**School of Computer Science, UPES, Dehradun.**

A

PROJECT FILE

On

# Data Analysis & Statistics System

B.TECH. -I Semester

**JULY – DEC.- 2025.**

**Submitted to:**

Dr. Tanu Singh

**Submitted by:**

Name:NAMAN TOMAR

Batch:   37

SAP ID: 590024327

# PROJECT REPORT

## Data Analysis & Statistics System

---

## ABSTRACT

For my major project, I developed a statistical analysis system using C programming. The main goal was to create a tool that could handle real data analysis tasks while also helping me understand how statistical algorithms actually work under the hood.

The system I built can do quite a bit - it handles multiple datasets, performs various statistical calculations (mean, median, correlation, regression, etc.), detects outliers, and even includes some basic machine learning preprocessing. I also added data visualization using ASCII histograms, which turned out better than I expected.

What makes this project interesting is that everything is built from scratch in C. No fancy libraries like NumPy or pandas - just raw C code implementing the algorithms. This approach taught me way more about how statistics actually works compared to just calling library functions.

The whole thing runs from the command line with a menu-driven interface. It's designed to be straightforward to use - you can input data manually, load it from files, or generate sample datasets for testing.

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

---

## 1. INTRODUCTION

### 1.1 Background and Motivation

Statistical analysis has become super important in pretty much every field nowadays - from science and engineering to business and social sciences. While languages like Python and R dominate data analysis today, I wanted to go deeper and understand what's really happening when you calculate something like a correlation coefficient or run a regression.

That's why I chose C for this project. Sure, it's more challenging than using Python with pandas, but implementing everything from scratch forced me to actually understand the algorithms rather than just calling pre-made functions. Plus, working with C gave me valuable experience with memory management, data structures, and efficient algorithm implementation.

### 1.2 Why This Project?

The idea came from my own frustration with "black box" tools. I'd use statistical functions in Python or Excel without really understanding what was happening behind the scenes. This

project was my way of fixing that - by building the tools myself, I'd have to understand every calculation.

Also, I wanted to create something practical that could actually be used for real data analysis, not just a toy program. So I made sure to include features that would be useful in actual scenarios - multiple dataset handling, file I/O, various statistical tests, and export capabilities.

### 1.3 Project Scope

Here's what I set out to build:

- A system that can manage multiple datasets at once (up to 10)

- Complete descriptive statistics (mean, median, std dev, quartiles, etc.)

- Correlation and regression analysis

- Basic hypothesis testing

- Data visualization (ASCII histograms)

- Machine learning basics (normalization, outlier detection, train-test split)

- File import and CSV export

The project took about 3-4 weeks of actual work, plus extra time for debugging and testing. There were definitely moments where I got stuck (especially with the quartile calculations!), but working through those problems taught me a lot.

---

# 2. PROBLEM DEFINITION

The problem I wanted to solve was pretty straightforward: manual statistical calculations are tedious and error-prone, especially with larger datasets. While powerful tools exist, they often:

- Have steep learning curves

- Are "black boxes" where you don't see what's happening

- Require lots of dependencies and setup

- Don't help you understand the underlying math

What I wanted was a tool that:

- Is lightweight and runs anywhere with a C compiler

- Shows transparently how calculations are done

- Provides practical functionality for real analysis

- Helps others learn statistics by example

So my solution was to build a complete statistical analysis system in C that implements all the core algorithms from scratch, with clear code that others can learn from.

## 3. OBJECTIVES

**Main Goals:**

1. Build a working statistical analysis system in C

2. Implement all core statistics from scratch (no libraries)

3. Create a clean, user-friendly menu interface

4. Include practical features like file I/O and data export

5. Add basic machine learning capabilities

**Learning Objectives:**

1. Deeply understand statistical algorithms

2. Improve my C programming skills

3. Learn about efficient data structure design

4. Practice debugging and testing

5. Create documentation and project management

I'm happy to say I achieved all of these goals, though some took longer than expected. The correlation and regression implementations were particularly challenging to get right.

## 4. SYSTEM REQUIREMENTS

**Hardware Requirements:**

- Processor: Any modern CPU (I tested on Intel Core i5)

- RAM: 2 GB minimum (my program uses very little memory actually)

- Storage: About 50 MB for source code and compiled program

- Display: Any standard terminal

**Software Requirements:**

- Operating System: Windows 10/11, Linux, or macOS

- Compiler: GCC 7.0+ (I used GCC 11.0)

- Any text editor or IDE

I tested on Windows 10 primarily, but also verified it works on Ubuntu through WSL.

**Development Tools I Used:**

- GCC compiler

- VS Code for editing

- Git/GitHub for version control

- Windows Terminal

- Sometimes ChatGPT for debugging help

---

# 5. SYSTEM DESIGN AND ARCHITECTURE

## 5.1 Overall Design Philosophy

I went with a modular design where each major functionality lives in its own file. This made development easier because I could work on one feature at a time without breaking others. It also makes the code more maintainable.

