

Design and Implementation of Data Cleansing System

Founders:

Mohammed Ahmed Rashad

Mohammed Osama Abdelghany

Mohamed Adly Mohammed

Toka Saied Nafie

Mohammed Abdelkarim Mohammed

Eyad Ahmed Ali

Mohammed Abdelmotaleb Mohammed

Abdelrahman Ghoniem Ghoniem

Supervised By:

Dr. Assem El-Ashaar

Dr. Omar Khalid

Acknowledgement

In the Name of **Allah**, the Most Merciful, the Most Compassionate,

Alhamdulillah all praises belong to Almighty Allah, peace be upon **Muhammad** His servant and messenger.

Our greatest gratitude to Modern University for Technology and Information (**M.T.I.**) for the role it played in forming our personalities, providing us with the means necessary for our educational development.

We would like to seize this opportunity to express **thanks to Dr. Assem Elshaar, Dr. Omar Khalid** for their efforts and support during the course of this project.

Furthermore, we would like to show our gratitude to all our Professors and University Staff whose efforts guided us through our academic endeavors. Special thanks to our Families for their continuous support without which we would have never been able to make it.

This would not have been possible unless **Prof.Dr. Mohamed Elgazzar**, President of the MTI University has established such a learning environment, providing the most effective facilities for our job to be done.

It is a pleasure to thank those who maneuvered this Possible, **Prof. Dr. Mohamed Taher El-Mayah**, and Dean of the faculty of Computers and Information for his unlimited help, support, encouragement, and guidance. Special thanks to **Prof.Dr. Hanafy Ismail**, head of Computer Science Department , We are Honored and got the pleasure of working with **Prof. Dr. Hesham El Deeb**, Head of AI department, As well as the Department staff, **Prof.Dr Alaa Abd El-Raheem** for their effort during our university studies. Also, Special thanks to **Prof. Dr. Hafez Abdel-Wahab**, head of Information System Department, as well as, for their effort along our university studies and **Prof. Dr. Mahamoud El-Shishtawy** for their effort along our university studies.

Special thanks to **Prof. Dr. Eman Taha**, head of Basic Science Department. As well as, the department staff, and for their effort teaching us the basic for our computer science degree, **Prof. Dr. Sayed Bakar, Prof.Dr. Rasha Saeed** and **Dr. Rania Ahmed** Special thanks to our families who always supported us during our academic life and we couldn't do it without their help.

Abstract

In the current landscape, one of the paramount challenges across various sectors is maintaining data integrity and quality. Data cleansing systems emerge as a vital solution, especially in domains like finance, retail, and healthcare, where accurate data is crucial for decision-making and compliance. Our project aims to design a robust data cleansing system. This system operates through a user-friendly graphical user interface, allowing seamless interaction with **datasets & databases** of varying complexities. The development process began with extensive research into existing data cleansing methodologies and tools, identifying key requirements and approaches. Subsequently, the design phase delineated the system's functionalities via comprehensive reports, outlining the expected outcomes. Implementation involved the utilization of cutting-edge frameworks and programming languages tailored to the specific needs of data cleansing. Rigorous testing ensued, ensuring the system's efficacy in identifying and rectifying anomalies within datasets & databases. Ultimately, the decision to develop the data cleansing system as a versatile application ensures accessibility and usability across diverse organizational contexts, empowering users to maintain data integrity effortlessly.

Table of Contents

Chapter-1-Introduction.....	11
1.1Introduction	11
1.2 Problem Definition.....	12
1.3 Objective	14
1.4 Project Outline	15
Chapter-2-system planning.....	18
2.1 introduction	18
2.2 Stages Involved in Creating a Project Plan.....	19
2.3 Define the Work Breakdown Structure	20
2.4 Identify the Required Resources.....	21
2.5 Construct a Project Schedule	24
2.6 Assumptions.....	26
2.7 Constraints	26
2.8 functional view:-	27
1 st Maintain Basic Data:	27
2 nd generate reports:.....	28
2.8 proposed approach solution	30
2.9 Agile framework selection	30
2.10 Board management	31
2.11 workflow	38
Chapter 3 – Survey & Literature Review.....	41
3.1 Introduction	41
3.2 Related Work	41
3.3 Related Papers	43
Chapter-4-Analysis phase	53
4.1 Introduction	53
4.2 System Scenario	53
4.3 DFD diagram.....	56
DFD Level (0) Process:-.....	57
DFD Level (0) Data stores:.....	58
DFD Level (0) Data entities:-	58
4.4 Context diagram:-	59
Context diagram entities:-	59

4.5 Required Database Type	60
4.6 Site map	61
Chapter 5 – design phase	64
5.1 Introduction	64
5.2 Flowcharts and pseudocodes.....	65
5.3 Block diagram.....	75
5.4 Figma designs.....	76
Chapter 6– system implementation	90
6.1 Database implementation	90
6.2 Raw backend implementation:-	106
6.2.1 Handling invalid formats	106
6.2.2 Handling missing values	109
6.2.3 Removing outliers	110
6.2.4 Handling duplicates.....	118
6.3 System Modules Implementation.....	122
6.3.1 Missing values	122
6.3.2 Duplicates.....	124
6.3.3 Outliers.....	126
6.3.4 Invalid formats	129
6.4 system integration	132
Introduction to System Integration with Flask and Jinja in an MVC Architecture:.....	132
Flask Framework:	132
Jinja:	133
Template Rendering with Jinja:	133
Flask Structure Files:	134
6.5 used softwares	135
System shot screens:-	139
1-missing values page	139
2-missing values results	139
3-duplicates page	140
4-duplicates results	140
5-outliers page	141
6-outliers results	141
7-Invalid formats page	142
8-invalid formats results	142

Chapter-7-System Testing.....	145
7.1 Introduction	145
7.2 Types of Testing	145
7.3 Test cases	146
Chapter-8- Conclusion	153
8.1 introduction	153
8.2 Conclusion.....	153
8.3 Future work.....	154
8.4 Appendix:-	155

Table of figures

Figure 1 stages involved in creating a project plan	19
Figure 2 agile lifecycle	20
Figure 3 task assignment and gantt chart	24
Figure 4 network diagram	25
Figure 5 agile life cycle	30
Figure 6 Scrum meetings lifecycle	30
Figure 7 trello board management tool.....	32
Figure 8 scrum poker estimation diagram	33
Figure 9 pair programming workflow	38
Figure 10 how data powers business opportunities	48
Figure 11 wh users give false data to brands databases.....	49
Figure 12 DFD diagram.....	56
Figure 13 context diagram	59
Figure 14 NOSQL Structure	60
Figure 15 Sitemap	61
Figure 16 check email pseudocode.....	65
Figure 17 check email flowchart	65
Figure 18 check phone number pseudocode.....	66
Figure 19 check phone number flowchart.....	66
Figure 20 uppercase pseudocode	67
Figure 21 uppercase flowchart.....	67
Figure 22 lowercase pseudocode.....	67
Figure 23 lowercase flowchart	67
Figure 24 remove double spaces pseudocode.....	67
Figure 25 remove double spaces flowchart	67
Figure 26 propercase pseudocode	68
Figure 27 propercase flowchart	68
Figure 28 deal with duplicates pseudocode.....	69
Figure 29 deal with duplicates flowchart	69
Figure 30 calculate outlier pseudocode	70
Figure 31 calculate outliers flowchart	70
Figure 32 find outliers flowchart	70
Figure 33 find outliers pseudocode	70
Figure 34 missing value flowchart.....	71
Figure 35 missing value pseudocode	72
Figure 36 remove dots flowchart	73
Figure 37 remove commas flowchart	73
Figure 38 remove dots pseudocode.....	73
Figure 39 remove commas pseudocode	73
Figure 40 remove letters flowchart.....	74
Figure 41 remove numbers flowchart.....	74
Figure 42 remove numbers pseudocode	74

Figure 43 remove letters pseudocode	74
Figure 44 Block Diagram	75
Figure 45 Sign up page	77
Figure 46 Sign in page	77
Figure 47 profile popup.....	77
Figure 48 verify email.....	77
Figure 49 Homepage.....	78
Figure 50 user account slider	79
Figure 51 Footer	80
Figure 52 header	80
Figure 53 Import page.....	81
Figure 54 export page	82
Figure 55 Dashboard	83
Figure 56 Missing values page	84
Figure 57 Handling duplicates page	85
Figure 58 invalid formats page.....	86
Figure 59 Remove outliers page.....	87
Figure 60 basic structure of mongo db	90
Figure 61 flask structure files.....	134
Figure 62 missing values system shotscreen.....	139
Figure 63 missing values results system shotscreen.....	139
Figure 64 duplicates page system shotscreen.....	140
Figure 65 duplicates results system shotscreen.....	140
Figure 66 outliers page system shot screen	141
Figure 67 outliers results system shotscreen	141
Figure 68 invalid formats page system shotscreen	142
Figure 69 invalid formats results system shotscreen	142
Figure 70 testcase 1 Database upload	146
Figure 71 testcase 2 Login.....	147
Figure 72 testcase 3 register	147
Figure 73 testcase 4 remove duplicates.....	148
Figure 74 test case 5 missing values	148
Figure 75 test case 6 delete outlier.....	149
Figure 76 testcase 7 invalid formats.....	150

Chapter 1

Introduction

Chapter-1-Introduction

1.1 Introduction

A data cleansing website system serves as a fundamental tool in ensuring the integrity and reliability of organizational data. It embodies a sophisticated application designed to interact with users in a seamless and intuitive manner, akin to human conversation. At its core, this system leverages advanced technologies to identify and rectify inconsistencies, errors, and redundancies within datasets, thereby enhancing data quality and accuracy. By harnessing the power of machine learning algorithms and data processing techniques, it offers comprehensive solutions for data validation, standardization, and enrichment. From financial institutions to e-commerce platforms, data cleansing systems play a pivotal role in optimizing decision-making processes and ensuring regulatory compliance. The development journey commences with thorough research into data cleansing methodologies and user requirements, followed by meticulous design and implementation phases tailored to the unique needs of the target audience. Rigorous testing ensures the robustness and efficacy of the system, ultimately culminating in a reliable and efficient data cleansing solution. With the potential to revolutionize data management practices across industries, a data cleansing website system represents a significant advancement towards achieving data-driven excellence and operational efficiency.

1.2 Problem Definition

The primary challenge at hand revolves around the need for efficient data management solutions that mitigate the risks associated with erroneous or inconsistent datasets & databases. Organizations grapple with the arduous task of sifting through vast amounts of data plagued by inaccuracies, duplicates, and inconsistencies, leading to inefficiencies in decision-making processes and regulatory non-compliance. Moreover, the conventional methods of data cleansing often entail significant time and resource investments, hindering operational agility and scalability. Additionally, the proliferation of data across disparate systems exacerbates the challenge of maintaining data integrity and consistency. Consequently, there arises a pressing need for a comprehensive data cleansing website system that offers streamlined processes for identifying, correcting, and standardizing data anomalies. Such a solution must cater to diverse organizational needs while ensuring user-friendly interfaces and robust security measures to safeguard sensitive information. By addressing these challenges, organizations can enhance data quality, optimize operational efficiency, and foster informed decision-making, thereby unlocking new opportunities for growth and innovation.

Given that

Current data is plagued by inconsistencies, inaccuracies, and duplicates. This leads to unreliable analysis, poor decision-making, and inefficient operations.

It is required to have a system to

Data analysis: Identify inconsistencies, inaccuracies, and duplicates.

Correction: Achieve high accuracy in fixing errors automatically.

User-friendliness: Monitor, manage, and generate reports with ease.

Security & Compliance: Ensure data privacy and regulatory adherence.

Such that

Data-driven decisions are based on reliable and accurate information.

Analysis becomes more efficient and insightful due to cleaner data.

Operational costs are reduced by eliminating errors and redundancies.

Compliance and security risks are minimized through data governance.

1.3 Objective

The main objectives are:-

- In General:-

- the main goal is to make users able to clean their datasets & databases for mainly 4 problems

- decide whether they want to clean their datasets & database by their own selves or hand us theirs to clean it for them

- save time and money by cleaning their own datasets & databases instead of paying highly amounts of money to offices specialized in data cleansing services

- In the Application:-

- develop a robust query subsystem to handle our 4 database problems

1-invalid formats

2-duplicates

3-outliers

4- Missing values

- develop dashboard that tracks changes that happen after and before performing a query subsystem action

1.4 Project Outline

Chapter 2: System Planning involves outlining project team tasks and presenting planning diagrams like Gantt charts and network diagrams.

Chapter 3: Literature Review provides an overview of related systems, including examples with screenshots to illustrate their advantages, disadvantages, and a comparison of pros and cons.

Chapter 4: System Analysis details user interactions, requirements, and uses diagrams like use case diagrams, application interaction diagrams (class, sequence, and activity diagrams) to depict functionalities.

Chapter 5: System Design presents the intended system design, starting with interface prototypes and followed by system flowcharts illustrating functionalities.

Chapter 6: System Implementation discusses Android development tools with detailed screenshots of usage, showcasing application screenshots and interface descriptions.

Chapter 7: System Testing covers the testing process during and after system implementation, including user feedback evaluation through messages.

Chapter 8: Conclusion and Future Work summarizes the implemented application and outlines potential future enhancements.

Chapter 2

System Planning

Chapter-2-system planning

2.1 introduction

The purpose of the Planning stage is to analyze the project in terms of work breakdown, cost, resources, and timing. At the end of this stage all team members should be clear on the sub tasks and deliverables with the project, the time constraints they are working on too and the roles and responsibilities that are expected from them.

Software project plan defines what the work is, and how this work can be completed.

In This Chapter, the project planning is presented using the project 356 software to manage the project activities and usually with defined stages, and with designated resources the project planning contains three parts: Project task, Gantt chart and Network diagram

- Creating a comprehensive Project Plan is a critical step in the Project Lifecycle, as it is used to :
 - Monitor and control the overall progress of the project.
 - Determine whether the project activities are complete.
 - Determine if the project is ready for closure and assess the level of success of the project after it has been closed.

2.2 Stages Involved in Creating a Project Plan

First we need to define the work breakdown structure (WBS).

This means listing the processes, activities and tasks required to undertake the project, as well as the key project milestones.

Then you'll need to quantify the human resources required to carry out each activity listed.

The initial step involves building a project schedule which describes the flow of project activities and the timeframes involved, as well as any plan in assumptions and constraints.



Figure 1 stages involved in creating a project plan

2.3 Define the Work Breakdown Structure

The First step to creating a detailed project plan for your project is to develop a comprehensive WBS, listing all of the activities and tasks required to undertake the project. The life cycle defines a methodology for improving the quality of software and the overall development process

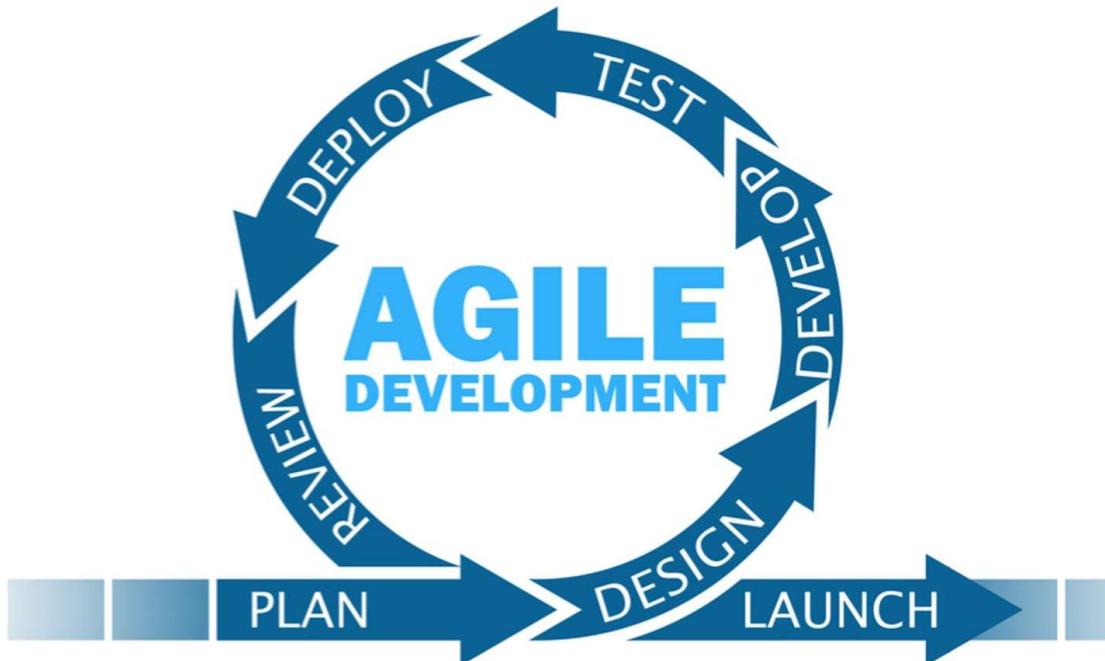


Figure 2 agile lifecycle

The project includes the following tasks:

- 1) Analysis and Data Collection
- 2) System Design
- 3) System Implementation
- 4) Testing
- 5) Documentation
- 6) Demo
- 7) Final Presentation

2.4 Identify the Required Resources

Having listed all the tasks required to undertake the project, you now need to identify the generic resources required to complete each task, as described in the table below. You will create a detailed list of the

Resources required for the project when you create a Resource Plan in the next stage

Tasks Description:-

A task is an item of work to be completed within a project. The table below will help us list each task, describing the tasks and identify the sequencing order that is appropriate. The completed activity and task lists will comprise the WBS for the project but we also need to specify any critical project milestones

• Data Collection Task

Involves actions and methods performed on data that help

Describe facts, detect patterns, develop explanations, and test

Hypotheses. This includes data quality assurance, statistical data

Analysis, modeling, and interpretation of analysis results.

• Involved to:-

- 1) Create new data
- 2) using your own previously collected data
- 3) collecting other data
- 4) purchasing data

- **System Design**

Is the process of defining the architecture, components, modules, interfaces, and data for a system to satisfy specified requirements?

Systems design could be seen as the application of systems theory to product development

- **System Implementation**

Describe the database environment where the software system and the database(s), if any is present, will be installed. Include a description of the diverse types of database and library environments (such as production, test, and training databases).

- **Documentation**

Documentation is a set of documents provided on the paper used as a quick reference or a guidebook to help the reader understand what is going on.

- **System Testing**

Is the type of testing to check the behavior of a complete and fully integrated software product based on the software requirements?

- **Demo**

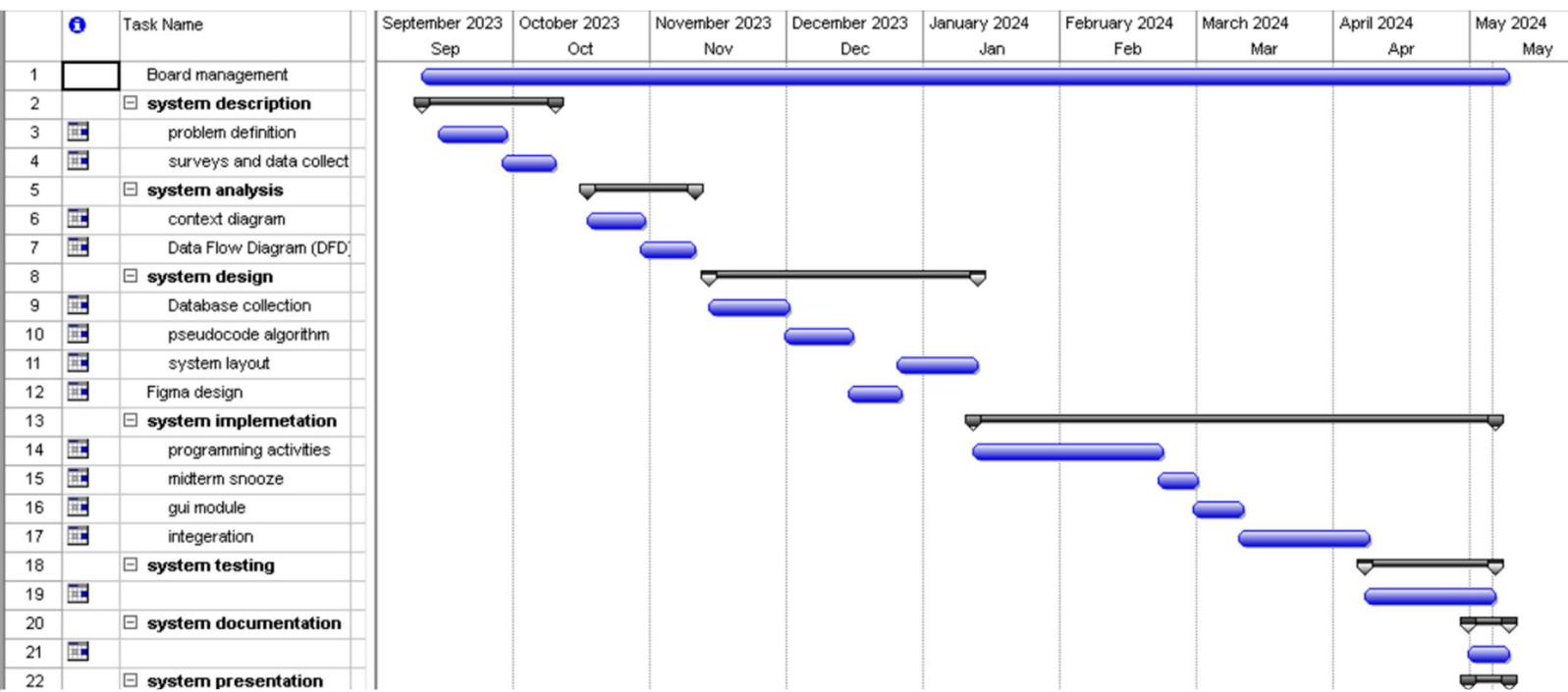
A system demo is a freely distributed piece of an upcoming or recently released app. Demos are typically released by the system publisher to help consumers get a feel of the system before deciding whether to buy the full version.

- **Presentation**

Representing the final work and everything that was done project by explaining what was learned and what was done by whom and how long it took to be done.

2.5 Construct a Project Schedule

Project Schedule Create a detailed project schedule by listing the processes, activities and tasks required to complete the project, as well as the dependencies, and sequence.



Network diagram:-

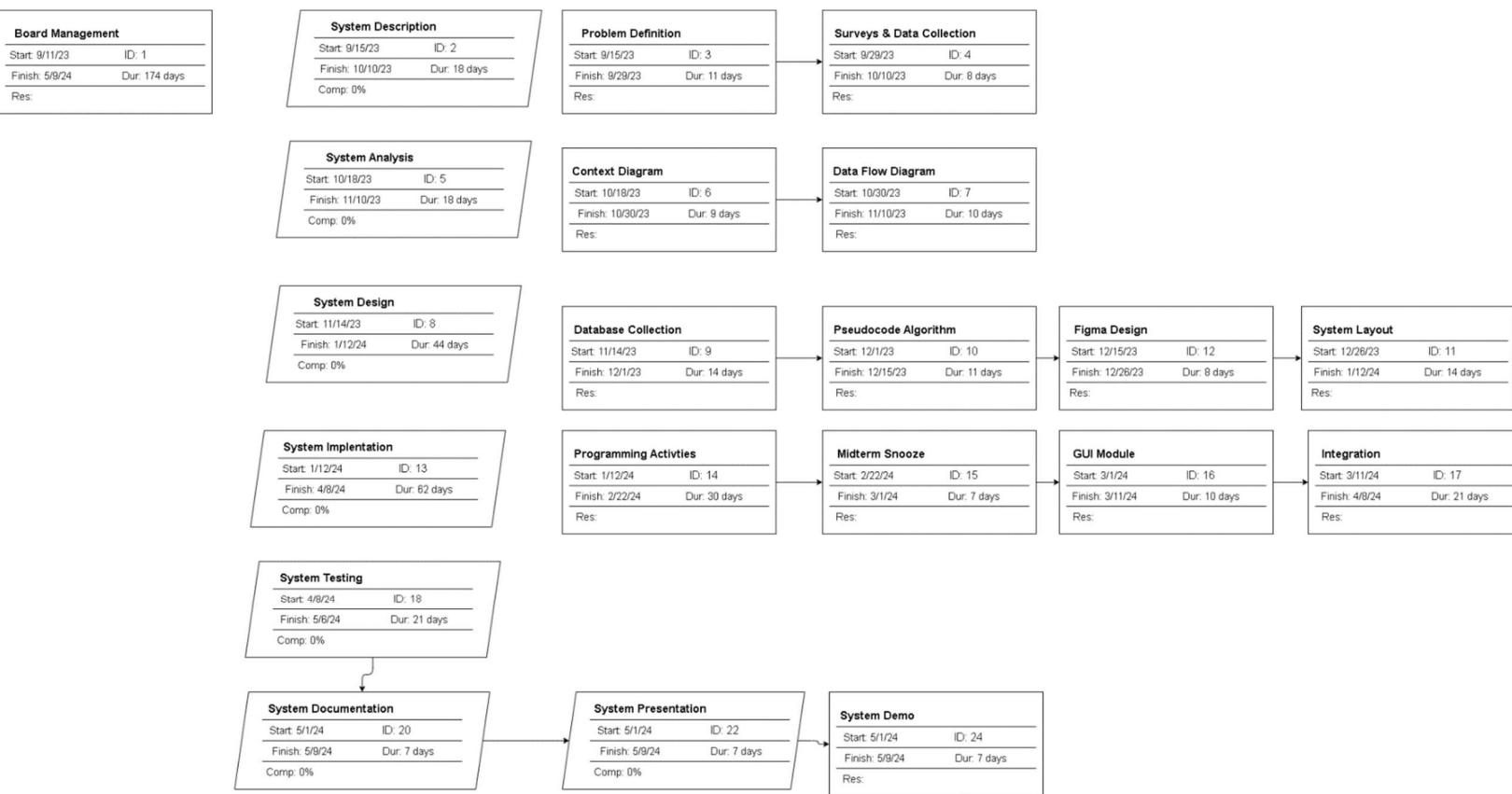


Figure 4 network diagram

2.6 Assumptions

List any assumptions made during this planning process. For example, it may be assumed that:

- The project will not change in scope.
- The resources identified will be available upon request.
- Approved funding will be available upon request.

2.7 Constraints

List any constraints identified during this planning process. For example:

- The project team must create all the physical deliverables within the allocated budget.
- Work must be undertaken within normal working hours only

2.8 functional view:-

1st Maintain Basic Data:

User:

First Name, Last Name, Email, Password, Phone, Address

Admin:

Username, Email, First Name, Last Name, Password

Log:

Username, Database Name, Log Entries

Reports:

Timestamp, Cleansing_Techniques_Used, Summary_Statistics, Missing_Values,
Duplicates, Invalid Formats, Statistical_Measures, user_inputs,
Data_Cleansing_Status

2nd generate reports:

Summary Report:

1. Total number of records processed.
2. Overall cleansing status (success/failure).
3. Timestamp of when the cleansing process occurred.
4. User responsible for the cleansing process.

Summary Statistics:

1. Negative Numbers Found: Number of negative values encountered.

Outliers:

1. Outliers Found: Number of outliers detected.
2. Outliers Removed: Number of outliers removed.
3. Percentage of outliers before and after removal.

Missing Values:

1. String: Number of missing string values.
2. Integer: Number of missing integer values.
3. Missing Values Found: Total number of missing values encountered.
4. Missing Values Removed: Number of missing values removed.
5. Percentage of missing values before and after removal.

Duplicates:

1. Duplicates Found: Number of duplicate records found.
2. Duplicates Removed: Number of duplicate records removed.
3. Percentage of duplicates before and after removal.

Invalid Formats:

1. Various types of invalid formats encountered (e.g., email, phone number, cases).
2. Invalid Formats Found: Total number of invalid format instances.
3. Invalid Formats Removed: Number of invalid format instances removed.
4. Percentage of invalid formats before and after removal.

Statistical Measures:

1. Mean: Average value of the dataset.
2. Median: Middle value of the dataset.
3. Mode: Most frequent value(s) in the dataset.
4. Standard Deviation: Measure of the dispersion of values from the mean.
5. Variance: Measure of how spread out the values are.
6. Minimum Value: Lowest value in the dataset.
7. Maximum Value: Highest value in the dataset.

Data Source Information:

1. File Name: Name of the file from which the data was sourced (if applicable).
2. Database Name: Name of the database from which the data was sourced (if applicable).
3. Table Name: Name of the table within the database from which the data was sourced (if applicable).

User Inputs:

1. Any additional information provided by the user during the cleansing process.

Data Cleansing Status:

1. Final status of the data cleansing process (e.g., completed, in progress, failed).

2.8 proposed approach solution

We chose to work as an agile self-developing team as we focused through the process in presenting deliverable items at the end of each sprint

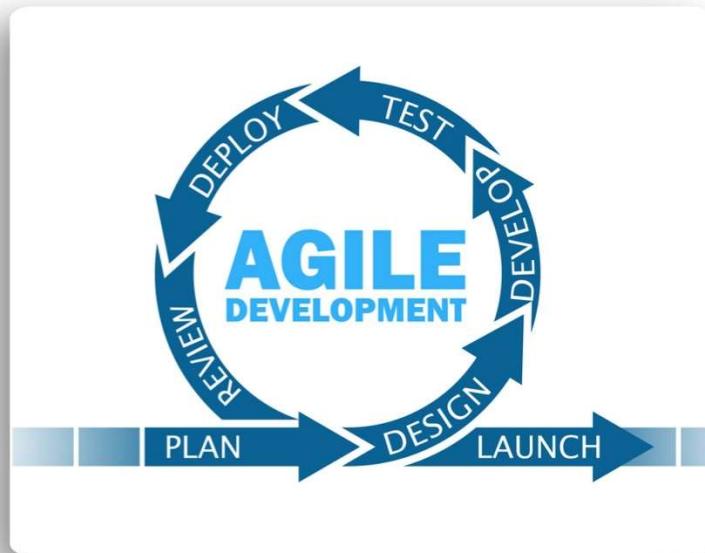


Figure 5 agile life cycle

2.9 Agile framework selection

We chose to work with a **Scrumban** framework that depended on many thing we did in the last sprint like

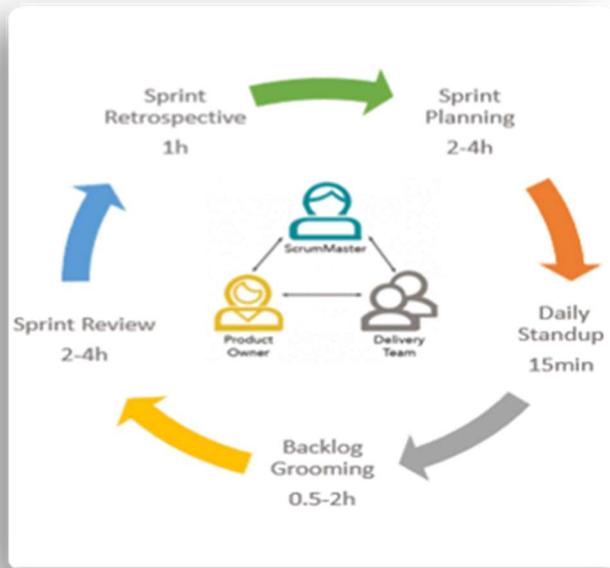


Figure 6 Scrum meetings lifecycle

2.10 Board management

We chose **Trello** software to manage our tasks where trello is a visual tool that empowers your team to manage any type of project, workflow, or task tracking. Add files, checklists, or even automation: Customize it all for how your team works best.

What is a board?

A board (A) represents a place to keep track of information — often for large projects, teams, or workflows. Whether you are launching a new website, tracking sales, or planning your next office party, a Trello board is the place to organize tasks, all the little details, and most importantly—collaborate with your colleagues.

What is a list?

Lists (B) keep cards, or specific tasks or pieces of information, organized in their various stages of progress. Lists can be used to create a workflow where cards are moved across each step in the process from start to finish, or simply act as a place to keep track of ideas and information. There's no limit to the number of lists you can add to a board, and they can be arranged and titled however you'd like.

What is a card?

The smallest, but most detailed unit of a board is a card (C). Cards are used to represent tasks and ideas. A card can be something that needs to get done, like a blog post to be written, or something that needs to be remembered, like company vacation policies. Just click “Add a card...” at the bottom of any list to create a new card, and give it a name like “Hire a new marketing manager” or “Write a blog post.”

What is the board menu?

On the right side of your Trello board is the menu (D)—the mission control center for your board. The menu is where you manage members' board permissions, control settings, search cards, enable Power-Ups, and create automations. You can also see all of the activity that has taken place on a board in the menu's activity feed. Take some time to check out everything the menu has to offer.

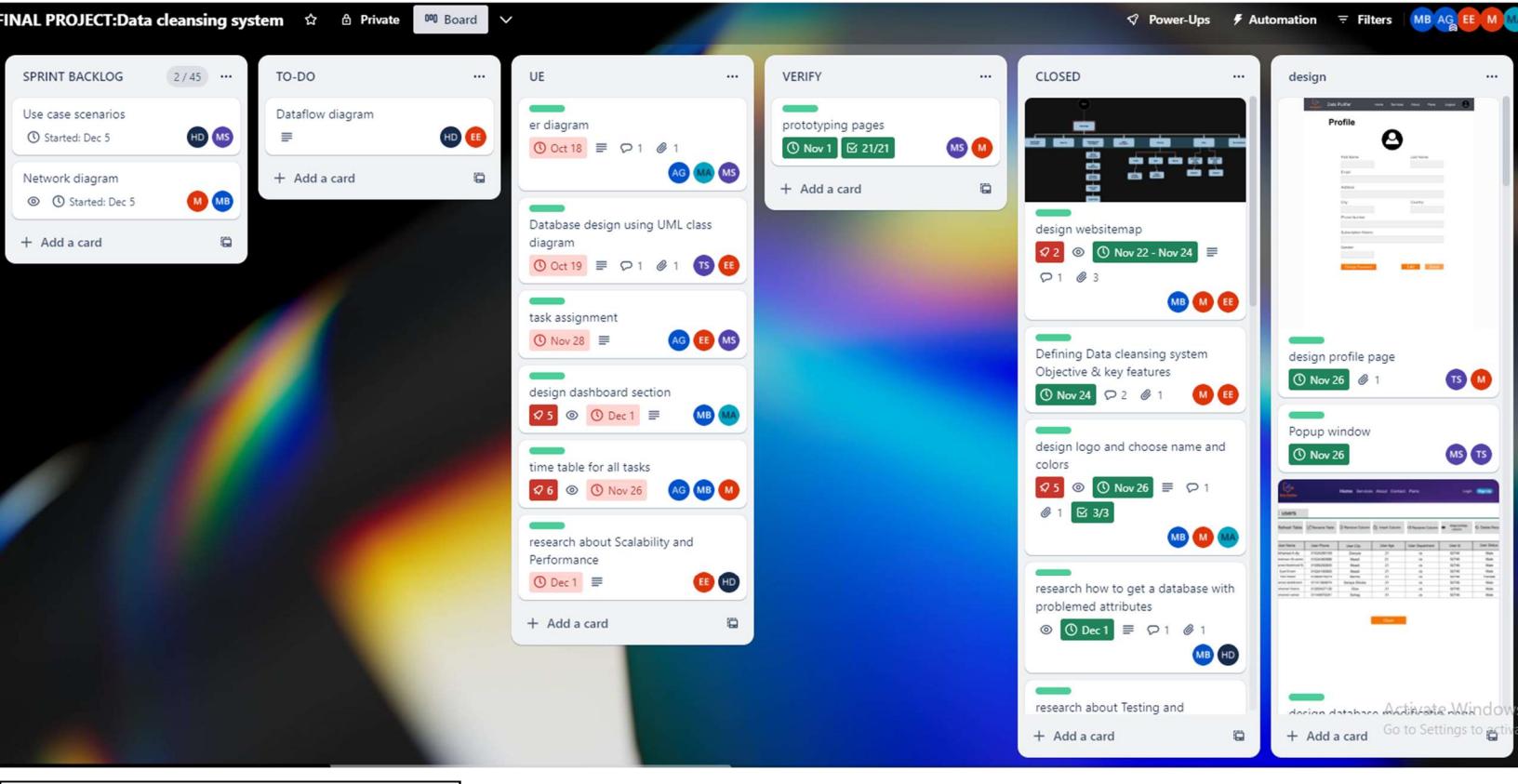


Figure 7 trello board management tool

<https://trello.com/b/5xonGb74/final-projectdata-cleansing-system>

How we setup deadlines in our tasks?



