

ANALYZING QUANTITATIVE DATA USING AI



VIRTUAL WORKSHOP
ORGANISED BY LEAD CITY UNIVERSITY



WORKSHOP DETAILS



MODE
Virtual



ORGANISED BY
Lead City University



FOCUS
Leveraging AI to
analyze quantitative
data for research
and decision-making.



FACILITATOR



**DR. BAMIRO
Nurudeen Babatunde**

PhD, MSc, M.Ed, B.Sc (Ed), NCE



PROFESSIONAL ROLE

✓ Senior Lecturer



DEPARTMENT

✓ Department of
Economics



FACULTY

✓ Faculty of Management
and Economics



**UNIVERSITI PENDIDIKAN
SULTAN IDRIS**

Perak, Malaysia



WHY THIS WORKSHOP MATTERS



Enhance data
analysis accuracy



Extract deeper
insights



Save time and
improve efficiency



Support evidence-
based decisions



Stay ahead with
AI-driven tools

“ Empowering researchers and professionals with AI to transform
quantitative data into **meaningful insights**. ”

PREAMBLE & APPRECIATION



OPENING REMARK

I am truly honoured to be invited as a facilitator for the workshop organised by Lead City University, to engage with the highly relevant and timely topic:

“ ANALYSING QUANTITATIVE DATA USING CHATGPT. ”



ACKNOWLEDGEMENT TO PARTICIPANTS

- ✓ To all participants, I commend your presence and commitment to ongoing knowledge sharing session.
- ✓ Your participation reflects dedication to continuous learning and collective growth.



SIGNIFICANCE OF YOUR ATTENDANCE

- ✓ Your attendance underscores a shared dedication to harnessing innovative tools.
- ✓ Together, we are advancing scholarship through smart, data-driven solutions.



**TOGETHER, WE ARE EMBRACING INNOVATION
TO TRANSFORM RESEARCH AND ACADEMIC PRACTICE.**



LEARNING OUTCOMES

By the end of this workshop, participants should be able to:

01

Design and develop **research questionnaire**



02

Choose appropriate method of **data analysis** in research investigation



03

Use ChatGPT to run and simplify **quantitative data analysis**.



04

Apply ChatGPT AI-assisted techniques for **interpreting** and **presenting** research data effectively.



Better Questions.

Better Analysis.

Better Decisions.



INTRODUCTION



The **rapid advancement** in artificial intelligence and open-source AI software (ChatGPT, Deepseek) has created **unhindered access to analytical tools** and **cost-effective data analysis opportunities** for educators to enhance their research.



“Knowledge without application is simply knowledge. Applying the knowledge to one’s life is wisdom — and that is the ultimate virtue”
— **Kasi Kaye Iliopoulos**



**EMPOWERING
EDUCATORS**

Equipping educators with powerful tools for better research.



**SMARTER
ANALYSIS**

Transforming data into meaningful insights.



**COST-
EFFECTIVE**

Leveraging free and open-source AI tools for impact.

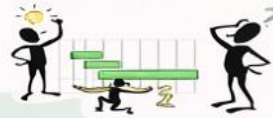


**RESEARCH
EXCELLENCE**

Driving innovation and academic excellence.

Concept of Variables in Research

What is a Variable?



What is a Variable?

Variables are **measurable characteristics** or **attributes** of a **concept, event, or entity** that can assume **different values** and can be **examined**, quantified, described, and interpreted to understand relationships or patterns.



A Variable is...

- ★ **something that can change or take on different values.** Examples:
 - weight, height, anxiety levels,
 - income, temperature.



Importance of Identifying Variables

- Identifying variables** is critical to define relationships and measure them accurately.
- More variables** increase complexity in analysis and data collection.



Example of Variables in Research

- “**Age**” and “**Height**” are the variables if examining in a sample of children; their values vary from child to child.

⌚ Age → 📏 Height



Example of Variables in Research
“**Age**” and “**Height**” are the variables if examining in a sample of children; their values vary from child to child.



Types of Variables



Independent Variable (IV)

The variable you **manipulate or categorize** to examine its effect on another variable.

Purpose: Acts as the cause or predictor.



Type of teaching method (lecture, group discussion, e-learning) in an education study.



Dependent Variable (DV)

The variable **measured** to assess the **effect** of the independent variable.

Purpose: Acts as the outcome or effect.



Students' test scores after applying different teaching methods.



Extraneous Variable

Variables not of primary interest but may affect the outcome if not controlled.

Example: Room temperature during an experiment on concentration.



Types of Variables



4. Confounding Variable



A type of **extraneous variable** that influences both the independent and dependent variables, potentially misleading the results.

Example: In a study on exercise and weight loss, diet could be a confounding variable.



5. Intervening Variable



A broader term for any variable that occurs between the IV and DV in the causal chain. It **may** act as a **mediator**, but it can also be simply a step in the sequence without a full explanatory role.

Example: Learning environment in a study linking a new teaching method to student engagement.



6. Mediating Variable



A variable that explains the **mechanism or process** through which the independent variable (IV) influences the dependent variable (DV).

Example: Self-efficacy mediates the relationship between training (IV) and job performance.



7. Moderating Variable



A variable that changes the strength or **direction** of the relationship between the independent and dependent variables.

Example: Age moderates the effect of training (IV) on job performance (DV): younger employees may respond differently compared to older employees.



Variable Operationalization



Operationalization is the process of translating abstract concepts into measurable variables.

- ✓ It involves defining the procedures or steps to observe and measure the concept in a study.

Example: “Socioeconomic *Status*” operationalized as a combination of annual income, highest education level, and job category

- ✓ Converting concepts or variables in Research Question and Research Hypothesis into measurable characteristic
- ✓ Identifying IV, DV and Other Variable in RQ and RH



Steps in Variable Operationalization

- 1 Identify the abstract concept (e.g., “Job Satisfaction”)
- 2 Define the dimensions of the concept (e.g., work environment, salary satisfaction, career growth)
- 3 Choose indicators for each dimension (e.g., Likert scale questions, performance ratings)
- 4 Specify measurement tools and scales (e.g., surveys, tests, observational checklists)
- 5 Ensure variables are measurable, valid, and reliable



Steps in Variable Operationalization



1 Step 1: Identify the Abstract Concept

- An abstract concept is a broad, theoretical idea that forms the basis of your research.
- This concept is often derived from literature, theory, or observation.

Example: “Job Satisfaction” in organizational psychology



2 Step 2: Define the Dimensions of the Concept

Dimensions are the specific aspects or facets of a concept.

- Breaking a concept into dimensions helps capture its complexity.

Example: For “Job Satisfaction”, possible dimensions include work environment, salary satisfaction, and career growth.



3 Step 3: Choose Indicators for Each Dimension

Indicators are observable signs or measures that reflect the presence or absence of each dimension.

- **Example:** For “salary satisfaction,” an indicator could be self-reported agreement with “I am fairly paid for my work”.



4 Step 4: Specify Measurement Tools and Scales

Measurement tools and scales translate indicators into quantifiable or recordable data.