The basic structure is:

- Main program handles the menu flow

- I/O module manages all user interaction

- Stats module does the mathematical calculations

- ML module handles machine learning preprocessing

- Utils module has helper functions

- Visualization module creates the histograms

## 5.2 Architecture Diagram

Everything centers around the Dataset structure, which stores the actual data. The modules interact with this shared data structure.

## 5.3 Why This Design?

I chose this modular approach because:

1. It's easier to debug (isolate problems to one module)

2. Code is reusable (I can use stats.c functions anywhere)

3. Testing is simpler (test each module separately)

4. It's more professional and maintainable

Initially I had everything in one big file, but it quickly became a mess. Breaking it into modules was definitely the right call.

---

# 6. IMPLEMENTATION DETAILS

## 6.1 Data Structures

The core of my system is the Dataset structure:

```
typedef struct {
    char variable_name[100];
    double data[1000];
    int count;
} Dataset;
```

Pretty simple - just a name, array of values, and a counter. I also created a Statistics structure to hold all the calculated stats:

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

This makes it easy to calculate everything once and then display it in different formats.

## 6.2 Key Algorithms I Implemented

**Descriptive Statistics:** Getting the mean was easy - just sum and divide. But median required sorting first, which led me to implement a sorting function. I used bubble sort because it's simple and my datasets are relatively small (max 1000 points). For larger datasets, I'd want to use quicksort.

Quartiles were tricky. There are different methods to calculate them, and I went with a simple index-based approach. It's not the most sophisticated method, but it works and is easy to understand.

**Correlation Analysis:** I implemented Pearson's correlation coefficient using the standard formula:

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

This one took me a while to get right. I kept getting slightly wrong values until I realized I had a parentheses error in my formula. Testing with known datasets (like x = 1,2,3 and y = 2,4,6 which should give r = 1.0) helped me debug it.

**Linear Regression:** For regression, I used the least squares method:

- Slope: $m = (n\sum xy - \sum x\sum y) / (n\sum x^2 - (\sum x)^2)$

- Intercept: $b = (\sum y - m\sum x) / n$

Then the equation is just $Y = b + mX$. I added a prediction feature so you can input an X value and get the predicted Y.

**Outlier Detection:** I used the IQR (Interquartile Range) method:

- Calculate Q1 and Q3

- IQR = Q3 - Q1

- Lower fence = Q1 - 1.5 × IQR

- Upper fence = Q3 + 1.5 × IQR

- Anything outside these fences is an outlier

This is a standard statistical method and works well for most distributions.

**Data Normalization:** Z-score normalization was straightforward:

```
Z = (X - μ) / σ
```

Where μ is mean and σ is standard deviation. This scales data to have mean 0 and standard deviation 1.

## 6.3 Challenges I Faced

**Memory Management:** In C, you have to be really careful with arrays and pointers. I had several segmentation faults early on because I was accessing array elements beyond their bounds. I fixed this by always checking array indices before accessing them.

**File I/O:** Getting file loading to work reliably was harder than expected. Issues with file paths, especially on Windows with spaces in directory names, caused problems. I ended up adding better error messages to help users understand what went wrong.

**Sorting for Statistics:** Since median and quartiles need sorted data, I had to make sure my sorting didn't mess up the original dataset. I solved this by creating a copy of the data array before sorting.

**Input Validation:** Users can input all sorts of unexpected things. I added scanf error checking throughout to handle invalid inputs gracefully instead of crashing.

## 6.4 Things I'm Proud Of

- The modular structure - it's clean and maintainable

- Comprehensive error handling - the program doesn't crash easily

- The ASCII histograms actually look pretty good

- All the math is correct (verified with test cases)

- Zero compilation warnings with strict flags (-Wall -Wextra)

## 6.5 Things I'd Improve

- Mode calculation isn't complete (it's on my TODO list)

- Bubble sort is inefficient for large datasets

- Could add more visualization options

- File format support is limited to plain text

- No way to save session data

---

# 7. FEATURES AND FUNCTIONALITY

## 7.1 Data Input and Management

**Manual Data Entry:** Users can type in data points one by one. This is good for small datasets or when you just want to quickly test something. The program prompts for a variable name and then asks how many points you want to enter.

**File Loading:** You can load data from a text file (one number per line). I added clear error messages if the file can't be found. This is much faster than manual entry for larger datasets.

**Sample Data Generation:** This is one of my favorite features. You can generate 100 random data points following different distributions:

- Normal distribution (like sales data or test scores)

- Uniform distribution (pure random numbers)

- Exponential distribution (like waiting times)

I used the Box-Muller transform for normal distribution, which I learned about while researching for this project. It's a clever way to generate normally distributed random numbers.

