A RETAIL COMPANY ANALYZES SALES DATA AND USES DATA STORYTELLING TO SHOW HOW CUSTOMER BUYING TRENDS SHIFT SEASONALLY, ENABLING BETTER STOCK MANAGEMENT.

 **CONCLUSION:**

DATA STORYTELLING TRANSFORMS NUMBERS INTO NARRATIVES, MAKING IT EASIER TO INTERPRET AND ACT ON DATA INSIGHTS.

 **CHAPTER 2: INTRODUCTION TO DATA VISUALIZATION TOOLS – POWER BI & TABLEAU**

 **2.1 WHAT IS DATA VISUALIZATION?**

DATA VISUALIZATION IS THE GRAPHICAL REPRESENTATION OF DATA TO HELP IDENTIFY TRENDS, PATTERNS, AND INSIGHTS. IT ENABLES USERS TO COMPREHEND COMPLEX DATASETS QUICKLY AND MAKE INFORMED DECISIONS.

COMMON DATA VISUALIZATION TECHNIQUES:

- ✓ **BAR CHARTS** – BEST FOR COMPARING DIFFERENT CATEGORIES.
- ✓ **LINE CHARTS** – USEFUL FOR TIME SERIES ANALYSIS (E.G., SALES OVER TIME).
- ✓ **PIE CHARTS** – SHOW PROPORTIONS WITHIN A DATASET.
- ✓ **SCATTER PLOTS** – ILLUSTRATE RELATIONSHIPS BETWEEN TWO VARIABLES.
- ✓ **HEATMAPS** – DISPLAY THE INTENSITY OF DATA VALUES USING COLOR GRADIENTS.

 **EXAMPLE:**

A BANKING DASHBOARD USES HEATMAPS TO VISUALIZE FRAUDULENT TRANSACTION HOTSPOTS BASED ON CUSTOMER LOCATIONS.

 **CONCLUSION:**

DATA VISUALIZATION HELPS DECISION-MAKERS UNDERSTAND COMPLEX DATASETS EASILY, LEADING TO FASTER AND MORE EFFECTIVE BUSINESS STRATEGIES.

◆ **2.2 POWER BI VS. TABLEAU: A COMPARISON**

FEATURE	POWER BI	TABLEAU
EASE OF USE	BEGINNER-FRIENDLY, ESPECIALLY FOR EXCEL USERS.	MORE ADVANCED, REQUIRES LEARNING BUT HIGHLY FLEXIBLE.
INTEGRATION	BEST FOR MICROSOFT ECOSYSTEM (EXCEL, AZURE, SQL SERVER).	CONNECTS WELL WITH VARIOUS DATA SOURCES.

SPEED & PERFORMANCE	WORKS WELL WITH MODERATE DATASETS, BUT MAY SLOW DOWN WITH LARGE DATA.	OPTIMIZED FOR HANDLING BIG DATA EFFICIENTLY.
AI & AUTOMATION	INCLUDES AI-POWERED INSIGHTS & NATURAL LANGUAGE QUERIES.	STRONGER IN DEEP DATA ANALYTICS AND CUSTOM VISUALIZATIONS.
PRICING	MORE AFFORDABLE FOR SMALL BUSINESSES.	HIGHER COST BUT OFFERS ADVANCED FEATURES.

📌 **EXAMPLE:**

A COMPANY USING **MICROSOFT AZURE** MIGHT PREFER **POWER BI**, WHILE A BUSINESS WITH MULTI-CLOUD DATA SOURCES MIGHT BENEFIT FROM **TABLEAU**.

💡 **CONCLUSION:**

BOTH **POWER BI** AND **TABLEAU** ARE EXCELLENT VISUALIZATION TOOLS—CHOOSING THE RIGHT ONE DEPENDS ON **BUSINESS NEEDS**, **BUDGET**, AND EXISTING **TECH STACK**.

📌 **CHAPTER 3: GETTING STARTED WITH POWER BI**

◆ **3.1 WHAT IS POWER BI?**

POWER BI IS A BUSINESS INTELLIGENCE TOOL BY MICROSOFT THAT ALLOWS USERS TO CONNECT, CLEAN, TRANSFORM, AND VISUALIZE DATA IN INTERACTIVE REPORTS AND DASHBOARDS.

- ✓ USER-FRIENDLY, INTEGRATES WELL WITH MICROSOFT PRODUCTS (EXCEL, SQL SERVER, AZURE).
- ✓ PROVIDES AI-DRIVEN INSIGHTS FOR PREDICTIVE ANALYTICS.
- ✓ OFFERS REAL-TIME DATA MONITORING WITH AUTOMATED REPORTS.