- **Examples:** Surveys, structured interviews, observational checklists, standardized tests
- **From the guide:** This step is part of “operationalization” – defining exactly how we will observe and measure a concept.
- **Scales:** Likert scales (1–5), frequency counts, percentage scores



5 Ensure Variables are Measurable, Valid, and Reliable

- A measurable variable can be observed and recorded using clear procedures based on literature.
- **Validity** ensures the variable truly **measures** what it is supposed to measure (Validation process).
- **Reliability** ensures consistent results under consistent conditions.



SMART Tips for Items Formation in Questionnaire Design

S SPECIFIC

- Write clear & precise questions
- Focus on one idea
- Avoid double-barreled questions



M MEASURABLE

- Ensure questions allow for quantifiable responses
- Use Likert scales or numeric options
- Ask respondents to rate, rank, or choose from a list



A ATTAINABLE

- Match the questions to the target audience
- Ensure the topic is familiar and easy to understand
- Avoid overly complex or difficult questions



R RELEVANT

- Ensure questions are related to the research objective
- Remove unrelated or off-topic questions
- Keep items focused on key constructs of interest



T TRANSPARENT

- Use simple and direct language
- Avoid jargon or technical terms
- Ensure response options are clear & distinct



Additional Tips:

- ✓ Start with easy, non-sensitive questions
- ✓ Order items logically and group by topic

Additional Tips

- ✓ Start with easy, non-sensitive questions
- ✓ Order items logically and group by topic
- ✓ Randomize or rotate to avoid order bias
- ✓ Pilot test to refine questions & identify issues

How to Choose the Right Likert Scale

What is a Likert Scale?

A Likert scale is used to measure:



Attitudes Opinions Perceptions Satisfaction



Frequency Intensity Intensity

Why Use Likert Scales?

- ✓ Supports quantitative analysis
- ✓ Easy for respondents
- ✓ Simple to analyze
- ✓ Measures attitudes & perceptions
- ✓ Useful in Surveys & SEM Studies

Key Idea: Converts subjective opinions into measurable data.

Most Common Likert Scale Types



Agreement
Strongly Disagree → Agree



Frequency
Never → Always



Frequency
Never → Always



Satisfaction
Very Dissatisfied → Very Satisfied



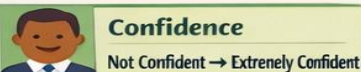
Importance
Not Important → Very Important



Likelihood
Very Unlikely → Very Likely



Quality
Very Poor → Excellent



Confidence
Not Confident → Extremely Confident



Usefulness
Not Useful → Extremely Useful

How to Choose a Likert Scale

Step 1: Identify what you want to measure

Construct	Recommended Scale
Opinion	Agreement
Behavior	Frequency
Feeling	Satisfaction / Comfort
Judgment	Importance / Quality
Probability	Likelihood

Mistakes to Avoid

- ✓ Mixing response directions
- ✓ Unclear or inconsistent labels
- ✓ Too many response options
- ✓ Double-barreled questions

Additional Tips

- ✓ Keep scales consistent
- ✓ Pilot test first
- ✓ Use 5 or 7-point scales
- ✓ Allow for neutral options

5-Point Likert Scale Example:

Strongly Disagree | Disagree | Neutral | Agree | Strongly Agree

Ethical Issues in Questionnaire Formation

Sources of Questionnaire Design

Self-Construction

- Creating items independently.
- Requires piloting to check validity and reliability.
- Risk of bias and incomplete item coverage.

Literature-Based Item Generation

- Generating items using theory or previous studies.
- Increases item validity and reliability.
- Risk of misrepresenting original intent.

Adoption and Adaptation

- Using existing questionnaires for your study.
- Adoption: Using entire questionnaires in original form.
 - Adaptation: Modifying existing items to fit new context.
 - Requires permission for ethical use.

Common Ethical Issues



Informed Consent

Ensuring participants understand the purpose, risks, and their rights to withdraw.



Privacy & Confidentiality

Protecting respondent identities and ensuring that responses remain confidential.



Informed Consent

Risk of misrepresenting "if amrurient".



Sensitive & Triggering Questions

Avoiding questions that cause discomfort or distress.

Common Ethical Issues



Pilot Test Your Questionnaire

Conduct pilot testing to assess clarity and identify ethical concerns.



Create a Clear Informed Consent Form

Clearly state the study's purpose, procedures, risks, and their rights.



Protect Respondent Privacy

Ensure data is anonymized and kept confidential.



Handle Sensitive Topics Carefully

Be cautious with sensitive questions and provide support resources.



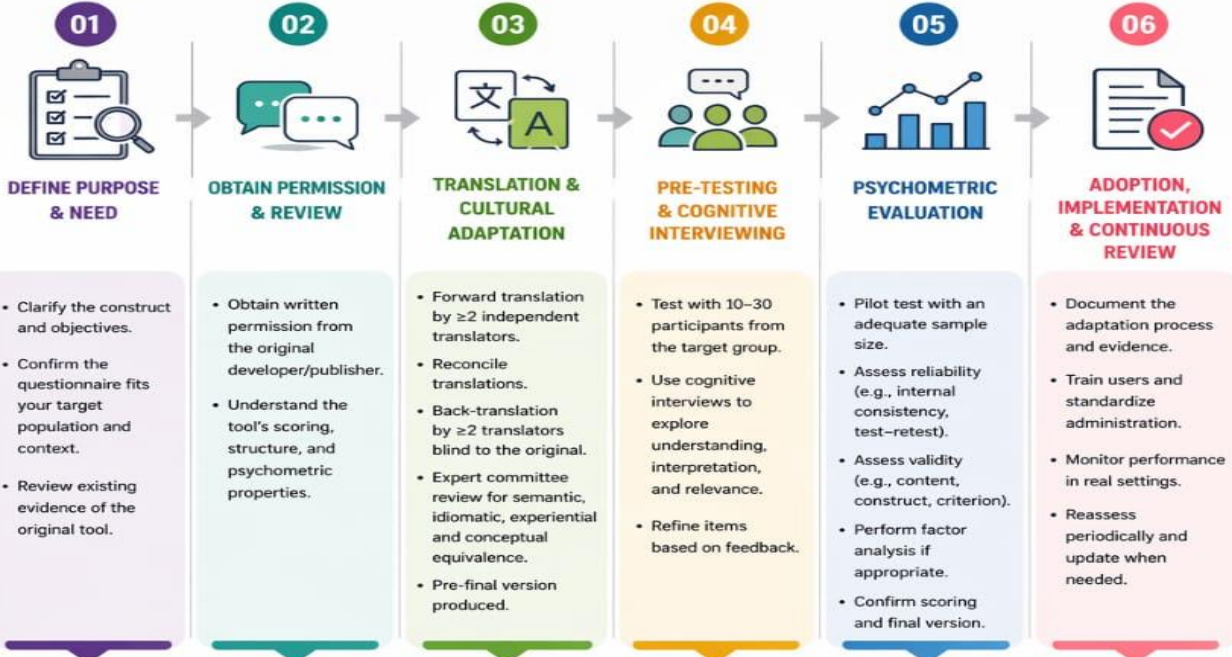
Maintain Honesty and Integrity

Report findings accurately and ensure transparency throughout the research process.