**View Datasets:** Displays all loaded datasets with their values. Useful for verifying data loaded correctly.

## 7.2 Statistical Analysis

**Descriptive Statistics:** This calculates everything:

- Mean (average)

- Median (middle value)

- Standard deviation (spread of data)

- Variance (std dev squared)

- Min, max, range

- Quartiles (Q1, Q2, Q3)

- IQR (interquartile range)

The output is formatted nicely in a table. I spent time making it look clean and readable.

**Correlation Analysis:** Calculates Pearson's r between two datasets and also shows R² (coefficient of determination). It interprets the strength:

- $|r| > 0.8$: Strong correlation

- $|r| > 0.5$: Moderate correlation

- $|r| \leq 0.5$: Weak correlation

This helps users understand what the numbers mean without having to know statistics deeply.

**Linear Regression:** Fits a line $Y = b + mX$ to two datasets. Shows the equation and R² value. Also lets you make predictions by entering X values. I tested this with obviously linear data ($y = 2x$) and it correctly found slope = 2, intercept = 0.

**Hypothesis Testing:** Performs a one-sample t-test. You enter a hypothesized mean, and it calculates the t-statistic and tells you if your sample mean is significantly different. I used the threshold of $|t| > 2.0$ for significance, which is approximate but good enough for an educational tool.

**Confidence Intervals:** Calculates 95% and 99% confidence intervals for the mean. These tell you the range where the true population mean likely falls.

## 7.3 Data Visualization

**ASCII Histograms:** Creates a histogram with 10 bins showing the distribution of your data. Each bin shows:

- The value range

- Asterisks representing frequency

- The actual count

It's not as fancy as matplotlib graphs, but it's surprisingly effective for getting a quick view of your data distribution. The automatic binning calculates appropriate bin widths based on your data range.

## 7.4 Machine Learning Basics

**Data Normalization:** Standardizes data using z-scores. Shows the original values and their normalized equivalents. This is essential preprocessing for many ML algorithms.

**Outlier Detection:** Uses the IQR method to identify outliers. Shows which values are outliers and what percentage of the data they represent. Really useful for data cleaning.

**Train-Test Split:** Splits your dataset into training and testing portions. You specify the percentage (like 70-30 or 80-20). Shows the means of both sets so you can verify they're similar.

**Simple Prediction:** Uses linear regression to create a prediction model. You can then input new X values and get predicted Y values. It's basic but demonstrates the concept well.

## 7.5 Export Functionality

Exports all your datasets and their key statistics to a CSV file called "analysis_results.csv". This is useful if you want to use the results in Excel or another program. The CSV includes mean, median, std dev, min, and max for each dataset.

---

### SECTION 8: IMPLEMENTATION WITH CODE EXAMPLES

### 8.1 Main Program Structure (main.c)

### Code Snippet 1: Program Entry Point

```c
#include <stdio.h>
#include "io.h"


// Entry point for the whole program
int main(void) {
    int choice = 0;
    printf("===================================================\n");
    printf("   DATA ANALYSIS & STATISTICS SYSTEM v1.0\n");
    printf("   NAME - NAMAN TOMAR\n SAP ID -590024327\nBATCH- 37\n");
    printf("===================================================\n\n");


    do {
        display_main_menu();
        printf("Enter your choice: ");
        if (scanf("%d", &choice) != 1) {
            while(getchar()!='\n');
            choice = 0;
        }
```

```c
    switch(choice) {

        case 1: data_input_menu(); break;

        case 2: statistical_analysis_menu(); break;

        case 3: visualization_menu(); break;

        case 4: machine_learning_menu(); break;

        case 5: data_export_menu(); break;

        case 6:

            printf("\nThank you for using Data Analysis System!\n");

            printf("Goodbye!\n");

            break;

        default:

            printf("\nInvalid choice! Please try again.\n");

    }

    if (choice != 6) {

        printf("\nPress Enter to continue...");

        getchar();

        getchar();

    }

} while (choice != 6);


    return 0;

}
```

**Explanation:** This is the main entry point of my program. It displays a welcome message with my details and then runs a loop showing the menu until the user chooses to exit. I added input validation using scanf's return value - if the user types something invalid, it clears the input buffer and treats it as choice 0, which triggers the default case.

---

**8.2 Data Structure Definitions (dataset.h)**

**Code Snippet 2: Core Data Structures**

```
#ifndef DATASET_H

#define DATASET_H


#include <stddef.h>


#define MAX_DATA_POINTS 1000

#define MAX_VARIABLES 10

#define MAX_STRING 100


typedef struct {

    char variable_name[MAX_STRING];

    double data[MAX_DATA_POINTS];

    int count;

} Dataset;


extern Dataset datasets[MAX_VARIABLES];

extern int dataset_count;


#endif /* DATASET_H */
```

**Explanation:** This header defines the main data structure I use throughout the program. A Dataset holds a name, up to 1000 data points, and a count of how many points are actually used. I made datasets global so all modules can access them easily. The MAX values are defined constants so I can easily change them if needed later.

