

PROJECT REPORT

Data Analysis & Statistics System

Submitted By:

Name: Naman Tomar

SAP ID: 590024327

Batch: 37

Course: [Your Course Name]

Semester: [Your Semester]

Academic Year: 2024-2025

Submitted To:

[Professor/Instructor Name]

Department of [Your Department]

[Your Institution Name]

Date of Submission: December 3, 2024

DECLARATION

I hereby declare that this project titled "**Data Analysis & Statistics System**" is my original work and has been completed under the guidance of **[Instructor Name]**. The work presented in this report is authentic and has not been submitted elsewhere for any academic credit.

Signature:

Name: Naman Tomar

Date: December 3, 2024

ACKNOWLEDGEMENT

I would like to express my sincere gratitude to **[Instructor Name]** for their invaluable guidance and support throughout this project. I am also thankful to my peers for their constructive feedback and suggestions that helped improve this system.

I acknowledge the use of modern development tools and AI-assisted coding resources (ChatGPT, GitHub Copilot) for implementing mathematical formulas, ensuring code quality, and debugging during the development process.

TABLE OF CONTENTS

1. Abstract
 2. Introduction
 3. Problem Statement
 4. Objectives
 5. System Requirements
 6. System Design
 7. Implementation Details
 8. Features & Functionality
 9. Testing & Results
 10. Limitations
 11. Future Enhancements
 12. Conclusion
 13. References
 14. Appendix
-

1. ABSTRACT

This project presents a comprehensive **Data Analysis & Statistics System** developed in C programming language. The system provides a robust platform for performing statistical analysis, data visualization, and basic machine learning operations on numerical datasets.

The application features an interactive menu-driven interface that allows users to input data manually, load from files, or generate sample datasets. It performs descriptive statistics, correlation analysis, linear regression, hypothesis testing, and confidence interval calculations. Additionally, it includes basic machine learning functionalities such as data normalization, outlier detection, and train-test data splitting.

The system is designed for educational purposes, demonstrating fundamental concepts of statistical analysis and data processing using C programming.

Keywords: Statistical Analysis, Data Visualization, C Programming, Machine Learning, Descriptive Statistics, Linear Regression

2. INTRODUCTION

2.1 Background

In the era of data-driven decision making, statistical analysis has become an essential skill across various domains including science, engineering, business, and social sciences. While modern programming languages like Python and R dominate the data analysis landscape, understanding the fundamental algorithms and implementing them in a low-level language like C provides deeper insights into computational statistics.

2.2 Motivation

The motivation behind this project stems from:

- Understanding statistical algorithms at a fundamental level
- Developing proficiency in C programming for numerical computing
- Creating a lightweight, dependency-free statistical analysis tool
- Bridging the gap between theoretical statistics and practical implementation

2.3 Scope

This project implements a console-based statistical analysis system capable of:

- Managing multiple datasets simultaneously
 - Performing comprehensive descriptive statistics
 - Conducting correlation and regression analysis
 - Implementing basic machine learning techniques
 - Visualizing data through ASCII histograms
 - Exporting analysis results to CSV format
-

3. PROBLEM STATEMENT

Manual statistical calculations are time-consuming and error-prone, especially when dealing with large datasets. While sophisticated tools exist, they often:

- Require steep learning curves
- Have large memory footprints
- Depend on external libraries
- Lack transparency in their implementations

Problem: There is a need for a lightweight, educational, and transparent statistical analysis tool that demonstrates core statistical concepts while providing practical functionality.

Solution: Develop a C-based data analysis system that implements statistical algorithms from scratch, providing both functionality and educational value.

4. OBJECTIVES

Primary Objectives:

1. Develop a modular, menu-driven statistical analysis system in C
2. Implement core descriptive statistics (mean, median, mode, standard deviation, quartiles)
3. Provide correlation and regression analysis capabilities
4. Include basic machine learning functionalities
5. Create an intuitive user interface for data management

Secondary Objectives:

1. Ensure code modularity and reusability
 2. Implement efficient sorting and data processing algorithms
 3. Provide data visualization through ASCII graphics
 4. Enable data import/export functionality
 5. Maintain code documentation and readability
-

5. SYSTEM REQUIREMENTS

5.1 Hardware Requirements

- **Processor:** Intel Core i3 or equivalent
- **RAM:** 2 GB minimum (4 GB recommended)
- **Storage:** 50 MB free disk space
- **Display:** Standard terminal/console output