QUESTIONNAIRE ADAPTATION & ADOPTION PROCESS



Adapting or adopting a questionnaire means more than translation. It requires a systematic, rigorous process to ensure the tool is **relevant, understandable, and valid** for the new context.



CAUTION & TIPS



REMEMBER: A well-adapted questionnaire improves data quality, supports valid conclusions, and respects the people and context it is meant to serve.

STRATEGIES TO SEARCH ONLINE QUESTIONNAIRES

FOR ADAPTATION OR ADOPTION



Finding the right questionnaire is the first step to good research. Use a **systematic and critical** approach to identify **valid, reliable, and context-appropriate** instruments.



CAUTIONS – AVOID THESE

- Using unvalidated Google Forms or informal questionnaires
- Copy-pasting items without citing original authors
- Ignoring cultural and contextual differences
- Removing or changing items without strong justification
- Relying only on convenience or popularity of the scale

ALWAYS DO THESE

- Use validated instruments first
- Document the adaptation/adoption process transparently
- Pilot test in the target population
- Assess and report reliability and validity
- Maintain original scale structure (unless justified)
- Cite the original source properly

PRO STRATEGY (For High-Impact Research)

- Use a systematic search (PRISMA-style) across multiple databases
- Conduct bibliometric mapping (Scopus/WoS) to identify key scales and trends
- Screen instruments based on psychometric quality and relevance
- This strengthens methodological rigor and increases publication acceptance in Q1 journals.

REMEMBER: The best questionnaire is not the most popular one, but the most valid, reliable, and suitable for your context. Good instruments lead to good data. Good data lead to good science.



WHAT IS DATA ANALYSIS?



1



DATA ANALYSIS DEFINITION

Data analysis is the process of **collecting, organizing, interpreting, and presenting** data to uncover patterns, draw conclusions, and support decision-making.

It involves using **statistical, logical, and analytical techniques** to understand what the data reveals about a phenomenon or research question.



2



WHY DATA ANALYSIS MATTERS

It plays a critical role in fulfilling the maxim **'Publish or Perish'** or, more recently, **'Be Visible or Vanish,'** by enabling the production of impactful, data-driven research.



DATA

+

INSIGHT

=

BETTER DECISIONS



Unlock the power of data.
Analyze today, impact tomorrow!





CAN CHATGPT CONDUCT DATA ANALYSIS?



Functionally, **ChatGPT** is designed to assist in data analysis with human guidance and specification of statistical tools, model specifications and variable definitions. Specifically, ChatGPT can be valuable for the following:

1



Interpret statistical results

Interpret statistical results (e.g., regression outputs, correlation matrices).

2



Help clean and organize data

Help clean and organize data (when combined with code or through structured instructions).

3



Generate scripts or syntax

Generate Python, R, or Stata scripts or syntax to run specific analyses.

4



Summarize and explain findings

Summarize and explain findings in plain language.

5



Assist in qualitative data analysis

Assist in qualitative data analysis (e.g., coding themes from interview transcripts).

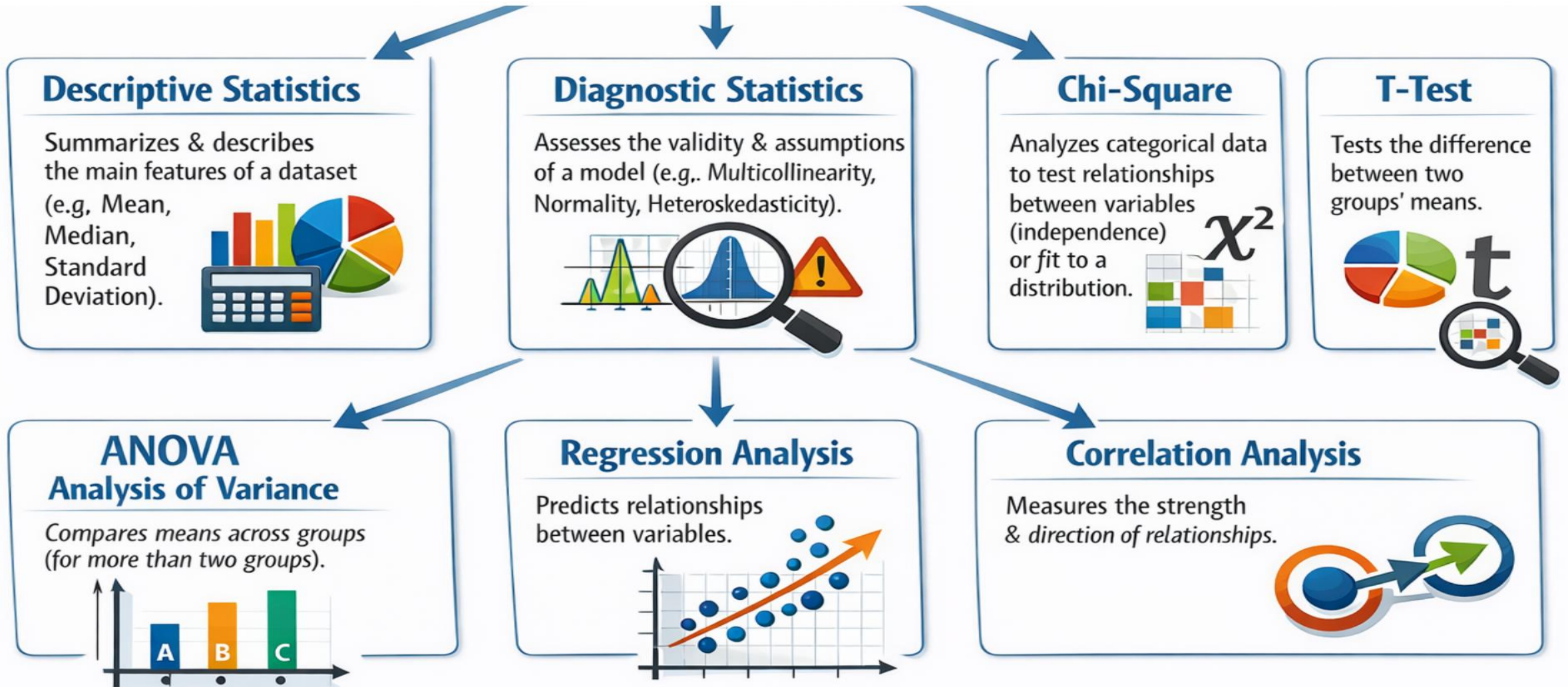


NB:

You need **premium version ChatGPT 4.0** to run reliable data analysis.



Why use statistical analysis?



The Meaning of Descriptive Statistics

Types:

Mean Median Mode



Standard Deviation



Frequency Tables



Interpretation Matters!

What does a High Mean Imply?

A High Mean indicates overall high performance!



What Does High Variability Mean in Research?

High Variability means **INCONSISTENCY** and **UNPREDICTABILITY!**



It's NOT Just About Calculations... It's About UNDERSTANDING!