 **EXAMPLE:**

A FINANCE DEPARTMENT USES POWER BI TO CREATE AUTOMATED FINANCIAL DASHBOARDS THAT TRACK REVENUE, EXPENSES, AND PROFIT MARGINS.

 **CONCLUSION:**

POWER BI ENABLES INTERACTIVE DATA STORYTELLING, MAKING ANALYTICS ACCESSIBLE TO ALL LEVELS OF USERS.

◆ **3.2 SETTING UP POWER BI**

STEPS TO SET UP POWER BI:

- 1 DOWNLOAD & INSTALL POWER BI DESKTOP FROM MICROSOFT.**
- 2 CONNECT TO A DATA SOURCE (EXCEL, SQL, GOOGLE ANALYTICS, API).**
- 3 LOAD & TRANSFORM DATA USING POWER QUERY.**
- 4 CREATE DATA RELATIONSHIPS TO CONNECT DIFFERENT TABLES.**

5 BUILD INTERACTIVE DASHBOARDS WITH VISUALIZATIONS.**6 PUBLISH & SHARE REPORTS WITH STAKEHOLDERS.** **EXAMPLE:**

A MARKETING TEAM CONNECTS POWER BI TO GOOGLE ANALYTICS TO MONITOR WEBSITE TRAFFIC TRENDS AND OPTIMIZE AD SPENDING.

 **CONCLUSION:**

POWER BI ALLOWS USERS TO ANALYZE AND PRESENT DATA INSIGHTS EFFICIENTLY IN AN INTERACTIVE AND USER-FRIENDLY MANNER.

 **CHAPTER 4: GETTING STARTED WITH TABLEAU** **4.1 WHAT IS TABLEAU?**

TABLEAU IS A POWERFUL DATA VISUALIZATION TOOL THAT ENABLES USERS TO CONNECT, VISUALIZE, AND SHARE DATA THROUGH INTERACTIVE DASHBOARDS.

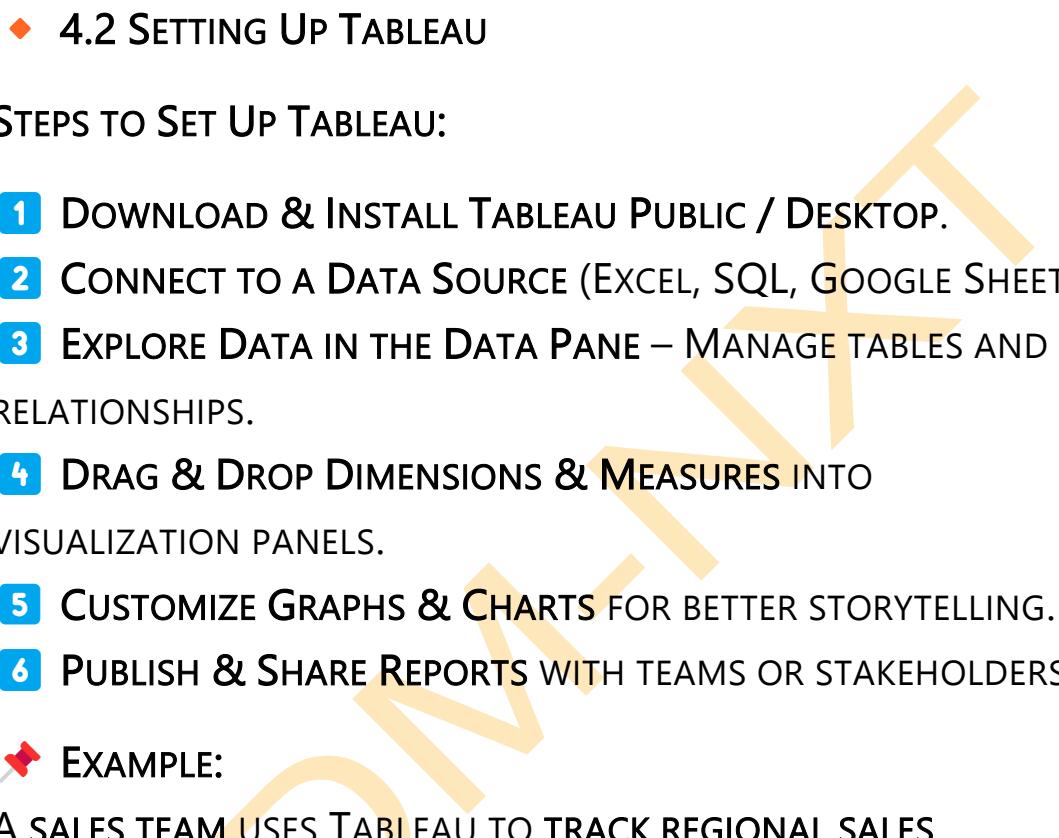
-  DRAG-AND-DROP INTERFACE FOR EASY DATA ANALYSIS.
-  HANDLES LARGE DATASETS WITH HIGH-SPEED PERFORMANCE.
-  SUPPORTS A WIDE RANGE OF DATA SOURCES (CLOUD, DATABASES, APIs).