Figure 8 scrum poker estimation diagram

we used a technique called **scrum poker estimation**

What is scrum poker?

Scrum poker, also known as “planning poker” and “pointing poker”, is a gamified technique that development teams use to guess the effort of project management tasks. These estimations are based on the entire group’s input and consensus, making them more engaging and accurate than other methods. To help gauge the number of story points for the relevant tasks, teams use planning poker cards, which are similar to poker cards.

How does planning poker work?

At the beginning of a poker planning session, the product owner or customer reviews an agile user story and reads it aloud. A user story is a general and informal explanation of a software feature that describes how it will offer value to the end-user (i.e. the customer).

Step 1: Hand out the cards to participants

Distribute an identical deck of cards to everybody. Each one has a number that the team has agreed to use as their estimate. Each player should have a deck consisting of different numbers. Cohn recommended a sequence of 0, 1, 2, 3, 5, 8, 13, 20, 40, and 100.

Other common sequences include doubling the next number (e.g. 1, 2, 4, 8, and so on). These values can represent a number of things: the number of story points, ideal days, or other units that the team uses for estimation.

The decks are intentionally kept minimal with considerable number-jumps. Doing this ensures that for each story, everyone can reach a consensus number. Otherwise, if they have a card for each number from one to 50, the process would be painfully slow.

Step 2: Read the story out loud

The moderator (either the product owner or product manager) narrates the story to the group. If participants have any questions, the moderator answers them.

Step 3: Discuss the story

Once the group finishes listening to the story, everyone shares their views on it. Some of these discussion points will likely include:

How should we handle the work?

How many people are expected to get involved?

What skills will be needed to work on the story?

How should we tackle any roadblocks that delay progress?

The group will also try to learn more about the story and ask questions to understand it better.

Step 4: Select and share

After the discussion, each person will privately select a card from the deck. Usually, it's used to show an estimate of story points (but can also be used to represent the number of ideal days). Once everyone selects a card, they show them at the same time.

If a player shows a higher card, it conveys that the story will be completed with greater difficulty and take a longer period to complete. Keep in mind that it's common for estimates to vary a lot.

Step 5: Reach a consensus

When team members show the same card, that number turns into a consensus. Now, the group can move forward and work on the next story.

However, if the cards continue to vary, then further discussions on the story will follow. Participants with higher or lower estimates than others will communicate their points of view. Then, they'll attempt to convince their teammates to understand their differing numbers.

Once this new deliberation ends, everyone will go through their deck and show them again. If a participant continues to agree with their last choice, then they will repeat the card or eventually choose a new one.

Usually, the estimates start to converge after the second round. If not, then the process repeats itself until the team agrees on a single number.

The benefits of scrum poker

According to one study, estimates from planning poker are statistically higher than individual ones. It was also noted that for the same tasks, planning poker estimates were more accurate than individual ones.

Other benefits include:

Estimating tasks relative to each other. Often, it can be tricky to gauge the time required to complete a project, especially when it's your first time. Planning poker familiarizes teams with assessing them. After playing the game for a while, you end up building a series of tasks that act as a future reference for comparison.

Lending an equal voice to everyone on the team. It can encourage newer employees to speak up by playing a card and explaining their logic. For example, imagine making a food reservation app. You and your colleague might give a smaller estimate, like 10 or 15. However, a new employee might give one of 100. Maybe they had experience with creating a similar app at their last job and knows that such an app takes a lot of time, especially compared to other ones.

Identifying gaps in requirement and implementation. When participants disclose their estimates, they will have to back them up with reasoning on why they are high or low. This can open up questions on the requirement and implementation – a feedback loop that can detect the gaps.

Who to include in scrum poker meetings

The right people should join the meeting or it becomes difficult to reap the benefits outlined above. These crucial roles include:

Scrum team members: scrum members deliver the items from the product backlog – a list of deliverables (e.g. new features). They will also provide their input for the discussions of story points.

Scrum master: a scrum leader is the facilitator in agile meetings. They should take part in all standard meetings.

Product owner: the owner or manager will describe all user stories to the team and answer their questions.

When to hold a planning poker session

Usually, teams arrange a session after creating the initial backlog. Although sessions can sometimes take more than one day, it leads to the development of initial estimates that are helpful in sizing or scoping the project.

Items are added to the product backlog incrementally throughout the project's lifespan. That's why it's usually more convenient for teams to conduct sessions once per iteration. In most cases, this happens some days after the iteration ends. Similarly, it also occurs right after a daily standup (a type of agile meeting) because the entire team is present.

2.11 workflow

We chose to work with **pair programming technique** where a task can never be hold by only one member

-pair programming is an integrated practice From **XP (extreme programming)** framework

-this can help the team to improve their teamwork skills and team management skills

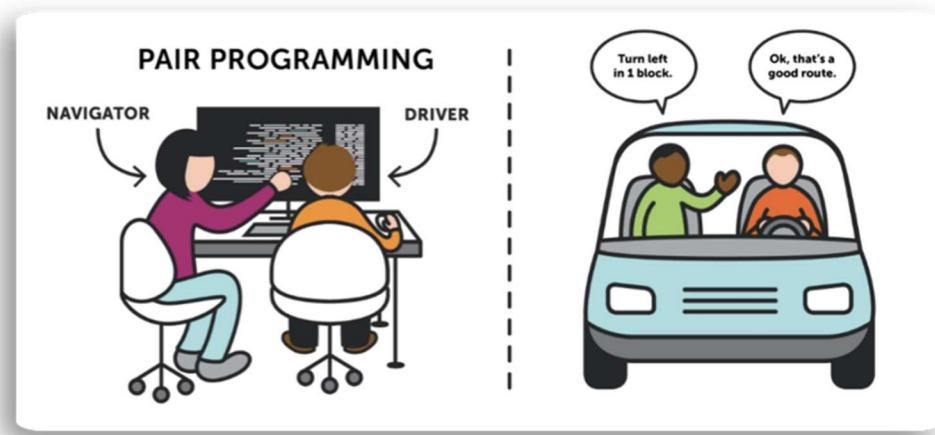


Figure 9 pair programming workflow

Chapter 3

Survey & Literature Review

Chapter 3 – Survey & Literature Review

3.1 Introduction

We will review some of the data cleansing systems listed in the top review before and build a comparison to show developments about the use of data cleansing systems reviewing previous works is about technology used at present technical possibilities. By addressing some of the techniques that modern data cleansing systems can understand the language user through it

We will explain:-

An introduction to previous works that revolve around our topic make a comparison between some works to clarify the best aspect of either of them in terms of technology, purpose, advantages, and disadvantages clarification of some topics on which our topic depends

3.2 Related Work

We found similar apps and similar research papers giving an idea about our desired system. In this section we made a comparison between a number of applications, each with its year, objective, technology, results, advantages and disadvantages, let's show the most interactive applications using a specific technology.

#	Name	Link	Advantages	Disadvantages
1	Winpure 	https://winpure.com/	<ul style="list-style-type: none">Industry-leading Data Matching Tool, quick, intelligent and easy-to-use	<ul style="list-style-type: none">Complexity: Some users have found the software to be complex, Requiring training to fully utilize all its features.Exporting Results: Exporting the Results of matching can be cumbersome and not user-friendly.
2	OpenRefine 	https://openrefine.org/	<ul style="list-style-type: none">It has an easy to use interfaceOpenRefine supports a wide variety of data	<ul style="list-style-type: none">Some users complain that the user interface looks a bit dated

3	<p>Knowledge</p>  <p>Net</p>	<p>https://kn-it.com/</p>	<ul style="list-style-type: none"> • Connected view of relationships beyond 1st & 2nd degree connections. • ease of use and the insights that KN gives on people 	<ul style="list-style-type: none"> • Not having a tool prior to Knowledge Net and not fully leveraging the power of the tool on clients
4	<p>Demand</p>  <p>tools</p>	<p>https://www.validity.com/demandtools/</p>	<ul style="list-style-type: none"> • Straightforward UI that does what it needs to do. • Makes after-the-fact de-duplication a breeze. • Easy mass updates to data sets - making cleaning data a breeze. 	<ul style="list-style-type: none"> • UI is outdated and could use a bit of a facelift. • Some competitors in this space are more all-in-one tools where DemandTools. • Lack of automation for regular data cleansing - everything we use DemandTools for is manual.
5	<p>Sage master</p>  <p>SageMaster</p>	<p>https://sagemaster.io/</p>	<ul style="list-style-type: none"> • Streamlined Financial Management • Enhanced accuracy • Time and cost savings • Scalability • Reporting and analytics 	<ul style="list-style-type: none"> • Initial setup complexity • High Learning curve • Cost of ownership • Limited mobility • Compatibility issues

3.3 Related Papers

A survey was done on 12 research papers which includes the required technologies needed to create the desired project. Those papers will be briefly discussed and the tables following will include the final results that are derived from the paper.

#	<u>title</u>	<u>objective</u>	<u>result</u>	<u>Advantage</u>	<u>Disadvantage</u>	<u>limitation</u>
1	"Exploring Data Profiling Techniques for Enhanced Data Cleansing Systems"	To investigate various data profiling methods for improving the effectiveness of data cleansing systems.	Identification of key data profiling techniques and their applicability in optimizing data cleansing processes.	Enables data cleansing systems to gain deeper insights into data characteristics and anomalies.	May require substantial computational resources and expertise to implement sophisticated profiling techniques.	Profiling techniques may not be universally applicable and may need customization based on data types and sources.
2	"Automation Strategies for Streamlining Data Cleansing Processes"	To develop automation strategies that reduce manual intervention in data cleansing tasks.	Implementation of automated workflows and algorithms for efficient and scalable data cleansing.	Increases operational efficiency and reduces the time required for data cleansing.	Over-reliance on automation may overlook nuanced data issues that require human intervention.	Automation may be challenging to implement in complex data environments with diverse data formats.
3	"Implementing Effective Data Governance Practices in Data Cleansing Systems"	To integrate robust data governance principles into data cleansing workflows.	Establishment of policies, procedures, and controls to ensure data quality and compliance during cleansing processes.	Enhances data integrity, security, and regulatory compliance.	Implementation of stringent governance practices may introduce additional overhead and complexity.	Data governance frameworks may vary in effectiveness across different organizational contexts.
4	"Anomaly Detection Methods for Proactive Data Cleansing and Quality Assurance"	To develop anomaly detection techniques for identifying irregularities in datasets and initiating proactive data cleansing.	Deployment of anomaly detection algorithms to detect and address data inconsistencies and outliers.	Enables early detection and mitigation of data anomalies, improving overall data quality.	False positives may lead to unnecessary data cleansing efforts and potential loss of valid information.	Anomaly detection algorithms may struggle with detecting subtle anomalies or patterns in noisy data.
5	"Developing Robust Data Quality Metrics for Assessing Data Cleansing Effectiveness"	To design comprehensive data quality metrics for evaluating the effectiveness of data cleansing efforts.	Creation of quantitative measures to assess data accuracy, completeness, consistency, and timeliness post-cleansing.	Facilitates objective evaluation of data cleansing outcomes and performance improvement over time.	Defining universal data quality metrics may be challenging due to the subjective nature of data quality.	Metrics may not capture qualitative aspects of data quality, such as relevance and interpretability.
6	"Strategies for Identifying and Handling"	To develop strategies for identifying and	Implementation of techniques for attribute profiling,	Improves the accuracy and reliability of data	Manual identification and handling of	Strategies may vary in effectiveness

	Problematic Attributes in Databases for Data Cleansing"	addressing problematic attributes that hinder data quality.	anomaly detection, and targeted cleansing.	by addressing specific issues at the attribute level.	problematic attributes may be time-consuming and error-prone.	depending on the complexity and diversity of data attributes.
7	"A Comparative Analysis of Data Cleaning Techniques for Improved Data Quality"	To conduct a comparative analysis of various data cleaning techniques to determine their effectiveness in improving data quality.	Evaluation of different data cleaning methods, such as deduplication, standardization, and error correction, in terms of their impact on data quality.	Provides insights into the strengths and weaknesses of different data cleaning approaches, aiding in informed decision-making.	Comparative analysis may overlook contextual factors that influence the suitability of cleaning techniques for specific datasets.	Results may not generalize to all datasets and may require validation in diverse data environments.
8	"Scalability and Performance Optimization in Large-Scale Data Cleansing Systems"	To address scalability and performance challenges in data cleansing systems when dealing with large-scale datasets.	Development of optimization strategies, parallel processing techniques, and distributed computing architectures to enhance scalability and performance.	Enables data cleansing systems to handle massive volumes of data efficiently, reducing processing time and resource consumption.	Implementation of scalable solutions may require significant upfront investment in infrastructure and technology.	Scalability enhancements may introduce complexity and overhead, impacting system maintainability and agility.
9	"Efficient Data Transformation Approaches for Seamless Integration in Data Cleansing Pipelines"	To explore efficient data transformation techniques that facilitate seamless integration into data cleansing pipelines.	Identification and implementation of data transformation methods, such as normalization, aggregation, and enrichment, to prepare data for cleansing.	Streamlines data preprocessing tasks and enhances compatibility with downstream cleansing processes.	Complex transformation requirements may necessitate custom development and configuration, increasing implementation effort.	Data transformation approaches may not adequately address all data compatibility and integration challenges, requiring ongoing refinement.
10	"Testing and Validation Frameworks for Ensuring the Accuracy and Reliability of Data Cleansing Processes"	To establish testing and validation frameworks for assessing the accuracy and reliability of data cleansing processes.	Development of systematic approaches, test suites, and validation metrics to verify the effectiveness of cleansing techniques.	Enhances confidence in data cleansing outcomes and validates adherence to predefined quality standards.	Testing and validation processes may introduce additional overhead and delay project timelines.	Validation frameworks may not capture all aspects of data quality and may require continuous refinement to align with evolving requirements.
11	"Training and Documentation Strategies for Empowering Users"	To design training and documentation strategies that	Creation of user-friendly training materials, tutorials, and documentation	Promotes user engagement, knowledge transfer, and self-	Developing comprehensive training and documentation	User training and documentation may not fully address

	"in Data Cleansing Systems"	empower users to effectively utilize data cleansing systems.	resources to facilitate user adoption and proficiency.	sufficiency in data cleansing operations.	resources may require significant time and effort.	individual learning preferences and skill levels, necessitating ongoing support and education.
12	"Data Integration Techniques for Seamlessly Incorporating Cleansed Data into Organizational Workflows"	To explore data integration techniques that enable seamless incorporation of cleansed data into organizational workflows.	Implementation of data integration mechanisms, such as APIs, ETL processes, and data warehouses, to facilitate data flow and interoperability.	Facilitates the integration of cleansed data with existing systems, applications, and business processes, enhancing data accessibility and utilization.	Data integration efforts may encounter compatibility issues, data format mismatches, and synchronization challenges.	Integration techniques may require ongoing maintenance and updates to adapt to evolving data requirements and technological landscapes.

1. Exploring Data Profiling Techniques for Enhanced Data Cleaning Systems

- This research investigates various data profiling methods to improve data cleaning systems by gaining deeper insights into data characteristics and anomalies.

2. Automation Strategies for Streamlining Data Cleaning Processes

- This study focuses on developing automation strategies to reduce manual intervention in data cleaning tasks, thereby enhancing operational efficiency and scalability.

3. Implementing Effective Data Governance Practices in Data Cleaning Systems

- The objective is to integrate robust data governance principles into data cleaning workflows, ensuring data quality, security, and regulatory compliance.

4. Anomaly Detection Methods for Proactive Data Cleaning and Quality Assurance

- This research aims to develop anomaly detection techniques to identify irregularities in datasets and initiate proactive data cleansing, improving overall data quality.

5. Developing Robust Data Quality Metrics for Assessing Data Cleansing Effectiveness

- The goal is to design comprehensive data quality metrics to evaluate the effectiveness of data cleansing efforts objectively and improve data quality over time.

6. Strategies for Identifying and Handling Problematic Attributes in Databases for Data Cleansing

- This study focuses on developing strategies to identify and address problematic attributes hindering data quality, thereby improving overall data reliability.

7. A Comparative Analysis of Data Cleaning Techniques for Improved Data Quality

- This research conducts a comparative analysis of different data cleaning techniques to determine their effectiveness in improving data quality, aiding in informed decision-making.

8. Scalability and Performance Optimization in Large-Scale Data Cleansing Systems

- The objective is to address scalability and performance challenges in data cleansing systems when dealing with large-scale datasets, reducing processing time and resource consumption.

9. Efficient Data Transformation Approaches for Seamless Integration in Data Cleansing Pipelines

- This study explores efficient data transformation techniques to prepare data for cleansing and ensure compatibility with downstream processes, streamlining data preprocessing tasks.

10. Testing and Validation Frameworks for Ensuring the Accuracy and Reliability of Data Cleansing Processes

- This research establishes testing and validation frameworks to verify the accuracy and reliability of data cleansing processes, enhancing confidence in data quality outcomes.

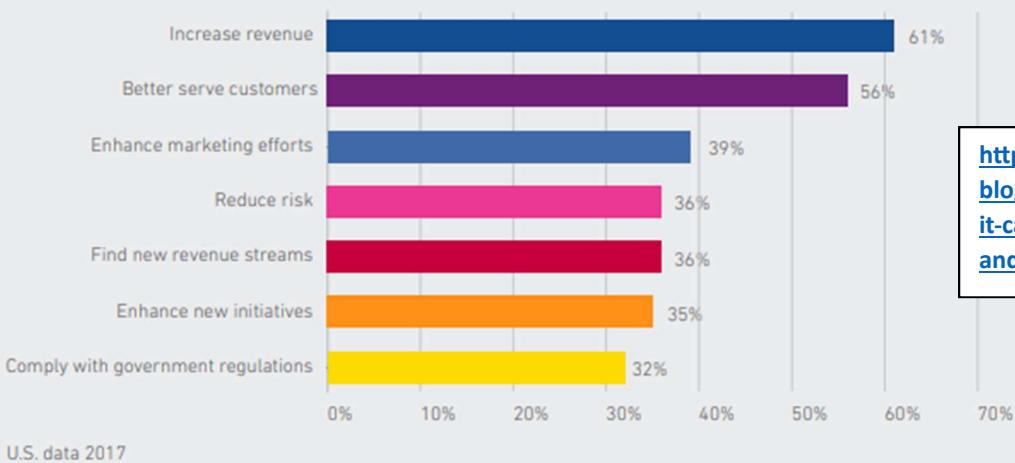
11. Training and Documentation Strategies for Empowering Users in Data Cleansing Systems

- The goal is to design training and documentation strategies to empower users in effectively utilizing data cleansing systems, promoting user engagement and proficiency.

12. Data Integration Techniques for Seamlessly Incorporating Cleansed Data into Organizational Workflows

- This study explores data integration techniques to seamlessly incorporate cleansed data into organizational workflows, enhancing data accessibility and utilization.

Chart 1
How data powers business opportunities



[https://www.cazoomi.com/
blog/dirty-data-how-much-
it-can-cost-your-business-
and-how-to-get-rid-of-it/](https://www.cazoomi.com/blog/dirty-data-how-much-it-can-cost-your-business-and-how-to-get-rid-of-it/)

What Do You Use Data for?

Figure 10 how data powers business opportunities

Typically, data is analyzed to help businesses predict future or emerging consumer behavior. This process involves understanding the data collected and knowing how to organize it, analyze it, and apply it to improve marketing efforts.

Most marketers use the same data to improve and personalize their customer experience. Seventy-eight percent of organizations say that the data they collect helps them increase customer acquisitions and lead conversions.

Data-driven marketing has, therefore, become a crucial strategy for marketers who want to succeed in the current hypercompetitive economy. Ninety-five percent of organizations are using data to power their business opportunities.

However, this can only occur if your data is clean. Dirty data can negatively affect your bottom line. Eighty-nine percent of business executives say that inaccurate data affects their ability to provide great customer experiences.

When your company is swimming in dirty data, the sales and marketing teams cannot make well-informed and accurate, data-driven decisions. On average, organizations believe that 25 percent of their data is inaccurate, a factor that impacts on the bottom line.

Dirty data can be costly in the long run. It can eventually lead to lower productivity, unnecessary spending, and unreliable decision-making.

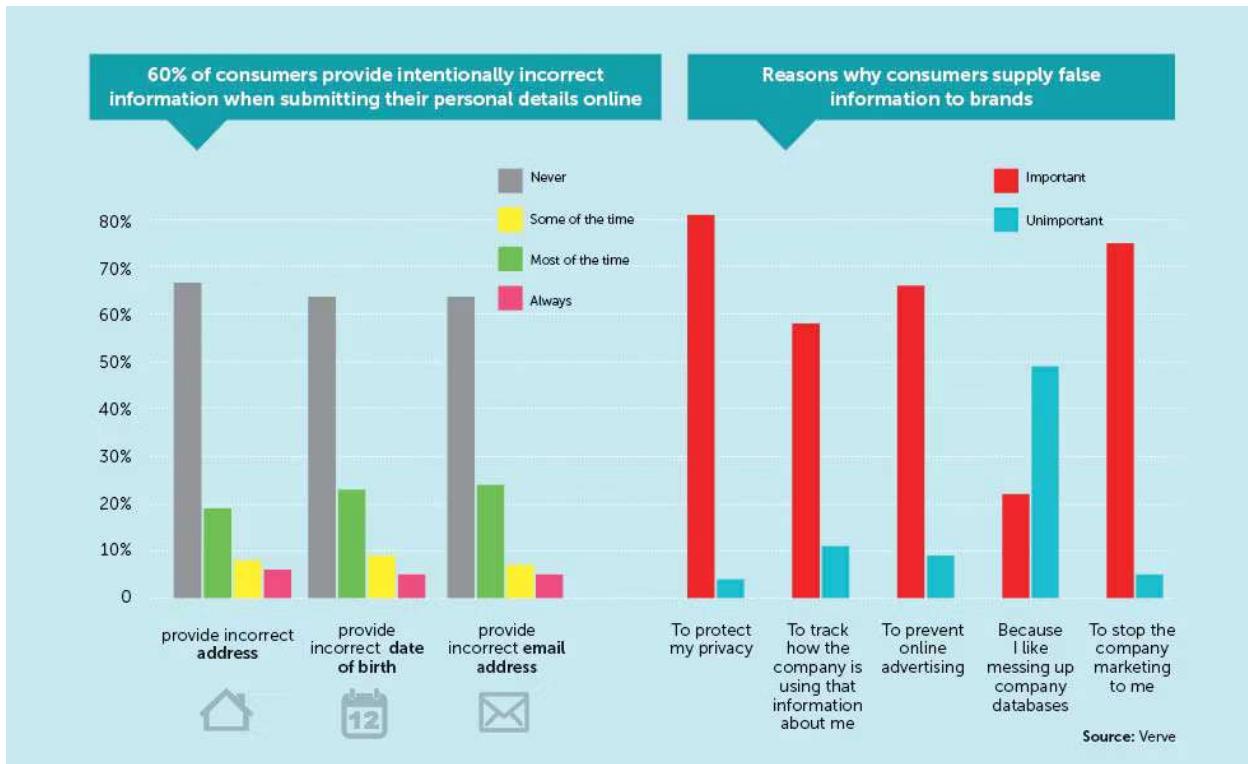


Figure 11 wh users give false data to brands databases

Why do we give false data to brands databases?

People are deliberately giving brands false data about themselves to protect their privacy, and are ignoring brands' efforts to empower them to take control of their data, according to a study of more than 2,400 UK consumers by research company Verve.

The research shows that 60% of consumers intentionally provide incorrect information when submitting their personal details online. Broken down by the types of data provided, birth dates are the most commonly falsified, as almost a quarter of consumers (23%) give the wrong date of birth to companies 'some of the time', 9% do this 'most of the time' and 5% 'always' give the wrong date.

The research also shows that nearly one-third of people give a fake email address and a made-up name at least some of the time. It is a similar story for incorrect information given about home addresses, phone numbers, job titles and company names.

"The upside of providing information has not been articulated," says managing director at Verve Colin Strong. "The case is not always made by companies about what consumers

are going to get in return for providing information, but people see the immediate effects of being put on more marketing lists and being pursued by online advertising and email spam."

The study therefore looks at why people give false data. Privacy concerns top the list with 81% of people saying this is important, while 77% believe that the information being requested is intrusive and 76% say it is unnecessary.

Also high on the list of reasons for lying on forms is to stop companies marketing to them – 75% of consumers say this – and two thirds say their motivation for giving false information is to prevent online advertising. Malice towards companies and embarrassment about having internet browsing behavior tracked are less common reasons.

Vanda Candeias

Digital & E-Commerce Transformation Director

<https://www.linkedin.com/pulse/why-do-we-give-false-data-brands-bases-vanda-candeias/>

Chapter 4

Analysis Phase

Chapter-4-Analysis phase

4.1 Introduction

The output gained from chapter three is an overview of the system's functionality and operations. Based on this output, the system analysis phase will take place, which involves going through all expected scenarios in order to create our system in terms of what the system will do, how each part of the system will be linked with other parts, and how it will be implemented from scratch.

4.2 System Scenario

1- User Login: Mohammed, a Data Analyst, starts his day by accessing the data cleansing system. He navigates to the login page and enters his username and password.

2- Multi-Factor Authentication: The system prompts Mohammed for multi-factor authentication. He receives a code on his registered mobile device, enters it, and gains access to the system.

3- Dashboard Overview: Once logged in, Mohammed is presented with a dashboard displaying recent activities, status of ongoing data cleansing processes, and notifications.

4- Upload Raw Data: Mohammed initiates a new data cleansing task by selecting the "Upload Data" option. He uploads raw data files from his local machine, including CSV and Excel files.

5- Select Cleansing Techniques: Mohammed reviews the uploaded data and selects specific cleansing rules based on the data quality requirements. She chooses rules for handling missing values, removing duplicates, and standardizing formats.

6- Initiate Cleansing Process: With the selected rules in place, Mohammed initiates the data cleansing process. The system begins processing the data according to the defined rules.

7- Validation and Sanitization: The system validates and sanitizes the data, identifying and removing any anomalies, errors, or potential threats. It logs the cleansing activities for auditing purposes.

8- Data Review: Once the cleansing process is complete, Mohammed reviews the cleansed data using interactive visualization tools and data preview features provided by the system.

9- Approval of Cleansed Data: Satisfied with the quality of the cleansed data, Mohammed approves it for further use. He acknowledges the system-generated summary of changes made during the cleansing process.

10- Export Cleansed Data: Mohammed exports the cleansed data in CSV format to share with other team members for analysis. The system generates a download link, and Mohammed saves the file to his local machine.

11- Generate Data Cleansing Report: Before logging out, Mohammed generates a comprehensive report summarizing the data cleansing activities, including the number of records processed, errors identified, and the effectiveness of applied rules.

12- Logout: Mohammed logs out of the data cleansing system, ensuring the termination of her session to maintain security.

13-Missing Values Detection: The system identifies columns with missing values and provides options for handling them, such as imputation or removal.

14-Imputation Strategy Selection: Mohammed selects the imputation strategy for handling missing values, opting for mean imputation for numerical columns and mode imputation for categorical columns.

15-Invalid Format Detection: The system detects columns with invalid formats, such as dates or phone numbers, and suggests standardization methods.

16-Standardization Method Selection: Mohammed chooses to standardize date formats to YYYY-MM-DD and phone numbers to a specific pattern.

17-Duplicate Records Identification: The system identifies duplicate records based on specified criteria, such as unique identifiers or combinations of attributes.

18-Duplicate Removal Strategy Selection: Mohammed selects a strategy for removing duplicate records, opting for keeping the first occurrence and discarding subsequent duplicates.

19-Outlier Detection: The system detects outliers in numerical columns using statistical methods like Z-score or IQR.

20-Outlier Removal Strategy Selection: Mohammed selects a strategy for handling outliers, choosing to remove them based on a predefined threshold.

21-Progress Monitoring: Mohammed monitors the progress of the cleansing process through the dashboard, observing the completion status and any encountered errors.

22-Error Handling: If errors occur during cleansing, Mohammed receives notifications and accesses detailed error logs to troubleshoot and resolve issues.

23-Data Quality Assessment: After cleansing, Mohammed assesses the data quality metrics provided by the system, including completeness, accuracy, and consistency.

24-Preview Cleansed Data: Mohammed previews the cleansed SQL data to ensure it meets the desired quality standards before proceeding.

25-User Access Management: Mohammed manages user access to the cleansing system, granting permissions based on roles and responsibilities.

26- Audit Trail Review: Mohammed reviews the audit trail of cleansing activities to track changes made to the data and maintain data governance compliance.

27- Documentation Generation: Mohammed generates documentation summarizing the cleansing process, techniques applied, and results achieved for record-keeping and compliance purposes.

28- Performance Optimization: Mohammed optimizes the cleansing process by fine-tuning parameters, such as batch size or parallel processing, to enhance efficiency and reduce processing time.

29 - Feedback Submission: Mohammed provides feedback on the cleansing system's performance and suggests improvements for future enhancements.

30 - Training and Support: Mohammed accesses training materials and support resources provided by the system vendor to enhance his skills and troubleshoot issues effectively.

31- System Updates: Mohammed receives notifications about system updates and new features, ensuring he stays informed about the latest improvements and enhancements.

4.3 DFD diagram

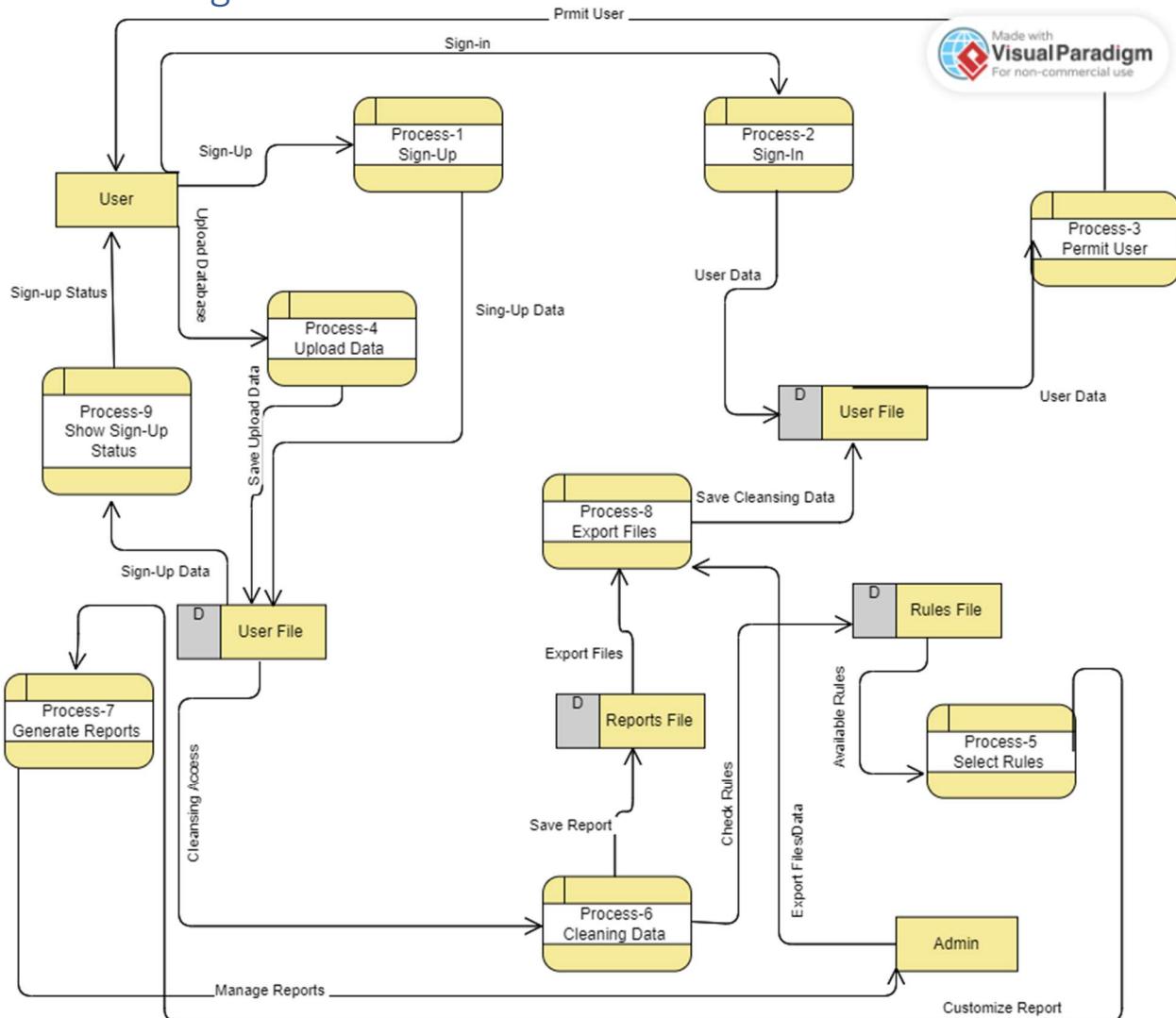


Figure 12 DFD diagram

DFD Level (0) Process:-

1.0: Sign-Up ->

The user initiates the sign-up process by providing their user data.

2.0: Sign-In ->

The user signs in to the system using their credentials.

3.0: Permit User ->

The admin permits the user to access the system.

4.0: Upload Data ->

The user uploads their data to the system.

5.0: Select Rules -> The user selects the rules for data cleansing.

6.0: Cleaning Data -> The system cleanses the uploaded data based on the selected rules.

7.0: Generate Reports -> The system generates reports based on the cleansed data.

8.0: Export Files -> The user exports the generated reports in a file format.

9.0: Show Sign-Up Status -> The system displays the sign-up status to the user.