Focus on Meaning!

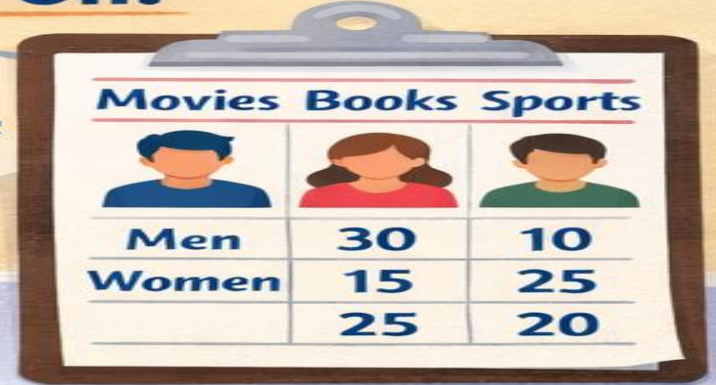
What is Chi-Square?

Types:

- Goodness-of-fit
- Test of Independence

Hands-On:

Gender vs Preference Dataset



	Movies	Books	Sports
Men	30	10	15
Women	25	20	25

Interpretation:

p-value

Low p-value?



✓ Likely significant!



Association vs Causation



Reporting:



“There is a **significant association** between gender and entertainment preferences.”

T-TEST

- Independent vs Paired T-test

- Hands-On:

- Compare:

- Male vs Female Performance

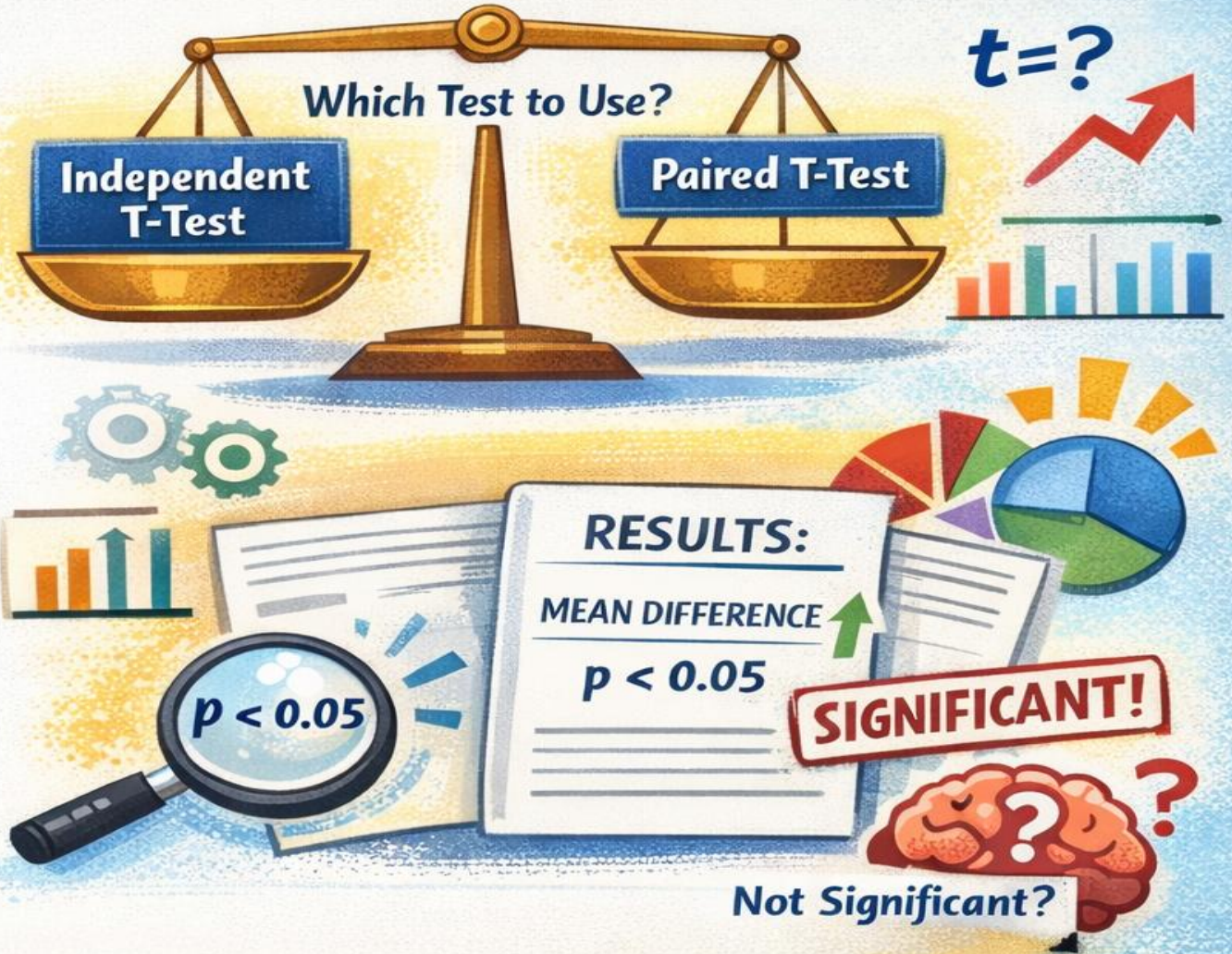


- Before vs After Intervention



- Interpreting:

- Mean Difference
- Significance



“Difference \neq meaningful unless statistically significant”


ANOVA (Analysis of Variance)

Comparing means of “three or more groups”

Example: *Income levels by education: Bachelor’s, Master’s, PhD*




One-Way ANOVA
One factor, 3+ groups



e.g. Teaching Method
(Online, Hybrid, Face-to-Face)

Two-Way ANOVA
Two factors + Interaction



e.g. Method & Gender
(Male vs Female)

Post Hoc Tests

- Tukey HSD,
- Bonferroni
- Scheffè

Find where the difference is.



Is there a difference?
ANOVA: YES... but **WHERE?**

Post Hoc Tests: Show **WHERE** the difference is!



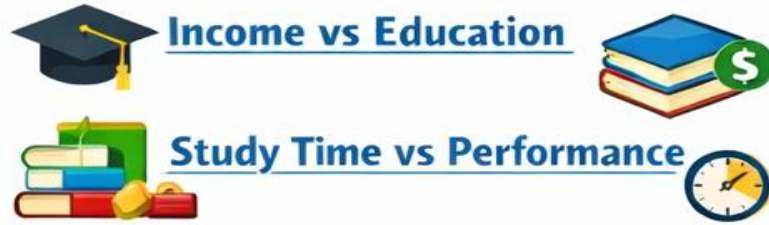
CORRELATION

Meaning of Correlation

Correlation measures the relationship between two variables.

- Do students who study more hours get better grades?
- Do people with higher education levels earn more income?

Hands-On Examples



Correlation \neq Causation!