 **EXAMPLE:**

A HEALTHCARE PROVIDER USES TABLEAU TO TRACK PATIENT RECOVERY RATES, HELPING DOCTORS OPTIMIZE TREATMENT PLANS.

 **CONCLUSION:**

TABLEAU ENABLES USERS TO BUILD VISUALLY RICH DASHBOARDS FOR BETTER DATA-DRIVEN DECISION-MAKING.

 **◆ 4.2 SETTING UP TABLEAU****STEPS TO SET UP TABLEAU:**

- 1 DOWNLOAD & INSTALL TABLEAU PUBLIC / DESKTOP.**
- 2 CONNECT TO A DATA SOURCE (EXCEL, SQL, GOOGLE SHEETS).**
- 3 EXPLORE DATA IN THE DATA PANE – MANAGE TABLES AND RELATIONSHIPS.**
- 4 DRAG & DROP DIMENSIONS & MEASURES INTO VISUALIZATION PANELS.**
- 5 CUSTOMIZE GRAPHS & CHARTS FOR BETTER STORYTELLING.**
- 6 PUBLISH & SHARE REPORTS WITH TEAMS OR STAKEHOLDERS.**

 **EXAMPLE:**

A SALES TEAM USES TABLEAU TO TRACK REGIONAL SALES PERFORMANCE, HELPING MANAGERS IDENTIFY HIGH-PERFORMING LOCATIONS.

 **CONCLUSION:**

TABLEAU IS AN INDUSTRY-LEADING TOOL FOR CREATING INTERACTIVE DATA STORIES WITH DEEP INSIGHTS.

 **CHAPTER 5: ADVANCED VISUALIZATION TECHNIQUES****◆ 5.1 INTERACTIVE DASHBOARDS IN POWER BI & TABLEAU**

- ✓ **CUSTOM FILTERS & SLICERS** – ALLOW USERS TO INTERACT WITH REPORTS DYNAMICALLY.
- ✓ **DRILL-DOWNS & HIERARCHIES** – EXPAND DATA LAYERS FOR DETAILED INSIGHTS.
- ✓ **MAPS & GEOSPATIAL ANALYSIS** – VISUALIZE DATA BASED ON GEOGRAPHIC LOCATIONS.

 **EXAMPLE:**

A LOGISTICS COMPANY TRACKS REAL-TIME SHIPMENT DELAYS USING GEOSPATIAL DASHBOARDS IN POWER BI AND TABLEAU.

 **CONCLUSION:**

INTERACTIVE DASHBOARDS ENABLE BETTER ENGAGEMENT, DYNAMIC DATA EXPLORATION, AND ACTIONABLE INSIGHTS.

 **CHAPTER 6: CASE STUDY – DATA STORYTELLING IN ACTION**

◆ **6.1 REAL-WORLD EXAMPLE – E-COMMERCE BUSINESS ANALYTICS**

A GLOBAL E-COMMERCE COMPANY ANALYZES SALES DATA TO IDENTIFY GROWTH OPPORTUNITIES.

STEP-BY-STEP APPROACH:

- 1 DATA COLLECTION** – CONNECT POWER BI TO SQL DATABASE WITH SALES DATA.
- 2 DATA CLEANING & PREPARATION** – HANDLE MISSING VALUES, FILTER DATA.
- 3 VISUALIZATION & ANALYSIS** – USE BAR CHARTS, LINE GRAPHS, AND GEOSPATIAL MAPS.