DFD Level (0) Data stores:-

1. User File
2. Rules File
3. Reports File

DFD Level (0) Data entities:-

1. User
2. Admin

4.4 Context diagram:-

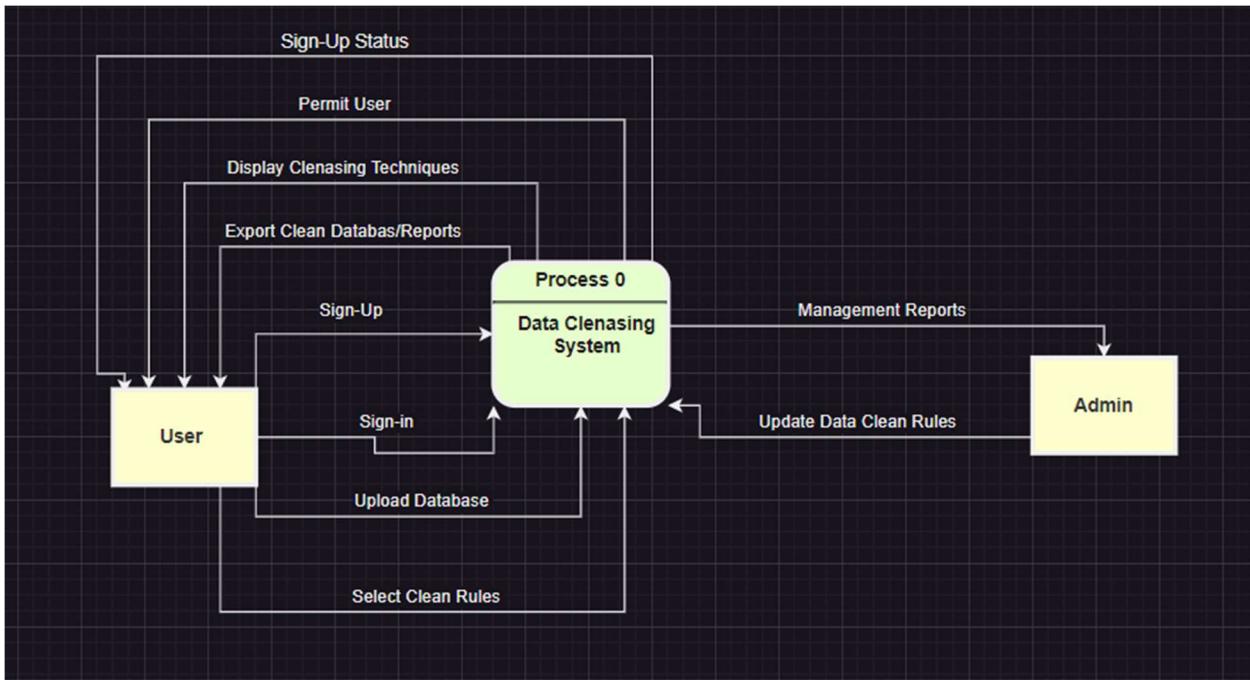


Figure 13 context diagram

Context diagram entities:-

1. User File
2. Admin

4.5 Required Database Type

By having large amounts of data, it's required to have a suitable database to be able to store this amount of data. In our case, using NoSQL database type enables the storage of the system to be simpler having data in a JSON-like tree structure and storing data in documents rather than tables. It handles unstructured data better than relational database types. One of these databases is called Firebase, which is a cloud-hosted NoSQL database, meaning that data is always available for retrieval.

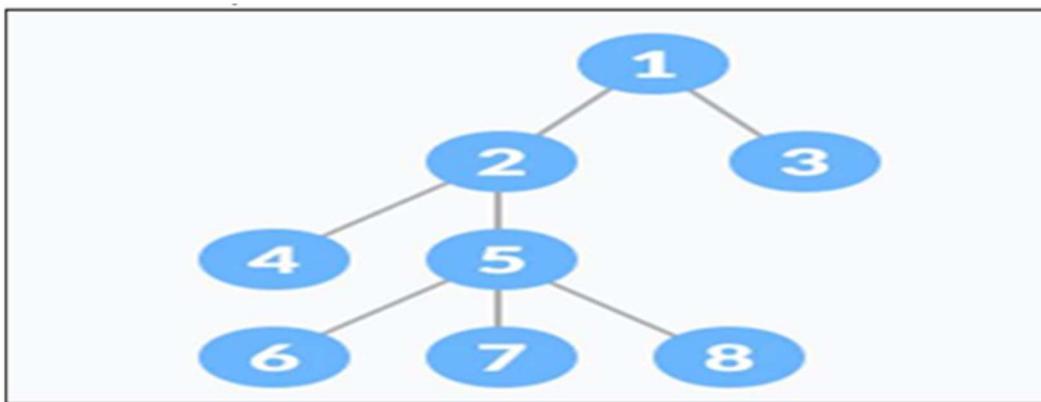


Figure 14 NOSQL Structure

4.6 Site map

A hierarchical diagram that represents the different pages of the system by visualizing a map that provides the flow of the system, the relationship between each page, and information about each.

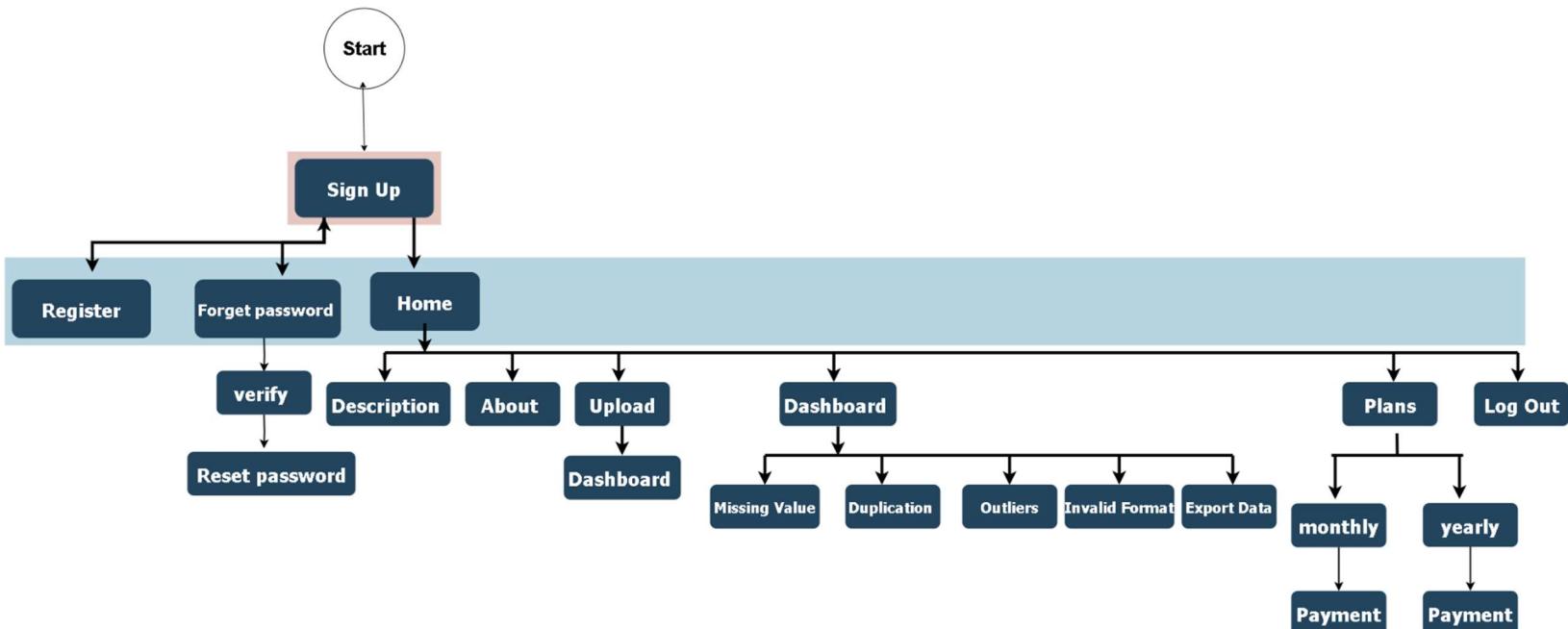


Figure 15 Sitemap

Chapter FIVE

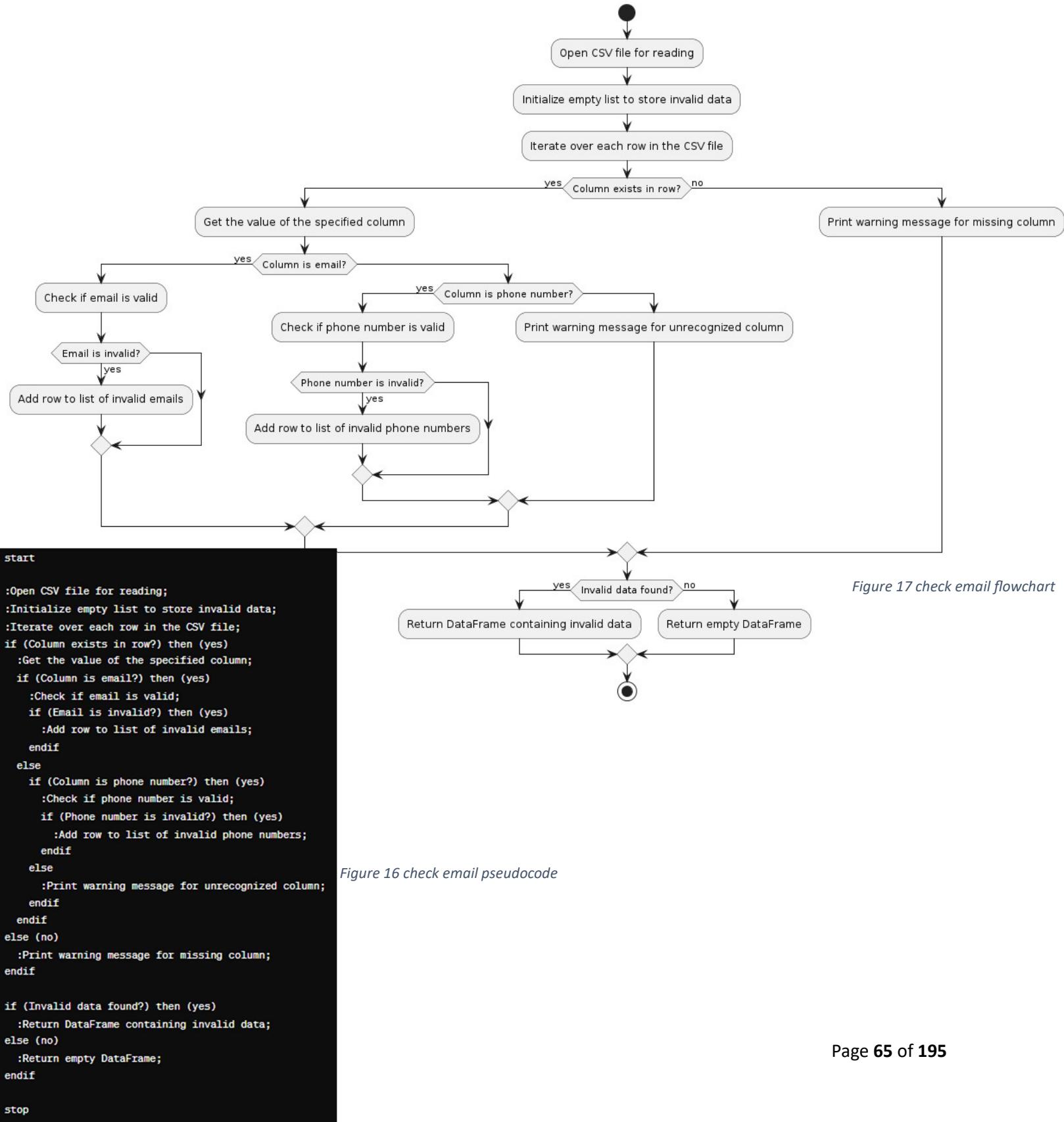
Design Phase

Chapter 5 – design phase

5.1 Introduction

The design phase aims to construct the system by planning the development of the system, which is achieved by having all the basic concepts of the system understood, gathering all the data needed, and clarifying how each part will be connected with other parts in order to achieve the full functionality of the system. This chapter aims to provide flowcharts and diagrams that describe the functionality of the different parts of the system.

5.2 Flowcharts and pseudocodes



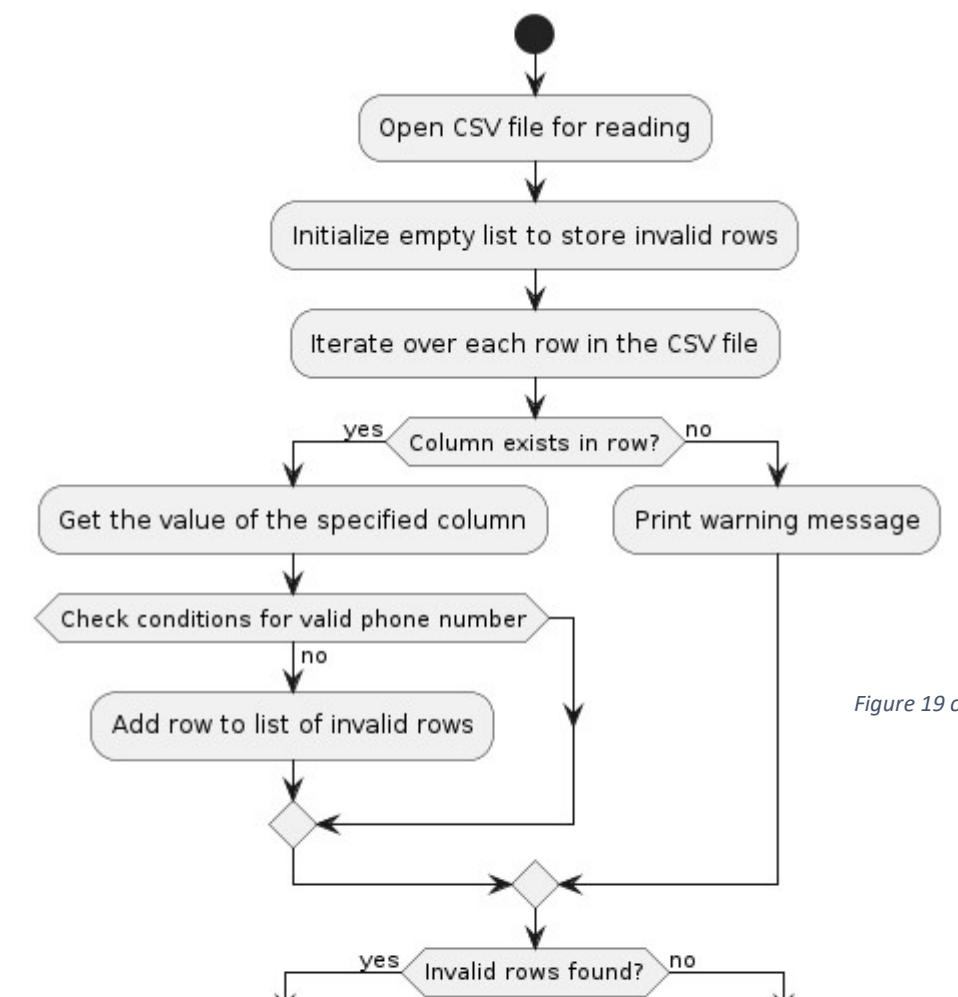


Figure 19 check phone number flowchart

```

start

:Open CSV file for reading;
:Initialize empty list to store invalid rows;
:Iterate over each row in the CSV file;
if (Column exists in row?) then (yes)
  :Get the value of the specified column;
  if (Check conditions for valid phone number) then (no)
    :Add row to list of invalid rows;
  endif
else (no)
  :Print warning message;
endif

if (Invalid rows found?) then (yes)
  :Return DataFrame containing invalid rows;
else (no)
  :Return empty DataFrame;
endif

stop

```

Figure 18 check phone number pseudocode

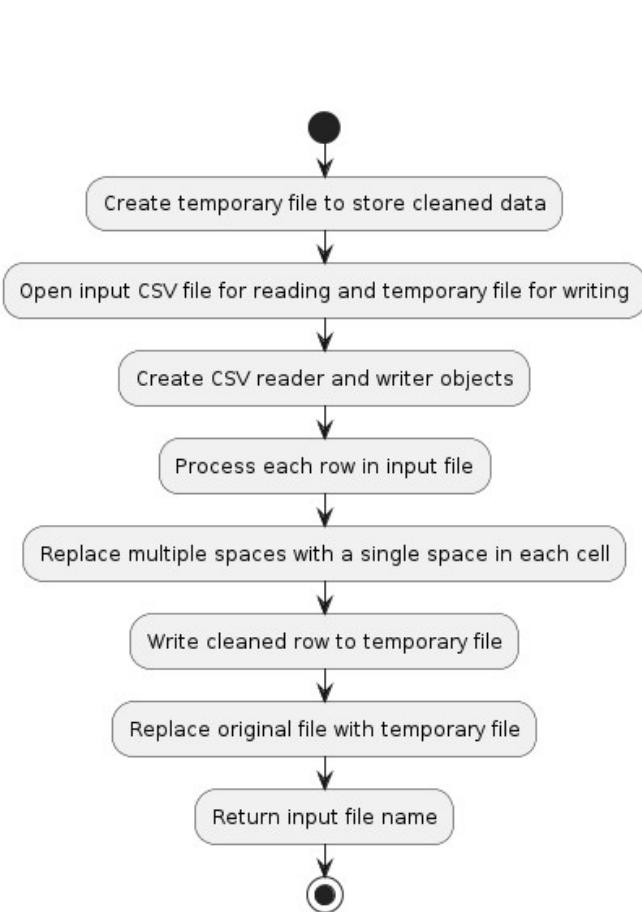


Figure 25 remove double spaces flowchart

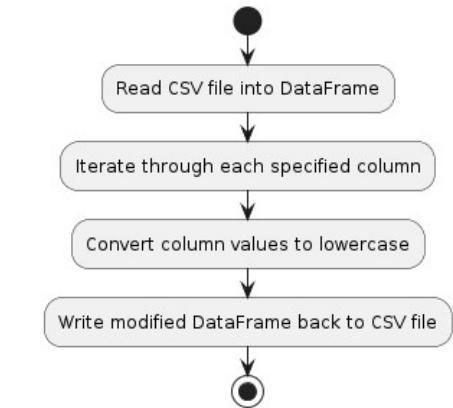


Figure 23 lowercase flowchart

```

start
:Read CSV file into DataFrame;
:Iterate through each specified column;
:Convert column values to lowercase;
:Write modified DataFrame back to CSV file;

stop

```

Figure 22 lowercase pseudocode

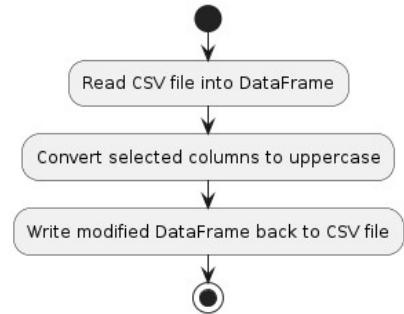


Figure 21 uppercase flowchart

```

start
:Read CSV file into DataFrame;
:Convert selected columns to uppercase;
:Write modified DataFrame back to CSV file;

stop

```

Figure 20 uppercase pseudocode

```

start
:Create temporary file to store cleaned data;
:Open input file for reading and temporary file for writing;
:Create CSV reader and writer objects;
:Process each row in input file;
:Remove leading and trailing spaces from each cell;
:Write cleaned row to temporary file;
:Replace original file with temporary file;
:Return input file name;

stop

```

Figure 24 remove double spaces pseudocode

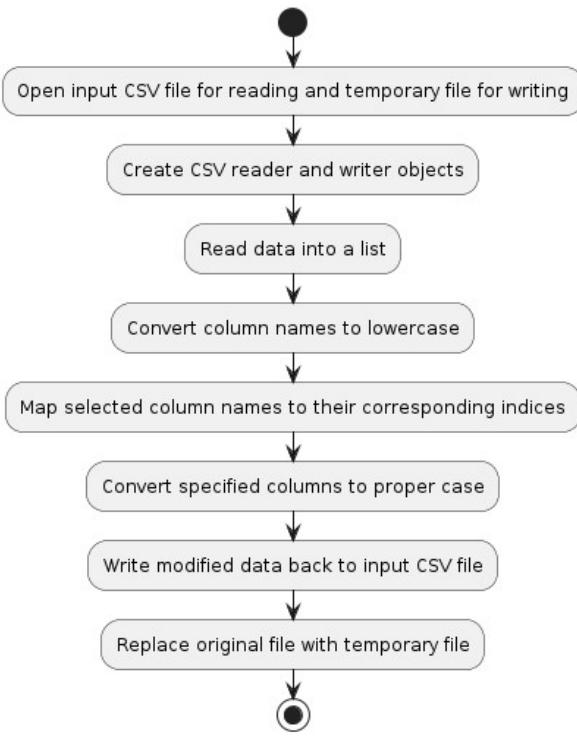


Figure 27 propercase flowchart

```

start

:Open input CSV file for reading and temporary file for writing;
:Create CSV reader and writer objects;
:Read data into a list;
:Convert column names to lowercase;
:Map selected column names to their corresponding indices;
:Convert specified columns to proper case;
:Write modified data back to input CSV file;
:Replace original file with temporary file;

stop

```

Figure 26 propercase pseudocode

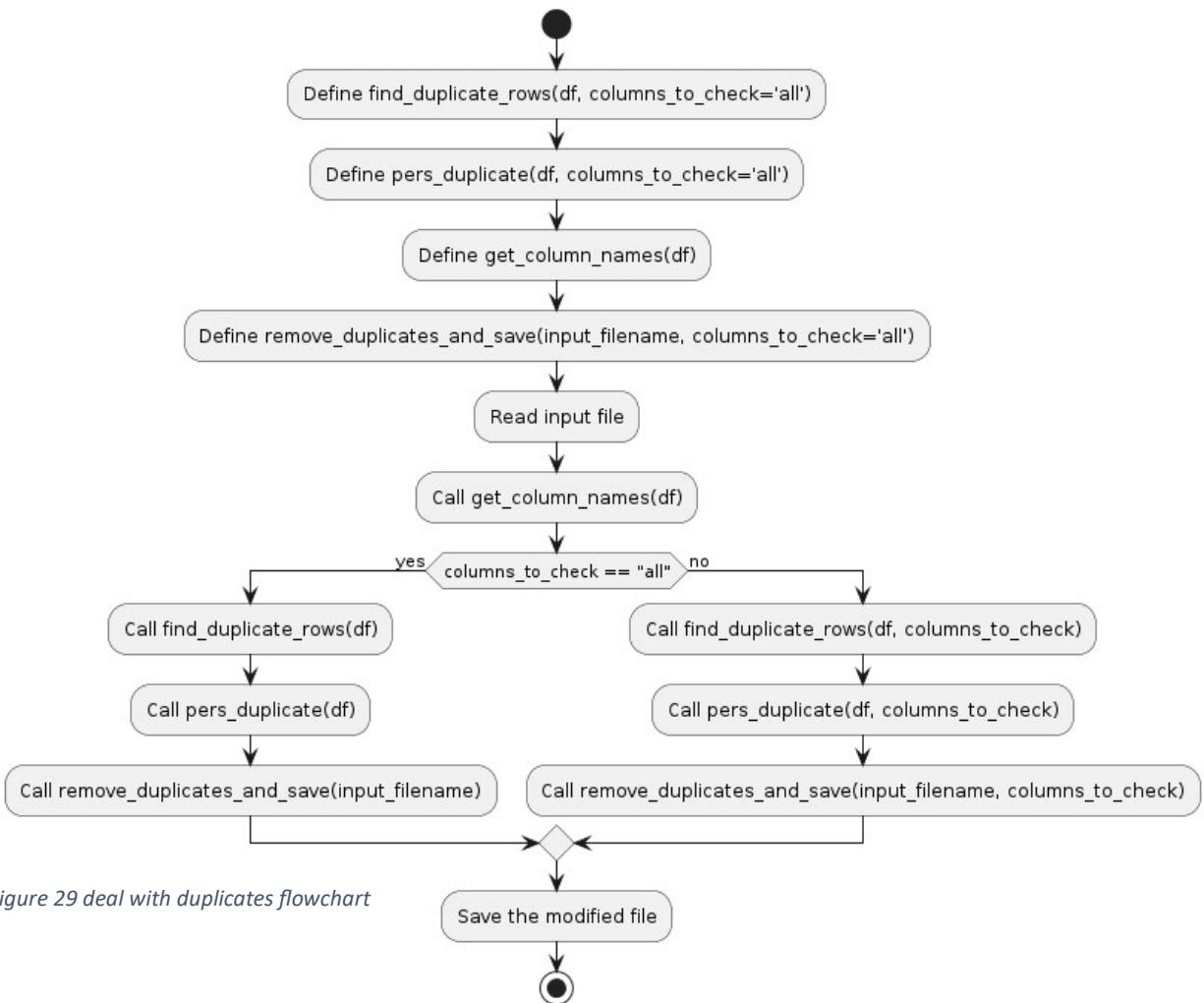


Figure 29 deal with duplicates flowchart

```

start

:Define find_duplicate_rows(df, columns_to_check='all');
:Define pers_duplicate(df, columns_to_check='all');
:Define get_column_names(df);
:Define remove_duplicates_and_save(input_filename, columns_to_check='all');

:Read input file;

:Call get_column_names(df);
if (columns_to_check == "all") then (yes)
  :Call find_duplicate_rows(df);
  :Call pers_duplicate(df);
  :Call remove_duplicates_and_save(input_filename);
else (no)
  :Call find_duplicate_rows(df, columns_to_check);
  :Call pers_duplicate(df, columns_to_check);
  :Call remove_duplicates_and_save(input_filename, columns_to_check);
endif

:Save the modified file;

stop

```

Figure 28 deal with duplicates pseudocode

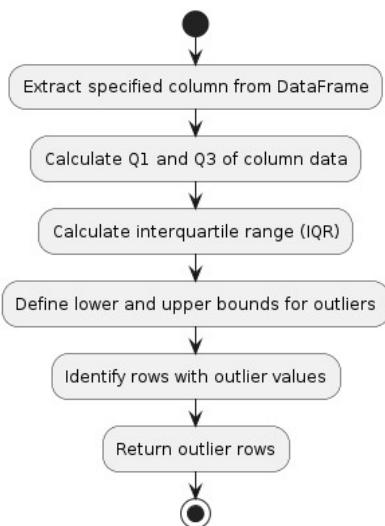


Figure 32 find outliers flowchart

```

start
:Extract specified column from DataFrame;
:Calculate Q1 and Q3 of column data;
:Calculate interquartile range (IQR);
:Define lower and upper bounds for outliers;
:Identify rows with outlier values;
:Return outlier rows;

stop

```

Figure 33 find outliers pseudocode

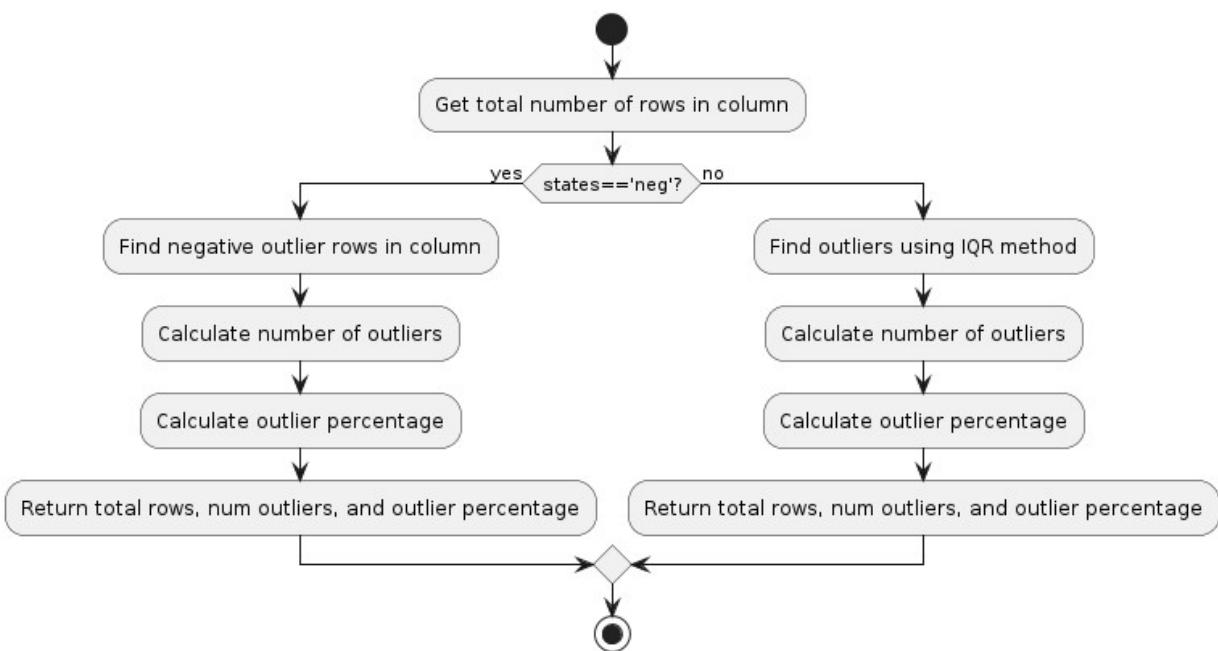


Figure 31 calculate outliers flowchart

```

start
:Get total number of rows in column;
if (states=='neg') then (yes)
    :Find negative outlier rows in column;
    :Calculate number of outliers;
    :Calculate outlier percentage;
    :Return total rows, num outliers, and outlier percentage;
else (no)
    :Find outliers using IQR method;
    :Calculate number of outliers;
    :Calculate outlier percentage;
    :Return total rows, num outliers, and outlier percentage;
endif

stop

```

Figure 30 calculate outlier pseudocode

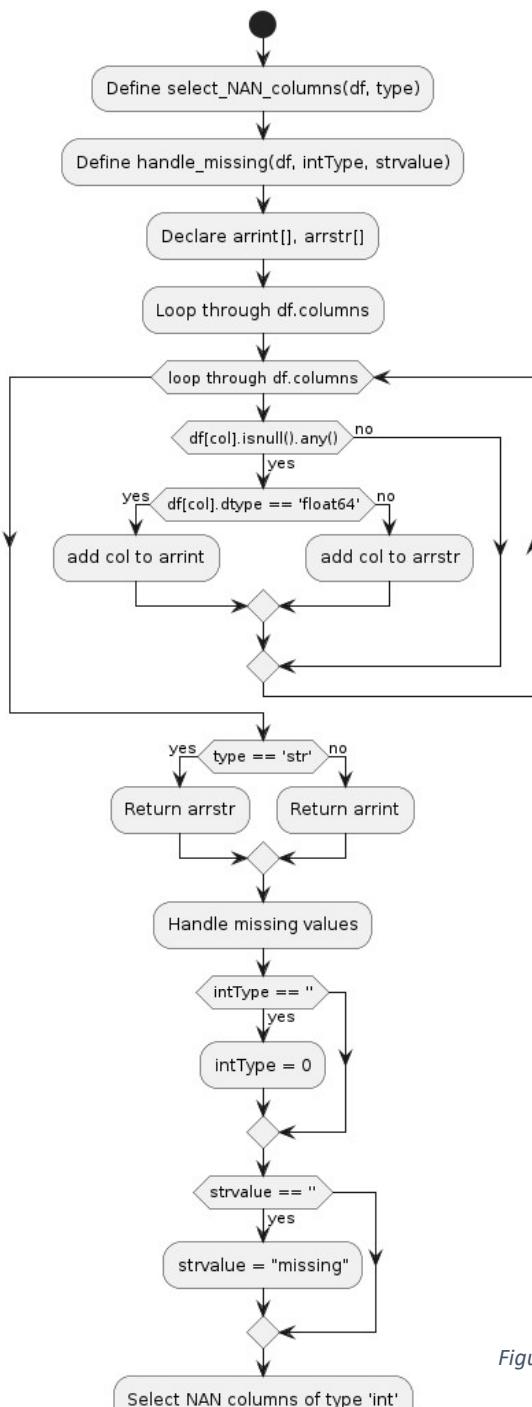
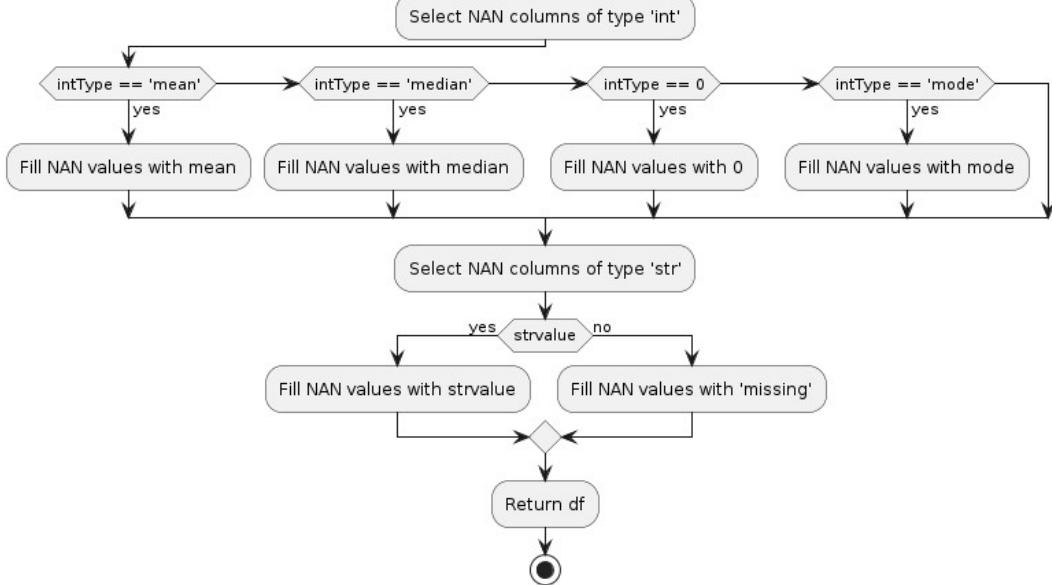


Figure 34 missing value flowchart



```

start
:Define select_NAN_columns(df, type);
:Define handle_missing(df, intType, strvalue);
:Declare arrint[], arrstr[];
:Loop through df.columns;
while (loop through df.columns)
if (df[col].isnull().any()) then (yes)
if (df[col].dtype == 'float64') then (yes)
:add col to arrint;
else (no)
:add col to arrstr;
endif
else (no)
endif
endwhile
if (type == 'str') then (yes)
:Return arrstr;
else (no)
:Return arrint;
endif
:Handle missing values;
if (intType == '') then (yes)
:intType = 0;
endif
if (strvalue == '') then (yes)
:strvalue = "missing";
endif
:Select NAN columns of type 'int';
if (intType == 'mean') then (yes)
:Fill NAN values with mean;
elseif (intType == 'median') then (yes)
:Fill NAN values with median;
elseif (intType == 0) then (yes)
:Fill NAN values with 0;
elseif (intType == 'mode') then (yes)
:Fill NAN values with mode;
endif
:Select NAN columns of type 'str';
if (strvalue) then (yes)
:Fill NAN values with strvalue;
else (no)
:Fill NAN values with 'missing';
endif
:Return df;
stop