Correlation does **NOT** imply **Causation!**

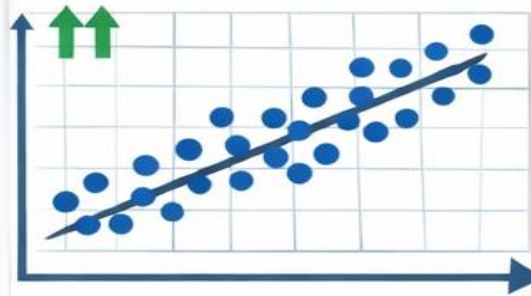


Third Variable: **Hot Weather**

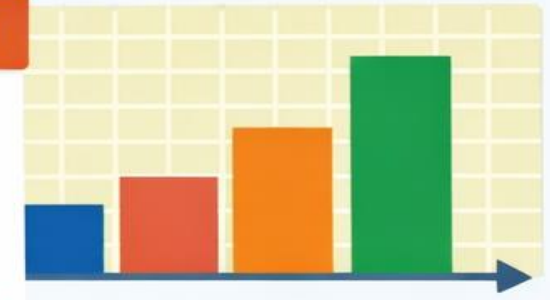
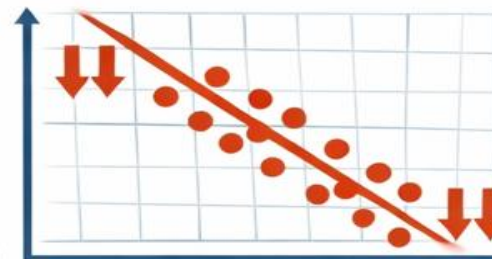
Interpreting r Values

Interpreting r Values

r Value	Meaning
0.1	Weak
0.5	Moderate
0.8	Strong



Negative Correlation



Why it Matters

- Predict Trends
- Analyze Relationships
- Make Decisions

Key Points

- **Correlation \neq Causation**
- **$r \neq$ Cause & Effect**

REGRESSION ANALYSIS OVERVIEW

HOW DOES X INFLUENCE Y? AND CAN WE PREDICT Y?

SIMPLE LINEAR REGRESSION

One Predictor (X)

$$Y = \beta_0 + \beta_1 X$$

MULTIPLE REGRESSION

Two or More Predictors (X_1, X_2, \dots)

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots$$

HANDS-ON EXAMPLE

$$\text{Performance} = 50 + 2(\text{Study Time}) + 5(\text{Motivation})$$

β COEFFICIENTS

$$\beta_1 = 2$$

1 hr study ++2 points
in performance

R^2 (R-SQUARED)

$$R^2 = 0.70$$

70% of variation
explained

SIGNIFICANCE (p-value)

$$p < 0.05$$

Significant Effect

• PREDICTION



Forecast Student Scores

• EXPLANATION



Understand Relationships

KEY CONCEPTS IN REGRESSION

• UNIT CHANGE MEANING

A 1-unit increase in X leads to a β -unit change in Y
(holding other variables constant)

1 hr Study Time \uparrow \rightarrow +2 Points in Performance

• PREDICTION vs EXPLANATION

PREDICTION

Goal: Accurate Estimates

Focus: R^2 , Model Fit



EXPLANATION

Goal: Understand Causes

Focus: β , Significance



CRITICAL WARNING: HIGH R^2 \neq CAUSATION!

High R^2 only means "Good Fit" – Not Proof of Cause



CHATGPT PROMPTS FOR DESCRIPTIVE STATISTICS & CHI-SQUARE TEST



1

DESCRIPTIVE STATISTICS

Summarize and describe the main features of your data

CHATGPT PROMPT EXAMPLE

Provide a descriptive statistical analysis for the following dataset.

Dataset: (e.g., Exam scores of 100 students)

Tasks:

- ✓ Calculate Mean, Median, Mode
- ✓ Calculate Standard Deviation and Variance
- ✓ Find Minimum, Maximum, Range
- ✓ Provide Frequency Distribution (for categorical data)
- ✓ Present results in a clear table
- ✓ Interpret the results in simple terms

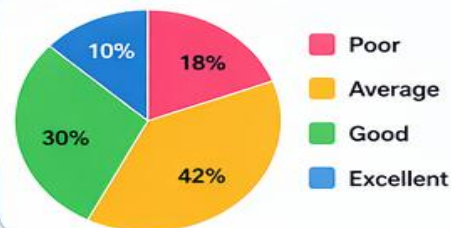
VISUAL EXAMPLE (NUMERICAL DATA)

Statistic	Exam Scores (out of 100)
Mean	72.45
Median	74.00
Mode	78
Standard Deviation	8.67
Variance	75.20
Minimum	48
Maximum	95
Range	47
Count (N)	100

FREQUENCY TABLE EXAMPLE (CATEGORICAL DATA)

Study Habit	Frequency	Percentage (%)
Poor	18	18.0%
Average	42	42.0%
Good	30	30.0%
Excellent	10	10.0%
Total	100	100.0%

VISUAL: PIE CHART EXAMPLE



INTERPRETATION EXAMPLE



The average exam score is 72.45 with a standard deviation of 8.67, which indicates that most students scored between 63.78 and 81.12.

The most common score is 78. A total of 42% of students have an average study habit.

TIPS FOR BETTER RESULTS



- ✓ Provide clear variable names and descriptions
- ✓ Mention the scale of measurement
- ✓ Specify the type of output you want (table, chart, interpretation)
- ✓ Mention the context or objective

2

CHI-SQUARE TEST

Test the association between categorical variables

χ^2

CHATGPT PROMPT EXAMPLE

Conduct a Chi-Square test of independence for the following data.

Variables:

- Gender (Male, Female)
- Preference for Online Learning (Yes, No)

Tasks:

- ✓ Create a contingency table
- ✓ Calculate Chi-Square statistic
- ✓ Determine degrees of freedom
- ✓ Compute p-value
- ✓ Test significance at 5% level ($\alpha = 0.05$)
- ✓ State the decision and interpret the result

CONTINGENCY TABLE EXAMPLE

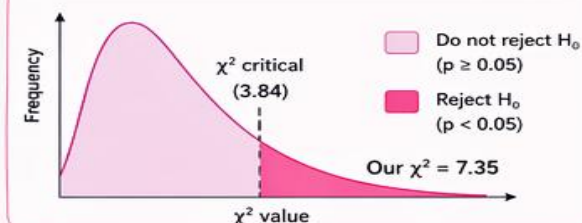
Gender vs. Preference for Online Learning

Gender	Preference for Online Learning		Total
	Yes	No	
Male	45	25	70
Female	35	45	80
Total	80	70	150

CHI-SQUARE TEST OUTPUT EXAMPLE

Statistic	Value	Interpretation
Chi-Square (χ^2)	7.35	Test statistic
Degrees of Freedom (df)	1	(rows-1) x (cols-1)
p-value	0.0067	Probability value
Significance Level (α)	0.05	Threshold
Decision	Reject H_0	Since $p < 0.05$
Conclusion	Significant association exists between gender and preference for online learning.	

VISUAL: CHI-SQUARE DISTRIBUTION



INTERPRETATION EXAMPLE

Since p-value (0.0067) is less than 0.05, we reject the null hypothesis. There is a significant association between gender and preference for online learning.