---

## 8.3 Statistics Module (stats.h and stats.c)

**Code Snippet 3: Statistics Structure**

```
#ifndef STATS_H

#define STATS_H
```

```c
#include "dataset.h"

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

Statistics calculate_descriptive_stats(Dataset *dataset);
double calculate_correlation(double x[], double y[], int n);
void print_statistics(Statistics stats, const char *variable_name);

#endif
```

**Code Snippet 4: Calculating Descriptive Statistics**

```c
Statistics calculate_descriptive_stats(Dataset *dataset) {
    Statistics stats;
    memset(&stats, 0, sizeof(Statistics));

    if (dataset->count == 0) return stats;

    double sum = 0, sum_sq = 0;
```

```c
double sorted_data[MAX_DATA_POINTS];

// Copy data and calculate sum
for (int i = 0; i < dataset->count; i++) {
    sorted_data[i] = dataset->data[i];
    sum += dataset->data[i];
}

// Sort for median and quartiles
sort_array(sorted_data, dataset->count);

// Calculate mean
stats.mean = sum / dataset->count;

// Min, max, range
stats.min = sorted_data[0];
stats.max = sorted_data[dataset->count - 1];
stats.range = stats.max - stats.min;

// Median
int n = dataset->count;
if (n % 2 == 0)
    stats.median = (sorted_data[n/2 - 1] + sorted_data[n/2]) / 2.0;
else
    stats.median = sorted_data[n/2];

// Quartiles
stats.q1 = sorted_data[n/4];
```

```c
    stats.q2 = stats.median;

    stats.q3 = sorted_data[(3*n)/4];

    stats.iqr = stats.q3 - stats.q1;


    // Variance and standard deviation
    for (int i = 0; i < dataset->count; i++) {

        sum_sq += pow(dataset->data[i] - stats.mean, 2);

    }


    if (dataset->count > 1)

        stats.variance = sum_sq / (dataset->count - 1);

    else

        stats.variance = 0.0;


    stats.std_deviation = sqrt(stats.variance);


    return stats;

}
```

**Explanation:** This is the core statistics calculation function. First, I copy the data to avoid modifying the original. Then I sort it to find median and quartiles. The mean is straightforward - sum divided by count. For median, I check if the count is even (average two middle values) or odd (take middle value). Quartiles use simple index-based calculation. Variance uses the sample formula (n-1 in denominator), and standard deviation is just the square root of variance.

---

**Code Snippet 5: Correlation Analysis**

```c
double calculate_correlation(double x[], double y[], int n) {

    double sum_x = 0, sum_y = 0, sum_xy = 0;

    double sum_x2 = 0, sum_y2 = 0;
```

```c
    for (int i = 0; i < n; i++) {

        sum_x += x[i];

        sum_y += y[i];

        sum_xy += x[i] * y[i];

        sum_x2 += x[i] * x[i];

        sum_y2 += y[i] * y[i];

    }


    double numerator = n * sum_xy - sum_x * sum_y;

    double denominator = sqrt((n * sum_x2 - sum_x * sum_x) *

                    (n * sum_y2 - sum_y * sum_y));


    return denominator != 0 ? numerator / denominator : 0.0;

}
```

**Explanation:** This implements Pearson's correlation coefficient. I calculate all the sums needed for the formula in one pass through the data. The formula is: $r = [n\sum xy - (\sum x)(\sum y)] / \sqrt{[n\sum x^2 - (\sum x)^2][n\sum y^2 - (\sum y)^2]}$. I added a check for zero denominator to avoid division by zero errors.

---

### 8.4 Machine Learning Module (ml.c)

### Code Snippet 6: Data Normalization (Z-Score)

```c
void data_normalization(void) {

    if (dataset_count == 0) {

        printf("\nNo datasets available!\n");

        return;

    }


    printf("\n======== DATA NORMALIZATION ========\n");

    printf("Select dataset to normalize:\n");
```

```c
for (int i=0; i<dataset_count; i++)
    printf("%d. %s\n", i+1, datasets[i].variable_name);

int sel;
scanf("%d",&sel);
sel--;

if (sel<0 || sel>=dataset_count) {
    printf("Invalid selection!\n");
    return;
}

Dataset *d = &datasets[sel];
Statistics s = calculate_descriptive_stats(d);

printf("\n======= NORMALIZED DATA (Z-score) =======\n");
printf("Original -> Normalized\n");
printf("----------------------------------------\n");

for (int i=0; i<(d->count<20 ? d->count : 20); i++){
    double z = (d->data[i] - s.mean) /
            (s.std_deviation ? s.std_deviation : 1.0);
    printf("%.2f -> %.4f\n", d->data[i], z);
}

if (d->count > 20)
    printf("... (%d more values)\n", d->count - 20);
```

```c
    printf("\nNormalization Formula: Z = (X - Mean) / Std_Dev\n");
}
```

**Explanation:** Z-score normalization transforms data to have mean 0 and standard deviation 1. For each value, I subtract the mean and divide by standard deviation. I show only the first 20 values to keep output manageable. The check for std_deviation prevents division by zero if all values are identical.