```

Figure 35 missing value pseudocode

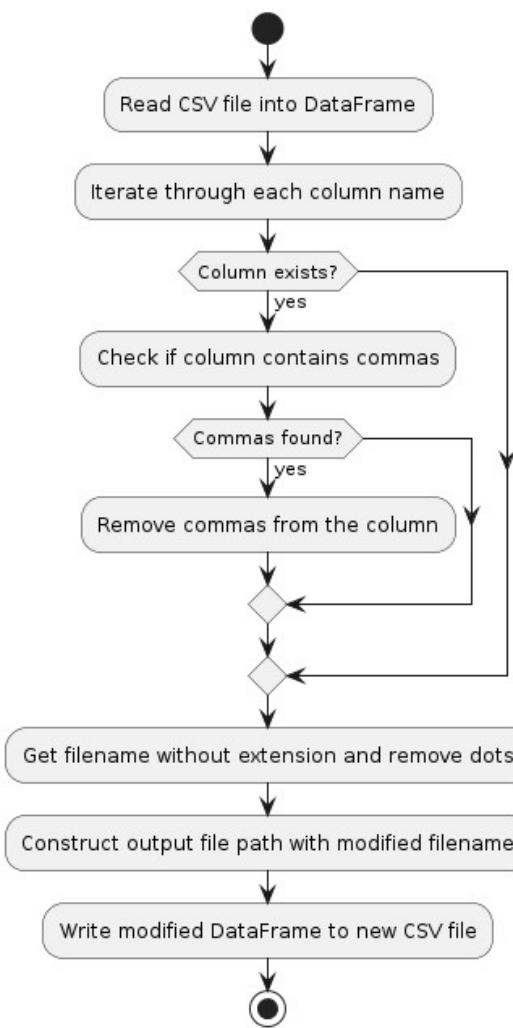


Figure 37 remove commas flowchart

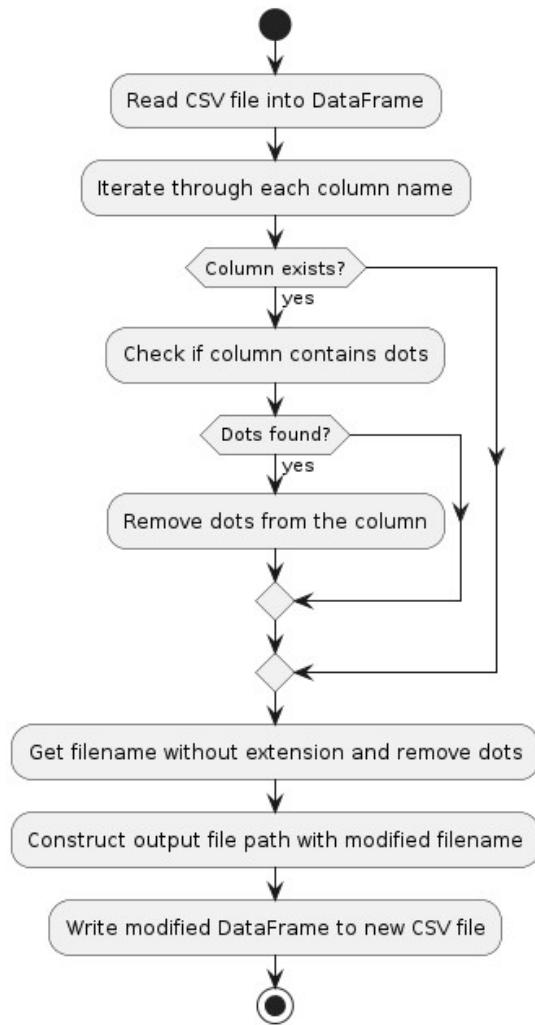


Figure 36 remove dots flowchart

```

start

:Read CSV file into DataFrame;
:Iterate through each column name;
if (Column exists?) then (yes)
    :Check if column contains commas;
    if (Commas found?) then (yes)
        :Remove commas from the column;
    endif
endif
:Get filename without extension and remove dots;
:Construct output file path with modified filename;
:Write modified DataFrame to new CSV file;

stop

```

Figure 39 remove commas pseudocode

```

start

:Read CSV file into DataFrame;
:Iterate through each column name;
if (Column exists?) then (yes)
    :Check if column contains dots;
    if (Dots found?) then (yes)
        :Remove dots from the column;
    endif
endif
:Get filename without extension and remove dots;
:Construct output file path with modified filename;
:Write modified DataFrame to new CSV file;

stop

```

Figure 38 remove dots pseudocode

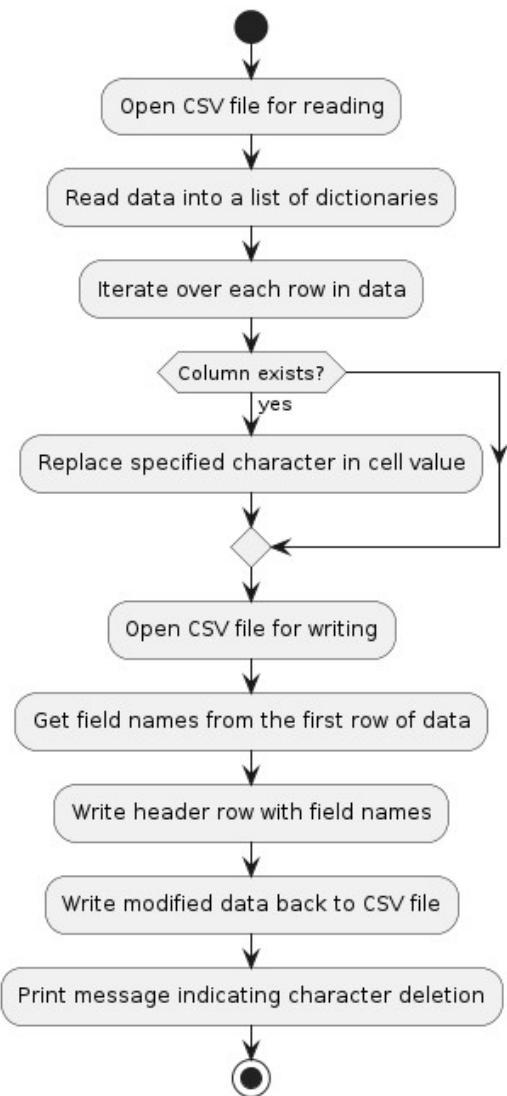


Figure 40 remove letters flowchart

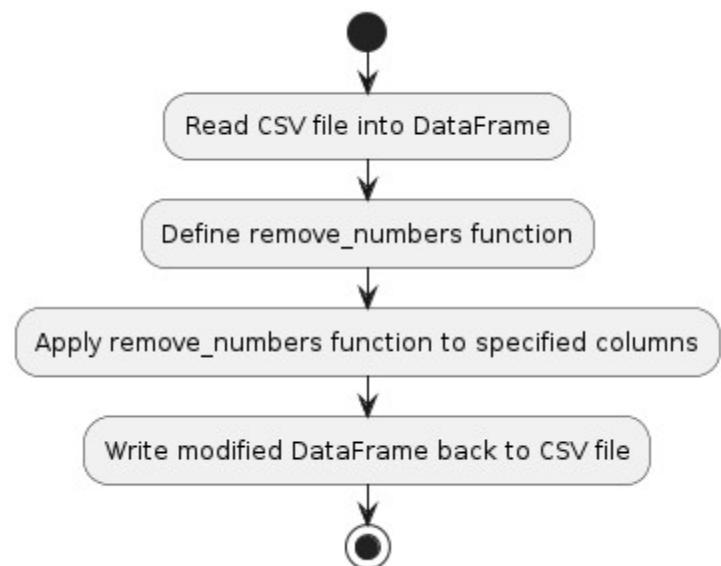


Figure 41 remove numbers flowchart

```

start
:Read CSV file into DataFrame;
:Define remove_numbers function;
:Apply remove_numbers function to specified columns;
:Write modified DataFrame back to CSV file;

stop

```

Figure 42 remove numbers pseudocode

```

start
:Open CSV file for reading;
:Read data into a list of dictionaries;
:Iterate over each row in data;
if (Column exists?) then (yes)
    :Replace specified character in cell value;
endif
:Open CSV file for writing;
:Get field names from the first row of data;
:Write header row with field names;
:Write modified data back to CSV file;
:Print message indicating character deletion;

stop

```

Figure 43 remove letters pseudocode

5.3 Block diagram

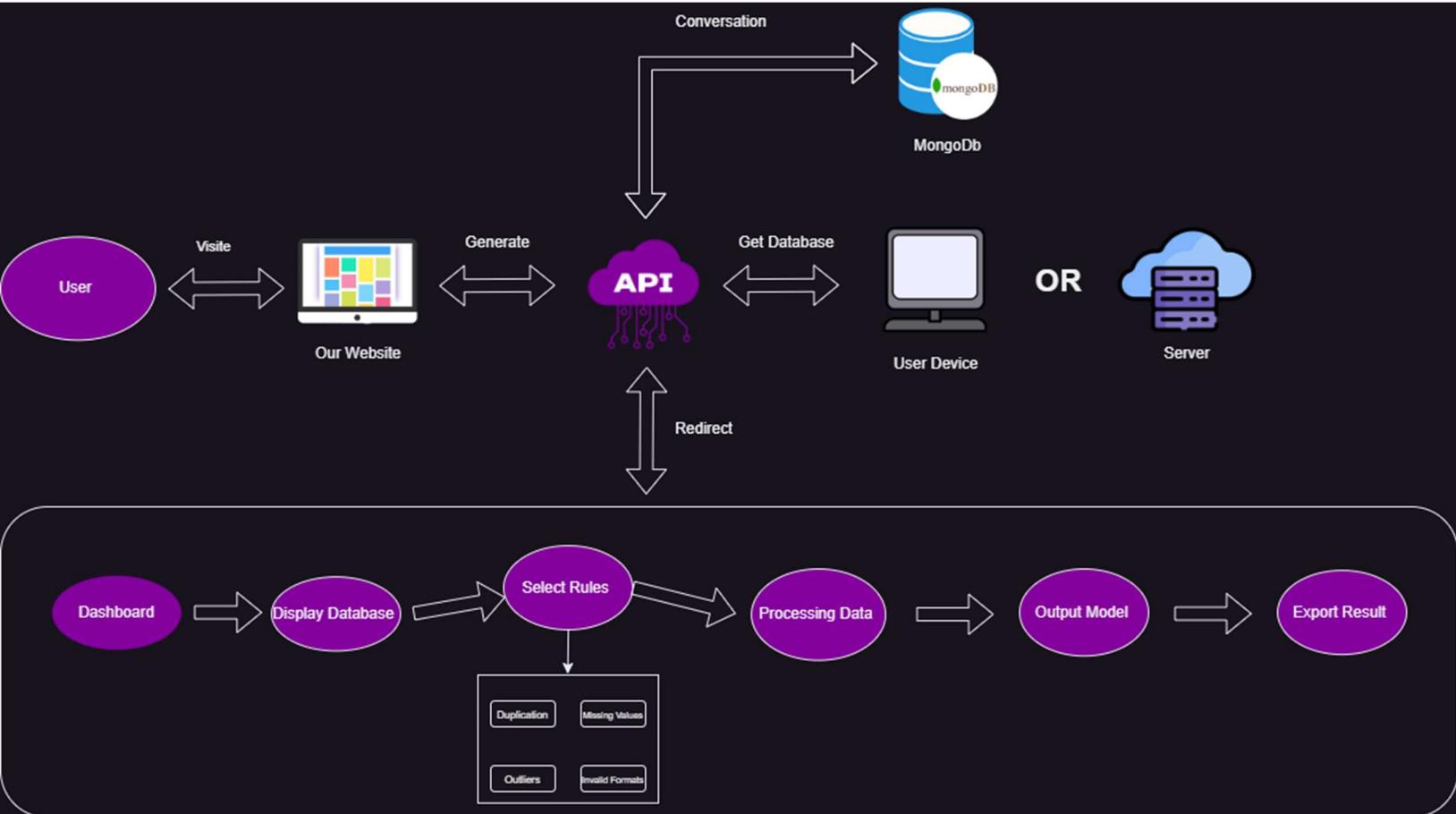


Figure 44 Block Diagram

5.4 Figma designs

At this part, prototypes of the system interface are created to simulate the initial design of the system. :-

SIGN IN

TO CONTINUE TO FREELANCE.COM

Continue with Facebook

Continue with Google

OR

USE PHONE

Email or phone number

password

CONTINUE

NO ACCOUNT SIGN UP



Figure 46 Sign in page

SIGN UP

TO CONTINUE TO FREELANCE.COM

Continue with Facebook

Continue with Google

OR

USE PHONE

First Name(optional)

First Name(optional)

Username

Email address

phone number (optional)

password

YOUR PASSWORD MEETS ALL THE NECESSARY REQUIREMENTS

CONTINUE

HAVE AN ACCOUNT? SIGN IN



Figure 45 Sign up page

VERIFY YOUR EMAIL

TO CONTINUE TO FREELANCE.COM

'example@ex.com'

Verification code
ENTER THE VERIFICATION CODE SENT TO YOUR EMAIL ADDRESS

DIDN'T RECEIVE A CODE? RESEND



Figure 48 verify email

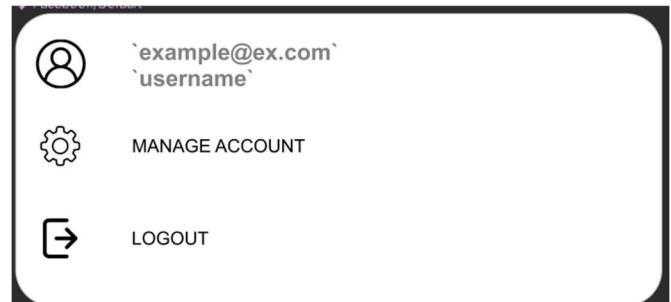


Figure 47 profile popup



Clean Up Your Data Now!

With **Data Cleanser**

Improve Your Data Quality



Data Cleanser

Improves your Data Quality

What's Data Cleanser ?

Data Cleanser is a tool that process of correcting poor quality or inaccurate customer data within your CRM system, database, record set or table.

Data decays over time and very quickly you can find that you have data in your business application that is duplicated, incomplete or poorly formatted. Data Cleansing is important for business performance and also more recently for regulatory compliance such as GDPR.

Our Data Cleanser Software will enable you to have duplicate free, accurate and trusted data which will drive out waste, reduce risk and help you to make better informed business decisions. If you are looking for a Single Customer View, our software can help.

We Provide Best Services For Customers



Cleansing Data

Lorum ipsum dolor sit amet etc
adipiscing elitipsum.



Export Data

Lorum ipsum dolor sit amet etc
adipiscing elitipsum.



Reports

Lorum ipsum dolor sit amet etc
adipiscing elitipsum.

Why Others Prefer Data Cleanser ?

Figure 49 Homepage

 ACCOUNT
 SECURITY

ACCOUNT

Manage your account information

PROFILE

 'username' →

USERNAME

'username'

 change username

EMAIL ADDRESSES

email address

(Primary)



+ Add an email address

PHONE NUMBERS

+ Add a phone number

CONNECTED ACCOUNTS

+ Add an account

SECURITY

Manage your security preferences

PASSWORD



 change password

TWO-STEP VERIFICATION

+ add two-step verification

DANGER

DELETE ACCOUNT

Delete your account and all its associated data

DELETE ACCOUNT

Figure 50 user account slider

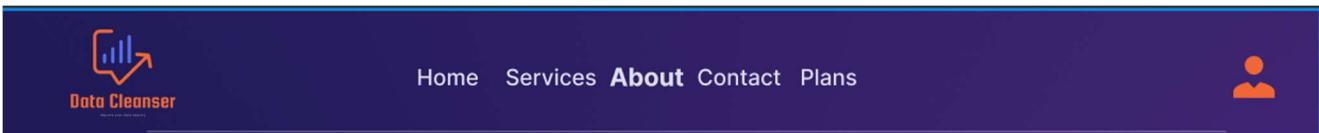


Figure 52 header



Figure 51 Footer

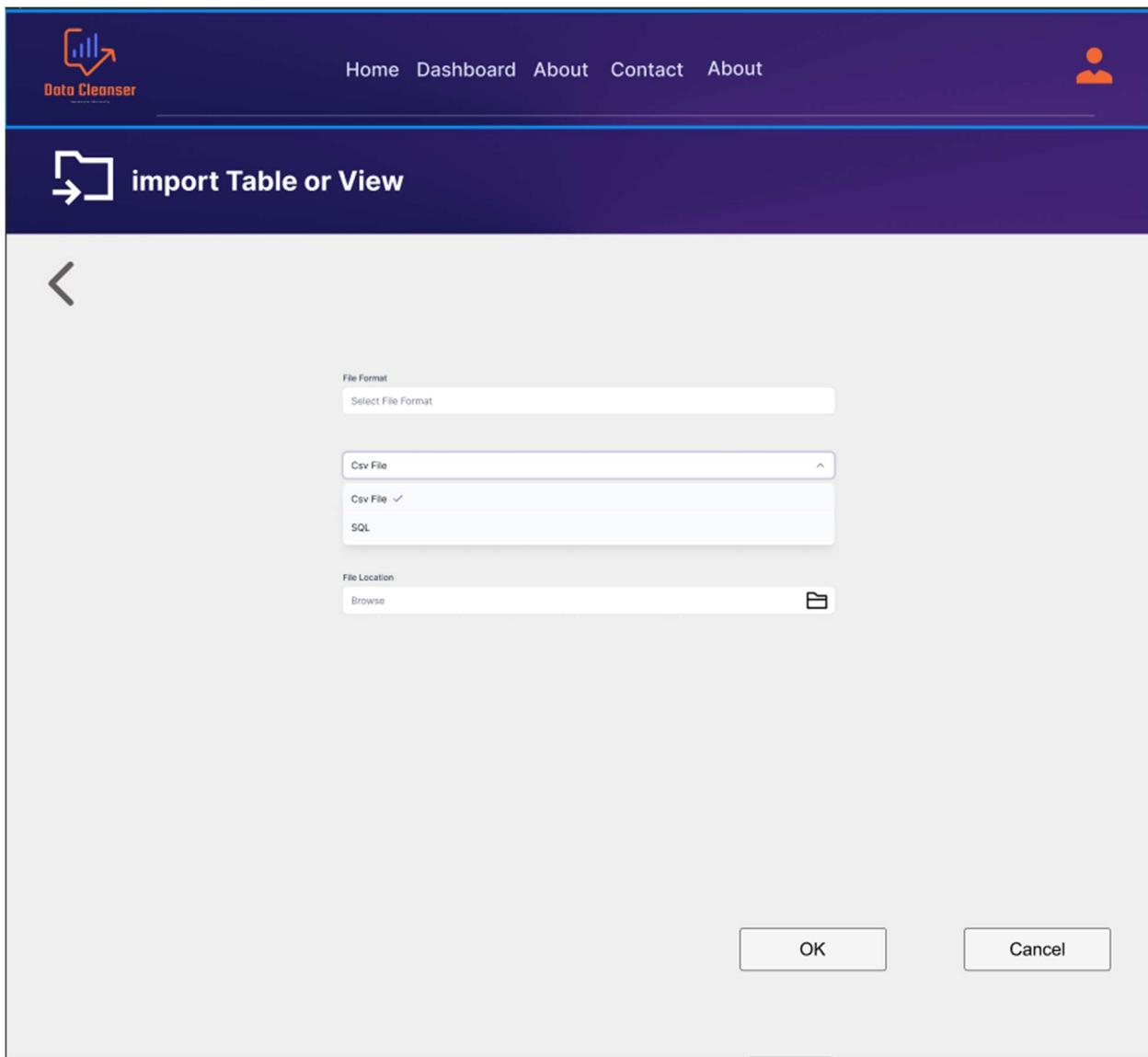


Figure 53 Import page

Export Options

Select Export Preferences

File Format

Sql Server

Sql Server

CSV Files

File Name

Enter File Name

File Location

Browse

Connection information

Server Name / IP Address

Port

Log on to the server

- Use Windows Authentication
 Use SQL Server Authentication

Login :

Password :

Database

Table Name

Export

Figure 54 export page

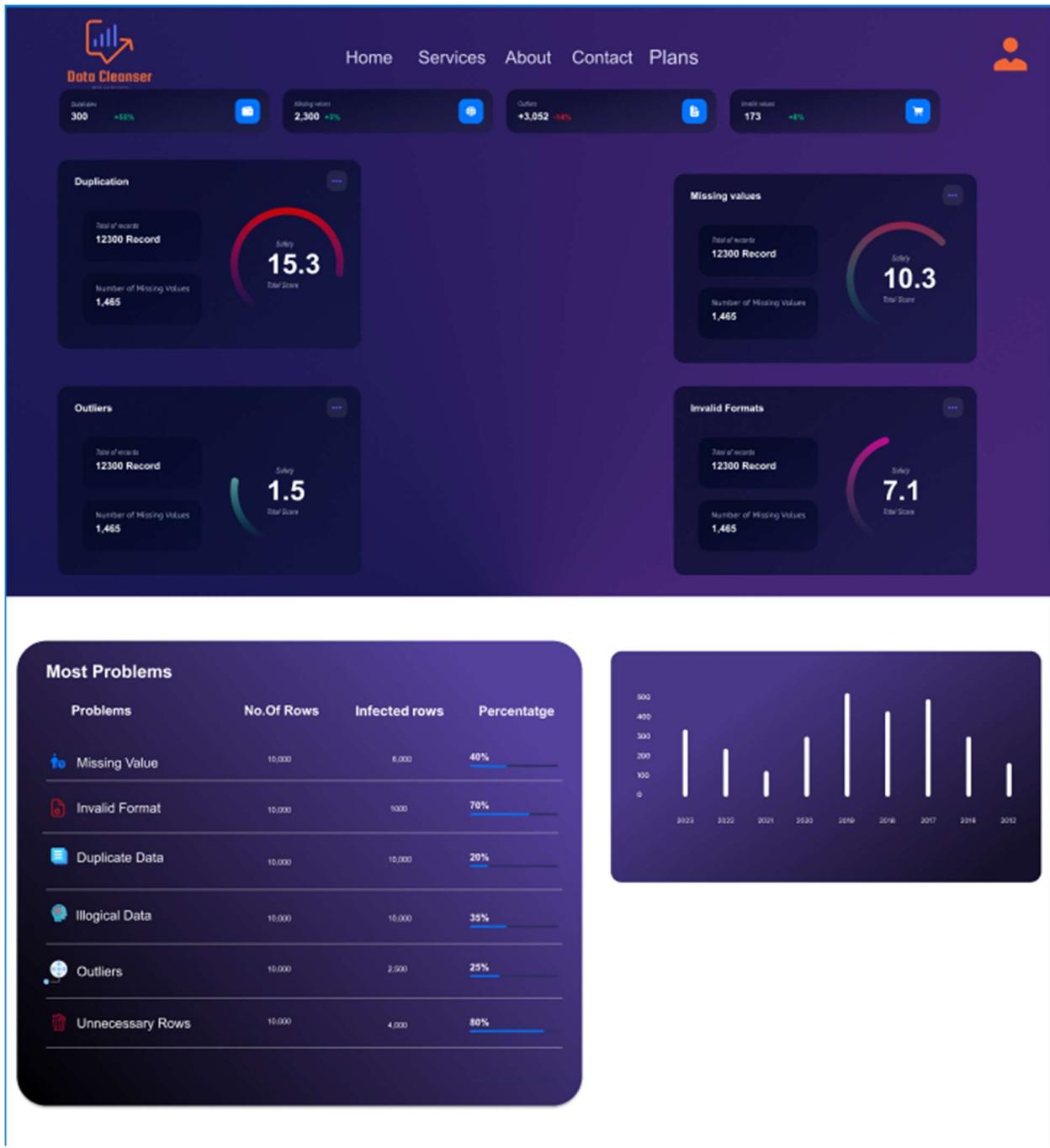


Figure 55 Dashboard

The screenshot shows the Data Cleanser application interface. At the top, there is a navigation bar with links: Home, Services, About, Contact, Plans, and a user profile icon. A large circular progress bar indicates 30% completion. Below the progress bar are two main sections: "Fill In" and "DELETE". Under "Fill In", there are four sub-sections: STRING, NUMBER, ROW, and COLUMN. Each section has several buttons: for STRING, FIXED VALUE, FORWARD, MEDIAN, FIXED NUM, FORWARD, and BACKWARD; for NUMBER, MOST REPEATED; and for ROW and COLUMN, various selection and processing buttons. To the right of these sections is a table titled "Table" with columns for Patient, Nurse, Doctor, and Room. Below this table is a data grid with columns: Action, Employees Name, Email Address, and Role. The data grid contains eight rows of employee information, each with edit and delete icons. At the bottom of the page is a footer with the Data Purifier logo, social media links (Facebook, X, Instagram, LinkedIn), a services menu, terms & conditions, and a newsletter sign-up form.

Action	Employees Name	Email Address	Role
	Lindsey Stroud	lindsey.stroud@gmail.com	Head of Technology
	Sarah brown	sarah.brown@gmail.com	Head of Technology
	Micheal Owen	michael.owen@gmail.com	Head of Technology
	Mary Jane	mary.jane@gmail.com	Head of Technology
	Peter doodle	peter.doodle@gmail.com	Head of Technology
	Peter doodle	peter.doodle@gmail.com	Head of Technology
	Peter doodle	peter.doodle@gmail.com	Head of Technology
	Peter doodle	peter.doodle@gmail.com	Head of Technology

Figure 56 Missing values page



15.3

Data Score

Show

DELETE

Employee Data Overview				
	Employee Name	Email Address	Role	Action
<input type="checkbox"/>	Lindsey Stroud	lindsey.stroud@gmail.com	Head of Technology	
<input type="checkbox"/>	Sarah brown	sarah.brown@gmail.com	Head of Technology	
<input type="checkbox"/>	Michael Owen	michael.owen@gmail.com	Head of Technology	
<input type="checkbox"/>	Mary Jane	mary.jane@gmail.com	Head of Technology	
<input type="checkbox"/>	Peter doodle	peter.doodle@gmail.com	Head of Technology	
<input type="checkbox"/>	Peter doodle	peter.doodle@gmail.com	Head of Technology	
<input type="checkbox"/>	Peter doodle	peter.doodle@gmail.com	Head of Technology	
<input type="checkbox"/>	Peter doodle	peter.doodle@gmail.com	Head of Technology	

Figure 57 Handling duplicates page



Invalid Formats

15.3
Data Score**Remove**

- Remove?
- Carriage
- dots
- white Spaces
- Letters
- Numbers

Standardize

- Standardize?
- upper case
- Lower Case
- Proper Case

Check

- Check?
- Phone
- email

Split

- Split?
- /
- *
-
- :
- ,

- Table
- Patient
- Murke
- Doctor
- Kids
- Treatment

Refresh Table**Rename****Remove****Insert Column****Merge Column****Hide/UnHide column**

Type	Type	Type	Action
<input type="checkbox"/> General	<input type="checkbox"/> Email Address	<input type="checkbox"/> Role	<input checked="" type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/> Text	lindsey.stroud@gmail.com	Head of Technology	<input checked="" type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/> Integer	sarah.brown@gmail.com	Head of Technology	<input checked="" type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/> Float	michael.owen@gmail.com	Head of Technology	<input checked="" type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/> Date	mary.jane@gmail.com	Head of Technology	<input checked="" type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/> Micheal Owen	peter.doodle@gmail.com	Head of Technology	<input checked="" type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/> Mary Jane	peter.doodle@gmail.com	Head of Technology	<input checked="" type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/> Peter doodle	peter.doodle@gmail.com	Head of Technology	<input checked="" type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/> Peter doodle	peter.doodle@gmail.com	Head of Technology	<input checked="" type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/> Peter doodle	peter.doodle@gmail.com	Head of Technology	<input checked="" type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/> Peter doodle	peter.doodle@gmail.com	Head of Technology	<input checked="" type="checkbox"/> <input type="checkbox"/>

Apply**Save**

Figure 58 invalid formats page

Remove Outliers

outliers percentage

20%

Total Score

No. of missing records

10/50

records

No. of suspected records

20/100

records

Table
Patient
Nurse
Doctor
Room
Treatment

Select preferred table you want to run
the test on

Select preferred data visualization
diagram
from the drop list beside

data visualization
line chart
scatter plot
histogram
box plot
basic pie charts

 Show all

 edit

 Delete selected

 Delete all

 Hide/UnHide column

User Name	User Phone	User City	User Age	User Department	User Id	User Status
Mohamed A dly	01024280109	Zawyaa	21	cs	92746	Male
Abdelrahman Gh.oenim	01024360988	Maadi	21	cs	92747	Male
Mohamed Abdelmote'lb	01095282605	Maadi	21	cs	92748	Male
Eyad Ernam	01024165905	Maadi	21	cs	92749	Male
Toka Saeed	01065479374	Banha	21	cs	92750	Female
Mohamed abdelkreem	01141364674	Saraya Elkoba	21	cs	92751	Male
mohamed Osama	01200427130	Giza	21	cs	92752	Male
mohamed rashad	01144970241	Sohag	21	cs	92753	Male


 Apply

 Save

Figure 59 Remove outliers page

Chapter SIX

System

Implementation

Chapter 6– system implementation

6.1 Database implementation

BASIC STRUCTURE OF MONGO DB

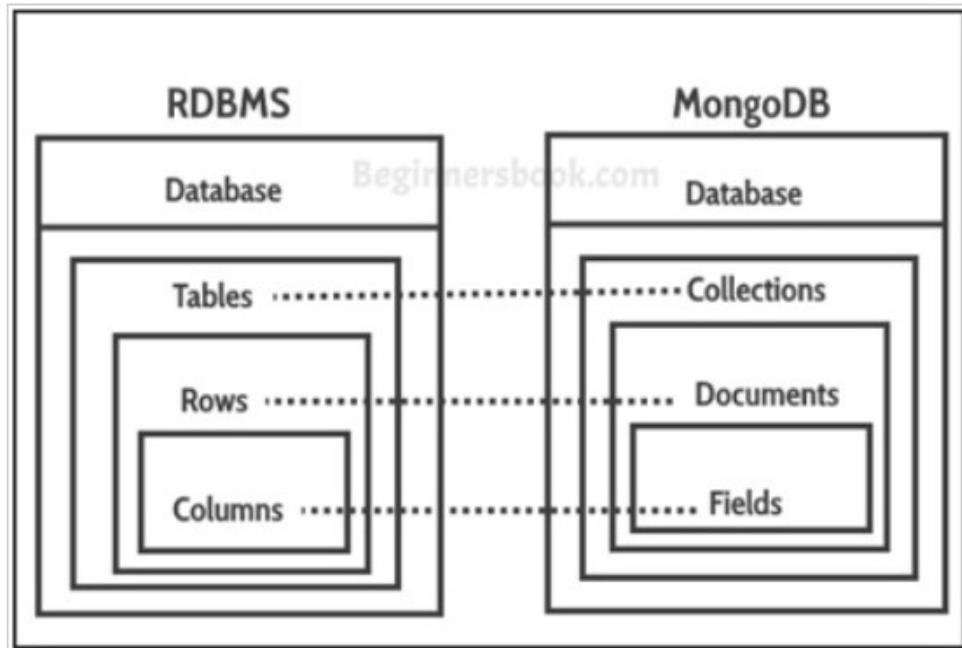


Figure 60 basic structure of mongo db

We chose mongo compass to implement our database as a non relational one in order to achieve the highest retrieval speed and quality

We wanted to create a containerized database in order to achieve the concept of **lightweight database on a virtual machine**

This means that you can run and retrieve fast your database on all types of machines without making assumptions for the operating system working right now or the database version your working with

So how do these containers actually look like?

Database server container is defined by elements called Dockerfiles, which basically define a sequence of steps required to build it. Dockerfiles begin with a “base image,” followed by databases and SQL scripts. It’s also possible to use snapshots and database clones when required.

Databases can be copied and run within the container file system, or mounted to the container directly by using the MOUNTDB command.

Pros:

1 – CI/CD Friendly – Database Containers are extremely consistent, which is very important while creating and maintaining an agile environment. Furthermore, these containers are proving to be very effective with approaches like Continuous Integration and Continuous Delivery (CI/CD).

2 – Multi-Cloud Compatibility – More and more companies today are leaving the on-premise format or creating hybrid pipelines to enable the seamless scaling up of their operations during peak-times. Not only can containers operate smoothly on the cloud, but they can do so on multiple cloud platforms.

The widespread popularity of the Docker image format further helps with portability. You can now use containers wherever you wish to operate.

3 – Cost Efficiency – Containers are very cost-efficient. In spite of the initial investment required for memory, CPU, and storage, it is possible to run many containers on the same infrastructure. They also integrate better with third party solutions (i.e – replication, mirroring, etc.), hence saving time and money.

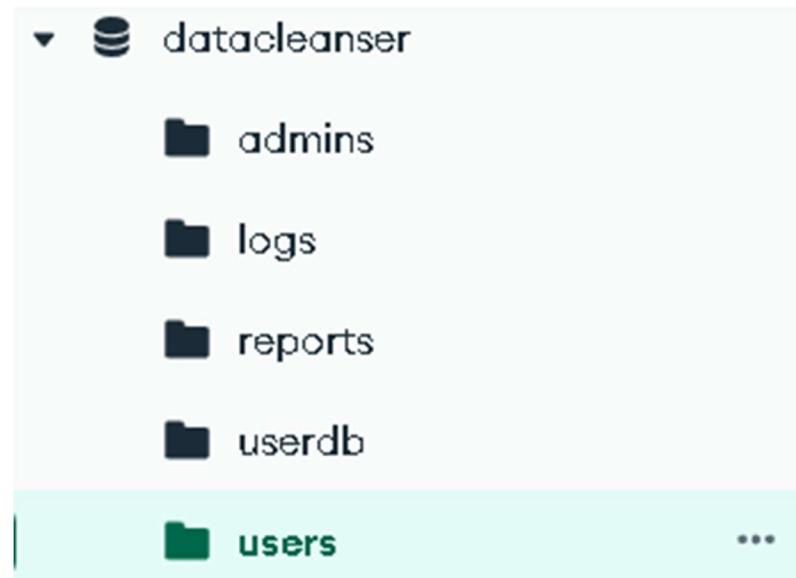
Cons:

1 – Security – Since all of the database containers usually share a common host, it becomes potentially easy to penetrate to the system due to the lack of isolation. Traditional security methods can be helpful in detecting such breaches, but this is by no means an air-tight methodology, something that needs to be clear.

2 – Lack of Isolation – Another issue with a container is the lack of OS flexibility which means containers must be of the same OS as of the Base OS. It can also become challenging to monitor database activity when hundreds (and potentially thousands) of containers are running on the same server.

3 – New Technology – As mentioned earlier, the use of Database Containers is a new practice and the methodology has not yet matured. There are still many aspects that need to be developed and polished before database containerization becomes a must-have component in every ecosystem.

Database Collection Implementation:



admins

Storage size: 20.48 kB
Documents: 1
Avg. document size: 97.00 B
Indexes: 1
Total index size: 20.48 kB

logs

Storage size: 20.48 kB
Documents: 1
Avg. document size: 94.00 B
Indexes: 1
Total index size: 20.48 kB

reports

Storage size: 20.48 kB
Documents: 1
Avg. document size: 1.15 kB
Indexes: 1
Total index size: 20.48 kB

userdb

Storage size: 4.10 kB
Documents: 0
Avg. document size: 0 B
Indexes: 1
Total index size: 4.10 kB

users

Storage size: 20.48 kB
Documents: 1
Avg. document size: 97.00 B
Indexes: 1
Total index size: 20.48 kB

1-User db



```
_id: ObjectId('667771824e24d291c9cdfacb')
Username : ***
Email : ***
First Name : ***
Last Name : ***
Password : ***
```

- **_id:** This field represents the unique identifier for the document in Integer data type.
- **Username:** This field stores the username of the user in Integer data type.
- **Email:** This field stores the email address of the user in varchar data type.
- **Name:** This field stores the name of the user in varchar.
- **Password:** This field stores the hashed password of the user for authentication purposes in varchar data type.