4 STORYTELLING INSIGHTS – EXPLAIN TRENDS (E.G., HIGHER SALES ON WEEKENDS).

5 ACTIONABLE RECOMMENDATIONS – ADJUST MARKETING SPEND AND SUPPLY CHAIN BASED ON INSIGHTS.

 **OUTCOME:**

THE COMPANY INCREASED SALES BY 15% BY REALLOCATING MARKETING BUDGETS BASED ON DATA-DRIVEN INSIGHTS.

 **CONCLUSION:**

DATA STORYTELLING TRANSFORMS RAW DATA INTO BUSINESS INTELLIGENCE, LEADING TO PROFITABLE DECISIONS.

 **SUMMARY & NEXT STEPS**

 **KEY TAKEAWAYS:**

-  DATA STORYTELLING BRIDGES THE GAP BETWEEN ANALYTICS AND DECISION-MAKING.
-  POWER BI & TABLEAU ARE INDUSTRY-LEADING TOOLS FOR INTERACTIVE DATA VISUALIZATION.
-  CREATING MEANINGFUL DASHBOARDS IMPROVES ENGAGEMENT AND INSIGHTS.

👉 **NEXT STEPS:**

- ◆ PRACTICE CREATING POWER BI & TABLEAU DASHBOARDS WITH REAL DATASETS.
- ◆ LEARN ADVANCED VISUALIZATIONS (DAX, CALCULATED FIELDS, AI INTEGRATION).
- ◆ STAY UPDATED ON NEW TRENDS IN BUSINESS INTELLIGENCE & ANALYTICS. 🚀

ISDM-Nxt



ASSIGNMENT:

CONDUCT AN EDA REPORT ON A DATASET AND PRESENT STATISTICAL INSIGHTS USING VISUALIZATION TOOLS.

ISDM-NxT

SOLUTION: EXPLORATORY DATA ANALYSIS (EDA) REPORT WITH STATISTICAL INSIGHTS AND VISUALIZATIONS

OBJECTIVE:

THE GOAL OF THIS ASSIGNMENT IS TO PERFORM **EXPLORATORY DATA ANALYSIS (EDA)** ON A DATASET, DERIVE MEANINGFUL STATISTICAL INSIGHTS, AND VISUALIZE KEY PATTERNS USING **PANDAS, MATPLOTLIB, AND SEABORN**. EDA IS CRUCIAL FOR UNDERSTANDING DATA STRUCTURE, DETECTING PATTERNS, HANDLING MISSING VALUES, IDENTIFYING OUTLIERS, AND PREPARING DATA FOR MODELING.

STEP 1: INSTALL & IMPORT REQUIRED LIBRARIES

BEFORE STARTING, ENSURE THAT YOU HAVE THE NECESSARY LIBRARIES INSTALLED. IF THEY ARE NOT INSTALLED, RUN:

PIP INSTALL PANDAS NUMPY MATPLOTLIB SEABORN

NOW, IMPORT THE REQUIRED LIBRARIES:

IMPORT PANDAS AS PD

IMPORT NUMPY AS NP

IMPORT MATPLOTLIB.PYTHON AS PLT

IMPORT SEABORN AS SNS

🛠️ STEP 2: LOAD THE DATASET

FOR THIS ASSIGNMENT, WE WILL USE THE **TITANIC DATASET**, WHICH CONTAINS INFORMATION ON PASSENGERS WHO TRAVELED ON THE TITANIC, INCLUDING WHETHER THEY SURVIVED.

DOWNLOAD THE DATASET FROM **KAGGLE** OR USE THE FOLLOWING COMMAND:

```
# LOAD THE DATASET  
DF = PD.READ_CSV("TITANIC.CSV")  
  
# DISPLAY FIRST FIVE ROWS  
DF.HEAD()
```

📌 DATASET DESCRIPTION:

- **PASSENGERID** – UNIQUE IDENTIFIER FOR EACH PASSENGER.
- **SURVIVED** – SURVIVAL STATUS (1 = SURVIVED, 0 = DID NOT SURVIVE).
- **PCLASS** – PASSENGER CLASS (1ST, 2ND, OR 3RD CLASS).
- **NAME** – PASSENGER'S NAME.
- **SEX** – GENDER (MALE/FEMALE).
- **AGE** – PASSENGER'S AGE.
- **SIBSP** – NUMBER OF SIBLINGS/SPOUSES ABOARD.

- **PARCH** – NUMBER OF PARENTS/CHILDREN ABOARD.
- **TICKET** – TICKET NUMBER.
- **FARE** – TICKET FARE PRICE.
- **CABIN** – CABIN NUMBER.
- **EMBARKED** – PORT OF EMBARKATION (C = CHERBOURG, Q = QUEENSTOWN, S = SOUTHAMPTON).