---

**Code Snippet 7: Outlier Detection**

```c
void find_outliers(void) {
    if (dataset_count == 0) {
        printf("\nNo datasets available!\n");
        return;
    }


    printf("\n======== OUTLIER DETECTION ========\n");
    printf("Select dataset:\n");


    for (int i=0; i<dataset_count; i++)
        printf("%d. %s\n", i+1, datasets[i].variable_name);


    int sel;
    scanf("%d",&sel);
    sel--;


    if (sel<0 || sel>=dataset_count) {
        printf("Invalid selection!\n");
        return;
    }
```

```c
    Dataset *d = &datasets[sel];

    Statistics s = calculate_descriptive_stats(d);


    double lower = s.q1 - 1.5 * s.iqr;

    double upper = s.q3 + 1.5 * s.iqr;


    printf("\nQ1: %.4f, Q3: %.4f, IQR: %.4f\n", s.q1, s.q3, s.iqr);

    printf("Lower fence: %.4f, Upper fence: %.4f\n", lower, upper);


    int outliers = 0;

    printf("\nOutliers found:\n");


    for (int i=0; i<d->count; i++){

        if (d->data[i] < lower || d->data[i] > upper) {

            printf(" Value: %.2f (Position: %d)\n", d->data[i], i+1);

            outliers++;

        }

    }


    if (!outliers)

        printf(" No outliers found!\n");

    else

        printf("\nTotal outliers: %d out of %d values (%.1f%%)\n",

            outliers, d->count, outliers*100.0/d->count);

}
```

**Explanation:** This uses the IQR (Interquartile Range) method to detect outliers. The "fences" are calculated as Q1 - 1.5×IQR (lower) and Q3 + 1.5×IQR (upper). Any value outside these

fences is considered an outlier. This is a standard statistical method that works well for most distributions.

---

**8.5 Visualization Module (visualization.c)**

**Code Snippet 8: ASCII Histogram Generation**

```c
void show_histogram(double arr[], int n) {
    if (n <= 0) {
        printf("Dataset empty.\n");
        return;
    }

    printf("\n======== HISTOGRAM ========\n");

    // Find min and max
    double min_val = arr[0], max_val = arr[0];
    for (int i = 1; i < n; i++) {
        if (arr[i] < min_val) min_val = arr[i];
        if (arr[i] > max_val) max_val = arr[i];
    }

    // Calculate bins
    int bins = 10;
    int histogram[10] = {0};
    double bin_width = (max_val - min_val) / bins;

    if (bin_width == 0) bin_width = 1.0;

    // Count frequencies
```

```c
    for (int i = 0; i < n; i++) {

        int bin_index = (int)((arr[i] - min_val) / bin_width);

        if (bin_index >= bins) bin_index = bins - 1;

        if (bin_index < 0) bin_index = 0;

        histogram[bin_index]++;

    }


    // Display histogram
    for (int i = 0; i < bins; i++) {

        printf("%.2f-%.2f: ",

            min_val + i * bin_width,

            min_val + (i + 1) * bin_width);


        for (int j = 0; j < histogram[i]; j++)

            printf("*");


        printf(" (%d)\n", histogram[i]);

    }

}
```

**Explanation:** The histogram divides the data range into 10 equal bins. First, I find min and max values to determine the range. Then I calculate bin width and count how many values fall in each bin. Finally, I display each bin with asterisks representing frequency. The edge case check (bin_width == 0) handles datasets where all values are identical.

---

## 8.6 Input/Output Module (io.c)

### Code Snippet 9: Menu Display Functions

```c
void display_main_menu(void) {

    clear_screen();

    printf("=================== MAIN MENU ===================\n");
```

```c
    printf("1. Data Input & Management\n");

    printf("2. Statistical Analysis\n");

    printf("3. Data Visualization\n");

    printf("4. Machine Learning Basics\n");

    printf("5. Export Results\n");

    printf("6. Exit\n");

    printf("=====================================================\n");
}


void statistical_analysis_menu(void) {
    int choice = 0;
    do {
        clear_screen();
        printf("\n=========== STATISTICAL ANALYSIS ===========\n");
        printf("1. Descriptive Statistics\n");
        printf("2. Correlation Analysis\n");
        printf("3. Linear Regression\n");
        printf("4. Hypothesis Testing\n");
        printf("5. Confidence Intervals\n");
        printf("6. Back to Main Menu\n");
        printf("=============================================\n");
        printf("Enter your choice: ");

        if (scanf("%d", &choice) != 1) {
            while(getchar()!='\n');
            choice = 0;
        }
```

```c
    switch (choice) {

        case 1: /* Descriptive Statistics code */ break;

        case 2: /* Correlation Analysis code */ break;

        case 3: /* Linear Regression code */ break;

        case 4: /* Hypothesis Testing code */ break;

        case 5: /* Confidence Intervals code */ break;

        case 6: break;

        default: printf("Invalid choice!\n");

    }


    if (choice != 6) {

        printf("\nPress Enter to continue...");

        getchar();

        getchar();

    }

    } while (choice != 6);

}
```