2- Log db



```
_id: ObjectId('667771114e24d291c9cdfac4')
Username : ""
Database_Name : ""
LogEntries : []
timestamp : Timestamp({ t: 0, i: 0 })
```

- **_id:** This field represents the unique identifier for the document in Integer data type.
- **Username:** The username of the user associated with the logged activity.
- **Database_Name:** The name of the database on which the operation was performed.
- **LogEntries:** An array of objects representing individual log entries, capturing details such as the operation performed and the timestamp.

3- Reports db

```
1  _id: ObjectId('6677713b4e24d291c9cdfac6')
2  Timestamp : Timestamp({ t: 0, i: 0 })
3  User : "JohnDoe"
4  ▼ Cleansing_Techniques_Used : Array (empty)
5  ▼ Summary_Statistics : Object
6      Negative_Numbers_Found : 0
7      ▼ Outliers : Object
8          Outliers_Found : 0
9          Outliers_Removed : 0
10         percentage_before_out : 0
11         percentage_after_out : 0
12         ▼ Missing_Values : Object
13             String : 0
14             Integer : 0
15             Missing_Values_Found : 0
16             Missing_Values_Removed : 0
17             percentage_before_ms : 0
18             percentage_after_ms : 0
19             ▼ Duplicates : Object
20                 Duplicates_Found : 0
21                 Duplicates_Removed : 0
22                 percentage_before_dup : 0
23                 percentage_after_dup : 0
24             ▼ Invalid_Formats : Object
25                 Email : 0
26                 Phone : 0
27                 Lower_Cases : 0
28                 Upper_Cases : 0
29                 Proper_Cases : 0
30                 Merge : 0
31                 Split : 0
32                 Proper_Case : 0
33                 Remove_Dots : 0
34                 Remove_Letters : 0
35                 Remove_Numbers : 0
36                 Remove_Selected_Value_From_Column : 0
37                 Remove_Double_Spaces : 0
38                 Strip : 0
39                 Invalid_formats_found : 0
40                 Invalid_formats_removed : 0
41                 percentage_before_Inv : 0
42                 percentage_after_Inv : 0
43             ▼ Statistical_Measures : Object
44                 Mean : 0
45                 Median : 0
46                 Mode : 0
47                 Standard_Deviation : 0
48                 Variance : 0
49                 Minimum_Value : 0
50                 Maximum_Value : 0
51                 user_inputs : " "
52             ▼ Data_Source : Object
53                 File_Name : " "
54                 Database_Name : " "
55                 Table_name : " "
```

_id: This field represents the unique identifier for the document in Object ID data type.

User: The username of the user associated with the report.

Timestamp: The date and time when the report was generated.

Summary_Statistics: Contains various summary statistics calculated during data cleansing, including:

Negative_Numbers_Found: Count of negative numbers found in the data.

Outliers: Count of outliers detected in the data.

Missing_Values: An object containing counts of missing values identified as string and integer types.

Duplicates: Count of duplicate records found in the data.

Invalid Formats: An object containing counts for various invalid format validations and manipulations, including email and phone number formats, conversions to lower case, upper case, proper case, merge and split operations, removal of dots, letters, numbers, selected values from columns, double spaces, and leading and trailing whitespaces.

Statistical_Measures: An object containing statistical measures calculated from the data, including mean, median, mode, standard deviation, variance, minimum value, and maximum value.

user_inputs: Additional inputs provided by the user.

Data_Source: Information about the source of the data being analyzed, including the name of the file, database, and table containing the data.

Data_Cleansing_Status: Indicates the status of the data cleansing process.

Cleansing_Techniques_Used: An array containing the cleansing techniques applied during data processing.

4-user db

```
{ WebDb.User_db.json X
  Click here to ask Blackbox to help you code faster
1  [
2    {
3      "_id": {
4        "$oid": "6635198d4760878753526692"
5      },
6      "Username": "",
7      "Databases": {
8        "database_1_info": {
9          "database_name": "Hospital_Management_System",
10         "database_type": "SQLiteDatabase",
11         "table_1_info": {
12           "table_name": "Patient",
13           "table_type": "varchar",
14           "Columns": {
15             "ColumnsName1": "SID_P",
16             "ColumnsType1": "Integer",
17             "ColumnsName2": "PatientID_P",
18             "ColumnsType2": "Integer",
19             "ColumnsName3": "PatientName",
20             "ColumnsType3": "Char",
21             "ColumnsName4": "Gender",
22             "ColumnsType4": "Char",
23             "ColumnsName5": "Address",
24             "ColumnsType5": "Char",
25             "ColumnsName6": "PhoneNumber_P",
26             "ColumnsType6": "Integer"
27           }
28         },
29         "table_2_info": {
30           "table_name": "Room",
31           "table_type": "varchar",
32           "Columns": {
33             "ColumnsName1": "RoomID",
34             "ColumnsType1": "Integer",
35             "ColumnsName2": "Type_of_room",
36             "ColumnsType2": "Char",
37             "ColumnsName3": "NoOfDays",
38             "ColumnsType3": "Char"
39           }
40         }
41       }
42     }
43   ]
44 }
```

```
47     "table_3_info":  
48     {  
49         "table_name": "Staff",  
50         "table_type": "varchar",  
51         "Columns":  
52         {  
53             "ColumnName1": "SID_S",  
54             "ColumnType1": "Integer",  
55             "ColumnName2": "AreaHandle",  
56             "ColumnType2": "Char",  
57             "ColumnName3": "Name_S",  
58             "ColumnType3": "Char",  
59             "ColumnName4": "PhoneNumber_S",  
60             "ColumnType4": "Integer"  
61         }  
62     },  
63     "table_4_info":  
64     {  
65         "table_name": "Nurse_Belongs_To",  
66         "table_type": "varchar",  
67         "Columns":  
68         {  
69             "ColumnName1": "SID_S",  
70             "ColumnType1": "Integer",  
71             "ColumnName2": "Area",  
72             "ColumnType2": "Integer",  
73             "ColumnName3": "NO_OF_DAYS_WORKED",  
74             "ColumnType3": "Integer"  
75         }  
76     },  
77     "table_5_info":  
78     {  
79         "table_name": "Ward_boy_Care_Belongs_To",  
80         "table_type": "varchar",  
81         "Columns":  
82         {  
83             "ColumnName1": "SID_S",  
84             "ColumnType1": "Integer",  
85             "ColumnName2": "Age",  
86             "ColumnType2": "Integer",  
87             "ColumnName3": "NO_OF_DAYS_WORKED",  
88             "ColumnType3": "Integer"  
89         }  
90     }.
```

```

91     "table_6_info":
92     {
93         "table_name": "Treatment",
94         "table_type": "varchar",
95         "Columns":
96         {
97             "ColumnName1": "Treatmentcode",
98             "ColumnType1": "Integer",
99             "ColumnName2": "DiseaseName_T",
100            "ColumnType2": "Char",
101            "ColumnName3": "Description_T",
102            "ColumnType3": "Char"
103        }
104    },
105    "table_7_info":
106    {
107        "table_name": "Bill_Payed_By",
108        "table_type": "varchar",
109        "Columns":
110        {
111            "ColumnName1": "PatientID_P",
112            "ColumnType1": "Integer",
113            "ColumnName2": "Bill_ID",
114            "ColumnType2": "Integer",
115            "ColumnName3": "NoOfDays",
116            "ColumnType3": "Char",
117            "ColumnName4": "fees",
118            "ColumnType4": "REAL"
119        }
120    },
121    "table_8_info":
122    {
123        "table_name": "Disease",
124        "table_type": "varchar",
125        "Columns":
126        {
127            "ColumnName1": "DiseaseName",
128            "ColumnType1": "Integer",
129            "ColumnName2": "Symptoms",
130            "ColumnType2": "Char"
131        }
132    },

```

```

133     "table_9_info":
134     {
135         "table_name": "Doctor",
136         "table_type": "varchar",
137         "Columns":
138         {
139             "ColumnName1": "DoctorID_D",
140             "ColumnType1": "Integer",
141             "ColumnName2": "Speciallization",
142             "ColumnType2": "Char",
143             "ColumnName3": "DoctorName",
144             "ColumnType3": "Char",
145             "ColumnName4": "Qualification",
146             "ColumnType4": "Char",
147             "ColumnName5": "Phone_Number",
148             "ColumnType5": "Integer"
149         }
150     },
151     "table_10_info":
152     {
153         "table_name": "Assigned",
154         "table_type": "varchar",
155         "Columns":
156         {
157             "ColumnName1": "RoomID",
158             "ColumnType1": "Integer",
159             "ColumnName2": "PatientID_P",
160             "ColumnType2": "Integer"
161         }
162     },
163     "table_11_info":
164     {
165         "table_name": "Care",
166         "table_type": "varchar",
167         "Columns":
168         {
169             "ColumnName1": "SID_S",
170             "ColumnType1": "Integer",
171             "ColumnName2": "PatientID_P",
172             "ColumnType2": "Integer"
173         }
174     },

```

```
175     "table_12_info":  
176     {  
177         "table_name": "Maintain",  
178         "table_type": "varchar",  
179         "Columns":  
180         {  
181             "ColumnName1": "SID_S",  
182             "ColumnType1": "Integer",  
183             "ColumnName2": "RoomID",  
184             "ColumnType2": "Integer"  
185         }  
186     },  
187     "table_13_info":  
188     {  
189         "table_name": "Treats",  
190         "table_type": "varchar",  
191         "Columns":  
192         {  
193             "ColumnName1": "DoctorID_D",  
194             "ColumnType1": "Integer",  
195             "ColumnName2": "PatientID_P",  
196             "ColumnType2": "Integer"  
197         }  
198     },  
199     "table_14_info":  
200     {  
201         "table_name": "Super_Form",  
202         "table_type": "varchar",  
203         "Columns":  
204         {  
205             "ColumnName1": "DiseaseName",  
206             "ColumnType1": "Integer",  
207             "ColumnName2": "PatientID_P",  
208             "ColumnType2": "Integer"  
209         }  
210     },  
211     "table_15_info":  
212     {  
213         "table_name": "Gets",  
214         "table_type": "varchar",  
215         "Columns":  
216         {  
217             "ColumnName1": "Treatmentcode",  
218             "ColumnType1": "Integer",  
219             "ColumnName2": "PatientID_P",  
220             "ColumnType2": "Integer"  
221         }  
222     }  
223 }
```

```
224     "database_2_info":  
225     {  
226         "database_name": "",  
227         "database_type": "",  
228         "table_1_info":  
229         {  
230             "table_name": "",  
231             "table_type": "",  
232             "Columns":  
233             {  
234                 "ColumnName1": "",  
235                 "ColumnType1": "",  
236                 "ColumnName2": "",  
237                 "ColumnType2": "",  
238                 "ColumnName3": "",  
239                 "ColumnType3": ""  
240             }  
241         },  
242         "table_2_info":  
243         {  
244             "table_name": "",  
245             "table_type": "",  
246             "Columns":  
247             {  
248                 "ColumnName1": "",  
249                 "ColumnType1": "",  
250                 "ColumnName2": "",  
251                 "ColumnType2": "",  
252                 "ColumnName3": "",  
253                 "ColumnType3": ""  
254             }  
255         }  
256     }  
257 }  
258 }]
```

- **_id**: This field represents the unique identifier for the document in Integer data type.
- **Username**: This field stores the username of the user to whom the database configuration belongs and it's a foreign key to the "Users" collection in varchar data type.
- **Databases**: This field contains information about various databases associated with the user. The "User_db" collection appears to provide a mechanism for users to define and manage their own databases and associated column configurations. This kind of flexibility can be useful in applications where users have varying data storage needs or where custom database structures are required for different purposes. By storing database configurations in the "User_db" collection, the application can dynamically create, modify, and access databases and their associated columns based on user-defined settings.

6.2 Raw backend implementation:-

6.2.1 Handling invalid formats

"Data consistency is paramount in any system or database. However, ensuring that all entries adhere to a specific format or standard can be challenging. From email addresses and phone numbers to whitespace discrepancies and case sensitivity, inconsistencies can arise in various forms. In this endeavor, we aim to categorize and tackle these invalid formats systematically, employing techniques such as merging, splitting, converting to proper case, and removing unwanted characters or values. By addressing these issues, we can enhance data integrity and optimize the functionality of our systems."

1. Format Validation:

- Email
- Phone

2. Case Handling:

- Lower cases
- Upper cases
- Proper cases

3. Data Manipulation:

- Merge
- Split
- Proper case
- Remove dots
- Remove letters
- Remove numbers
- Remove Selected Value From Column
- Remove double spaces
- Strip

4. Data Analysis:

- Show Repeated in all tables
- Show Repeated in one column

1st check emails

```
1 import csv
2 import re
3 import dns.resolver
4 def is_valid_email(email):
5     # Adjusted regular expression pattern for email validation
6     email_pattern = re.compile(r'^[a-zA-Z0-9._%+-]+@[a-zA-Z0-9.-]+\.[a-zA-Z]{2,}\$')
7
8     # Match the email pattern
9     if not email_pattern.match(email):
10         return False # Email pattern does not match
11
12     # Extract domain from the email address
13     domain = email.split('@')[1]
14     try:
15         # Resolve MX records for the domain
16         mx_records = dns.resolver.resolve(domain, 'MX')
17         return True # Domain has valid MX records
18     except dns.resolver.NoAnswer:
19         print(f"No MX records found for the domain: {domain}")
20         return False # No MX records found for the domain
21     except dns.resolver.NXDOMAIN:
22         print(f"Domain does not exist: {domain}")
23         return False # Domain does not exist
24     except dns.resolver.Timeout:
25         print(f"Timeout occurred while querying DNS servers for domain: {domain}")
26         return False # Timeout occurred while querying DNS servers
27     except Exception as e:
28         print(f"An error occurred while validating email {email}: {e}")
29         return False # Other errors
30
31 def show_valid_emails_only(input_file, email_column_index):
32     with open(input_file, 'r', newline='') as csv_in:
33         reader = csv.reader(csv_in)
34         header = next(reader) # Skip the header
35         print("Email,Is_Valid_Email") # Print header for email and validation status
36
37         for row in reader:
38             email = row[email_column_index] # Choose email address from the specified column index
39             is_valid = is_valid_email(email)
40             print(f"{email},{is_valid}") # Print the email address and its validation status
41
42 # Example usage:
43 input_file = 'Example.csv'
44 chosen_column_index = 2 # Specify the column index for email addresses
45 show_valid_emails_only(input_file, chosen_column_index)
```

2nd check phone

```
1 import csv
2
3 def check_phone_numbers_in_column(csv_file, column_index):
4     invalid_numbers = []
5     with open(csv_file, 'r') as file:
6         reader = csv.reader(file)
7         next(reader) # Skip the first row
8         for row_number, row in enumerate(reader, start=2):
9             if len(row) >= column_index:
10                 item = row[column_index - 1] # Adjust index to 0-based
11                 if len(item) != 11 or not item.startswith(('012', '010', '011', '015')) or not item.isdigit():
12                     invalid_numbers.append(f"Invalid phone number at row {row_number}, column {column_index}: {item}")
13             else:
14                 print(f"Warning: Row {row_number} does not have enough columns.")
15
16 return invalid_numbers
17
18 # Example usage:
19 column_index = 4 # Adjust this to the index of the column you want to check (1-based index)
20 invalid_numbers = check_phone_numbers_in_column('comma.csv', column_index)
21 if invalid_numbers:
22     print("Invalid phone numbers:")
23     for number in invalid_numbers:
24         print(number)
25 else:
26     print("All phone numbers in the specified column are valid.")
```

Rest in Appendex chapter...

6.2.2 Handling missing values

```
import pandas as pd
def select_NAN_columns(df,type):
    arrint=[]
    arrstr=[]
    for col in df.columns:
        if df[col].isnull().any():
            if df[col].dtype == 'float64':
                arrint.append(col)
            else:
                arrstr.append(col)
    if(type=='str') :
        return arrstr
    else:
        return arrint

def handle_missing(df,intType,strvalue):
    # handling missing int values
    if (intType==''):
        intType=0
    if (strvalue==''):
        strvalue="missing"

    arrint=select_NAN_columns(df,'int')
    if intType=='mean':
        for i in arrint:
            df[i]=df[i].fillna(df[i].mean())
    elif intType=='median':
        for i in arrint:
            df[i] = df[i].fillna(df[i].median())
    elif intType==0:
        for i in arrint:
            df[i] = df[i].fillna(0)
    elif intType=='mode':
        for i in arrint:
            mode=df[i].value_counts().idxmax()
            df[i] = df[i].fillna(mode)

    # handling missing str values
    arrstr=select_NAN_columns(df,'str')
    if (strvalue):
        for i in arrstr:
            df[i] = df[i].fillna(strvalue)
    else:
        for i in arrstr:
            df[i] = df[i].fillna('missing')
```

6.2.3 Removing outliers

Outliers, although sometimes disruptive, are pivotal in data cleaning. They can distort results and mislead interpretations but may also reveal important insights or errors. Researchers employ various methods to manage outliers, including removal, replacement, or transformation. The choice depends on the data's nature, analysis goals, and domain knowledge. Transparent reporting of outlier handling methods is crucial for research integrity. Ultimately, understanding and appropriately managing outliers enhance the reliability and validity of data-driven insights.

Detecting and handling outliers in a dataset is crucial for several reasons:

1. **Data Quality:** Outliers can skew statistical analyses and machine learning models, leading to incorrect conclusions or predictions. Removing or adjusting outliers helps maintain the quality and integrity of the dataset, ensuring more accurate results.
2. **Model Performance:** Outliers can influence the parameters of predictive models, reducing their effectiveness. By removing or replacing outliers, models can better capture the underlying patterns in the data and make more reliable predictions.
3. **Assumption Violation:** Many statistical methods and machine learning algorithms assume that the data is normally distributed or follows a specific distribution. Outliers can violate these assumptions, leading to biased estimates or incorrect inferences. Handling outliers helps ensure that the assumptions of the chosen analysis method are met.
4. **Data Interpretation:** Outliers can sometimes represent genuine but rare events or errors in data collection. Understanding the nature of outliers is essential for accurately interpreting the data and drawing meaningful conclusions. Handling outliers appropriately allows analysts to differentiate between anomalous observations and true insights.
5. **Data Visualization:** Outliers can distort visualizations such as scatter plots, histograms, and box plots, making it challenging to visualize the underlying patterns in the data. Removing or adjusting outliers improves the interpretability of visualizations and facilitates better insights into the data.

Methods for handling outliers include:

- **Removing:** Simply removing outliers from the dataset can be appropriate if they are likely to be errors or do not represent meaningful information. However, this approach should be used cautiously, as it can lead to loss of valuable data and potentially bias the analysis.
- **Replacing:** Outliers can be replaced with more typical values, such as the mean, median, or mode of the dataset. This approach can help mitigate the impact of outliers without removing them entirely. However, it may also introduce some distortion to the data if the outliers are indicative of genuine variability.

Ultimately, the choice of how to handle outliers depends on the specific characteristics of the dataset, the objectives of the analysis, and the domain knowledge of the analyst. It's essential to carefully consider the implications of different outlier handling strategies and choose the approach that best suits the particular context and goals of the analysis.

1. Detecting outliers in a numerical dataset using the Interquartile Range (IQR) method and Plot them:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

def find_outliers_iqr(column_data):
    """
    Detect outliers in a numerical dataset using the Interquartile Range (IQR) method.

    Args:
    - column_data: A pandas Series representing a column of numerical values.

    Returns:
    - outliers: A list of tuples containing outlier values found in the column and their indices.
    """
    # Calculate the first quartile (Q1) and third quartile (Q3) of the column data
    q1 = np.percentile(column_data, 25)
    q3 = np.percentile(column_data, 75)

    # Calculate the interquartile range (IQR)
    iqr = q3 - q1

    # Define the lower and upper bounds for outlier detection
    lower_bound = q1 - 1.5 * iqr
    upper_bound = q3 + 1.5 * iqr

    # Identify outliers and their indices
    outliers = [(value, idx) for idx, value in enumerate(column_data) if value < lower_bound or value > upper_bound]

    return outliers

# Load the CSV file into a pandas DataFrame
df = pd.read_csv("example11.csv")

# Choose the column index you want to perform outlier detection on
column_index = 2 # Change this to the desired column index

# Extract the specified column from the DataFrame as a pandas Series
column_data = df.iloc[:, column_index]

# Get the column name from the DataFrame based on the provided column index
column_name = df.columns[column_index]

# Find outliers in the specified column using the modified function
outliers = find_outliers_iqr(column_data)

# Calculate the total number of data points
total_data_points = len(column_data)

# Calculate the number of outliers
num_outliers = len(outliers)

# Calculate the number of non-outliers
num_non_outliers = total_data_points - num_outliers

# Calculate the percentage of outliers
percentage_outliers = (num_outliers / total_data_points) * 100
```

```

# Print the outliers and their locations
print("Outliers found in column '{}'".format(column_name))
for outlier in outliers:
    row_index = outlier[1] # Assuming the index corresponds to row number
    print("Value:", outlier[0], "at Row:", row_index)

# Print the number of outliers in the dataset
print("Number of outliers in the dataset: {} out of {}".format(num_outliers,
total_data_points))

# Print the percentage of outliers in the dataset
print("Percentage of outliers in the dataset: {:.2f}%".format(percentage_outliers))

# Scatter Plot
plt.figure(figsize=(8, 6))
plt.scatter(range(len(column_data)), column_data, color='blue', label='Data Points')
plt.scatter(*zip(*outliers), color='red', label='Outliers')
plt.title('Scatter Plot of ' + column_name)
plt.xlabel('Index')
plt.ylabel('Values')
plt.legend()
plt.grid(True)
plt.savefig('scatter_plot.png') # Save scatter plot as image
plt.show()

# Box Plot
plt.figure(figsize=(8, 6))
plt.boxplot(column_data, vert=False)
plt.title('Box Plot of ' + column_name)
plt.xlabel('Values')
plt.grid(True)
plt.savefig('box_plot.png') # Save box plot as image
plt.show()

# Histogram
plt.figure(figsize=(8, 6))
plt.hist(column_data, bins=20, color='skyblue', edgecolor='black')
plt.title('Histogram of ' + column_name)
plt.xlabel('Values')
plt.ylabel('Frequency')
plt.grid(True)
plt.savefig('histogram.png') # Save histogram as image
plt.show()

# Line Chart
plt.figure(figsize=(8, 6))
plt.plot(range(len(column_data)), column_data, color='green')
plt.title('Line Chart of ' + column_name)
plt.xlabel('Index')
plt.ylabel('Values')
plt.grid(True)
plt.savefig('line_chart.png') # Save line chart as image
plt.show()

```

Here's how the function works:

- Input:** It takes a pandas Series (`column_data`), representing a column of numerical values, as input.
- Calculating Quartiles and IQR:** The function calculates the first quartile (Q1) and third quartile (Q3) of the column data, and then computes the Interquartile Range (IQR), which is the difference between Q3 and Q1.

3. **Defining Bounds:** Using the IQR, the function defines lower and upper bounds for outlier detection. Values outside these bounds are considered outliers. The bounds are typically set at 1.5 times the IQR away from the first and third quartiles.
4. **Identifying Outliers:** It then iterates through the column data, identifying values that fall below the lower bound or above the upper bound. For each outlier found, it stores the value and its index in a list of tuples.
5. **Output:** The function returns a list of tuples containing the outlier values and their indices.

In addition to outlier detection, the script also includes data visualization techniques such as scatter plots, box plots, histograms, and line charts to visually represent the data and highlight the detected outliers.

2. Remove outlier values from a specified column:

```
import numpy as np
import pandas as pd
import os

def delete_outliers(df, column_index):
    """
    Delete outlier values from a DataFrame.

    Args:
    - df: A pandas DataFrame.
    - column_index: The index of the column containing outliers.

    Returns:
    - df_cleaned: The DataFrame with outlier values removed.
    """
    # Find outliers in the specified column
    column_data = df.iloc[:, column_index]
    q1 = np.percentile(column_data, 25)
    q3 = np.percentile(column_data, 75)
    iqr = q3 - q1
    lower_bound = q1 - 1.5 * iqr
    upper_bound = q3 + 1.5 * iqr

    # Delete outlier values
    df_cleaned = df.copy() # Make a copy of the original DataFrame
    df_cleaned.loc[(df_cleaned.iloc[:, column_index] < lower_bound) |
    (df_cleaned.iloc[:, column_index] > upper_bound), df.columns[column_index]] = np.nan

    return df_cleaned

# Load the CSV file into a pandas DataFrame
input_filename = "example11.csv" # Change this to the desired input filename
df = pd.read_csv(input_filename)

# Choose the column index you want to perform outlier detection on
column_index = 1 # Change this to the desired column index

# Delete outlier values from the specified column
df_cleaned = delete_outliers(df, column_index)

# Extract the basename of the original filename without the extension
filename, extension = os.path.splitext(input_filename)
basename = os.path.basename(filename)
```

```
# Construct the output filename with the suffix "deleted_outliers"
output_file = basename + "_deleted_outliers" + extension

# Save the cleaned DataFrame to a new CSV file with the modified filename
df_cleaned.to_csv(output_file, index=False)

print("Cleaned DataFrame with outlier values removed has been saved to
'{}'.".format(output_file))
```

Here's how the function works:

1. **Input:** It takes two arguments: the Data Frame (**df**) and the index of the column containing outliers (**column_index**).
2. **Outlier Detection:** Using the Interquartile Range (IQR) method, the function identifies outliers in the specified column. It calculates the first quartile (Q1), third quartile (Q3), and the interquartile range (IQR) of the column data.
3. **Defining Bounds:** Lower and upper bounds for outlier detection are then determined based on the IQR. Values falling below the lower bound or above the upper bound are considered outliers.
4. **Deleting Outliers:** The function creates a copy of the original Data Frame (**df_cleaned**) and replaces outlier values in the specified column with NaN (missing values). This effectively removes the outliers from the Data Frame.
5. **Output:** The cleaned Data Frame (**df_cleaned**) with outlier values removed is returned.

After cleaning the Data Frame, the script saves the cleaned Data Frame to a new CSV file with the suffix "`_deleted_outliers`" appended to the original filename.

3. Replace outlier values in a specified column with the average of non-outlier values within that column:

```
import numpy as np
import pandas as pd
import os

def replace_outliers_with_average(df, column_index):
    """
    Replace outliers in a DataFrame with the average of values in the column.

    Args:
    - df: A pandas DataFrame.
    - column_index: The index of the column containing outliers.

    Returns:
    - df_corrected: The DataFrame with outliers replaced by the average of values in the column.
    """
    # Find outliers in the specified column
    column_data = df.iloc[:, column_index]
    q1 = np.percentile(column_data, 25)
    q3 = np.percentile(column_data, 75)
    iqr = q3 - q1
    lower_bound = q1 - 1.5 * iqr
    upper_bound = q3 + 1.5 * iqr

    # Calculate the average of values within the interquartile range
    avg = np.mean(df[(df.iloc[:, column_index] >= lower_bound) & (df.iloc[:, column_index] <= upper_bound)].iloc[:, column_index])

    # Ensure that the average value has the same data type as the column
    avg = df.iloc[:, column_index].dtype.type(avg)

    # Replace outliers with the average of values
    df_corrected = df.copy()
    df_corrected.loc[(df_corrected.iloc[:, column_index] < lower_bound) | (df_corrected.iloc[:, column_index] > upper_bound), df.columns[column_index]] = avg

    return df_corrected

# Load the CSV file into a pandas DataFrame
input_filename = "example11.csv" # Change this to the desired input filename
df = pd.read_csv(input_filename)

# Choose the column index you want to perform outlier detection on
column_index = 1 # Change this to the desired column index

# Replace outliers with the average of values in the specified column
df_corrected = replace_outliers_with_average(df, column_index)

# Extract the basename of the original filename without the extension
filename, extension = os.path.splitext(input_filename)
basename = os.path.basename(filename)

# Construct the output filename with the suffix "_average_to_outliers"
output_file = basename + "_average_to_outliers" + extension

# Save the corrected DataFrame to a new CSV file with the modified filename
df_corrected.to_csv(output_file, index=False)

print("Corrected DataFrame with outliers replaced by the average of values has been saved to '{}'.format(output_file))
```

Here's how the function works:

1. **Input:** The function takes two arguments: the Data Frame (**DF**) and the index of the column containing outliers (**column_index**).
2. **Outlier Detection:** Using the Interquartile Range (IQR) method, the function identifies outliers in the specified column. It calculates the first quartile (Q1), third quartile (Q3), and the interquartile range (IQR) of the column data.
3. **Calculating Average:** It then calculates the average of values within the interquartile range, excluding outliers. This average serves as the replacement value for outliers.
4. **Replacing Outliers:** Outliers are replaced with the calculated average of non-outlier values in the column.
5. **Output:** The corrected Data Frame (**df_corrected**) with outliers replaced by the average of values is returned.

After replacing outliers, the script saves the corrected Data Frame to a new CSV file with the suffix "_average_to_outliers" appended to the original filename.

4. Replace outlier values in a specified column of with a chosen value:

```
import numpy as np
import pandas as pd
import os

def replace_outliers_with_value(df, column_index, new_value):
    """
    Replace outliers in a DataFrame with a chosen value.

    Args:
    - df: A pandas DataFrame.
    - column_index: The index of the column containing outliers.
    - new_value: The new value to replace outliers with.

    Returns:
    - df_corrected: The DataFrame with outliers replaced by the chosen value.
    """
    # Find outliers in the specified column
    column_data = df.iloc[:, column_index]
    q1 = np.percentile(column_data, 25)
    q3 = np.percentile(column_data, 75)
    iqr = q3 - q1
    lower_bound = q1 - 1.5 * iqr
    upper_bound = q3 + 1.5 * iqr

    # Replace outliers with the chosen value
    df_corrected = df.copy()
    df_corrected.loc[(df_corrected.iloc[:, column_index] < lower_bound) |
    (df_corrected.iloc[:, column_index] > upper_bound), df.columns[column_index]] = new_value

    return df_corrected

# Load the CSV file into a pandas DataFrame
input_filename = "example11.csv" # Change this to the desired input filename
df = pd.read_csv(input_filename)

# Choose the column index you want to perform outlier detection on
column_index = 1 # Change this to the desired column index

# Choose the value to replace outliers with
```

```

chosen_value = 999 # Change this to the desired value

# Replace outliers with the chosen value in the specified column
df_corrected = replace_outliers_with_value(df, column_index, chosen_value)

# Extract the basename of the original filename without the extension
filename, extension = os.path.splitext(input_filename)
basename = os.path.basename(filename)

# Construct the output filename with the suffix "_outliers_replaced_with_{chosen_value}"
output_file = basename + "_outliers_replaced_with_" + str(chosen_value) + extension

# Save the corrected DataFrame to a new CSV file with the modified filename
df_corrected.to_csv(output_file, index=False)

print("Corrected DataFrame with outliers replaced by the chosen value has been saved to
'{}'.".format(output_file))

```

Here's how the function works:

- Input:** The function takes three arguments: the Data Frame (**DF**), the index of the column containing outliers (**column_index**), and the value to replace outliers with (**new_value**).
- Outlier Detection:** Using the Interquartile Range (IQR) method, the function identifies outliers in the specified column. It calculates the first quartile (Q1), third quartile (Q3), and the interquartile range (IQR) of the column data.
- Replacing Outliers:** Outliers are replaced with the chosen value (**new_value**). Any value falling below the lower bound or above the upper bound is replaced with the specified value.
- Output:** The corrected Data Frame (**df_corrected**) with outliers replaced by the chosen value is returned.

After replacing outliers, the script saves the corrected Data Frame to a new CSV file. The filename is modified to include information about the chosen value used for outlier replacement.

6.2.4 Handling duplicates

Detect Duplicates:

```
import csv

def find_duplicates(input_file, column_index):
    duplicates = {}
    with open(input_file, 'r', newline='') as infile:
        reader = csv.reader(infile)
        next(reader) # Skip header row if present
        for row_number, row in enumerate(reader, start=1):
            value = row[column_index]
            if value in duplicates:
                duplicates[value].append(row_number)
            else:
                duplicates[value] = [row_number]

    # Filter duplicates
    duplicate_entries = {value: rows for value, rows in duplicates.items() if len(rows) > 1}
    if not duplicate_entries:
        print("No duplicates found.")
        return

    print("Duplicates found:")
    for value, rows in duplicate_entries.items():
        print(f"Value '{value}' duplicated at rows: {', '.join(map(str, rows))}")

input_file = 'output Dots.csv'
column_index = 2 # Specify the index of the column to check for duplicates
find_duplicates(input_file, column_index)
```

The `find_duplicates` function is designed to identify and report duplicate entries within a specified column of a CSV file. Let's delve into how this function operates:

1. Imports: The function imports the `csv` module for handling CSV files.
2. Function Definition: `find_duplicates` is defined with two parameters:
 - `input_file`: The path to the CSV file to be analyzed.
 - `column_index`: The index of the column within the CSV file where duplicates are to be checked.

3. Dictionary Initialization: A dictionary named `duplicates` is initialized. This dictionary will hold the values from the specified column as keys, and the corresponding row numbers where these values occur as values.
4. File Reading and Enumeration: The function opens the input CSV file using a `with` statement and iterates over its rows using a `csv.reader` object. It skips the header row (if present) using `next (reader)` `.
5. Duplicate Identification: For each row in the CSV file, the function retrieves the value from the specified column and checks if it's already present in the `duplicates` dictionary. If it is, the row number is appended to the list associated with that value. Otherwise, a new list is created with the row number.
6. Filtering Duplicate Entries: After processing all rows, the function filters out duplicate entries by retaining only those values in the `duplicates` dictionary which have more than one associated row number. This filtered dictionary is named `duplicate_entries`.
7. Output: If no duplicates are found, the function prints a message indicating so. Otherwise, it prints each duplicate value along with the rows where it occurs.
8. Function Invocation: Finally, the function is invoked with the specified input file and column index.

This function offers a straightforward approach to identify and report duplicate entries within a CSV file, making it useful for data validation and cleaning tasks.

Remove Duplicates:

```
import csv
import os

def remove_duplicates(input_file, column_index):
    unique_values = set()
    output_file = os.path.splitext(input_file)[0] + '_unique.csv'

    with open(input_file, 'r', newline='') as infile:
        reader = csv.reader(infile)
        with open(output_file, 'w', newline='') as outfile:
            writer = csv.writer(outfile)
            for row in reader:
                value = row[column_index]
                if value not in unique_values:
                    writer.writerow(row)
                    unique_values.add(value)

input_file = 'output Dots.csv'
column_index = 0 # Specify the index of the column to check for duplicates
remove_duplicates(input_file, column_index)
```

This function, `remove_duplicates`, serves the purpose of removing duplicate rows from a CSV file based on a specified column index. Here's a breakdown of how it works:

1. Imports: The function imports the necessary modules `csv` and `os` for handling CSV files and file paths, respectively.
2. Function Definition: The function `remove_duplicates` takes two arguments:
 - `input_file`: The path to the CSV file from which duplicates are to be removed.
 - `column_index`: The index of the column in the CSV file that is used to identify duplicates.
3. Set Initialization: A set named `unique_values` is initialized to store unique values encountered in the specified column.
4. Output File Path Determination: The output file path is determined by appending '_unique.csv' to the input file's base name. This ensures that the original file is not overwritten.

5. File handling with `with` Statement: The function uses `with` statements to open the input and output files. This ensures that the files are properly closed after use, even if an exception occurs.

6. CSV Reading and Writing: Inside the first `with` block, the function reads the input CSV file using `csv.reader`. Inside the second `with` block, it writes to the output CSV file using `csv.writer`.

7. Duplicate Removal Logic: For each row in the input CSV file, the function checks the value in the specified column (determined by `column_index`). If the value is not already in `unique_values`, the row is written to the output CSV file, and the value is added to `unique_values`. This ensures that only unique rows are written to the output file.

8. Function Invocation: Finally, the function is invoked with the specified input file and column index.

This function provides a convenient way to remove duplicate rows from a CSV file while retaining the order of appearance and without requiring the entire file to fit into memory at once.

6.3 System Modules Implementation

6.3.1 Missing values