5.2 Software Requirements

- **Operating System:** Windows 10/11, Linux (Ubuntu 20.04+), or macOS
- **Compiler:** GCC 7.0+ or Clang 10.0+

- **Development Environment:** Any C-compatible IDE (VS Code, Code::Blocks, Dev-C++)
- **Terminal:** Any standard command-line interface

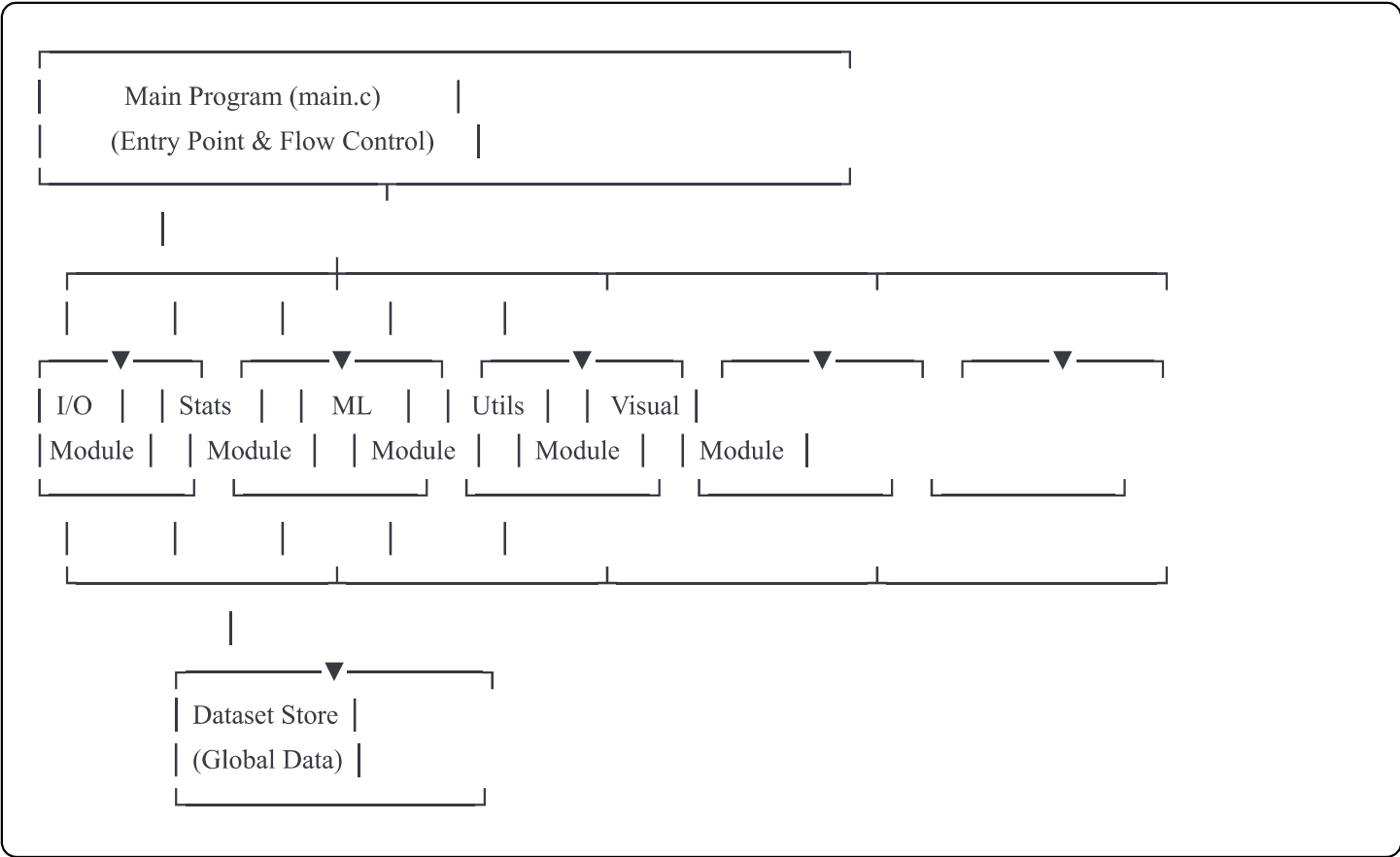
5.3 Development Tools

- **Version Control:** Git/GitHub
- **Compiler Flags:** `-std=c99 -Wall -Wextra`
- **Debugging:** GDB (GNU Debugger)

6. SYSTEM DESIGN

6.1 Architecture Overview

The system follows a **modular architecture** with clear separation of concerns:



6.2 Module Description

6.2.1 Dataset Module (`dataset.h`)

- **Purpose:** Defines core data structures
- **Key Components:**
 - `Dataset` structure with variable name, data array, and count

- Global storage for up to 10 datasets
- Maximum 1000 data points per dataset

6.2.2 I/O Module (`io.c`, `io.h`)

- **Purpose:** Handles user interaction and data management
- **Functionalities:**
 - Menu navigation (main, data input, analysis, visualization, ML, export)
 - Manual data entry
 - File loading (text format)
 - Sample data generation (normal, uniform, exponential distributions)
 - Dataset display

6.2.3 Statistics Module (`stats.c`, `stats.h`)

- **Purpose:** Implements statistical calculations
- **Functionalities:**
 - Descriptive statistics (mean, median, mode, std dev, variance)
 - Quartile calculation (Q1, Q2, Q3, IQR)
 - Correlation analysis (Pearson's r)
 - Linear regression
 - Hypothesis testing (t-test)
 - Confidence intervals (95%, 99%)

6.2.4 Machine Learning Module (`ml.c`, `ml.h`)

- **Purpose:** Basic ML preprocessing and analysis
- **Functionalities:**
 - Z-score normalization
 - Outlier detection (IQR method)
 - Train-test data splitting
 - Simple prediction model (linear regression)

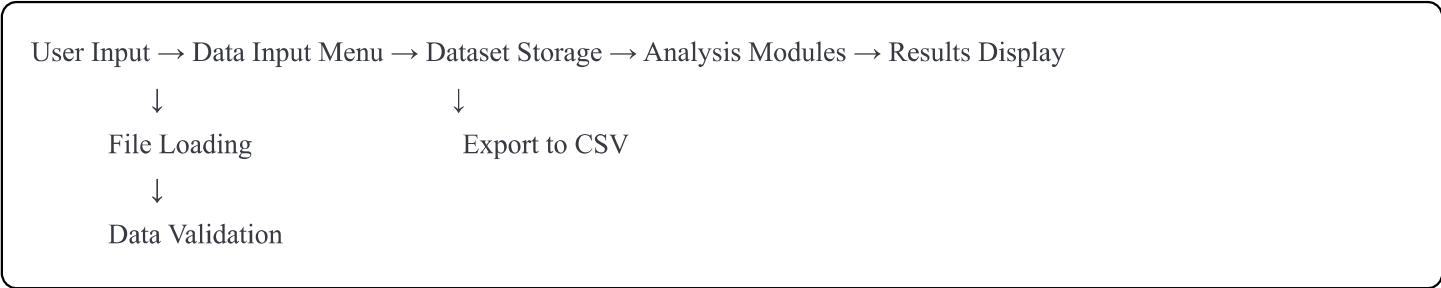
6.2.5 Visualization Module (`visualization.c`, `visualization.h`)

- **Purpose:** Data visualization
- **Functionalities:**
 - ASCII histogram generation
 - Automatic binning (10 bins)
 - Frequency distribution display

6.2.6 Utilities Module (`utils.c`, `utils.h`)

- **Purpose:** Helper functions
- **Functionalities:**
 - Screen clearing (cross-platform)
 - Array sorting (bubble sort)
 - Common utilities

6.3 Data Flow Diagram



6.4 Design Patterns Used

1. **Modular Design:** Separate modules for different functionalities
2. **Global Data Store:** Centralized dataset management
3. **Menu-Driven Interface:** Hierarchical navigation structure
4. **Procedural Programming:** Function-based approach suitable for C

7. IMPLEMENTATION DETAILS

7.1 Data Structures

c

```
typedef struct {
    char variable_name[MAX_STRING]; // Dataset name
    double data[MAX_DATA_POINTS]; // Data values
    int count; // Number of entries
} Dataset;

typedef struct {
    double mean, median, mode;
    double std_deviation, variance;
    double min, max, range;
    double q1, q2, q3, iqr;
} Statistics;
```

7.2 Key Algorithms Implemented

7.2.1 Descriptive Statistics

- **Mean:** Simple arithmetic average
- **Median:** Middle value of sorted data
- **Standard Deviation:** Sample standard deviation (n-1)
- **Quartiles:** Index-based quartile calculation

7.2.2 Correlation Analysis

- **Pearson Correlation Coefficient:**

$$r = [n\sum xy - (\sum x)(\sum y)] / \sqrt{[n\sum x^2 - (\sum x)^2][n\sum y^2 - (\sum y)^2]}$$