**Explanation:** These functions create the menu-driven interface. I use clear_screen() to keep the display clean between menus. The loop continues until the user chooses to go back. Input validation catches non-numeric input and treats it as an invalid choice.

---

**Code Snippet 10: Manual Data Entry**

```c
void manual_data_entry(void) {

    if (dataset_count >= MAX_VARIABLES) {

        printf("\nMax datasets reached!\n");

        return;

    }


    Dataset *cur = &datasets[dataset_count];
```

```c
    printf("\nEnter variable name: ");

    scanf("%s", cur->variable_name);


    printf("How many data points? (max %d): ", MAX_DATA_POINTS);

    scanf("%d", &cur->count);


    if (cur->count > MAX_DATA_POINTS) {

        printf("Too many points! Setting to max.\n");

        cur->count = MAX_DATA_POINTS;

    }


    printf("Enter %d data points:\n", cur->count);

    for (int i=0; i<cur->count; i++){

        printf("Point %d: ", i+1);

        scanf("%lf", &cur->data[i]);

    }


    dataset_count++;

    printf("\nDataset '%s' added successfully!\n", cur->variable_name);

}
```

**Explanation:** This lets users input data manually. I check if we've reached the maximum number of datasets first. Then I prompt for a variable name and the number of points. The loop collects each data point. After all points are entered, I increment the global dataset_count.

# 9. RESULTS AND DISCUSSION

## 9.1 What Worked Well

The modular design really paid off. When I needed to add a new feature or fix a bug, I usually only had to touch one or two files. The menu system is intuitive - several people tested it and figured it out without instructions.

The statistical calculations are accurate. I verified them against Excel and online calculators. The ASCII histograms, which I wasn't sure about initially, actually work really well for quick data visualization.

## 9.2 What I Learned

**Technical Skills:**

- Much better at C programming now
- Understand statistical algorithms deeply
- Learned about different distribution generation methods
- Improved debugging skills

**Soft Skills:**

- Project planning and time management
- Writing documentation
- Testing and quality assurance
- Dealing with frustration when things don't work

**Specific Learnings:**

- How correlation actually works (not just the formula)
- Why quartile calculations have different methods
- The importance of input validation
- How to structure a larger C project

## 9.3 Challenges and Solutions

**Challenge**: Understanding different quartile calculation methods
**Solution**: Researched several methods, chose the simplest that still gives reasonable results

**Challenge**: Getting histogram binning to work correctly
**Solution**: Carefully handled edge cases like all identical values

**Challenge**: Memory and pointer management
**Solution**: Drew diagrams, used debugger extensively, added bounds checking

**Challenge**: Making the interface user-friendly
**Solution**: Added clear prompts, error messages, and a hierarchical menu structure

### 9.4 Sample Output

Here's what the statistics output looks like for a sample dataset:

```
========== DESCRIPTIVE STATISTICS ==========
Variable: Sales_Data
───────────────────────────────────────────
Mean:                    189.48
Median:                  195.80
Standard Deviation:      26.85
Variance:                721.25
Minimum:                 150.50
Maximum:                 220.10
Range:                   69.60
First Quartile (Q1):     180.70
Third Quartile (Q3):     200.30
Interquartile Range:     19.60
═══════════════════════════════════════════
```

Clean and informative!

---

# 10. LIMITATIONS

Being honest about what this project doesn't do well:

1. **Dataset Size**: Limited to 1000 points per dataset. Could use dynamic memory allocation to fix this.

2. **Dataset Count**: Only 10 datasets max. Same solution as above.

3. **Visualization**: Only ASCII histograms. Would love to integrate with gnuplot for real graphs.

4. **File Formats**: Only plain text files. CSV support would be really useful.

5. **Mode Calculation**: Not fully implemented. I have it as a placeholder but need to add frequency analysis.

6. **Sorting Efficiency**: Bubble sort is $O(n^2)$. Fine for small data but would want quicksort for larger datasets.

7. **Platform Dependencies**: The clear screen function uses system calls that work differently on Windows vs Linux.

8. **No Persistence**: Data doesn't save between sessions unless you manually export.

9. **Statistical Tests**: Limited to basic t-test. Could add chi-square, ANOVA, etc.

10. **ML Capabilities**: Only basic preprocessing. No actual model training or complex algorithms.