```
@app.route('/Missing', methods=['GET', 'POST'])
def miss():
    filenames=get_uploaded_files()
    if request.method == 'POST':
        operation = request.form['operation']
        filename = request.form['table']
        if filename=="Not select":
            return render_template('Missing_value.html', filenames=filenames,table_original="You haven't selected any Table!",missingPers=0,missingCount=0,totalRows=0)
        else:
            df = pd.read_csv('uploads/'+filename)
            def show():
                missing_Table =find_missing_values_rows(df)
                if missing_Table.count().max()==0:
                    missing_Html_Table="No Missing Value Found.!"
                else:
                    missing_Html_Table = missing_Table.to_html(classes='table table-bordered table-striped sizee', index=False)
                    table_html_original = df.to_html(classes='table table-bordered table-striped sizee', index=False)
                    missingPers,missingCount,totalRows=pers_missing(df)
                return missingPers,missingCount,totalRows,missing_Html_Table,table_html_original
            if (operation=='show'):
                missingPers,missingCount,totalRows,missing_Html_Table,table_html_original=show()
                return render_template('Missing_value.html', table_original=table_html_original, table_cleaned=missing_Html_Table,filenames=filenames,missingPers=missingPers,missingCount=missingCount,totalRows=totalRows)
            elif (operation=='fill'):
                intype = request.form['fill-int-type']
                if intype=="Not select int":
                    return render_template('Missing_value.html',intmsg='Select Option To Fill Integer Missing Value',missingPers=0,missingCount=0,totalRows=0)
                else:
                    strValue = request.form['fill-str-value']
                    if strValue=='':
                        return render_template('Missing_value.html',strmsg='Select Option To Fill String Missing Value',missingPers=0,missingCount=0,totalRows=0)
                    else:
                        newdata=handle_missing(df,intType,strValue)
                        newdata.to_csv("uploads/"+filename,index=False)
                        table_html_cleaned = newdata.to_html(classes='table table-bordered table-striped ', index=False)
                        missingPers,missingCount,totalRows,missing_Html_Table,table_html_original=show()
                        return render_template('Missing_value.html', table_original=table_html_original, table_cleaned=table_html_cleaned,filenames=filenames,missingPers=missingPers,missingCount=missingCount,totalRows=totalRows)
            elif operation=='delete' :
                df=df.dropna()
                df.to_csv('uploads/'+filename)
                missingPers,missingCount,totalRows=pers_missing(df)
                table_html_original = df.to_html(classes='table table-bordered table-striped ', index=False)
                missingPers,missingCount,totalRows,missing_Html_Table,table_html_original=show()
                return render_template('Missing_value.html', table_original=table_html_original,table_cleaned="No Missing Value Found.!", filenames=filenames,missingPers=missingPers,missingCount=missingCount,totalRows=totalRows)
            else:
                return render_template('Missing_value.html', filenames=filenames,missingPers=0,missingCount=0,totalRows=0)
```

1. Routing and Request Handling:

- The Flask application defines a route `/Missing` to handle requests related to missing value operations using both `GET` and `POST` methods.

2. View Function miss():

- This function serves as the controller for the `/Missing` route.
- For `POST` requests:
 - It retrieves the list of uploaded filenames using `get_uploaded_files()`.
 - Depending on the operation specified by the user (`show`, `fill`, `delete`), it performs corresponding actions on the uploaded CSV files:
 - If the operation is `show`, it displays missing values in the CSV file along with basic statistics.
 - If the operation is `fill`, it handles missing values based on user-defined criteria and updates the CSV file accordingly.

- If the operation is `delete`, it removes rows with missing values from the CSV file.
- For `GET` requests, it renders the initial page for uploading files.

3. Data Processing Functions:

- The function renders the HTML template `Missing_value.html`, passing relevant data such as filenames, original and modified table HTML, and missing value statistics.
- The template dynamically adjusts its content based on user interactions and the operation performed.

4. Data Processing Functions:

- The code utilizes several helper functions such as `get_uploaded_files()`, `find_missing_values_rows()`, `pers_missing()`, and `handle_missing()` to process data and perform operations related to missing values in CSV files.

5. Input Validation and Error Handling:

- Basic input validation is implemented to ensure that users select a file and provide necessary inputs for filling missing values.
- Error messages are displayed when required inputs are missing or invalid.

6. File Operations:

- CSV files uploaded by users are read into pandas Data Frames for processing.
- Modified Data Frames are saved back to CSV files in the `uploads` directory.

7. User Interface:

- The system provides a user-friendly interface for uploading files, selecting operations, and viewing results.
- Feedback messages and error alerts are displayed to guide users through the process.

Overall, this implementation facilitates the exploration and manipulation of missing values in CSV files through a web-based interface, enhancing usability and accessibility for users.

6.3.2 Duplicates

```
@app.route('/duplicate', methods=['POST', 'GET'])
def duplicate():
    filenames=get_uploaded_files()
    if request.method == 'POST':
        operation = request.form['operation']
        filename = request.form['table']
        if filename=="Not select":
            return render_template('duplicate.html', filenames=filenames,table_original="You haven't selected any Table!",pers=0,effectedRows=0,totalRows=0)
        else:
            def show_dupl():
                selected_columns = request.form.getlist('columns')
                filename = request.form["table"]
                df = pd.read_csv('uploads/'+filename)
                pers,effectedRows,totalRows=pers_duplicate(df)
                columns=get_column_names(df)
                if selected_columns:
                    show_dup=find_duplicate_rows(df,selected_columns)
                    pers,effectedRows,totalRows=pers_duplicate(df,selected_columns)
                else:
                    show_dup=find_duplicate_rows(df,'all')

                if show_dup.count().max() ==0:
                    table_html_duplicate='No Duplicates Found.'
                else:
                    table_html_duplicate = show_dup.to_html(classes='table table-striped table-bordered sizee ', index=False)
                    table_html_original = df.to_html(classes='table table-striped table-bordered sizee ', index=False)
                return pers,effectedRows,totalRows,table_html_original,table_html_duplicate,columns,selected_columns
            if operation == 'show':
                pers,effectedRows,totalRows,table_html_original,table_html_duplicate,columns,selected_columns=show_dupl()
                return render_template('duplicate.html',filenames=filenames,table_original=table_html_original,table_duplicate=table_html_duplicate,pers=pers,effectedRows=effectedRows,totalRows=totalRows,columns=columns,selected_columns=selected_columns)
            else:
                try:
                    selected_columns = request.form.getlist('columns')
                    if selected_columns:
                        remove_duplicates_and_save(filename,selected_columns)
                    pers,effectedRows,totalRows,table_html_original,table_html_duplicate,columns,selected_columns=show_dupl()
                    return render_template('duplicate.html',filenames=filenames,table_original=table_html_original,table_duplicate=table_html_duplicate,pers=pers,effectedRows=effectedRows,totalRows=totalRows)
                except:
                    pers,effectedRows,totalRows,table_html_original,table_html_duplicate=show_dupl()
                    return render_template('duplicate.html',filenames=filenames,table_original=table_html_original,table_duplicate=table_html_duplicate,pers=pers,effectedRows=effectedRows,totalRows=totalRows)

    return render_template('duplicate.html',filenames=filenames,pers=0,effectedRows=0,totalRows=0)
```

1. Route Definition:

- The Flask application defines a route `/duplicate` to handle requests related to duplicate value operations using both `GET` and `POST` methods

2. View Function `duplicate ()`:

- This function acts as the controller for the `/duplicate` route.
- For `POST` requests:
 - It retrieves the list of uploaded filenames using `get_uploaded_files()`.
 - Depending on the operation specified by the user (`show`, `remove`), it carries out corresponding actions on the selected CSV file.
- For `GET` requests, it renders the initial page for selecting files and operations.

3. Data Processing Function `show_dupl()`:

- This function, nested within `duplicate()`, processes data related to duplicate values.
- It reads the selected CSV file into a Data Frame (`DF`) and calculates statistics such as the percentage of duplicate rows and the number of affected rows.

- Depending on the selected columns, it identifies duplicate rows using `find_duplicate_rows()` function.
- It generates HTML tables for displaying the original and duplicate rows.

4. Template Rendering:

- The function renders the HTML template `duplicate.html`, passing relevant data such as filenames, original and duplicate table HTML, statistics, columns, and selected columns.
- The template dynamically adjusts its content based on user interactions and the operation performed.

5. Input Validation and Error Handling:

- Basic input validation ensures that users select a file and provide necessary inputs for identifying and removing duplicate values.
- Error handling manages exceptions that may occur during file processing.

6. File Operations:

- CSV files uploaded by users are read into pandas Data Frames for processing.
- Modified Data Frames are saved back to CSV files in the `uploads` directory after removing duplicate rows.

7. User Interface:

- The system offers a user-friendly interface for selecting files, columns, and operations related to duplicate values.
- It provides feedback messages and error alerts to assist users throughout the process.

This implementation provides a convenient web-based interface for users to explore and manage duplicate values in CSV files effectively, enhancing usability and accessibility.

6.3.3 Outliers

```
@app.route('/outlier',methods=['GET','POST'])
def outlier():
    filenames=get_uploaded_files()
    if request.method == 'POST':
        filename = request.form['table']
        operation = request.form['operation']
        def show_outlier(filename):
            df = pd.read_csv('uploads/'+filename)
            columns=get_column_names(df)
            table_html_original = df.to_html(classes='table table-striped table-bordered ', index=False)
            table_html_original = df.to_html(classes='table table-striped table-bordered ', index=False)
            outliers,total_rows,num_outliers,outlier_percentage=getOutlierDetails(selected_columns[0],df)
            if outliers.count().max() ==0:
                return render_template('outlier.html',filenames=filenames,table_original=table_html_original,table_outliers="No Outlier Found!",total_rows=total_rows,num_outliers=num_outliers,outlier_percentage=outlier_percentage,columns=columns)
            else:
                table_html_outliers = outliers.to_html(classes='table table-striped table-bordered ', index=False)
                return render_template('outlier.html',filenames=filenames,table_original=table_html_original,table_outliers=table_html_outliers,total_rows=total_rows,num_outliers=num_outliers,outlier_percentage=outlier_percentage,columns=columns)
        if filename=='Select':
            return render_template('outlier.html',filenames=filenames,table_original="You Should Select Table !",total_rows=0,num_outliers=0,outlier_percentage=0)
        else:
            if operation == 'show':
                df = pd.read_csv('uploads/'+filename)
                columns=get_column_names(df)
                table_html_original = df.to_html(classes='table table-striped table-bordered ', index=False)
                return render_template('outlier.html',filenames=filenames,table_original=table_html_original,total_rows=0,num_outliers=0,outlier_percentage=0,columns=columns)

            elif operation == 'find_outlier':
                df = pd.read_csv('uploads/'+filename)
                columns=get_column_names(df)
                chart = request.form['chart']
                selected_columns = request.form.getlist('columns')
                save_and_display_plot(df,selected_columns[0],chart)

                table_html_original = df.to_html(classes='table table-striped table-bordered ', index=False)
                if selected_columns:
                    outliers,total_rows,num_outliers,outlier_percentage=getOutliersDetails(selected_columns[0],df)

                if outliers.count().max() ==0:
                    return render_template('outlier.html',filenames=filenames,table_original=table_html_original,table_outliers="No Outlier Found!",total_rows=total_rows,num_outliers=num_outliers,outlier_percentage=outlier_percentage,columns=columns)
                else:
                    table_html_outliers = outliers.to_html(classes='table table-striped table-bordered ', index=False)
                    return render_template('outlier.html',filenames=filenames,table_original=table_html_original,table_outliers=table_html_outliers,total_rows=total_rows,num_outliers=num_outliers,outlier_percentage=outlier_percentage,columns=columns)
            else:
                return render_template('outlier.html',filenames=filenames,table_original="You Don't Select Any Row!",total_rows=0,num_outliers=0,outlier_percentage=0,columns=columns)

        elif operation == 'negative':
            df = pd.read_csv('uploads/'+filename)
            columns=get_column_names(df)
            selected_columns = request.form.getlist('columns')
            table_html_original = df.to_html(classes='table table-striped table-bordered ', index=False)
            if selected_columns:
                outliers,total_rows,num_outliers,outlier_percentage=getNegative(selected_columns[0],df)
            if outliers.count().max() ==0:
                return render_template('outlier.html',filenames=filenames,table_original=table_html_original,table_outliers="No Negative Value Found!",total_rows=total_rows,num_outliers=num_outliers,outlier_percentage=outlier_percentage,columns=columns)
            else:
                return render_template('outlier.html',filenames=filenames,table_original=table_html_original,table_outliers="You Don't Select Any Row!",total_rows=0,num_outliers=0,outlier_percentage=0,columns=columns)

        elif operation == 'delete_outlier':
            selected_columns = request.form.getlist('columns')
            df = pd.read_csv('uploads/'+filename)
            columns=get_column_names(df)
            if selected_columns:
                remove_outliers_and_save(selected_columns[0],'uploads/'+filename)
            return show_outlier(filename)

        elif operation == 'delete_negative':
            selected_columns = request.form.getlist('columns')
            df = pd.read_csv('uploads/'+filename)
            columns=get_column_names(df)
            if selected_columns:
                remove_negative_rows(selected_columns[0],'uploads/'+filename)
            return show_outlier(filename)

    else:
        return render_template('outlier.html',filenames=filenames,table_original="You Should Select Table !",total_rows=0,num_outliers=0,outlier_percentage=0)
```

1. Route Definition:

- The route `/outlier` is defined using the `@app.route()` decorator in Flask.
- The route handles both `GET` and `POST` requests, allowing users to interact with the system through form submissions.

2. View Function `outlier()`:

- This function is responsible for processing user requests related to outlier detection.
- Upon receiving a `POST` request, it retrieves the list of uploaded filenames using the `get_uploaded_files()` function.
- It extracts the selected filename and operation from the form data.
- Depending on the operation specified (`show`, `find_outlier`, `negative`, `delete_outlier`, `delete_negative`), it performs corresponding actions.
- For `GET` requests, it renders the initial page for selecting files and operations.

3. Data Processing Function `show_outlier():`

- This nested function handles the core data processing tasks related to outlier detection.
- It reads the selected CSV file into a pandas Data Frame (`DF`).
- It generates HTML tables for displaying the original data and outlier details.
- Depending on the operation, it calculates outlier details such as the number of outliers, total rows, and outlier percentage using appropriate functions (`getOutliersDetails()`, `getNegative()`).
-

4. Template Rendering:

- The function renders the HTML template `outlier.html` to display the user interface.
- It passes relevant data such as filenames, original and outlier table HTML, statistics, and column names to the template for rendering.
- The template is designed to dynamically update its content based on user interactions and the operation performed.

5. Input Validation and Error Handling:

- Basic input validation is implemented to ensure that users select a file and provide necessary inputs for outlier detection.
- Error handling is incorporated to manage exceptions that may occur during file processing or data manipulation.

6. File Operations:

- Uploaded CSV files are read into pandas Data Frames for processing.
- Modified Data Frames, after outlier removal, are saved back to CSV files in the `uploads` directory.

- Functions like `remove_outliers_and_save()` and `remove_negative_rows()` handle the removal of outliers or negative values from the Data Frame and save the modified data back to files.

7. User Interface:

- The system provides a user-friendly interface for users to select files, columns, and operations related to outlier detection.
- It includes interactive elements such as dropdown menus, checkboxes, and buttons to facilitate user interactions.
- Feedback messages and error alerts are displayed to guide users through the process and provide status updates.

By combining these components, the system offers a comprehensive and intuitive web-based solution for users to detect, analyze, and manage outliers in their CSV datasets.

6.3.4 Invalid formats

```
#app.route('/invalid', methods=['POST', 'GET'])
def invalid():
    filename = get_uploaded_file()
    if request.method == 'POST':
        operation = request.form['operation']
        filename = request.form['table']
        def show_dupl():
            filename = request.form['table']
            df = pd.read_csv('uploads/' + filename)
            columnnames_get_column_names(df)
            table_html_original = df.to_html(classes='table table-striped table-bordered size', index=False)
            return render_template('invalid_format.html', filenames=filenames, table_original=table_html_original, columns=columns)
        if filename=='Not select':
            return render_template('invalid_format.html', filenames=filenames, table_original="You haven't selected any Table!", pern=0, affectedRows=0, totalRows=0)
        else:
            if operation == 'show':
                return show_dupl()
            if operation == 'apply':
                Remove_operation = request.form['Remove']
                Standardize = request.form['Standardize']
                Split = request.form['Split']
                check = request.form['check']
                selected_columns = request.form.getlist('columns')
                if not (Remove_operation==''):
                    if Remove_operation=='Comma':
                        Remove.remove_commas_from_csv(filename,selected_columns)
                        return show_dupl()
                    elif Remove_operation=='Data':
                        Remove.remove_data_from_csv(filename,selected_columns)
                        return show_dupl()
                    elif Remove_operation=='Strip':
                        Remove.remove_strip(filename)
                        return show_dupl()
                    elif Remove_operation=='WhiteSpace':
                        Remove.remove_white_spaces_in_csv(filename)
                        return show_dupl()
                    elif Remove_operation=='char':
                        selected_columns = request.form.getlist('columns')
                        char = request.form['char']
                        Remove.delete_character_from_csv(filename,selected_columns[0],char)
                        return show_dupl()
                    elif Remove_operation=='num':
                        selected_columns = request.form.getlist('columns')
                        Remove.remove_numbers_and_writes_by_column_name(filename,selected_columns)
                        return show_dupl()
                if not (Standardize==''):
                    if Standardize=='lower':
                        xt.convert_csv_columns_to_lowercase(filename,selected_columns)
                        return show_dupl()
                    elif Standardize=='upper':
                        xt.convert_csv_to_uppercase(filename,selected_columns)
                        return show_dupl()
                    elif Standardize=='Proper':
                        xt.convert_csv_columns_to_propercase(filename,selected_columns)
                        return show_dupl()
                if not (Split==''):
                    first_col = request.form['first_col']
                    second_col = request.form['second_col']
                    selected_columns = request.form.getlist('columns')

                    if Split=='slash':
                        sp.split_name_column(filename,selected_columns[0],first_col,second_col,'/')
                        return show_dupl()
                    elif Split=='and':
                        sp.split_name_column(filename,selected_columns[0],first_col,second_col,'&')
                        return show_dupl()
                    elif Split=='dot':
                        sp.split_name_column(filename,selected_columns[0],first_col,second_col,'.')
                        return show_dupl()
                    elif Split=='commas':
                        sp.split_name_column(filename,selected_columns[0],first_col,second_col,',')
                        return show_dupl()
                    elif Split=='specific':
                        separator = request.form['separator']
                        sp.split_name_column(filename,selected_columns[0],first_col,second_col,sign=separator)
                        return show_dupl()

                if not (check==''):
                    if check=='Email':
                        filename = request.form['table']
                        df = pd.read_csv('uploads/' + filename)
                        columnnames_get_column_names(df)
                        invalid_rows = get_invalid_emails(filename,selected_columns[0])
                        table_html_original = df.to_html(classes='table table-striped table-bordered size', index=False)
                        if invalid_rows.count().max() != 0:
                            table_cleaned = invalid_rows.to_html(classes='table table-striped table-bordered size', index=False)
                            return render_template('invalid_format.html', filenames=filenames, table_original=table_html_original, table_cleaned=table_cleaned, columns=columns)
                        else:
                            return render_template('invalid_format.html', filenames=filenames, table_original=table_html_original, table_cleaned="No Invalid Rows Found!", columns=columns)
                    elif check=='Phone':
                        filename = request.form['table']
                        df = pd.read_csv('uploads/' + filename)
                        columnnames_get_column_names(df)
                        invalid_rows = check_phone_numbers_in_column(filename,selected_columns[0])
                        table_html_original = df.to_html(classes='table table-striped table-bordered size', index=False)
                        if invalid_rows.count().max() != 0:
                            table_cleaned = invalid_rows.to_html(classes='table table-striped table-bordered size', index=False)
                            return render_template('invalid_format.html', filenames=filenames, table_original=table_html_original, table_cleaned=table_cleaned, columns=columns)
                        else:
                            return render_template('invalid_format.html', filenames=filenames, table_original=table_html_original, table_cleaned="No Invalid Rows Found!", columns=columns)
    return render_template('invalid_format.html', filenames=filenames)
```

1. Route Definition:

- The Flask application defines a route `/invalid` to handle requests related to data validation and formatting.

2. View Function `invalid()`:

- This function serves as the controller for the `/invalid` route.
- It retrieves the list of uploaded filenames using the `get_uploaded_files()` function.
- For POST requests:
 - It extracts the operation and filename from the form data submitted by the user.
 - Depending on the selected operation (`show`, `apply`), it performs corresponding actions.
- For GET requests, it renders the initial page for selecting files and operations.

3. Data Processing Function `show_dupl()`:

- This nested function handles the core data processing tasks related to data validation and formatting.
- It reads the selected CSV file into a DataFrame (`DF`).
- It generates HTML tables for displaying the original data.
- It is used to render the `invalid_format.html` template.

4. Template Rendering:

- The function renders the HTML template `invalid_format.html`, passing relevant data such as filenames, original table HTML, columns, and cleaned data (if applicable).
- The template dynamically adjusts its content based on user interactions and the operation performed.

5. Data Formatting Operations:

- Various data formatting operations are supported, including removing characters (e.g., commas, dots, whitespace), converting text case (lowercase, uppercase, proper case), and splitting columns based on separators (e.g., slash, dot, comma).
- The user can choose specific columns and apply formatting operations accordingly.
 1. Data Validation Operations:
- Data validation operations such as checking for invalid email addresses or phone numbers in specific columns are supported.
- If invalid data is found, the corresponding rows are displayed in a separate table to inform the user.

6. Input Validation and Error Handling:

- Basic input validation is implemented to ensure that users select a file and provide necessary inputs for data validation and formatting.
- Error handling manages exceptions that may occur during file processing or data manipulation.

7. User Interface:

- The system provides a user-friendly interface for users to select files, columns, and operations related to data validation and formatting.
- Interactive elements such as dropdown menus, checkboxes, and buttons facilitate user interactions.
- Feedback messages and error alerts guide users through the process and provide status updates.

This implementation offers a comprehensive web-based solution for users to validate and format data in CSV files, enhancing data quality and consistency.

6.4 system integration

Introduction to System Integration with Flask and Jinja in an MVC Architecture:

In the dynamic landscape of web development, effective system integration is fundamental to building robust and scalable applications. Harnessing the power of frameworks like Flask and templating engines such as Jinja, developers can seamlessly integrate various components within the Model-View-Controller (MVC) architecture. This introduction sets the stage for exploring how Flask, Jinja, and MVC intertwine to create cohesive web applications.

Flask Framework:

Flask provides a simple yet powerful framework for developing web applications in Python. Its minimalistic design allows developers to quickly set up and extend functionality as needed. Flask offers features such as URL routing, request handling, and template rendering, making it an ideal choice for building web applications of varying complexities.

```
app=Flask(__name__)

UPLOAD_FOLDER = 'uploads'
app.config['UPLOAD_FOLDER'] = UPLOAD_FOLDER
|
UPLOAD_FOLDER = 'uploads'
app.config['UPLOAD_FOLDER'] = UPLOAD_FOLDER

@app.route('/')
def Auth():
    return render_template('Auth.html')
```

Jinja:

Jinja is a modern and designer-friendly templating engine for Python, inspired by Django's template system. It allows developers to create dynamic web pages by embedding Python code within HTML templates. Jinja provides features like template inheritance, macros, filters, and control structures, enabling developers to create reusable and maintainable templates.

```
<select name="columns" class="form-control mt-3 fv py-5 Table-name" id="column-name" multiple>
  <option disabled value="">Select Which Column To Clean</option>
  <option value="all">all</option>

  {% for column in columns %}
  <option value={{column}}>{{ column }}</option>
  {% endfor %}
  |
</select>
```

```
{% if table_original=="You haven't selected any Table!" %}
<p style="margin-top: 3%;" class="alert alert-danger mx-auto w-100" role="alert">{{table_original}}</p>

{% elif table_original %}
<div class="card">
  <div class="card-body" style="max-height: 400px; overflow-y: scroll;">
    {{ table_original|safe }}
  </div>
</div>

{% endif %}
```

Template Rendering with Jinja:

Jinja templates enable developers to generate dynamic HTML content by embedding Python code within HTML markup. This allows for the seamless integration of data processed by Flask into the presentation layer of the application. Developers can use Jinja's features like template inheritance and macros to create modular and maintainable templates.