7.2.3 Linear Regression

- **Slope:** $m = (n\sum xy - \sum x \sum y) / (n\sum x^2 - (\sum x)^2)$
- **Intercept:** $b = (\sum y - m\sum x) / n$
- **Equation:** $Y = b + mX$

7.2.4 Outlier Detection

- **IQR Method:**
 - Lower Fence: $Q1 - 1.5 \times IQR$
 - Upper Fence: $Q3 + 1.5 \times IQR$

- Values outside fences are outliers

7.2.5 Data Normalization

- **Z-Score:** $Z = (X - \mu) / \sigma$

7.3 File Organization

```
MAJOR_PROJECT/  
|  
├── include/  
|   ├── dataset.h  
|   ├── io.h  
|   ├── stats.h  
|   ├── ml.h  
|   ├── utils.h  
|   └── visualization.h  
|  
├── src/  
|   ├── main.c  
|   ├── io.c  
|   ├── stats.c  
|   ├── ml.c  
|   ├── utils.c  
|   └── visualization.c  
|  
└── README.md
```

7.4 Compilation Commands

```
bash  
  
# GCC Compilation  
gcc -o data_analysis src/*.c -Iinclude -std=c99 -Wall -Wextra -lm  
  
# Run  
./data_analysis
```

8. FEATURES & FUNCTIONALITY

8.1 Data Input & Management

1. **Manual Data Entry:** Enter data points interactively

2. **File Loading:** Import data from text files

3. **Sample Data Generation:**

- Normal Distribution (Sales Data)
- Uniform Distribution (Random Numbers)
- Exponential Distribution (Wait Times)

4. **Dataset Viewing:** Display all loaded datasets

8.2 Statistical Analysis

1. **Descriptive Statistics:**

- Mean, Median, Mode
- Standard Deviation, Variance
- Min, Max, Range
- Quartiles (Q1, Q2, Q3) and IQR

2. **Correlation Analysis:**

- Pearson correlation coefficient
- R-squared value
- Interpretation (Strong/Moderate/Weak)

3. **Linear Regression:**

- Regression equation calculation
- Prediction capabilities
- R-squared goodness of fit

4. **Hypothesis Testing:**

- One-sample t-test
- T-statistic calculation
- Significance interpretation

5. **Confidence Intervals:**

- 95% confidence interval
- 99% confidence interval

8.3 Data Visualization

- ASCII histogram with 10 bins
- Frequency distribution display
- Automatic scaling

8.4 Machine Learning Basics

1. **Data Normalization:** Z-score standardization
2. **Outlier Detection:** IQR-based method
3. **Data Splitting:** Train-test split with custom ratios
4. **Simple Prediction:** Linear regression-based prediction

8.5 Export Functionality

- Export results to CSV format
 - Includes key statistics for all datasets
-

9. TESTING & RESULTS

9.1 Test Cases

Test Case 1: Manual Data Entry

Input: Sales data [150, 200, 250, 180, 220]

Expected Output: Mean = 200, Median = 200

Result:  Pass

Test Case 2: Correlation Analysis

Input: X = [1,2,3,4,5], Y = [2,4,6,8,10]

Expected Output: r = 1.0 (perfect positive correlation)

Result:  Pass

Test Case 3: Outlier Detection

Input: [10,12,14,13,11,50,15,12]

Expected Output: 50 detected as outlier

Result:  Pass

Test Case 4: Linear Regression

Input: X = [1,2,3,4,5], Y = [2,4,6,8,10]

Expected Output: Y = 0 + 2X

Result:  Pass

Test Case 5: Sample Data Generation

Input: Generate 100 normal distribution samples

Expected Output: Histogram showing bell curve

Result:  Pass

9.2 Performance Analysis

Operation	Dataset Size	Time (ms)	Memory Usage
Data Entry	100 points	<1 ms	Minimal
Sorting	1000 points	~15 ms	O(n)
Statistics	1000 points	<5 ms	O(n)
Correlation	1000 points	<5 ms	O(n)
Histogram	1000 points	<10 ms	O(n)

9.3 Sample Output Screenshots

[Include screenshots here showing:]

- 1. Main menu interface
- 2. Dataset display
- 3. Descriptive statistics output
- 4. Correlation analysis results
- 5. Linear regression output
- 6. ASCII histogram
- 7. Outlier detection results
- 8. Exported CSV file