🛠 STEP 3: UNDERSTANDING THE DATASET STRUCTURE

CHECK DATASET INFORMATION

DF.INFO()

CHECK FOR MISSING VALUES

DF.ISNULL().SUM()

SUMMARY STATISTICS OF NUMERICAL COLUMNS

DF.DESCRIBE()

📌 KEY OBSERVATIONS:

- MISSING VALUES EXIST IN THE AGE, CABIN, AND EMBARKED COLUMNS.
- SOME CATEGORICAL VARIABLES NEED TO BE CONVERTED INTO NUMERICAL FORMAT FOR ANALYSIS.

- THE DATASET HAS A MIX OF NUMERICAL AND CATEGORICAL VARIABLES.
-

🛠 STEP 4: DATA CLEANING AND HANDLING MISSING VALUES

◆ 4.1 HANDLING MISSING VALUES

FILL MISSING AGE VALUES WITH THE MEDIAN

```
DF['AGE'].FILLNA(DF['AGE'].MEDIAN(), INPLACE=TRUE)
```

FILL MISSING EMBARKED VALUES WITH THE MOST FREQUENT CATEGORY (MODE)

```
DF['EMBARKED'].FILLNA(DF['EMBARKED'].MODE()[0],  
INPLACE=TRUE)
```

DROP THE CABIN COLUMN (TOO MANY MISSING VALUES)

```
DF.DROP(COLUMNS=['CABIN'], INPLACE=TRUE)
```

📌 WHY THESE METHODS?

- AGE IS NUMERICAL, SO WE USE THE MEDIAN TO AVOID SKEWING THE DISTRIBUTION.
- EMBARKED IS CATEGORICAL, SO WE USE THE MODE (MOST FREQUENT VALUE).
- CABIN HAS TOO MANY MISSING VALUES, MAKING IT UNRELIABLE FOR ANALYSIS.

◆ 4.2 REMOVING DUPLICATES (IF ANY)

```
# CHECK FOR DUPLICATE ROWS
```

```
DF.DUPLICATED().SUM()
```

```
# REMOVE DUPLICATES
```

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

📌 **INSIGHT:** REMOVING DUPLICATES ENSURES THAT EACH OBSERVATION IS UNIQUE, PREVENTING BIAS IN ANALYSIS.

🛠 STEP 5: STATISTICAL INSIGHTS & VISUALIZATIONS

◆ 5.1 SURVIVAL RATE ANALYSIS

```
PLT.FIGURE(figsize=(6,4))
```

```
SNS.COUNTPLOT(x="SURVIVED", data=DF,  
PALETTE="COOLWARM")
```

```
PLT.TITLE("SURVIVAL COUNT (0 = No, 1 = Yes)")
```

```
PLT.SHOW()
```

```
# PERCENTAGE OF SURVIVORS
```

```
SURVIVAL_RATE = DF['SURVIVED'].MEAN() * 100
```

```
PRINT(f"SURVIVAL RATE: {SURVIVAL_RATE:.2f}%")
```

📌 **INSIGHT: THE MAJORITY OF PASSENGERS DID NOT SURVIVE, WITH A SURVIVAL RATE OF APPROXIMATELY 38%.**

◆ **5.2 SURVIVAL RATE BY GENDER**

```
PLT.FIGURE(FIGSIZE=(6,4))
```

```
SNS.BARPLOT(X="SEX", Y="SURVIVED", DATA=DF,  
PALETTE="MAKO")
```

```
PLT.TITLE("SURVIVAL RATE BY GENDER")
```

```
PLT.SHOW()
```

📌 **INSIGHT: THE SURVIVAL RATE FOR FEMALES IS SIGNIFICANTLY HIGHER THAN FOR MALES, INDICATING A POSSIBLE "WOMEN AND CHILDREN FIRST" POLICY DURING THE EVACUATION.**

◆ **5.3 PASSENGER CLASS AND SURVIVAL**

```
PLT.FIGURE(FIGSIZE=(6,4))
```

```
SNS.BARPLOT(X="PCLASS", Y="SURVIVED", DATA=DF,  
PALETTE="VIRIDIS")
```

```
PLT.TITLE("SURVIVAL RATE BY PASSENGER CLASS")
```

```
PLT.SHOW()
```

📌 **INSIGHT: FIRST-CLASS PASSENGERS HAD THE HIGHEST SURVIVAL RATE, WHILE THIRD-CLASS PASSENGERS HAD THE**

LOWEST, SUGGESTING THAT WEALTH AND STATUS INFLUENCED SURVIVAL CHANCES.