```
{% extends "layout.html" %}
```

Flask Structure Files:

1- `__pycache__`: This is a hidden folder that Python creates to store compiled bytecode files to speed up program startup.

2- `Img`: This folder likely holds images used in the project.

3- `Static`: This folder typically stores static files such as CSS, JavaScript, or images that are used by the web application.

4- `Templates`: This folder likely stores HTML templates used to generate the web application's user interface.

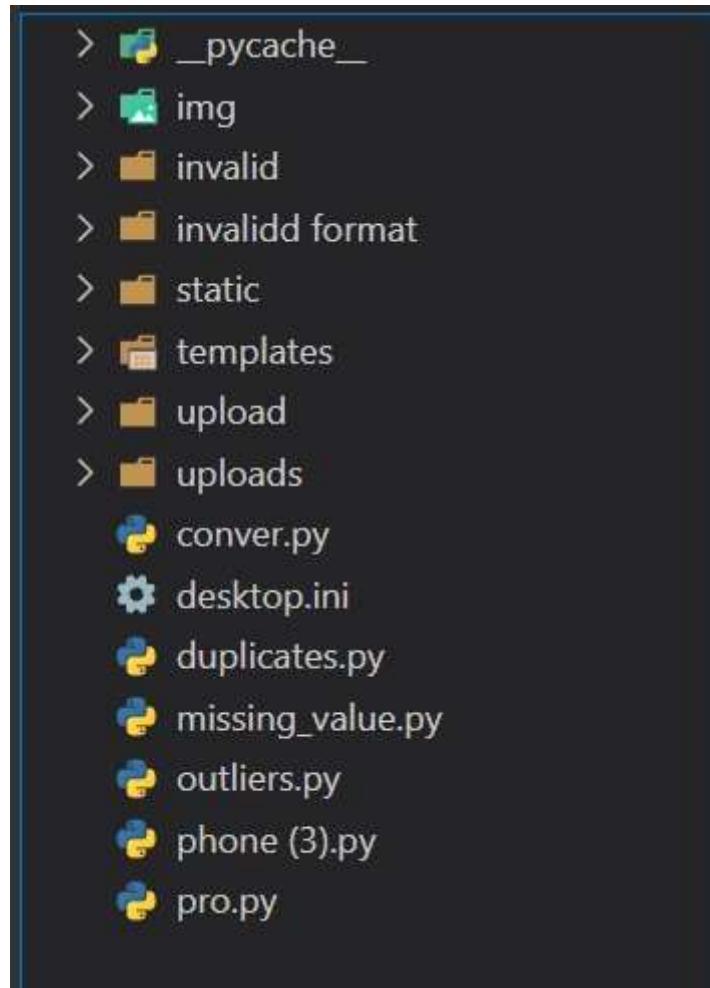
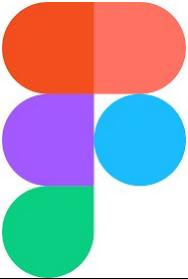


Figure 61 flask structure files

6.5 used softwares

They are the tools required for the system to be built. The usage of those tools is determined by previous studies performed during research on the appropriate tools to be used to achieve the desired results.

#	Name	Description
1	Trello 	Trello is the visual work management tool that empowers teams to ideate, plan, manage, and celebrate their work together in a collaborative, productive, and organized way.
2	Draw.io 	Draw.io is an online diagramming tool with integrations with Jira, google and confluence available free online or at cost depending on integration chosen
3	Postman 	Postman is an API platform for building and using APIs. Postman simplifies each step of the API lifecycle and streamlines collaboration so you can create better APIs—faster.
4	Visual studio code 	Visual Studio is an integrated development environment (IDE) and Visual Studio Code is a rich text editor like Sublime Text and Atom. But the difference between the tools is more than just IDE and text editor. An IDE is a robust tool for writing, editing, debugging, and running your code.
5	Notion 	Notion is a single space where you can think, write, and plan. Capture thoughts, manage projects, or even run an entire company — and do it exactly the way you want.
6	Discord 	Discord is a free communications app that lets you share voice, video, and text chat with friends, game communities, and developers. It has hundreds of millions of users, making it one of the most popular ways to connect with people online.
7	PowerPoint ai 	PowerPoint (PPT) is a powerful, easy-to-use presentation graphics software program that allows you to create professional-looking electronic slide shows.

8	Figma 	Figma design is for people to create, share, and test designs for websites, mobile apps, and other digital products and experiences. It is a popular tool for designers, product managers, writers and developers and helps anyone involved in the design process contribute, give feedback, and make better decisions, faster.
9	Mongo db 	MongoDB Compass is a user-friendly graphical interface for interacting with MongoDB databases. It simplifies tasks like data exploration, querying, and administration, offering features for performance monitoring, schema visualization, and index management. It's designed to streamline database interaction for developers, administrators, and analysts.

4.2.4 Used technologies

#	Name	description
1	Flask API 	Flask is a lightweight web framework for Python. Flask API allows developers to build web APIs using Flask, providing tools and extensions to streamline the process.
2	Fast API 	D is a modern, fast (high-performance) web framework for building APIs with Python 3.7+ based on standard Python type hints. It is designed for high performance and productivity, with automatic interactive API documentation.
3	REST API 	REST (Representational State Transfer) is an architectural style for designing networked applications. REST APIs allow systems to communicate over HTTP by defining a set of operations that correspond to the HTTP methods GET, POST, PUT, DELETE, etc.
4	jinja 2 API 	Jinja2 is a modern and designer-friendly templating language for Python, modelled after Django's templates. It is a text-based template language that renders templates to produce documents.

5	Python 3 	Python is a high-level programming language known for its simplicity and readability. Python 3 is the latest major release of the Python programming language, featuring improvements over Python 2.x.
6	Bootstrap 	Bootstrap is a popular open-source CSS framework for building responsive and mobile-first websites and web applications. It provides pre-designed templates and components for easier development.
7	CSS 	CSS is the acronym of “Cascading Style Sheets”. CSS is a computer language for laying out and structuring web pages (HTML or XML). This language contains coding elements and is composed of these “cascading style sheets” which are equally called CSS files.
8	Csv editor 	A CSV (Comma-Separated Values) editor is a tool used to view, edit, and manipulate CSV files. It typically provides features for importing, exporting, sorting, filtering, and formatting CSV data, making it easier to work with tabular data in spreadsheet-like formats.

System shot screens:-

1-missing values page

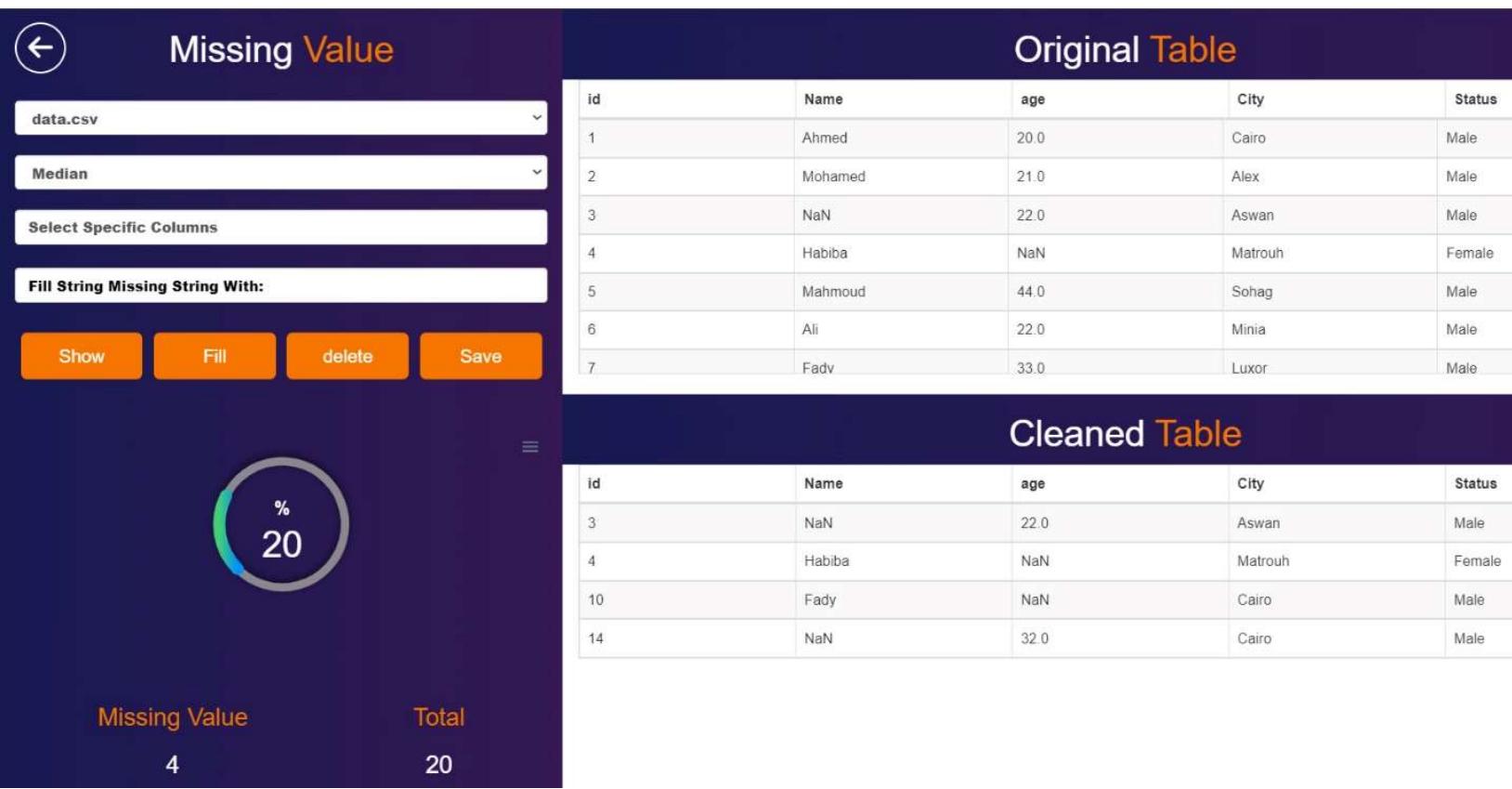


Figure 62 missing values system shotscreen

2-missing values results

id	Name	age	City	Status
1	Ahmed	20.0	Cairo	Male
2	Mohamed	21.0	Alex	Male
3	Hamza	22.0	Aswan	Male
4	Habiba	32.0	Matrouh	Female
5	Mahmoud	44.0	Sohag	Male
6	Ali	22.0	Minia	Male
7	Fady	33.0	Luxor	Male
8	Bassant	32.0	Aswan	Female
9	Nader	40.0	Cairo	Male
10	Fady	32.0	Cairo	Male

Figure 63 missing values results system shotscreen

3-duplicates page

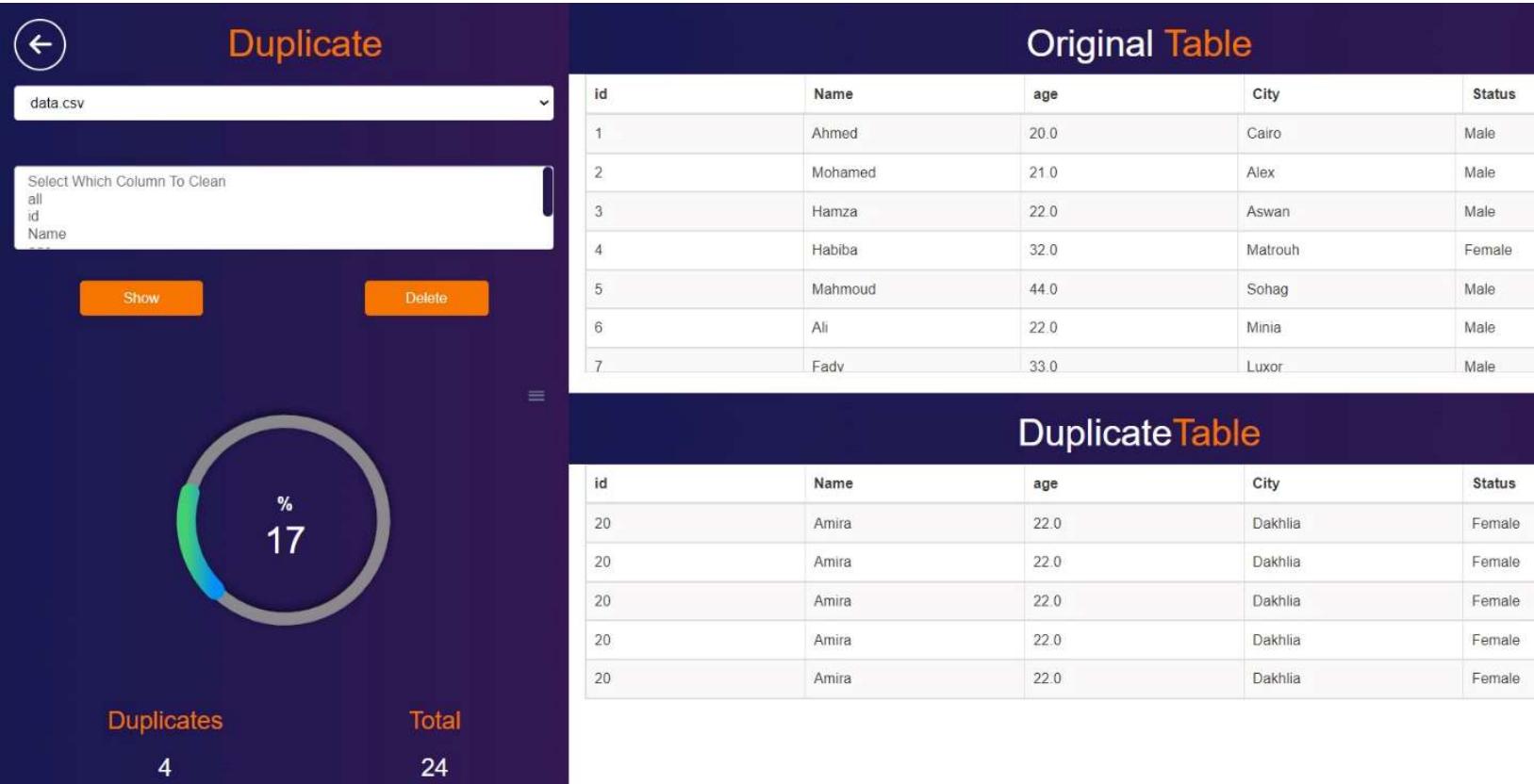


Figure 64 duplicates page system shotscreen

4-duplicates results

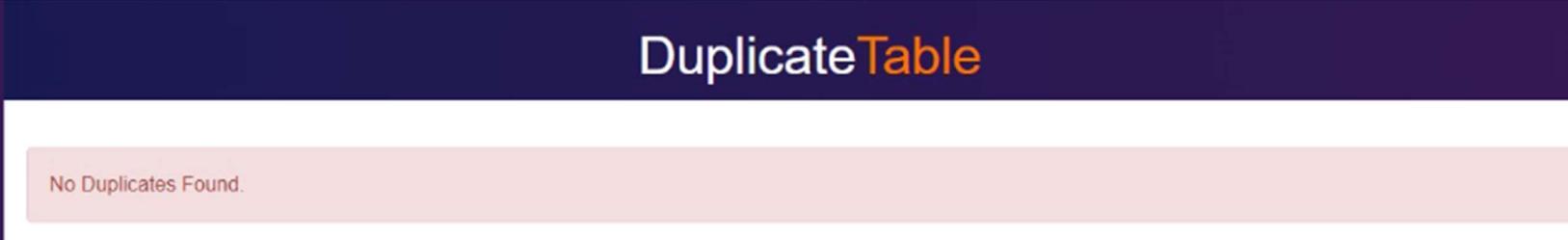


Figure 65 duplicates results system shotscreen

5-outliers page

The screenshot shows the Outliers page interface. At the top, there's a navigation bar with a back arrow and the title "Outliers". Below it are three input fields: "Select Table", "Select Specific Columns", and "Data Visualization". A prominent orange button labeled "Detect Negative Values" is centered. To the right of these controls is a circular progress bar indicating 14%. Below the progress bar, the text "Outlier" is followed by "1" and "Total" followed by "7". At the bottom are four buttons: "Show Table", "Find Outlier", "Delete Outlier", and "Delete Negative".

id	Name	age	City	address	email	phone
0	John mido	21233.0	New York	1.0	midoboy@gmail.com	1144970241
1	Adly ahmed	22.0	Cairo	1.0	midoboy5456@org.com	1120422328
2	Adly Hassan	24.0	Sohag	1.0	midoboy@gmail.com	1153623367
3	Toka Rashad	25.0	Unkown	1.0	midoboy@gmail.com	1157838235
4	Razan Saad	25.0	Elread	1.0	midoboy@gmail.eg	1129073337
5	Mido Hossam	24.0	Adly	1.0	@gmail.com	5588
6	Mido Hossam	29.0	Adly	1.0	@gmail.com	5588

Outlier Table

id	Name	age	City	address	email	phone
0	John mido	21233.0	New York	1.0	midoboy@gmail.com	1144970241

Figure 66 outliers page system shot screen

6-outliers results

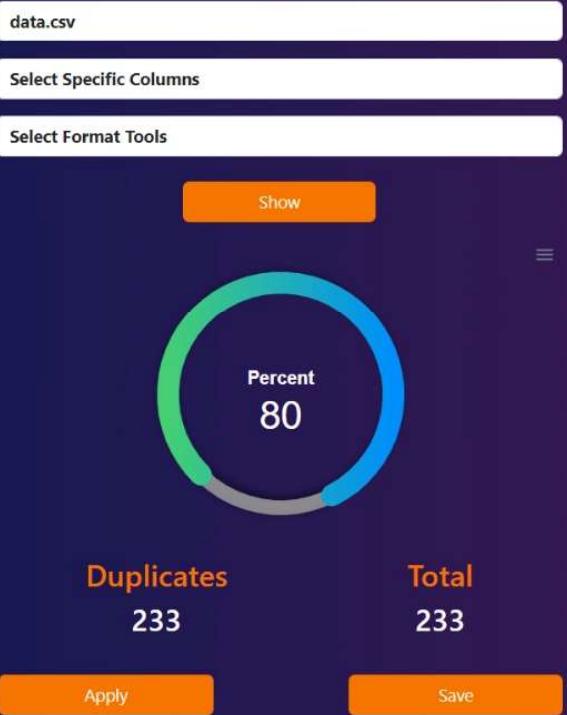
Outlier Table

No Outlier Found!

Figure 67 outliers results system shotscreen

7-Invalid formats page

Invalid Format



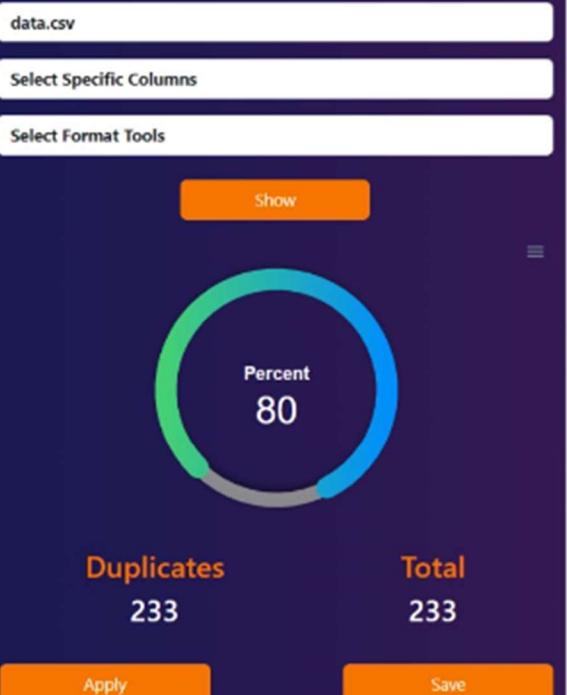
Original Table

	id	Name	age	City	
1	Ahmed mosa	20.0	Cairo	Male	
2	Mohamed Rashad	21.0	Alex	Male	
3	Hamza nady	22.0	Aswan	Male	
4	Habiba mohamed	500.0	Matrouh	Female	
5	Mahmoud hamza	44.0	Sohag	Male	
6	Ali rashad	22.0	Minia	Male	
7	Fady nido	33.0	Luxor	Male	
8	Bassant ahmed	432.0	Aswan	Female	

Figure 68 invalid formats page system shotscreen

8-invalid formats results

Invalid Format



Cleaned Table

	Id	First Name	Last Name	age	City	
1	Ahmed	mosa	20.0	Cairo	Male	
2	Mohamed	Rashad	21.0	Alex	Male	
3	Hamza	nady	22.0	Aswan	Male	
4	Habiba	mohamed	500.0	Matrouh	Female	
5	Mahmoud	hamza	44.0	Sohag	Male	
6	Ali	rashad	22.0	Minia	Male	
7	Fady	nido	33.0	Luxor	Male	
8	Bassant	ahmed	432.0	Aswan	Female	

Figure 69 invalid formats results system shotscreen

Chapter SEVEN

System Testing

Chapter-7-System Testing

7.1 Introduction

The main purpose of system testing is to evaluate and assess the quality of work performed at each step of the system development process. It is a crucial step in ensuring that the system functions as intended and meets the desired requirements. The primary goal of testing is to verify that the system performs its intended functions accurately and reliably. By conducting thorough testing, any defects, errors, or issues in the system can be identified and addressed, thereby improving the overall quality, reliability, and maintainability of the system. Testing helps in validating the system against the specified requirements, ensuring that it meets the expected standards and performs optimally. Additionally, through testing, areas of improvement can be identified, allowing for enhancements and refinements to be made to the system, further enhancing its quality and reliability.

7.2 Types of Testing

There are numbers of types of software testing, categorized by what is being tested and the purpose, or objective, of the test. The objectives range from usability to disaster recovery. Generally, the most common testing types are unit testing, interface testing, functionality was testing, compatibility testing, performance testing, scalability testing, and security testing. A simple definition of the previous testing types is as follows:

- Unit testing: Functional and reliability testing in an Engineering environment. Producing tests for the behavior of components of a product to ensure their correct behavior before system integration.
- Interface testing: done by the user.
- Functionality testing: Validating an application or website conforms to its specifications and correctly accomplishes all its required functions. It may include testing the mathematical and algorithm correctness of scientific and financial software, as well as testing GUI functionality.

- Compatibility testing: To ensure compatibility of an application or website with different browsers, operating systems, and hardware platforms. Compatibility testing can be performed manually or can be driven by an automated functional or regression test suite.
- Performance testing: Performance testing determines how well the software performs in terms of the speed of computations or responsiveness to the user.
- Scalability testing: Scalability testing is performed to ensure that the software will function well as the number of users, size of data sets or other factors change from small to large values.
- Security testing: Application security testing determines how well the software can defend against attacks, such as firewall software securing a computer against 99 Internet viruses and worms. In the case of the web application, it refers to the Testing of the site and web server configuration to eliminate any security or access loopholes.

7.3 Test cases

Test case 1: upload database:-

A	B	C	D
Step	Action	Expected Result	Status
1 Upload With Local Server			
1.1	Click on Local Server Button	Display Select option	Pass ▾
1.2	Select Data Base Name Which want to upload	Display upload button	Pass ▾
1.3	Click on Upload button	Message "Data Successfully Uploaded To MongoDB"	Pass ▾
1.4	Click On Dashboard button in nav	Go to Dashboard Page To start cleaning	Pass ▾
2 Upload With Another Server			
2.1	Click on Another Server Button	Display input field to enter server name	Pass ▾
2.2	Enter Server Name	Display submit button	Pass ▾
2.3	Click on Submit Button	Message "Connection server Success"	Pass ▾
2.4	Click On Dashboard button in nav	Go to Dashboard Page To start cleaning	Pass ▾

Figure 70 testcase 1 Database upload

Test case 2: login:-

A	B	C	D
Step	Action	Expected Result	Status
Preq.	1. User already registered with valid data 2. Internet connection is present		
1	Login with a valid username & a valid password		
1.1	Enter a valid username	Username is entered successfully and the placeholder is removed	Pass ▾
1.2	Enter a valid password	Password is entered successfully and encrypted	Ready to Test ▾
1.3	Click on Log in	User is logged in successfully and redirected to home page	Ready to Test ▾
2	Login with a valid username & a invalid password		
2.1	Enter a valid username	Username is entered successfully and the placeholder is removed	Ready to Test ▾
2.2	Enter an invalid password	Password is entered successfully and encrypted	Ready to Test ▾
2.3	Click on Log in	An error message appears telling the user that his username or password is not correct	Ready to Test ▾
3	Login with an invalid username		
3.1	Enter an invalid username	Username is entered successfully and the placeholder is removed	Ready to Test ▾
3.2	Enter any password	Password is entered successfully and encrypted	Ready to Test ▾
3.3	Click on Log in	An error message appears telling the user that his username or password is not correct	Ready to Test ▾
4	Login with leaving both username and password fields empty		
4.1	Leave username and password fields empty and click on login	An error message appears telling the user that he can't leave the fields empty	Ready to Test ▾

Figure 71 testcase 2 Login

Test case 3: register:-

A	B	C	D
Step	Action	Expected Result	Status
Preq.			
1	User Entered Valid Email, username and Password		
1.1	input username field	Username is entered successfully and the placeholder is removed	Pass ▾
1.2	input Email field	Email is entered successfully and the placeholder is removed	Pass ▾
1.3	input Password field	Password is entered successfully and encrypted	Pass ▾
2	User Entered invalid Password		
2.1	Enter a password	Message "Please Enter Valid Password"	Pass ▾
3	User Entered invalid Username		
3.1	Enter a password	Message "Please Enter Valid Username"	Pass ▾

Figure 72 testcase 3 register

Test case 4: remove duplicates:-

A	B	C	D
Step	Action	Expected Result	Status
Prereq. User Already Choosed a Table			
1 User Want to Delete Duplicates			
1.1 Click in multiple select Specific column and choose option	option is choosed	Pass	▼
1.2 Click on Show Button	Display original and Duplicates table with Duplicated Row rows	Pass	▼
1.3 Click on Delete Button	Display original and Duplicates table without Duplicated Row rows	Pass	▼

Figure 73 testcase 4 remove duplicates

Test case 5: missing values:-

A	B	C	D
Step	Action	Expected Result	Status
Prereq. User Already Choosed a Table			
1 User Want to Fill integer Missing Value			
1.1 Click in fill integer select option	Option is displayed (Mean, Median Mode)	Pass	▼
1.2 Choose Fill integer Option	option is choosed	Pass	▼
1.3 Click in multiple select Specific column and choose option	option is choosed	Pass	▼
1.4 Click on Show Button	Display original and Missing value table with effected rows	Pass	▼
1.5 Click on Fill Button	Display original and Cleaned Table after fill the integers	Pass	▼
1.6 Click on Save Button	Save Data after Cleaned	Pass	▼
2 User Want to Fill String Missing Value			
2.1 Click in multiple select Specific column and choose option	option is choosed	Pass	▼
2.2 Enter String to fill string missing value	String is entered successfully	Pass	▼
2.3 Click on Show Button	Display original and Missing value table with effected rows	Pass	▼
2.4 Click on Fill Button	Display original and Cleaned Table after fill the String	Pass	▼
2.5 Click on Save Button	Save Data after Cleaned	Pass	▼
3 User Want to Delete Missing Value			
3.1 Click in multiple select Specific column and choose option	option is choosed	Pass	▼
3.2 Click on Show Button	Display original and Missing value table with effected rows	Pass	▼
3.3 Click on Delete Button	Display Message "No Missing Value Found!"	Pass	▼

Figure 74 test case 5 missing values

Test case 6: delete outlier:-

A	B	C	D
Step	Action	Expected Result	Status
Prereq. User Already Choosed a Table			
1 User Want to Delete Outlier			
1.1 Click in multiple select Specific column and choose option		option is choosed	Pass ▼
1.2 Select Kind of Data Visualzaton		option is choosed	Pass ▼
1.3 Click on Show Button		Display origanal table	Pass ▼
1.4 Click on Find Outlier Button		Display origanal and Outlier table	Pass ▼
1.5 Click on Delete Outlier Button		Message "Outlier is Deleted"	Pass ▼
2 User Want to Delete Negative Values			
2.1 Click in multiple select Specific column and choose option		option is choosed	Pass ▼
2.2 Select Kind of Data Visualzaton		option is choosed	Pass ▼
2.3 Click on Show Button		Display origanal table	Pass ▼
2.4 Click on Detect Negative Values Button		Display origanal and Negative Values Table	Pass ▼
2.5 Click on Delete Negative Values Button		Message "Negative Values is Deleted"	Pass ▼

Figure 75 test case 6 delete outlier

Test case 7: Invalid Formats:--

A	B	C	D
Step	Action	Expected Result	Status
Prereq: User Already Choosed a Table			
1 User Want to Remove (Comma,Dots,White Space)			
1.1 Click in multiple select Specific column and choose option	option is choosed		Pass
1.2 Select option From Remove Options	option is choosed		Pass
1.3 Click On Apply Button	Message "Updated Format Done"		Pass
2 User Want to Remove (Character,Number)			
2.1 Click in multiple select Specific column and choose option	option is choosed		Pass
2.2 Select option From Remove Options (Character,Number)	option is choosed		Pass
2.3 input Character Field	Input Entered Successfully		Pass
2.3 input Number Field	Input Entered Successfully		Fail
2.4 Click On Apply Button	Message "Updated Format Done"		Pass
3 User Want to Standardize Column			
3.1 Click in multiple select Specific column and choose option	option is choosed		Pass
3.2 Select option From Standardize Options	option is choosed		Pass
3.3 Click On Apply Button	Display Message "Updated Format Done"		Pass
4 User Want to Check (Phone,Email)			
4.1 Click in multiple select Specific column and choose option	option is choosed		Pass
4.2 Select option From Check Options (Phone,Email)	option is choosed		Pass
4.3 Click On Apply Button	Display Valid Format Table		Pass
4.4 Click On Save Button	Display Message "Updated Format Done"		Pass
5 User Want to Split Column			
5.1 Click in multiple select Specific column and choose option	option is choosed		Pass
5.2 Select Specific Separator from Split options	option is choosed		Pass
5.3 input First Column Name Field	Input Entered Successfully		Pass
5.4 input Second Column Name Field	Input Entered Successfully		Pass
5.5 Click On Apply Button	Display Valid Format Table		Pass
5.6 Click On Save Button	Display Message "Updated Format Done"		Pass
6 User Want to Merge Column			
6.1 Click in multiple select Specific column and choose 2 columns To merge	option is choosed		Pass
6.2 Check Merge Column Input	Input is Checked		Pass
6.3 input New Column Name	Input Entered Successfully		Pass
6.4 Click On Apply Button	Display Valid Format Table		Pass
6.5 Click On Save Button	Display Message "Updated Format Done"		Pass
7 User Want to Convert Date			
7.1 Click in multiple select Specific column and choose option	option is choosed		Pass
7.2 Check Convert Date Column Input	Input is Checked		Pass
7.3 Click On Apply Button	Convert Select Column From String Format To Date Format		Pass
7.4 Click On Save Button	Display Message "Updated Format Done"		Pass

Figure 76 testcase 7 invalid formats

Chapter EIGHT

Conclusion & Future Work

Chapter-8- Conclusion

8.1 introduction

In this chapter the project has been created and tested in terms of its functionalities and the overall use of the application. The system achieved results will be displayed as what is concluded from the creation process. Also, a comparison of our project and other related work will be shown. In addition, some extra features will be given as future updates to the final project.

8.2 Conclusion

1. Our data cleansing website, integrated with Flask and developed using HTML, CSS, and Python, marks a significant leap forward in addressing data quality challenges comprehensively.
2. By prioritizing critical issues like missing values, invalid formats, outliers, and duplicates, our approach ensures a targeted cleansing process that effectively resolves key data quality challenges.
3. The platform boasts an intuitive user interface and detailed graphical representations, empowering users to effortlessly clean and visualize their databases while making informed decisions.
4. Offering flexibility in exporting data across various formats, the platform also provides comprehensive training modules to equip users with effective database management skills.
5. Our cleansing process preserves the original structure of the database, maintaining the format of the data throughout.
6. Guided by Agile methodologies, our development process is iterative and responsive, allowing us to continually refine and enhance the platform to meet evolving user needs effectively.
7. Looking ahead, the platform remains poised for future enhancements, promising even greater efficiency and responsiveness to evolving user needs.
8. User feedback plays a crucial role in shaping the evolution of the platform, ensuring that it continues to meet the diverse needs of its user base.

8.3 Future work

1-Advanced Machine Learning Techniques:

- **Automated Error Detection:** Implement machine learning algorithms to automatically detect and suggest corrections for common data errors, such as typos, duplicates, and inconsistencies.
- **Natural Language Processing (NLP):** Use NLP to clean and standardize text data, especially useful for processing unstructured datalike reviews or comments

2. Enhanced User Interface and Experience:

- **Drag-and-Drop Functionality:** Simplify the user experience by enabling drag-and-drop capabilities for data import and export
- **Customization Options:** Allow users to customize cleaning rules and workflows to better fit their specific needs.
-

3. Integration and Interoperability:

- **Cloud Storage Support:** Enable integration with popular cloud storage solutions like AWS S3, Google Drive, and Dropbox for easy data import and export.

4. Security:

- **Data Encryption:** Ensure that all data, both in transit and at rest, is encrypted to protect sensitive information.

5. Enhanced Reporting and Monitoring:

- **Audit Logs:** Maintain detailed logs of all data cleaning activities for accountability and tracking purposes.
- **Custom Reports:** Provide tools for generating custom reports that summarize the results of data cleaning operations.

5. Training and Support:

- **In-App Tutorials:** Develop comprehensive in-app tutorials and guides to help users understand how to use the application effectively
- **Community and Support Forum:** Provide tools for generating custom reports that summarize the results of data cleaning operations.

8.4 Appendix:-

1st raw backend implementation	157
1 st check emails	157
2 nd check phone	158
3 rd lower cases.....	159
4 th upper cases	159
5 th Proper case	160
6 th Merge.....	161
7 th split.....	162
8 th remove dots.....	163
9 th remove letters.....	164
10 th remove numbers.....	165
11 th remove selected from column.....	166
12 th remove double spaces	167
13 th strip	168
14 th show repeated in all tables.....	169
15 th show repeated on one column.....	170
16 th handling missing values	171
17 th <i>Detecting outliers</i>	172
18 th Remove outlier values from a specified column:	174
19 th Replace outlier values in a specified column with the average of non-outlier values within that column:	175
20 th Replace outlier values in a specified column of with a chosen value:.....	176
21 st Detect Duplicates:	177
22 nd Remove Duplicates:.....	178
2nd System Modules Implementation.....	179
1 st Missing values	179
2 nd Duplicates	179
3 rd Invalid formats.....	180
3rd Python:-	181
1 st __Init__	181
2 nd Configure	185
3 rd Routs	185
4 th Models	185

5 th Convert.....	185
6 th Duplicates	187
7 th Missing values.....	188
8 th Outliers.....	188
9 th Run	191
4th system integration.....	192
Flask Framework:	192
Jinja	192
Template Rendering with Jinja:	193
Flask Structure Files:	193
8.5 References:-	194

1st raw backend implementation

1st check emails

```
1 import csv
2 import re
3 import dns.resolver
4 def is_valid_email(email):
5     # Adjusted regular expression pattern for email validation
6     email_pattern = re.compile(r'^[a-zA-Z0-9._%+-]+@[a-zA-Z0-9.-]+\.[a-zA-Z]{2,}$')
7
8     # Match the email pattern
9     if not email_pattern.match(email):
10         return False # Email pattern does not match
11
12     # Extract domain from the email address
13     domain = email.split('@')[1]
14     try:
15         # Resolve MX records for the domain
16         mx_records = dns.resolver.resolve(domain, 'MX')
17         return True # Domain has valid MX records
18     except dns.resolver.NoAnswer:
19         print(f"No MX records found for the domain: {domain}")
20         return False # No MX records found for the domain
21     except dns.resolver.NXDOMAIN:
22         print(f"Domain does not exist: {domain}")
23         return False # Domain does not exist
24     except dns.resolver.Timeout:
25         print(f"Timeout occurred while querying DNS servers for domain: {domain}")
26         return False # Timeout occurred while querying DNS servers
27     except Exception as e:
28         print(f"An error occurred while validating email {email}: {e}")
29         return False # Other errors
30
31 def show_valid_emails_only(input_file, email_column_index):
32     with open(input_file, 'r', newline='') as csv_in:
33         reader = csv.reader(csv_in)
34         header = next(reader) # Skip the header
35         print("Email,Is_Valid_Email") # Print header for email and validation status
36
37         for row in reader:
38             email = row[email_column_index] # Choose email address from the specified column index
39             is_valid = is_valid_email(email)
40             print(f"{email},{is_valid}") # Print the email address and its validation status
41
42 # Example usage:
43 input_file = 'Example.csv'
44 chosen_column_index = 2 # Specify the column index for email addresses
45 show_valid_emails_only(input_file, chosen_column_index)
```

2nd check phone

```
1 import csv
2
3 def check_phone_numbers_in_column(csv_file, column_index):
4     invalid_numbers = []
5     with open(csv_file, 'r') as file:
6         reader = csv.reader(file)
7         next(reader) # Skip the first row
8         for row_number, row in enumerate(reader, start=2):
9             if len(row) >= column_index:
10                 item = row[column_index - 1] # Adjust index to 0-based
11                 if len(item) != 11 or not item.startswith(('012', '010', '011', '015')) or not item.isdigit():
12                     invalid_numbers.append(f"Invalid phone number at row {row_number}, column {column_index}: {item}")
13             else:
14                 print(f"Warning: Row {row_number} does not have enough columns.")
15
16     return invalid_numbers
17
18 # Example usage:
19 column_index = 4 # Adjust this to the index of the column you want to check (1-based index)
20 invalid_numbers = check_phone_numbers_in_column('comma.csv', column_index)
21 if invalid_numbers:
22     print("Invalid phone numbers:")
23     for number in invalid_numbers:
24         print(number)
25 else:
26     print("All phone numbers in the specified column are valid.")
```

3rd lower cases

```
1 import csv
2 import os
3
4 def convert_csv_columns_to_lowercase(input_file, columns):
5     output_file = os.path.splitext(input_file)[0] + '_lowercase.csv'
6     with open(input_file, 'r', newline='') as infile:
7         reader = csv.reader(infile)
8         with open(output_file, 'w', newline='') as outfile:
9             writer = csv.writer(outfile)
10            for row in reader:
11                lower_row = [cell.lower() if index in columns else cell for index, cell in enumerate(row)]
12                writer.writerow(lower_row)
13
14 # Example usage:
15 input_file = 'input.csv'
16 selected_columns = [0, 2] # Selecting columns 1 and 3 (0-based index)
17 convert_csv_columns_to_lowercase(input_file, selected_columns)
18
```

4th upper cases

```
1 import csv
2 import os
3
4 def convert_csv_to_uppercase(input_file, columns):
5     output_file = os.path.splitext(input_file)[0] + '_uppercase.csv'
6     with open(input_file, 'r', newline='') as infile:
7         reader = csv.reader(infile)
8         with open(output_file, 'w', newline='') as outfile:
9             writer = csv.writer(outfile)
10            for row in reader:
11                upper_row = [cell.upper() if index in columns else cell for index, cell in enumerate(row)]
12                writer.writerow(upper_row)
13
14 # Example usage:
15 input_file = 'input.csv'
16 selected_columns = [0, 2] # Selecting columns 1 and 3 (0-based index)
17 convert_csv_to_uppercase(input_file, selected_columns)
18
```

5th Proper case

```
import csv
import os

def convert_csv_columns_to_propercase(input_file, columns):
    output_file = os.path.splitext(input_file)[0] + '_propercase.csv'
    with open(input_file, 'r', newline='') as infile:
        reader = csv.reader(infile)
        with open(output_file, 'w', newline='') as outfile:
            writer = csv.writer(outfile)
            for row in reader:
                proper_row = [cell.title() if index in columns else cell for index, cell in enumerate(row)]
                writer.writerow(proper_row)

# Example usage:
input_file = 'input.csv'
selected_columns = [0, 2] # Selecting columns 1 and 3 (0-based index)
convert_csv_columns_to_propercase(input_file, selected_columns)
```

6th Merge

```
import csv
import os

def merge_selected_columns(input_file, column1, column2, new_header):
    # Extract the file name and extension
    file_name, file_extension = os.path.splitext(input_file)

    output_file = f"{file_name}_merged{file_extension}"

    with open(input_file, 'r', newline='') as infile, open(output_file, 'w', newline='') as outfile:
        reader = csv.reader(infile)
        writer = csv.writer(outfile)

        # Read the header row
        header = next(reader)

        # Modify the header to replace the merged column header
        header[column1] = new_header

        # Remove the second column header
        del header[column2]

        # Write modified header to the output file
        writer.writerow(header)

        # Write modified rows
        for row in reader:
            row[column1] = f"{row[column1]} {row[column2]}"
            del row[column2]  # Delete the second column after merging
            writer.writerow(row)

# Example usage:
input_file = "splitt_first_name_splitted.csv"
column1 = 1  # First column to merge (ID)
column2 = 2  # Second column to merge (Name)
new_header = "Name"  # New header name for the merged column
merge_selected_columns(input_file, column1, column2, new_header)
```

7th split

```
import csv
import os

def split_name_column(input_file, first_column_index, second_column_index, first_column_name, last_column_name, sign):
    output_file_name, extension = os.path.splitext(input_file)
    column_name = first_column_name.lower()
    output_file_name = f"split_{column_name}_splitted{extension}"

    with open(input_file, 'r') as csv_input, open(output_file_name, 'w', newline='') as csv_output:
        reader = csv.reader(csv_input)
        writer = csv.writer(csv_output)

        header = next(reader) # Read the header row
        header.insert(second_column_index, last_column_name) # Insert new column for last name
        header[first_column_index] = first_column_name # Update existing column name for first name
        writer.writerow(header) # Write the modified header row

        for row in reader:
            if sign in row[first_column_index]:
                first_name, last_name = row[first_column_index].split(sign, 1)
                row[first_column_index] = first_name.strip() # Update first name
                row.insert(second_column_index, last_name.strip()) # Insert last name
            writer.writerow(row)

# Example usage:
input_filename = 'split.csv'
name_of_first_column = "First Name"
name_of_second_column = "Last Name"
index_for_first_column = 1
index_for_second_column = 2
sign = "/"

split_name_column(input_filename, index_for_first_column, index_for_second_column, name_of_first_column, name_of_second_column, sign)
```

8th remove dots

```
import pandas as pd
import os

def remove_dots_from_csv(input_file, selected_columns=None):
    # Read CSV file into a DataFrame
    df = pd.read_csv(input_file)

    # Iterate through each column index
    for col_index in selected_columns:
        col = df.columns[col_index]
        # Check if column contains dots
        if df[col].dtype == 'object' and df[col].str.contains('.').any():
            # Remove dots from the column
            df[col] = df[col].str.replace('.', '')

    # Get the filename without extension and remove dots
    filename_without_dots = os.path.splitext(os.path.basename(input_file))[0].replace('.', '')

    # Construct the output file path with the modified filename
    output_file = filename_without_dots + '_without_dots.csv'

    # Write the modified DataFrame to a new CSV file
    df.to_csv(output_file, index=False)

# Example usage:
input_file = 'Bad Dots.csv'
selected_columns = [1, 2] # Add the column indices you want to select
remove_dots_from_csv(input_file, selected_columns)
```

9th remove letters

```
✓ import csv
  import re

✓ def remove_letters_from_column(input_csv, column_index):
    with open(input_csv, 'r', newline='') as infile:
        reader = csv.reader(infile)
        header = next(reader) # Read the header row
        column_name = header[column_index] # Get the name of the specified column
        output_csv = f"{column_name}_without_letters.csv" # Create new file name

    with open(output_csv, 'w', newline='') as outfile:
        writer = csv.writer(outfile)
        writer.writerow(header) # Write the header row to the output file

    for row in reader:
        if column_index < len(row):
            cell = row[column_index]
            processed_cell = re.sub(r'[a-zA-Z]', '', cell)
            row[column_index] = processed_cell
        writer.writerow(row)

input_csv = 'upper.csv'
column_index = 2 # Assuming the column index to process is 2 (0-indexed)
remove_letters_from_column(input_csv, column_index)
```

10th remove numbers

```
import csv
import re

def remove_numbers_and_write_by_index(input_file, column_index):
    output_file = ""
    # Extract the column name
    with open(input_file, 'r') as f:
        reader = csv.reader(f)
        headers = next(reader)
        column_name = headers[column_index]
        output_file = f"{column_name}_without_numbers.csv"

    with open(input_file, 'r') as f, open(output_file, 'w', newline='') as f_out:
        reader = csv.reader(f)
        writer = csv.writer(f_out)

        # Define a function to remove numbers from a string
        def remove_numbers(text):
            return re.sub(r'\d+', '', text)

        # Process each row
        for row in reader:
            row[column_index] = remove_numbers(row[column_index])
            writer.writerow(row)

# Example usage:
input_file = 'upper.csv'
column_index = 3 # Assuming you want to operate on the third column (index 2)
remove_numbers_and_write_by_index(input_file, column_index)
```

11th remove selected from column

```
import csv
import re

def remove_numbers_and_write_by_index(input_file, column_index):
    output_file = ""
    # Extract the column name
    with open(input_file, 'r') as f:
        reader = csv.reader(f)
        headers = next(reader)
        column_name = headers[column_index]
        output_file = f"{column_name}_without_numbers.csv"

    with open(input_file, 'r') as f, open(output_file, 'w', newline='') as f_out:
        reader = csv.reader(f)
        writer = csv.writer(f_out)

        # Define a function to remove numbers from a string
        def remove_numbers(text):
            return re.sub(r'\d+', '', text)

        # Process each row
        for row in reader:
            row[column_index] = remove_numbers(row[column_index])
            writer.writerow(row)