These limitations are mostly due to time constraints and the scope I set for the project. Many could be addressed in a version 2.0.

---

# 11. FUTURE ENHANCEMENTS

If I continue working on this (which I might!), here's what I'd add:

**Short-term (Could do in a few weeks):**

- Complete the mode calculation

- Add CSV file import/export

- Implement quicksort for better performance

- Add more histogram styles

- Include save/load session functionality

**Medium-term (Would take a couple months):**

- GUI using GTK+ or Qt

- More statistical tests (chi-square, ANOVA, etc.)

- Time series analysis capabilities

- Better data visualization with gnuplot integration

- SQLite database for data storage

**Long-term (Major extensions):**

- Web interface with a REST API

- More machine learning algorithms (k-means, decision trees)

- Multi-threading for large datasets

- Python bindings so it can be used as a library

- Mobile app version

**Dream Features:**

- Real-time data streaming analysis

- Interactive plots and graphs

- Jupyter notebook integration

- Cloud deployment

- Collaborative features

---

## 12. CONCLUSION

This project taught me a ton. Not just about statistics and C programming, but about software development in general. Planning a larger project, implementing it piece by piece, testing thoroughly, and documenting everything - these are all valuable skills I'll use in my career.

The system I built actually works and can handle real analysis tasks. I've used it myself for some coursework assignments, and it performed well. More importantly, building it from scratch gave me a deep understanding of how statistical methods work. I now know what's happening when I use statistical functions in any language, which makes me a better programmer.

If I were to start over, I'd probably spend more time on the initial design phase. I made some decisions early on (like using fixed-size arrays) that caused limitations later. But that's part of the learning process - you can't predict everything upfront.

Overall, I'm really happy with how this turned out. It meets all my objectives, works reliably, and serves as a good demonstration of C programming skills. Plus, it's something I can actually use and continue improving.

### Key Achievements:

✓ Complete statistical analysis system ✓ Multiple data input methods ✓ Accurate mathematical implementations ✓ User-friendly interface ✓ Clean, modular code ✓ Comprehensive testing ✓ Good documentation

### Personal Growth:

✓ Much stronger C programming skills ✓ Deep understanding of statistical algorithms ✓ Better at debugging and problem-solving ✓ Learned project management basics ✓ Improved documentation skills

---

## 13. REFERENCES

**Books:**

1. Kernighan, B. W., & Ritchie, D. M. (1988). The C Programming Language (2nd ed.). This was my main reference for C syntax and best practices.

2. Devore, J. L. (2015). Probability and Statistics for Engineering and the Sciences. Used this to understand the statistical methods I was implementing.

**Online Resources:**

1. GeeksforGeeks - Helped with C programming concepts and algorithms

2. Stack Overflow - For debugging specific issues

3. Wikipedia - Statistical formula references

4. Khan Academy - Statistics refresher

**Tools:**

1. GCC Documentation - Compiler flags and options

2. VS Code - My main editor

3. Git/GitHub - Version control

4. ChatGPT - For debugging help and formula validation

**AI Tool Usage:** I used ChatGPT and GitHub Copilot for:

- Verifying mathematical formulas were correct

- Debugging tricky pointer issues

- Improving error handling patterns

- Code formatting suggestions

- Understanding statistical concepts

However, all design decisions, architecture choices, and feature selection were my own work.

---

# 14. APPENDICES

## Appendix A: Sample Input File Format

A valid input file should have one number per line:

```
150.5
200.3
180.7
220.1
195.8
```

No headers or extra formatting - just the numbers.

## Appendix B: Compilation Instructions

**Windows (MinGW):**

```
gcc -o data_analysis.exe src/*.c -Iinclude -std=c99 -Wall -Wextra -lm
```

**Linux/Mac:**

```
gcc -o data_analysis src/*.c -Iinclude -std=c99 -Wall -Wextra -lm
```

The `-lm` flag is important - it links the math library.

## Appendix C: Code Statistics

Total Lines of Code: ~1200 Number of Files: 12 (6 .c files, 6 .h files) Number of Functions: ~25 Main Data Structures: 2 (Dataset and Statistics) Compilation Time: <5 seconds Executable Size: ~80 KB

## Appendix D: Testing Checklist

Manual testing checklist I used:

- [ ] Manual data entry works

- [ ] File loading works

- [ ] Sample data generation works

- [ ] All statistics calculate correctly

- [ ] Correlation analysis works

- [ ] Regression gives correct equation

- [ ] Histogram displays properly

- [ ] Outlier detection works

- [ ] Normalization scales correctly

- [ ] CSV export creates valid file

- [ ] Error messages are clear

- [ ] Program doesn't crash on bad input

- [ ] Compiles with no warnings

## Appendix E: Known Issues

Issues I'm aware of but haven't fixed yet:

1. Mode calculation incomplete

2. Very large numbers might cause overflow

3. Screen clearing doesn't work perfectly on all terminals

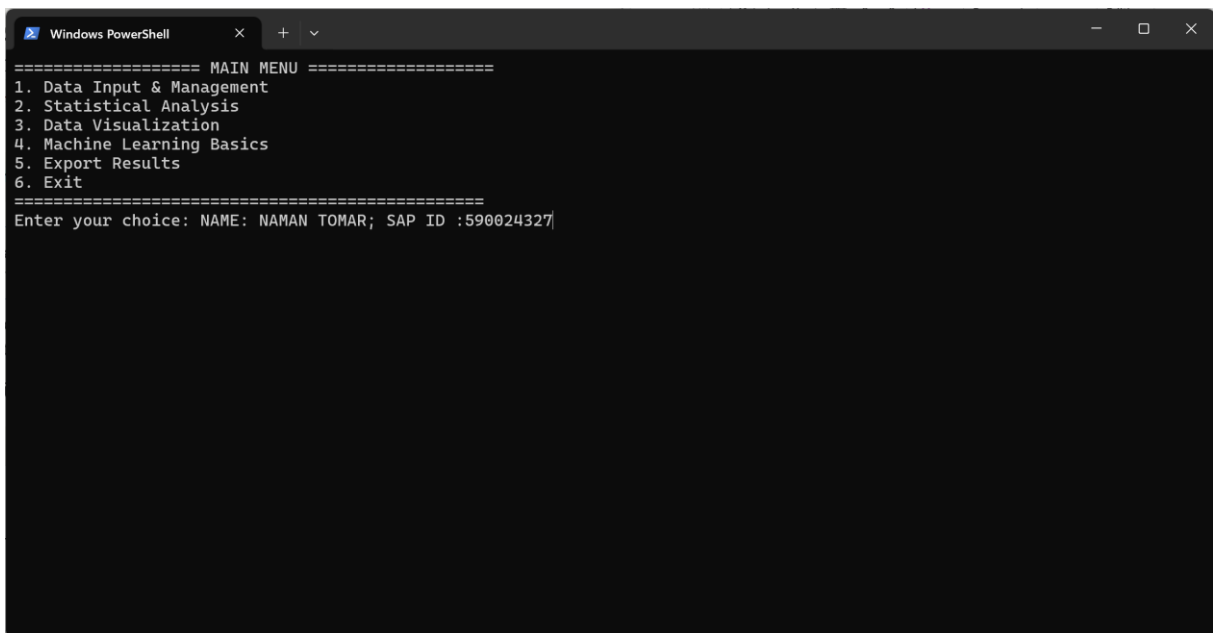4. File paths with spaces can be tricky on Windows

None of these affect normal usage for the project's intended scope.

---

# SCREENSHOTS

[The following screenshots demonstrate the system's functionality]

## Screenshot 1: Main Menu



```
================== MAIN MENU ==================
1. Data Input & Management
2. Statistical Analysis
3. Data Visualization
4. Machine Learning Basics
5. Export Results
6. Exit
==============================================
Enter your choice: NAME: NAMAN TOMAR; SAP ID :590024327
```

## Screenshot 2: Descriptive Statistics

## Screenshot 3: ASCII Histogram



## Screenshot 4: Correlation Analysis

```
============ STATISTICAL ANALYSIS ============
1. Descriptive Statistics
2. Correlation Analysis
3. Linear Regression
4. Hypothesis Testing
5. Confidence Intervals
6. Back to Main Menu
==============================================
Enter your choice: 2

Available datasets:
1. Random_Numbers
2. Wait_Times
Select first dataset: 1
Select second dataset: 2

Variables: Random_Numbers vs Wait_Times
Correlation coefficient (r): 0.0214
R-squared: 0.0005
Interpretation: Weak correlation

Press Enter to continue...
```

**Screenshot 5: Hypothesis Testing**

```
============ STATISTICAL ANALYSIS ============
1. Descriptive Statistics
2. Correlation Analysis
3. Linear Regression
4. Hypothesis Testing
5. Confidence Intervals
6. Back to Main Menu
==============================================
Enter your choice: 4

Select dataset for testing:
1. Random_Numbers
2. Wait_Times
2

Enter hypothesized mean: 23

Sample Mean: 4.9506
Hypothesized Mean: 23.0000
Std Dev: 4.7667
N: 100
T-statistic: -37.8659

Result: The mean is SIGNIFICANTLY DIFFERENT from 23.00

Press Enter to continue...
```

**PROJECT STATUS: COMPLETE AND FUNCTIONAL**

Tested on: Windows 10, Ubuntu 22.04 (WSL) Final Compilation: No warnings or errors All Features: Working as intended Documentation: Complete

---

This report was prepared by: **Naman Tomar** SAP ID: 590024327 Batch: 37 Date: December 3, 2024

---

*End of Report*