ASSUMPTIONS CHECK



- Data are in frequency counts
- Observations are independent
- Expected frequency in each cell ≥ 5 (If not, consider Fisher's Exact Test)

WHEN TO USE CHI-SQUARE TEST?



To test association between two categorical variables



To check if distribution of categories differ



Goodness-of-fit test for a single categorical variable



Widely used in survey analysis, marketing, health, social sciences

GENERAL TIP: Always provide clear data, variable names, and context in your prompt for accurate, meaningful, and well-interpreted results!





CHATGPT PROMPTS FOR DIAGNOSTIC STATISTICS



Use diagnostic statistics to check model/data assumptions, detect problems and improve analysis.

1 WHAT ARE DIAGNOSTIC STATISTICS?



Diagnostic statistics help evaluate the quality of data and the appropriateness of statistical models by checking assumptions, detecting outliers, heteroscedasticity, multicollinearity, non-normality and model fit.

2 CHATGPT PROMPT TEMPLATE

"Perform diagnostic tests for [TYPE OF MODEL/ANALYSIS].
Given the dataset with the following variables:
• [List of Variables]

Tasks:

- ✓ Check and report all relevant diagnostic statistics and tests
- ✓ Interpret the results in simple terms
- ✓ Indicate whether assumptions are met
- ✓ Suggest corrective actions if violations are found
- ✓ Provide summary and conclusions

3 EXAMPLE PROMPT (REGRESSION)

Perform diagnostic tests for a multiple linear regression model.
Dependent Variable: Sales
Independent Variables: Price, Advertising Spend, Store Size, Promotion

Tasks:

- ✓ Check normality of residuals
- ✓ Check multicollinearity
- ✓ Check heteroscedasticity
- ✓ Detect outliers and influential points
- ✓ Report all relevant diagnostic statistics and plots
- ✓ Interpret the results and give recommendations



KEY DIAGNOSTIC STATISTICS, TESTS & WHAT THEY TELL YOU

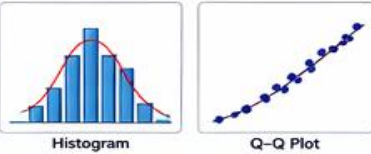
1 NORMALITY OF RESIDUALS

What it checks:
Are residuals normally distributed?

Tests / Statistics:

- Shapiro-Wilk Test
- Jarque-Bera Test
- Skewness & Kurtosis

Visuals:
Histogram, Q-Q Plot



Interpretation:

- p-value > 0.05 → Normality Assumption Met
- p-value ≤ 0.05 → Non-normal Residuals

If Violated:

- Transform variables (log, sqrt, Box-Cox)
- Use robust/bootstrapped methods

2 MULTICOLLINEARITY

What it checks:
Are independent variables highly correlated with each other?

Tests / Statistics:

- Variance Inflation Factor (VIF)
- Tolerance

Rule of Thumb:

- VIF > 5 (or 10) → Problem
- Tolerance < 0.2 → Problem

Example VIF Table

Variable	VIF	Tolerance
Price	2.45	0.41
Advertising	6.87	0.15
Store Size	1.98	0.50
Promotion	3.21	0.31

If Violated:

- Remove/Combine correlated variables
- Use PCA or Ridge Regression

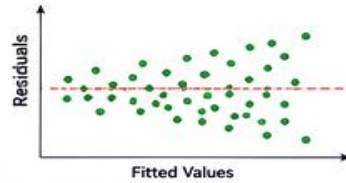
3 HETEROSCEDASTICITY

What it checks:
Is the variance of residuals constant across all levels of fitted values?

Tests / Statistics:

- Breusch-Pagan Test
- White Test

Visual:
Residuals vs Fitted Plot



Interpretation:

- p-value > 0.05 → Homoscedasticity (OK)
- p-value ≤ 0.05 → Heteroscedasticity

If Violated:

- Use robust standard errors (HC3)
- Transform dependent variable

4 OUTLIERS & INFLUENTIAL POINTS

What it checks:
Are there observations that unduly influence the model?

Measures / Statistics:

- Standardized Residuals
- Leverage (h)
- Cook's Distance (D)
- DFFITS, DFBETAS

Rules of Thumb:

- |Std. Residual| > 3 → Outlier
- Leverage > 2k/n → High leverage
- Cook's D > 1 → Influential



If Problematic:

- Investigate data
- Correct errors or consider robust methods

5 INDEPENDENCE (AUTOCORRELATION)

What it checks:
Are residuals independent from each other? (Common in time series data)

Tests / Statistics:

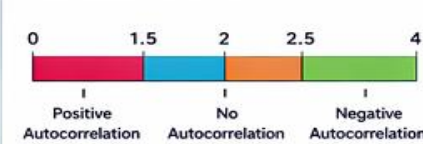
- Durbin-Watson (DW) Test

Rule of Thumb (DW Value):

- ≈ 2 → No autocorrelation
- < 1.5 → Positive autocorrelation
- > 2.5 → Negative autocorrelation

If Violated:

- Add missing variables
- Use ARIMA/GLS models



6 OVERALL MODEL FIT

What it checks:
How well does the model explain the data?

Statistics:

- R-squared (R²)
- Adjusted R²
- Root Mean Square Error (RMSE)
- AIC / BIC (for model comparison)

Interpretation:

- Higher R² (closer to 1) → Better fit
- Lower RMSE → Better accuracy
- Lower AIC/BIC → Better model



GENERAL TIPS FOR USING CHATGPT

- ✓ Provide clear variable names and definitions
- ✓ Specify the model or test you are using
- ✓ Share sample size and context
- ✓ Ask for tables, plots and interpretations
- ✓ Always ask for recommendations

COMMON MODELS & RELEVANT DIAGNOSTICS

Model Type	Key Diagnostics
Linear Regression	1, 2, 3, 4, 5, 6
Logistic Regression	2, 4, 6 (Hosmer-Lemeshow Test)
Time Series (ARIMA)	1, 5, 6 (ACF/PACF plots)
ANOVA/GLM	1, 3, 4, 6
Machine Learning	Residual Analysis, Overfitting, Feature Importance

BEST PRACTICES

- ✓ Check assumptions before interpreting results
- ✓ Do not ignore warning signs from diagnostics
- ✓ Report both statistical results and diagnostics
- ✓ Combine multiple diagnostics for a conclusion
- ✓ Document any corrective actions taken

OUTPUT YOU CAN ASK CHATGPT TO PROVIDE

- ✓ Summary table of all diagnostic tests
- ✓ Interpretation in simple language
- ✓ Diagnostic plots (with explanation)
- ✓ Whether assumptions are met or violated
- ✓ Actionable recommendations

★ **REMEMBER:** Good diagnostics lead to reliable results, valid inferences, and better decisions!



CHATGPT PROMPTS FOR T-TEST & ANOVA ANALYSIS

1

T-TEST

Compare the means of TWO groups to determine if the difference is statistically significant

CHATGPT PROMPT

Conduct an independent samples t-test.