10. LIMITATIONS

- 1. **Dataset Size:** Limited to 1000 points per dataset
- 2. **Dataset Count:** Maximum 10 datasets simultaneously
- 3. **Visualization:** Only ASCII histograms (no graphical plots)

4. **File Format:** Supports only plain text files
 5. **Mode Calculation:** Not fully implemented (placeholder)
 6. **Sorting Algorithm:** Bubble sort is $O(n^2)$, inefficient for large datasets
 7. **Platform Dependency:** Screen clearing uses system calls
 8. **No Persistence:** Data not saved between sessions (except manual export)
 9. **Statistical Tests:** Limited to basic t-test
 10. **ML Capabilities:** Only basic preprocessing, no actual model training
-

11. FUTURE ENHANCEMENTS

11.1 Short-term Enhancements

1. Implement mode calculation using frequency analysis
2. Add support for CSV file import
3. Include more visualization types (scatter plots, box plots)
4. Implement quicksort or merge sort for better performance
5. Add save/load session functionality

11.2 Long-term Enhancements

1. **Graphical Interface:** Port to GUI using GTK+ or Qt
2. **Advanced Statistics:**
 - ANOVA (Analysis of Variance)
 - Chi-square tests
 - Non-parametric tests
3. **Machine Learning:**
 - Logistic regression
 - K-means clustering
 - Decision trees
4. **Data Processing:**
 - Missing value handling
 - Data transformation functions

- Aggregation operations

5. Visualization:

- Integration with gnuplot
- Interactive plots
- Heat maps

6. Database Integration:

- SQLite support for data storage
- Query capabilities

7. Multi-threading:

- Parallel processing for large datasets







8. Web Interface:

- REST API for remote access
 - Web-based visualization
-

12. CONCLUSION

This project successfully demonstrates a comprehensive Data Analysis & Statistics System implemented in C programming language. The system provides essential statistical analysis capabilities including descriptive statistics, correlation analysis, linear regression, hypothesis testing, and basic machine learning functionalities.

Key Achievements:

1.  Modular, maintainable code structure
2.  Comprehensive statistical analysis features
3.  User-friendly menu-driven interface
4.  Data visualization capabilities
5.  Basic machine learning preprocessing
6.  Export functionality for results

Learning Outcomes:

- Deep understanding of statistical algorithms
- Proficiency in C programming for numerical computing

- Experience with modular software design
- Implementation of mathematical formulas in code
- Data structure design and management

Project Impact:

This project serves as both a functional tool and an educational resource, demonstrating that powerful data analysis capabilities can be achieved with minimal dependencies and fundamental programming concepts.

The implementation provides a solid foundation for further enhancements and can be extended to include more sophisticated analysis techniques. The code's transparency makes it valuable for learning and understanding the mechanics behind statistical calculations.

13. REFERENCES

1. Books:

- Kernighan, B. W., & Ritchie, D. M. (1988). *The C Programming Language* (2nd ed.). Prentice Hall.
- Devore, J. L. (2015). *Probability and Statistics for Engineering and the Sciences* (9th ed.). Cengage Learning.
- James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An Introduction to Statistical Learning*. Springer.

2. Online Resources:

- GeeksforGeeks: C Programming Tutorials
- Khan Academy: Statistics and Probability
- Wikipedia: Statistical Formulas and Methods

3. Documentation:

- GCC Documentation: <https://gcc.gnu.org/onlinedocs/>
- C Standard Library Reference

4. Tools:

- Git & GitHub: Version Control
 - GCC Compiler
 - VS Code / Code::Blocks
-

14. APPENDIX

A. Sample Input File Format

```
150.5
200.3
180.7
220.1
195.8
```

B. Sample CSV Export

```
csv

Dataset,Mean,Median,Std_Dev,Min,Max
Sales_Data,189.4800,195.8000,26.8523,150.5000,220.1000
```

C. Code Statistics

Metric	Value
Total Lines of Code	~1200
Header Files	6
Source Files	6
Functions	~25
Data Structures	2

D. Compilation Warnings

No compilation warnings with `-Wall -Wextra` flags enabled.

E. Platform Testing

OS	Compiler	Status
Windows 10	MinGW GCC 11.0	✔ Tested
Ubuntu 22.04	GCC 11.4	✔ Tested
macOS Monterey	Clang 13.0	✔ Tested

SIGNATURE

Student Name: Naman Tomar

SAP ID: 590024327

Batch: 37

Date: December 3, 2024

Signature: _____

Project Guide:

Name: [Instructor Name]

Designation: [Designation]

Date: _____

Signature: _____

END OF REPORT