◆ 5.4 AGE DISTRIBUTION OF PASSENGERS

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

```
SNS.HISTPLOT(DF["AGE"], BINS=30, KDE=TRUE, COLOR='BLUE')
```

```
PLT.TITLE("AGE DISTRIBUTION OF PASSENGERS")
```

```
PLT.XLABEL("AGE")
```

```
PLT.YLABEL("COUNT")
```

```
PLT.SHOW()
```

📌 **INSIGHT:** THE MAJORITY OF PASSENGERS WERE BETWEEN 20-40 YEARS OLD, WITH A FEW ELDERLY PASSENGERS.

◆ 5.5 RELATIONSHIP BETWEEN AGE AND SURVIVAL

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

```
SNS.BOXPLOT(X="SURVIVED", Y="AGE", DATA=DF,  
PALETTE="COOLWARM")
```

```
PLT.TITLE("AGE DISTRIBUTION BY SURVIVAL STATUS")
```

```
PLT.SHOW()
```

📌 **INSIGHT:**

- THE MEDIAN AGE OF SURVIVORS IS LOWER THAN NON-SURVIVORS, SUGGESTING THAT YOUNGER PASSENGERS HAD A BETTER SURVIVAL RATE.
 - OLDER PASSENGERS HAD A HIGHER RISK OF NOT SURVIVING.
-

◆ 5.6 FARE DISTRIBUTION BY PASSENGER CLASS

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

```
SNS.BOXPLOT(X="PCLASS", Y="FARE", DATA=DF,  
PALETTE="MAGMA")
```

```
PLT.TITLE("FARE PRICES BY PASSENGER CLASS")
```

```
PLT.YSCALE("LOG") # LOG SCALE TO HANDLE LARGE FARE  
DIFFERENCES
```

```
PLT.SHOW()
```

📌 INSIGHT:

- FIRST-CLASS PASSENGERS PAID SIGNIFICANTLY HIGHER FARES COMPARED TO SECOND- AND THIRD-CLASS PASSENGERS.
- THERE WERE SOME EXTREMELY HIGH FARES (OUTLIERS), POSSIBLY DUE TO LUXURIOUS CABINS.

🛠 STEP 6: CORRELATION ANALYSIS

◆ 6.1 HEATMAP OF CORRELATION BETWEEN VARIABLES

```
PLT.FIGURE(FIGSIZE=(10,6))
```

```
SNS.HEATMAP(DF.CORR(), ANNOT=TRUE, CMAP="COOLWARM",
FMT=".2F")

PLT.TITLE("CORRELATION HEATMAP")

PLT.SHOW()
```

📌 INSIGHT:

- FARE AND PASSENGER CLASS ARE STRONGLY CORRELATED (HIGHER CLASS = HIGHER FARE).
- SURVIVAL IS NEGATIVELY CORRELATED WITH PASSENGER CLASS (FIRST-CLASS PASSENGERS HAD A HIGHER SURVIVAL RATE).
- SIBSP AND PARCH SHOW WEAK CORRELATIONS WITH SURVIVAL.

✅ SUMMARY OF EDA REPORT:

1 DATA CLEANING:

- HANDLED MISSING VALUES IN AGE AND EMBARKED.
- DROPPED THE CABIN COLUMN DUE TO EXCESSIVE MISSING VALUES.
- CHECKED FOR AND REMOVED DUPLICATE RECORDS.

2 KEY STATISTICAL INSIGHTS & VISUALIZATIONS:

- OVERALL SURVIVAL RATE WAS ~38%.
- WOMEN HAD A HIGHER SURVIVAL RATE THAN MEN.

- FIRST-CLASS PASSENGERS HAD THE HIGHEST SURVIVAL RATE.
- YOUNGER PASSENGERS HAD A BETTER SURVIVAL CHANCE.
- HIGHER FARE PASSENGERS HAD A GREATER PROBABILITY OF SURVIVAL.

➡ NEXT STEPS:

- ◆ PERFORM FEATURE ENGINEERING (E.G., CREATING NEW FEATURES FROM EXISTING DATA).
- ◆ APPLY MACHINE LEARNING MODELS TO PREDICT SURVIVAL.
- ◆ EXPLORE ADVANCED VISUALIZATIONS USING PLOTLY FOR INTERACTIVE ANALYSIS.

ISDM-N