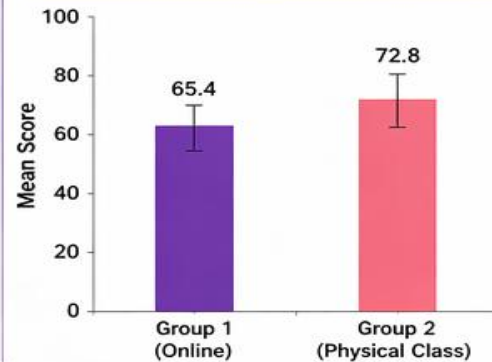
Dataset:

- Group 1: Online students' exam scores
- Group 2: Physical class students' exam scores

Tasks:

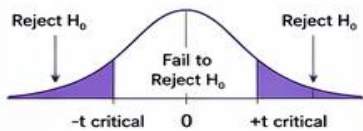
- ✓ Calculate t-statistic
- ✓ Compute p-value
- ✓ Test significance at 5% level ($\alpha = 0.05$)
- ✓ State decision (Reject or Fail to Reject H_0)
- ✓ Interpret the result in simple terms

VISUAL: TWO GROUPS COMPARISON



T-TEST OUTPUT – WHAT TO LOOK FOR

TEST STATISTIC



If $|t_{\text{calculated}}| > t_{\text{critical}} \rightarrow$ Reject H_0

EXAMPLE OUTPUT

- t(58) = -2.35
- p-value = 0.022
- $\alpha = 0.05$

DECISION RULE

- ✓ If p-value < 0.05 \rightarrow Reject H_0
- ✓ If p-value \geq 0.05 \rightarrow Fail to Reject H_0

INTERPRETATION EXAMPLE

Since p-value (0.022) < 0.05, we reject the null hypothesis. There is a significant difference in mean exam scores between the two groups.

TYPES OF T-TEST

- Independent (Two-Sample)**
Compare means between two independent groups.
- Paired (Dependent)**
Compare means from the same group at two times (e.g., before & after).

ASSUMPTIONS CHECK

- ✓ Independence of observations
- ✓ Normality of data in each group
- Homogeneity of variances (equal variances)



REAL WORLD EXAMPLE

A researcher wants to know if a new teaching method improves students' performance compared to the traditional method.

2

ANOVA (ANALYSIS OF VARIANCE)

Compare the means of THREE OR MORE groups to determine if at least one group mean is different



CHATGPT PROMPT

Perform a One-Way ANOVA.

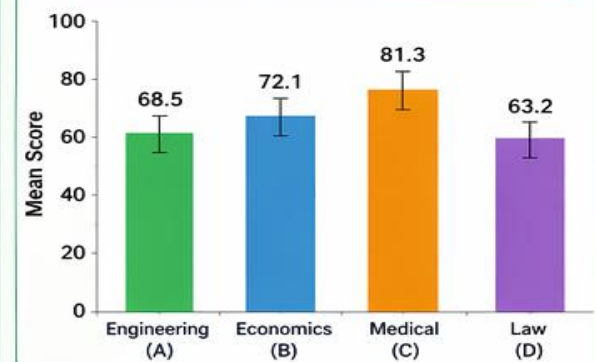
Dataset:

- Group A: Engineering students' scores
- Group B: Economics students' scores
- Group C: Medical students' scores
- Group D: Law students' scores

Tasks:

- ✓ Compute F-statistic
- ✓ Determine p-value
- ✓ Test significance at 5% level ($\alpha = 0.05$)
- ✓ State decision (Reject or Fail to Reject H_0)
- ✓ If significant, suggest post hoc test (e.g., Tukey HSD) to find which groups differ

VISUAL: THREE+ GROUPS COMPARISON



ANOVA OUTPUT – WHAT TO LOOK FOR

ANOVA TABLE (EXAMPLE)

Source	SS	df	MS	F	p-value
Between Groups	1257.6	3	419.2	5.67	0.002
Within Groups	3590.4	96	37.4	-	-
Total	4848.0	99	-	-	-

DECISION RULE

- ✓ If p-value < 0.05 \rightarrow Reject H_0 (at least one mean is different)
- ✓ If p-value \geq 0.05 \rightarrow Fail to Reject H_0 (all means are equal)

INTERPRETATION EXAMPLE

Since p-value (0.002) < 0.05, we reject the null hypothesis. There is a significant difference in mean scores among the groups. Use post hoc test to identify which groups differ.

POST HOC TEST (IF ANOVA IS SIGNIFICANT)

Use tests like Tukey HSD, Bonferroni, or Scheffé to find which specific groups are different.

Example (Tukey HSD Results)

A vs B	A vs C	A vs D	B vs C	B vs D	C vs D
No	Yes*	No	Yes*	Yes*	Yes*

* Significant difference ($p < 0.05$)

ASSUMPTIONS CHECK

- ✓ Independence of observations
- ✓ Normality within groups
- ✓ Homogeneity of variances (equal variances)



REAL WORLD EXAMPLE

A university wants to compare the average CGPA of students from multiple faculties to see if performance differs across faculties.

KEY INSIGHT

T-TEST compares TWO means, ANOVA compares THREE OR MORE means.

TIP: Provide clear data, group labels, and context for the best and most accurate analysis!

TIP: Always provide clear variable names, units, group labels, and context for accurate and meaningful results!

CHATGPT PROMPTS FOR CORRELATION & REGRESSION ANALYSIS

1

CORRELATION

Measure the strength & direction of relationship between two variables

CHATGPT PROMPT

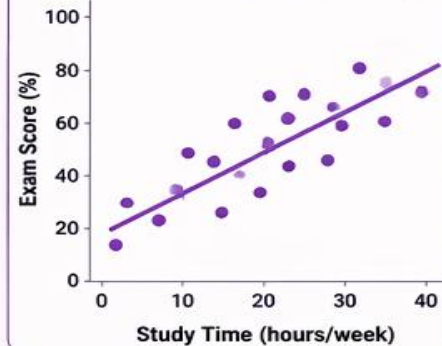
Analyze the correlation between the following variables:

- Study Time (hours per week)
- Exam Score (%)

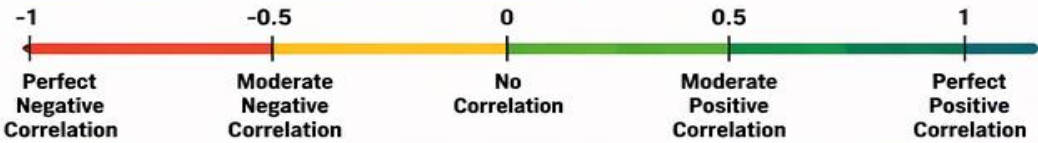
Tasks:

- ✓ Compute correlation coefficient (r)
- ✓ Interpret strength (weak / moderate / strong)
- ✓ Indicate direction (positive / negative)
- ✓ Test significance (p -value)
- ✓ Provide a brief conclusion in simple language

SCATTER PLOT EXAMPLE



INTERPRETING CORRELATION (r)



STRENGTH GUIDE ($|r|$)

0.00 – 0.19	Very Weak
0.20 – 0.39	Weak
0.40 – 0.69	Moderate
0.70 – 0.89	Strong
0.90 – 1.00	Very Strong

KEY TAKEAWAYS

- $r > 0 \rightarrow$ Positive relationship
- $r < 0 \rightarrow$ Negative relationship
- $|r|$ closer to 1 \rightarrow Stronger relationship
- Correlation does NOT imply causation.

REAL WORLD EXAMPLE

More study time is associated with higher exam scores. Strong positive correlation indicates students who study more tend to score higher.



2

REGRESSION

Predict and explain the effect of one or more independent variables on a dependent variable



CHATGPT PROMPT

Run a multiple regression analysis.

Dependent Variable:

- Exam Score (%)

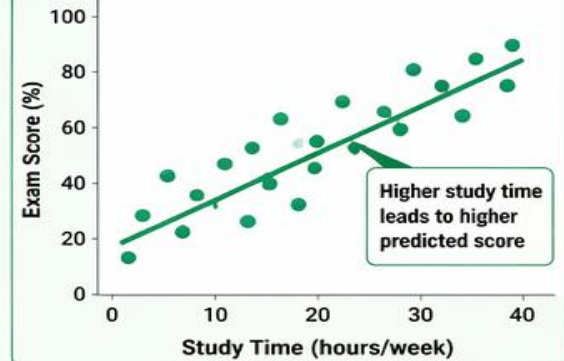
Independent Variables:

- Study Time (hours per week)
- Attendance (%)
- Motivation Level (1–10)

Tasks:

- ✓ Estimate the regression equation
- ✓ Interpret each coefficient
- ✓ Report R^2 and Adjusted R^2
- ✓ Test overall model significance
- ✓ Explain predictive meaning in simple terms

REGRESSION LINE EXAMPLE



REGRESSION OUTPUT – WHAT TO LOOK FOR

REGRESSION EQUATION

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

Example:

$$\text{Exam Score} = 20.5 + 1.8(\text{Study Time}) + 0.4(\text{Attendance}) + 2.3(\text{Motivation})$$

- **Intercept (β_0):** Expected score when all predictors = 0
- **Slopes (β):** Change in Y for 1 unit increase in X (holding others constant)

MODEL FIT (R^2)



R^2 shows % of variation in the dependent variable explained by the model.

Higher $R^2 \rightarrow$ Better fit

SIGNIFICANCE

- ✓ $p\text{-value} < 0.05 \rightarrow$ Model is statistically significant
- Check p -value for each predictor

ASSUMPTIONS CHECK

- ✓ Linearity
- ✓ Independence
- ✓ Homoscedasticity
- ✓ Normality

KEY TAKEAWAYS

Regression helps PREDICT and EXPLAIN. Coefficients show the direction and strength of effect. R^2 tells how well the model explains the data. Always check assumptions for reliable results.

REAL WORLD EXAMPLE



Using study time, attendance and motivation, we can predict a student's exam score and identify which factor has the greatest impact.



TIP: Always provide clear variable names, units, and context in your prompt for accurate and meaningful results!





RESEARCHER TASK WORKFLOW

A STEP-BY-STEP GUIDE TO DATA ANALYSIS USING AI



1



DESCRIPTIVE STATISTICS

Summarize and understand your data.



MEAN



MEDIAN



STANDARD
DEVIATION



DATA
SUMMARY

2



CHOOSE THE CORRECT TEST

Select the statistical test that best fits your research goal.



t-test

Compare
2 groups



ANOVA

Compare
3+ groups



CORRELATION

Relationship
between
variables



REGRESSION

Predict
outcome
(variable)



CHI-SQUARE

Association
(categorical
variables)

3



INTERPRET RESULTS

Give meaning to your findings.



State
Hypothesis
(H₀ & H₁)



Check
p-value



Make
Decision
(Reject / Fail to reject H₀)



Interpret & Explain
What the results
mean

4



WRITE REPORT

Communicate your results clearly and professionally.



Present
Tables & Figures



Explain
Clearly



Use
Academic Format
(APA/Chicago/Other)



KEY TAKEAWAY
AI makes analysis faster,
but understanding the method is your responsibility.



WORK SMARTER
WITH AI



ANALYZE BETTER
WITH KNOWLEDGE













RESEARCH DEEPER
WITH UNDERSTANDING



ACHIEVE IMPACT
WITH CONFIDENCE

CHOOSING THE RIGHT STATISTICAL TEST

Match Your Research Goal to the Appropriate Test

RESEARCH GOAL	TEST
 <p>COMPARE 2 GROUPS Determine whether there is a significant difference between the means of two groups.</p>	 <p>T-TEST Used to compare the means of two groups.</p>
 <p>COMPARE 3+ GROUPS Determine whether there is a significant difference among three or more group means.</p>	 <p>ANOVA Used to compare the means of three or more groups.</p>
 <p>RELATIONSHIP BETWEEN VARIABLES Determine the strength and direction of the relationship between two variables.</p>	 <p>CORRELATION Measures the strength and direction of the relationship between variables.</p>
 <p>PREDICT OUTCOME Predict the value of a dependent variable based on one or more independent variables.</p>	 <p>REGRESSION Used to model and predict outcomes based on one or more predictors.</p>
 <p>ASSOCIATION (CATEGORICAL) Determine whether there is an association between two categorical variables.</p>	 <p>CHI-SQUARE Tests whether there is a significant association between categorical variables.</p>



QUICK TIP: Choose the test that best aligns with your research question and data type.

Right Test → Valid Results → Strong Conclusions!



ADHERENCE TO THE CREP STRATEGIES



ChatGPT has the potential to **influence creative writing and research originality** when used ethically and productively.



The CREP strategies are a set of **principles** that I consistently advise researchers to follow while using ChatGPT to ensure that their work **aligns with academic values**.



C = USE ChatGPT Creatively

Use ChatGPT to brainstorm ideas, explore perspectives, and enhance the creativity of your research.



R = USE ChatGPT Responsibly

Verify information, fact-check outputs, and take ownership of the final content.



E = USE ChatGPT Ethically

Respect intellectual property, avoid plagiarism, and acknowledge AI assistance appropriately.



P = USE ChatGPT Professionally

Maintain academic integrity, use it to improve quality, and present work with professionalism.



CREP is your guide to using ChatGPT in a way that is **Creative, Responsible, Ethical, and Professional**.



Use ChatGPT wisely.
Elevate your research. Uphold your values.



Part 1

Conducting Quantitative Data Analysis with ChatGPT

using a real-world dataset (Practical Session 1)



+



OpenAI

Leverage the power of AI to analyze data,
discover insights, and make informed decisions.



Real-World
Dataset



Quantitative
Analysis



AI-Powered
Insights



Evidence-Based
Decisions



“

Technology will not
replace great teachers
but technology in
the hands of
great teachers
can be transformational.

”



George Couros



Technology
Empowers



Great Teachers
Inspire



Together
We Transform



Better Learning,
Better Future





**Questions and
Answers**