# Example usage:
input_file = 'upper.csv'
column_index = 3 # Assuming you want to operate on the third column (index 2)
remove_numbers_and_write_by_index(input_file, column_index)
```

12th remove double spaces

```
import csv
import re

def remove_numbers_and_write_by_index(input_file, column_index):
    output_file = ""
    # Extract the column name
    with open(input_file, 'r') as f:
        reader = csv.reader(f)
        headers = next(reader)
        column_name = headers[column_index]
        output_file = f"{column_name}_without_numbers.csv"

    with open(input_file, 'r') as f, open(output_file, 'w', newline='') as f_out:
        reader = csv.reader(f)
        writer = csv.writer(f_out)

        # Define a function to remove numbers from a string
        def remove_numbers(text):
            return re.sub(r'\d+', '', text)

        # Process each row
        for row in reader:
            row[column_index] = remove_numbers(row[column_index])
            writer.writerow(row)

# Example usage:
input_file = 'upper.csv'
column_index = 3 # Assuming you want to operate on the third column (index 2)
remove_numbers_and_write_by_index(input_file, column_index)
```

13th strip

```
import csv
import re
import os

def remove_spaces_in_csv(input_file):
    # Extract file name and extension
    file_name, file_extension = os.path.splitext(input_file)
    output_file = f"{file_name}_Striped{file_extension}"

    with open(input_file, 'r', newline='') as infile, \
        open(output_file, 'w', newline='') as outfile:
        reader = csv.reader(infile)
        writer = csv.writer(outfile)

        for row in reader:
            cleaned_row = [re.sub(r'^\s+|\s+$', '', cell) for cell in row] # Remove leading and trailing spaces
            writer.writerow(cleaned_row)

    return output_file

# Example usage:
input_filename = 'spaces.csv'
output_filename = remove_spaces_in_csv(input_filename)
print(f"Spaces removed successfully! Output saved to {output_filename}")
```

14th show repeated in all tables

```
import csv
from collections import Counter, defaultdict

def count_repeated_values(csv_file):
    repeated_values = defaultdict(list)

    with open(csv_file, 'r') as file:
        reader = csv.reader(file)
        headers = next(reader) # Read the header row
        for column_index, column_name in enumerate(headers):
            column_values = Counter()
            file.seek(0) # Reset file pointer to read again
            next(reader) # skip the header row
            for line_num, row in enumerate(reader, start=2): # start=2 because we skipped header
                try:
                    value = row[column_index]
                    column_values[value] += 1
                    if column_values[value] > 1:
                        repeated_values[(column_name, value)].append(line_num)
                except IndexError:
                    print(f"Error: Row {line_num} in column '{column_name}' has fewer elements.")
                    continue

    return repeated_values

# Example usage:
csv_file = 'upper_unrepeated.csv'
repeated_values = count_repeated_values(csv_file)

for (column_name, value), rows in repeated_values.items():
    print(f"Value '{value}' in column '{column_name}' is repeated {len(rows)} times in rows:", rows)
```

15th show repeated on one column

```
import csv
from collections import Counter

def find_repeated_values(csv_file, column_index):
    repeated_values = Counter()
    column_name = None
    repeated_rows = {}

    with open(csv_file, 'r') as file:
        reader = csv.reader(file)
        header = next(reader)
        if column_index < len(header):
            column_name = header[column_index]
            for row_num, row in enumerate(reader, start=2): # Start counting rows from 2 (1-indexed)
                if len(row) > column_index:
                    value = row[column_index]
                    if value in repeated_values:
                        repeated_values[value] += 1
                        repeated_rows[value].append(row_num)
                    else:
                        repeated_values[value] = 1
                        repeated_rows[value] = [row_num]

    return column_name, repeated_values, repeated_rows

# Example usage:
csv_file = 'upper_unrepeated.csv'
column_index = 3 # Replace 0 with the index of the column you want to analyze
column_name, repeated_values, repeated_rows = find_repeated_values(csv_file, column_index)

if repeated_values:
    print(f"Repeated values in column '{column_name}':")
    for value, count in repeated_values.items():
        if count > 1:
            print(f"  value '{value}' repeated {count} times in rows: {', '.join(map(str, repeated_rows[value]))}")
else:
    print(f"No repeated values found in column '{column_name}'")
```

16th handling missing values

```
import pandas as pd
def select_NAN_columns(df,type):
    arrint=[]
    arrstr=[]
    for col in df.columns:
        if df[col].isnull().any():
            if df[col].dtype == 'float64':
                arrint.append(col)
            else:
                arrstr.append(col)
    if(type=='str') :
        return arrstr
    else:
        return arrint

def handle_missing(df,intType,strvalue):
    # handling missing int values
    if (intType==''):
        intType=0
    if (strvalue==''):
        strvalue="missing"

    arrint=select_NAN_columns(df,'int')
    if intType=='mean':
        for i in arrint:
            df[i]=df[i].fillna(df[i].mean())
    elif intType=='median':
        for i in arrint:
            df[i] = df[i].fillna(df[i].median())
    elif intType==0:
        for i in arrint:
            df[i] = df[i].fillna(0)
    elif intType=='mode':
        for i in arrint:
            mode=df[i].value_counts().idxmax()
            df[i] = df[i].fillna(mode)
    # handling missing str values
    arrstr=select_NAN_columns(df,'str')
    if (strvalue):
        for i in arrstr:
            df[i] = df[i].fillna(strvalue)
    else:
        for i in arrstr:
            df[i] = df[i].fillna('missing')
    return df
```

17th Detecting outliers

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

def find_outliers_iqr(column_data):
    """
    Detect outliers in a numerical dataset using the Interquartile Range (IQR) method.

    Args:
    - column_data: A pandas Series representing a column of numerical values.

    Returns:
    - outliers: A list of tuples containing outlier values found in the column and their indices.
    """
    # Calculate the first quartile (Q1) and third quartile (Q3) of the column data
    q1 = np.percentile(column_data, 25)
    q3 = np.percentile(column_data, 75)

    # Calculate the interquartile range (IQR)
    iqr = q3 - q1

    # Define the lower and upper bounds for outlier detection
    lower_bound = q1 - 1.5 * iqr
    upper_bound = q3 + 1.5 * iqr

    # Identify outliers and their indices
    outliers = [(value, idx) for idx, value in enumerate(column_data) if value < lower_bound or value > upper_bound]

    return outliers

# Load the CSV file into a pandas DataFrame
df = pd.read_csv("example11.csv")

# Choose the column index you want to perform outlier detection on
column_index = 2 # Change this to the desired column index

# Extract the specified column from the DataFrame as a pandas Series
column_data = df.iloc[:, column_index]

# Get the column name from the DataFrame based on the provided column index
column_name = df.columns[column_index]

# Find outliers in the specified column using the modified function
outliers = find_outliers_iqr(column_data)

# Calculate the total number of data points
total_data_points = len(column_data)

# Calculate the number of outliers
num_outliers = len(outliers)

# Calculate the number of non-outliers
num_non_outliers = total_data_points - num_outliers

# Calculate the percentage of outliers
percentage_outliers = (num_outliers / total_data_points) * 100

# Print the outliers and their locations
print("Outliers found in column '{}'.format(column_name))\nfor outlier in outliers:
```

```

row_index = outlier[1] # Assuming the index corresponds to row number
print("Value:", outlier[0], "at Row:", row_index)

# Print the number of outliers in the dataset
print("Number of outliers in the dataset: {} out of {}".format(num_outliers,
total_data_points))

# Print the percentage of outliers in the dataset
print("Percentage of outliers in the dataset: {:.2f}%".format(percentage_outliers))

# Scatter Plot
plt.figure(figsize=(8, 6))
plt.scatter(range(len(column_data)), column_data, color='blue', label='Data Points')
plt.scatter(*zip(*outliers), color='red', label='Outliers')
plt.title('Scatter Plot of ' + column_name)
plt.xlabel('Index')
plt.ylabel('Values')
plt.legend()
plt.grid(True)
plt.savefig('scatter_plot.png') # Save scatter plot as image
plt.show()

# Box Plot
plt.figure(figsize=(8, 6))
plt.boxplot(column_data, vert=False)
plt.title('Box Plot of ' + column_name)
plt.xlabel('Values')
plt.grid(True)
plt.savefig('box_plot.png') # Save box plot as image
plt.show()

# Histogram
plt.figure(figsize=(8, 6))
plt.hist(column_data, bins=20, color='skyblue', edgecolor='black')
plt.title('Histogram of ' + column_name)
plt.xlabel('Values')
plt.ylabel('Frequency')
plt.grid(True)
plt.savefig('histogram.png') # Save histogram as image
plt.show()

# Line Chart
plt.figure(figsize=(8, 6))
plt.plot(range(len(column_data)), column_data, color='green')
plt.title('Line Chart of ' + column_name)
plt.xlabel('Index')
plt.ylabel('Values')
plt.grid(True)
plt.savefig('line_chart.png') # Save line chart as image
plt.show()

```

18th Remove outlier values from a specified column:

```
import numpy as np
import pandas as pd
import os

def delete_outliers(df, column_index):
    """
    Delete outlier values from a DataFrame.

    Args:
    - df: A pandas DataFrame.
    - column_index: The index of the column containing outliers.

    Returns:
    - df_cleaned: The DataFrame with outlier values removed.
    """
    # Find outliers in the specified column
    column_data = df.iloc[:, column_index]
    q1 = np.percentile(column_data, 25)
    q3 = np.percentile(column_data, 75)
    iqr = q3 - q1
    lower_bound = q1 - 1.5 * iqr
    upper_bound = q3 + 1.5 * iqr

    # Delete outlier values
    df_cleaned = df.copy() # Make a copy of the original DataFrame
    df_cleaned.loc[(df_cleaned.iloc[:, column_index] < lower_bound) | (df_cleaned.iloc[:, column_index] > upper_bound), df.columns[column_index]] = np.nan

    return df_cleaned

# Load the CSV file into a pandas DataFrame
input_filename = "example11.csv" # Change this to the desired input filename
df = pd.read_csv(input_filename)

# Choose the column index you want to perform outlier detection on
column_index = 1 # Change this to the desired column index

# Delete outlier values from the specified column
df_cleaned = delete_outliers(df, column_index)

# Extract the basename of the original filename without the extension
filename, extension = os.path.splitext(input_filename)
basename = os.path.basename(filename)

# Construct the output filename with the suffix "deleted_outliers"
output_file = basename + "_deleted_outliers" + extension

# Save the cleaned DataFrame to a new CSV file with the modified filename
df_cleaned.to_csv(output_file, index=False)

print("Cleaned DataFrame with outlier values removed has been saved to\n'{}'.".format(output_file))
```

19th Replace outlier values in a specified column with the average of non-outlier values within that column:

```
import numpy as np
import pandas as pd
import os

def replace_outliers_with_average(df, column_index):
    """
    Replace outliers in a DataFrame with the average of values in the column.

    Args:
    - df: A pandas DataFrame.
    - column_index: The index of the column containing outliers.

    Returns:
    - df_corrected: The DataFrame with outliers replaced by the average of values in the column.
    """
    # Find outliers in the specified column
    column_data = df.iloc[:, column_index]
    q1 = np.percentile(column_data, 25)
    q3 = np.percentile(column_data, 75)
    iqr = q3 - q1
    lower_bound = q1 - 1.5 * iqr
    upper_bound = q3 + 1.5 * iqr

    # Calculate the average of values within the interquartile range
    avg = np.mean(df[(df.iloc[:, column_index] >= lower_bound) & (df.iloc[:, column_index] <= upper_bound)].iloc[:, column_index])

    # Ensure that the average value has the same data type as the column
    avg = df.iloc[:, column_index].dtype.type(avg)

    # Replace outliers with the average of values
    df_corrected = df.copy()
    df_corrected.loc[(df_corrected.iloc[:, column_index] < lower_bound) | (df_corrected.iloc[:, column_index] > upper_bound), df.columns[column_index]] = avg

    return df_corrected

# Load the CSV file into a pandas DataFrame
input_filename = "example11.csv" # Change this to the desired input filename
df = pd.read_csv(input_filename)

# Choose the column index you want to perform outlier detection on
column_index = 1 # Change this to the desired column index

# Replace outliers with the average of values in the specified column
df_corrected = replace_outliers_with_average(df, column_index)

# Extract the basename of the original filename without the extension
filename, extension = os.path.splitext(input_filename)
basename = os.path.basename(filename)

# Construct the output filename with the suffix "_average_to_outliers"
output_file = basename + "_average_to_outliers" + extension

# Save the corrected DataFrame to a new CSV file with the modified filename
df_corrected.to_csv(output_file, index=False)

print("Corrected DataFrame with outliers replaced by the average of values has been saved to {}".format(output_file))
```

20th Replace outlier values in a specified column of with a chosen value:

```
import numpy as np
import pandas as pd
import os

def replace_outliers_with_value(df, column_index, new_value):
    """
    Replace outliers in a DataFrame with a chosen value.

    Args:
    - df: A pandas DataFrame.
    - column_index: The index of the column containing outliers.
    - new_value: The new value to replace outliers with.

    Returns:
    - df_corrected: The DataFrame with outliers replaced by the chosen value.
    """
    # Find outliers in the specified column
    column_data = df.iloc[:, column_index]
    q1 = np.percentile(column_data, 25)
    q3 = np.percentile(column_data, 75)
    iqr = q3 - q1
    lower_bound = q1 - 1.5 * iqr
    upper_bound = q3 + 1.5 * iqr

    # Replace outliers with the chosen value
    df_corrected = df.copy()
    df_corrected.loc[(df_corrected.iloc[:, column_index] < lower_bound) |
    (df_corrected.iloc[:, column_index] > upper_bound), df.columns[column_index]] = new_value

    return df_corrected

# Load the CSV file into a pandas DataFrame
input_filename = "example11.csv" # Change this to the desired input filename
df = pd.read_csv(input_filename)

# Choose the column index you want to perform outlier detection on
column_index = 1 # Change this to the desired column index

# Choose the value to replace outliers with
chosen_value = 999 # Change this to the desired value

# Replace outliers with the chosen value in the specified column
df_corrected = replace_outliers_with_value(df, column_index, chosen_value)

# Extract the basename of the original filename without the extension
filename, extension = os.path.splitext(input_filename)
basename = os.path.basename(filename)

# Construct the output filename with the suffix "_outliers_replaced_with_{chosen_value}"
output_file = basename + "_outliers_replaced_with_" + str(chosen_value) + extension

# Save the corrected DataFrame to a new CSV file with the modified filename
df_corrected.to_csv(output_file, index=False)

print("Corrected DataFrame with outliers replaced by the chosen value has been saved to\n'{}'.".format(output_file))
```

21st Detect Duplicates:

```
import csv

def find_duplicates(input_file, column_index):
    duplicates = {}
    with open(input_file, 'r', newline='') as infile:
        reader = csv.reader(infile)
        next(reader) # Skip header row if present
        for row_number, row in enumerate(reader, start=1):
            value = row[column_index]
            if value in duplicates:
                duplicates[value].append(row_number)
            else:
                duplicates[value] = [row_number]

    # Filter duplicates
    duplicate_entries = {value: rows for value, rows in duplicates.items() if
len(rows) > 1}
    if not duplicate_entries:
        print("No duplicates found.")
        return

    print("Duplicates found:")
    for value, rows in duplicate_entries.items():
        print(f"Value '{value}' duplicated at rows: {', '.join(map(str, rows))}")

input_file = 'output Dots.csv'
column_index = 2 # Specify the index of the column to check for duplicates
find_duplicates(input_file, column_index)
```

22nd Remove Duplicates:

```
import csv
import os

def remove_duplicates(input_file, column_index):
    unique_values = set()
    output_file = os.path.splitext(input_file)[0] + '_unique.csv'

    with open(input_file, 'r', newline='') as infile:
        reader = csv.reader(infile)
        with open(output_file, 'w', newline='') as outfile:
            writer = csv.writer(outfile)
            for row in reader:
                value = row[column_index]
                if value not in unique_values:
                    writer.writerow(row)
                    unique_values.add(value)

input_file = 'output Dots.csv'
column_index = 0 # Specify the index of the column to check for duplicates
remove_duplicates(input_file, column_index)
```

2nd System Modules Implementation

1st Missing values

```
@app.route('/Missing', methods=["GET", "POST"])
def miss():
    filenames=get_uploaded_files()
    if request.method == 'POST':
        operation = request.form['operation']
        filename = request.form['table']
        if filename=="Not select":
            return render_template('Missing_value.html', filenames=filenames,table_original="You haven't selected any Table!",missingPers=0,missingCount=0,totalRows=0)
        else:
            df = pd.read_csv('uploads/'+filename)
            show():
                missing_Table =find_missing_values_rows(df)
                if missing_Table.count().max()==0:
                    missing_Html_Table="No Missing Value Found.!"
                else:
                    missing_Html_Table = missing_Table.to_html(classes='table table-bordered table-striped sizee', index=False)
                    table_html_original = df.to_html(classes='table table-bordered table-striped sizee', index=False)
                    missingPers,missingCount,totalRows=pers_missing(df)
                    return missingPers,missingCount,totalRows,missing_Html_Table,table_html_original
            if (operation=="show"):
                missingPers,missingCount,totalRows,missing_Html_Table,table_html_original=show()
                return render_template('Missing_value.html', table_original=table_html_original, table_cleaned=missing_Html_Table,filenames=filenames,missingPers=missingPers,missingCount=missingCount,totalRows=totalRows)
            elif (operation=="fill"):
                intType = request.form['fill-int-type']
                if intType=="Not select int":
                    return render_template('Missing_value.html',intmsg='Select Option To Fill Integer Missing Value',missingPers=0,missingCount=0,totalRows=0)
                else:
                    strValue = request.form['fill-str-value']
                    if strValue== '':
                        return render_template('Missing_value.html',strmsg='Select Option To Fill String Missing Value',missingPers=0,missingCount=0,totalRows=0)
                    else:
                        newdata=handle_missing(df,intType,strValue)
                        newdata.to_csv('uploads/'+filename,index=False)
                        table_html_cleaned = newdata.to_html(classes='table table-bordered table-striped ', index=False)
                        missingPers,missingCount,totalRows,missing_Html_Table,table_html_original=show()
                        return render_template('Missing_value.html', table_original=table_html_original, table_cleaned=table_html_cleaned,filenames=filenames,missingPers=missingPers,missingCount=missingCount,totalRows=totalRows)
            elif operation=="delete" :
                df=df.dropna()
                df.to_csv('uploads/'+filename)
                missingPers,missingCount,totalRows=pers_missing(df)
                table_html_original = df.to_html(classes='table table-bordered table-striped ', index=False)
                missingPers,missingCount,totalRows,missing_Html_Table,table_html_original=show()
                return render_template('Missing_value.html', table_original=table_html_original,table_cleaned="No Missing Value Found.!", filenames=filenames,missingPers=missingPers,missingCount=missingCount,totalRows=totalRows)
    else:
        return render_template('Missing_value.html', filenames=filenames,missingPers=0,missingCount=0,totalRows=0)
```

2nd Duplicates

```
@app.route('/duplicate',methods=['POST', 'GET'])
def duplicate():
    filenames=get_uploaded_files()
    if request.method == 'POST':
        operation = request.form['operation']
        filename = request.form['table']
        if filename=="Not select":
            return render_template('duplicate.html', filenames=filenames,table_original="You haven't selected any Table!",pers=0,effectedRows=0,totalRows=0)
        else:
            def show_dup():
                selected_columns = request.form.getlist('columns')
                filename = request.form['table']
                df = pd.read_csv('uploads/'+filename)
                pers,effectedRows,totalRows=pers_duplicate(df)
                columns=get_column_names(df)
                if selected_columns:
                    show_dup=find_duplicate_rows(df,selected_columns)
                    pers,effectedRows,totalRows=pers_duplicate(df,selected_columns)
                else:
                    show_dup=find_duplicate_rows(df,'all')

                if show_dup.count().max() ==0:
                    table_html_duplicate="No Duplicates Found."
                else:
                    table_html_duplicate = show_dup.to_html(classes='table table-striped table-bordered sizee ', index=False)
                    table_html_original = df.to_html(classes='table table-striped table-bordered sizee ', index=False)
                    return pers,effectedRows,totalRows,table_html_original,table_html_duplicate,columns,selected_columns

            if operation == 'show':
                pers,effectedRows,totalRows,table_html_original,table_html_duplicate,columns,selected_columns=show_dup()
                return render_template("duplicate.html",filenames=filenames,table_original=table_html_original,table_duplicate=table_html_duplicate,pers=pers,effectedRows=effectedRows,totalRows=totalRows,columns=columns,selected_columns=selected_columns)
            else:
                try:
                    selected_columns = request.form.getlist('columns')
                    if selected_columns:
                        remove_duplicates_and_save(filename,selected_columns)
                        pers,effectedRows,totalRows,table_html_original,table_html_duplicate,columns,selected_columns=show_dup()
                        return render_template('duplicate.html',filenames=filenames,table_original=table_html_original,table_duplicate=table_html_duplicate,pers=pers,effectedRows=effectedRows,totalRows=totalRows)
                    else:
                        remove_duplicates_and_save(filename,'all')
                        pers,effectedRows,totalRows,table_html_original,table_html_duplicate,columns,selected_columns=show_dup()
                        return render_template('duplicate.html',filenames=filenames,table_original=table_html_original,table_duplicate=table_html_duplicate,pers=pers,effectedRows=effectedRows,totalRows=totalRows)
                except:
                    pers,effectedRows,totalRows,table_html_original,table_html_duplicate=show_dup()
                    return render_template('duplicate.html',filenames=filenames,table_original=table_html_original,table_duplicate=table_html_duplicate,pers=pers,effectedRows=effectedRows,totalRows=totalRows)

    return render_template('duplicate.html',filenames=filenames,pers=0,effectedRows=0,totalRows=0)
```

3rd Invalid formats

```
#app.route('/invalid', methods=['POST', 'GET'])
def invalid():
    filenames = get_uploaded_files()
    if request.method == 'POST':
        operation = request.form['operation']
        filename = request.form['table']
        def show_dupl():
            filename = request.form['table']
            df = pd.read_csv('uploads/' + filename)
            columns = get_column_names(df)
            table_html_original = df.to_html(classes='table table-striped table-bordered size', index=False)
            return render_template('invalid_format.html', filenames=filenames, table_original=table_html_original, columns=columns)
        if filename == "Not select":
            return render_template('invalid_format.html', filenames=filenames, table_original="You haven't selected any Table!", person=0, affectedRows=0, totalRows=0)
        else:
            if operation == 'show':
                return show_dupl()
            if operation == 'apply':
                Remove_operation = request.form['Remove']
                Standardize = request.form['Standardize']
                Split = request.form['Split']
                check = request.form['check']
                selected_columns = request.form.getlist('columns')
                if not (Remove_operation == ""):
                    if Remove_operation == 'Comma':
                        Remove.remove_commas_from_csv(filename, selected_columns)
                    return show_dupl()
                elif Remove_operation == 'Data':
                    Remove.remove_data_from_csv(filename, selected_columns)
                    return show_dupl()
                elif Remove_operation == 'Strip':
                    Remove.remove_strip(filename)
                    return show_dupl()
                elif Remove_operation == 'WhiteSpace':
                    Remove.remove_white_spaces_in_csv(filename)
                    return show_dupl()
                elif Remove_operation == 'char':
                    selected_columns = request.form.getlist('columns')
                    char = request.form['char']
                    Remove.delete_character_from_csv(filename, selected_columns[0], char)
                    return show_dupl()
                elif Remove_operation == 'num':
                    selected_columns = request.form.getlist('columns')
                    Remove.remove_numbers_and_write_by_column_name(filename, selected_columns)
                    return show_dupl()
                if not (Standardize == ""):
                    if Standardize == 'lower':
                        st.convert_csv_columns_to_lowercase(filename, selected_columns)
                    return show_dupl()
                elif Standardize == 'upper':
                    st.convert_csv_to_uppercase(filename, selected_columns)
                    return show_dupl()
                elif Standardize == 'Proper':
                    st.convert_csv_columns_to_propercase(filename, selected_columns)
                    return show_dupl()
            if not (Split == ""):
                first_col = request.form['first_col']
                second_col = request.form['second_col']
                selected_columns = request.form.getlist('columns')

                if Split == 'slash':
                    sp.split_name_column(filename, selected_columns[0], first_col, second_col, '/')
                    return show_dupl()
                elif Split == 'and':
                    sp.split_name_column(filename, selected_columns[0], first_col, second_col, '&')
                    return show_dupl()
                elif Split == 'dot':
                    sp.split_name_column(filename, selected_columns[0], first_col, second_col, '.')
                    return show_dupl()
                elif Split == 'comma':
                    sp.split_name_column(filename, selected_columns[0], first_col, second_col, ',')
                    return show_dupl()
                elif Split == 'specific':
                    separator = request.form['separator']
                    sp.split_name_column(filename, selected_columns[0], first_col, second_col, sign=separator)
                    return show_dupl()

            if not (check == ""):
                if check == 'Email':
                    filename = request.form['table']
                    df = pd.read_csv('uploads/' + filename)
                    columns = get_column_names(df)
                    invalid_rows_email = invalid_emails(filename, selected_columns[0])
                    table_html_original = df.to_html(classes='table table-striped table-bordered size', index=False)
                    if invalid_rows_email.count().max() != 0:
                        table_cleaned = invalid_rows_email.to_html(classes='table table-striped table-bordered size', index=False)
                        return render_template('invalid_format.html', filenames=filenames, table_original=table_html_original, table_cleaned=table_cleaned, columns=columns)
                    else:
                        return render_template('invalid_format.html', filenames=filenames, table_original=table_html_original, table_cleaned="No Invalid Rows Found!", columns=columns)
                elif check == 'Phone':
                    filename = request.form['table']
                    df = pd.read_csv('uploads/' + filename)
                    columns = get_column_names(df)
                    invalid_rows_phone = check_phone_numbers_in_column(filename, selected_columns[0])
                    table_html_original = df.to_html(classes='table table-striped table-bordered size', index=False)
                    if invalid_rows_phone.count().max() != 0:
                        table_cleaned = invalid_rows_phone.to_html(classes='table table-striped table-bordered size', index=False)
                        return render_template('invalid_format.html', filenames=filenames, table_original=table_html_original, table_cleaned=table_cleaned, columns=columns)
                    else:
                        return render_template('invalid_format.html', filenames=filenames, table_original=table_html_original, table_cleaned="No Invalid Rows Found!", columns=columns)
            return render_template('invalid_format.html', filenames=filenames)
```

3rd Python:-

1st _Init_

```
1 ✓ from flask import Flask ,request,render_template,url_for,redirect
2 import pandas as pd
3 import os
4 from missing_value import handle_missing,get_uploaded_files,pers_missing,find_missing_values_rows
5 from duplicates import find_duplicate_rows,get_column_names,pers_duplicate,remove_duplicates_and_save
6 # from outliers import getOutliersDetails,getNegative,remove_negative_rows
7 #remove_outliers_and_save,save_and_display_plot
8
9 # from conver import get_sql_server_databases,db_convert,get_server_name
10 # from dashboard import calculate_null_percentage_in_database
11 # from outliers import getOutliersDetails
12
13 app=Flask(__name__)
14
15 UPLOAD_FOLDER = 'uploads'
16 app.config['UPLOAD_FOLDER'] = UPLOAD_FOLDER
17
18
19
20 UPLOAD_FOLDER = 'uploads'
21 app.config['UPLOAD_FOLDER'] = UPLOAD_FOLDER
22
23
24 @app.route('/')
25 def Auth():
26     ||| return render_template('Auth.html')
27
28
29
30 @app.route('/home')
31 def home():
32     ||| return render_template('home.html')
33
34
35 @app.route('/About')
36 def About():
37     ||| return render_template('About.html')
38
39 @app.route('/upload')
40 def upload():
41     ||| return render_template('upload.html')
42
43 @app.route('/dashboard')
44 def dashboard():
45     ||| return render_template('dashboard.html')
46
47
48 #####Missing#####
49
50 @app.route('/Missing', )
51 def miss():
52     filenames=get_uploaded_files()
53     if request.method == 'POST':
54         operation = request.form['operation']
55         filename = request.form['table']
56         if filename=="Not select":
57             return render_template('Missing_value.html',filenames=filenames,table_original="You haven't selected any Table!",missingPers=0,missingCount=0,totalRows=0)
58         else:
59             df = pd.read_csv('uploads/'+filename)
60
61         if (operation=='show'):
62
63
64             return render_template('Missing_value.html')
65         elif (operation=='fill'):
66             intType = request.form['fill-int-type']
67             if intType=="Not select int":
68                 return render_template('Missing_value.html',intmsg='Select Option To Fill Integer Missing Value',missingPers=0,missingCount=0,totalRows=0)
69             else:
70                 strValue = request.form['fill-str-value']
71                 if strValue=='':
72                     return render_template('Missing_value.html',strmsg='Select Option To Fill String Missing Value',missingPers=0,missingCount=0,totalRows=0)
73                 else:
74                     newdata=handle_missing(df,intType,strValue)
75                     newdata.to_csv("uploads/"+filename,index=False)
76                     table_html_cleaned = newdata.to_html(classes='table table-bordered table-striped ', index=False)
77
78             return render_template('Missing_value.html', table_original=table_html_cleaned,table_cleaned=table_html_cleaned,
79             |||||filenames=filenames,missingPers=missingPers,missingCount=missingCount,totalRows=totalRows)
```

```

80    elif operation=='delete' :
81        df=df.dropna()
82        df.to_csv('uploads/'+filename)
83        missingPers,missingCount,totalRows=pers_missing(df)
84        table_html_original = df.to_html(classes='table table-bordered table-striped ', index=False)
85
86        return render_template('Missing_value.html', table_original=table_html_original,table_cleaned="No Missing Value Found.",filenames=filenames,missingPers=missingPers,missingCount=missingCount,totalRows=totalRows)
87        # if (operation=='Not select int'):
88
89        #     missingPers,missingCount,totalRows,missing_Html_Table,table_html_original=show()
90        #     return render_template('Missing_value.html', table_original=table_html_original, table_cleaned=missing_Html_Table,
91        #     filenames=filenames,missingPers=missingPers,missingCount=missingCount,totalRows=totalRows,intmsg='Select Option To Fill Integer Missing Value')
92
93    else:
94        return render_template('Missing_value.html', filenames=filenames,missingPers=0,missingCount=0,totalRows=0)
95
96
97 #####endMissing#####
98

```

```

106    @app.route('/duplicate',methods=['POST', 'GET'])
107    def duplicate():
108        filenames=get_uploaded_files()
109        if request.method == 'POST':
110            operation = request.form['operation']
111            filename = request.form['table']
112            if filename=="Not select":
113                return render_template('duplicate.html', filenames=filenames,table_original="You haven't selected any Table!",pers=0,effectedRows=0,totalRows=0)
114            else:
115                def show_dupl():
116                    selected_columns = request.form.getlist('columns')
117                    filename = request.form['table']
118                    df = pd.read_csv('uploads/'+filename)
119                    pers,effectedRows,totalRows=pers_duplicate(df)
120                    columns=get_column_names(df)
121                    if selected_columns:
122                        show_dup=find_duplicate_rows(df,selected_columns)
123                        pers,effectedRows,totalRows=pers_duplicate(df,selected_columns)
124                    else:
125                        show_dup=find_duplicate_rows(df,'all')
126
127
128                    if show_dup.count().max() ==0:
129                        table_html_duplicate='No Duplicates Found.'
130                    else:
131                        table_html_duplicate = show_dup.to_html(classes='table table-striped table-bordered sizee ', index=False)
132                        table_html_original = df.to_html(classes='table table-striped table-bordered sizee ', index=False)
133                        return pers,effectedRows,totalRows,table_html_original,table_html_duplicate,columns,selected_columns
134                if operation == 'show':
135                    pers,effectedRows,totalRows,table_html_original,table_html_duplicate,columns,selected_columns=show_dupl()
136                    return render_template('duplicate.html',filenames=filenames,table_original=table_html_original,table_duplicate=table_html_duplicate,pers=pers,
137                                         effectedRows=effectedRows,totalRows=totalRows,columns=columns,selected_columns=selected_columns)
138                else:
139                    try:
140                        selected_columns = request.form.getlist('columns')
141                        if selected_columns:
142                            remove_duplicates_and_save(filename,selected_columns)
143                            pers,effectedRows,totalRows,table_html_original,table_html_duplicate,columns,selected_columns=show_dupl()
144                            return render_template('duplicate.html',filenames=filenames,table_original=table_html_original,
145                                         table_duplicate=table_html_duplicate,pers=pers,effectedRows=effectedRows,totalRows=totalRows)
146
147                    else:
148                        remove_duplicates_and_save(filename,'all')
149                        pers,effectedRows,totalRows,table_html_original,table_html_duplicate,columns,selected_columns=show_dupl()
150                        return render_template('duplicate.html',filenames=filenames,table_original=table_html_original,
151                                         table_duplicate=table_html_duplicate,pers=pers,effectedRows=effectedRows,totalRows=totalRows)
152
153                except:
154                    pers,effectedRows,totalRows,table_html_original,table_html_duplicate=show_dupl()
155                    return render_template('duplicate.html',filenames=filenames,table_original=table_html_original,
156                                         table_duplicate=table_html_duplicate,pers=pers,effectedRows=effectedRows,totalRows=totalRows)
157
158
159    return render_template(['duplicate.html',filenames=filenames,pers=0,effectedRows=0,totalRows=0])
160
#####endDuplicate#####

```

```

4 #####outlier#####
5
6 @app.route('/outlier')
7 def outlier():
8     filenames=get_uploaded_files()
9     if request.method == 'POST':
10         filename = request.form['table']
11         operation = request.form['operation']
12     def show_outlier(filename):
13         df = pd.read_csv('uploads/'+filename)
14         columns=get_column_names(df)
15         table_html_original = df.to_html(classes='table table-striped table-bordered ', index=False)
16         table_html_outliers = df.to_html(classes='table table-striped table-bordered ', index=False)
17         if outliers.count().max() ==0:
18             return render_template('outlier.html',filenames=filenames,table_original=table_html_original,table_outliers="No Outlier Found!",total_rows=total_rows,num_outliers=num_outliers,outlier_percentage=outlier_percentage,columns=columns)
19         else:
20             table_html_outliers = outliers.to_html(classes='table table-striped table-bordered ', index=False)
21             return render_template('outlier.html',filenames=filenames,table_original=table_html_original,
22                                 table_outliers=table_html_outliers,total_rows=total_rows,num_outliers=num_outliers,
23                                 outlier_percentage=outlier_percentage,columns=columns)
24
25     if filename=='Select':
26         return render_template('outlier.html',filenames=filenames,table_original="You Should Select Table !",total_rows=0,num_outliers=0,outlier_percentage=0)
27     else:
28
29         if operation == 'show':
30             df = pd.read_csv('uploads/'+filename)
31             columns=get_column_names(df)
32             table_html_original = df.to_html(classes='table table-striped table-bordered ', index=False)
33             return render_template('outlier.html',filenames=filenames,table_original=table_html_original,total_rows=0,num_outliers=0,
34                                 outlier_percentage=0,columns=columns)
35
36     elif operation == 'find_outlier':
37         df = pd.read_csv('uploads/'+filename)
38         columns=get_column_names(df)
39         chart = request.form['chart']
40         selected_columns = request.form.getlist('columns')
41
42     elif operation == 'find_outlier':
43         df = pd.read_csv('uploads/'+filename)
44         columns=get_column_names(df)
45         chart = request.form['chart']
46         selected_columns = request.form.getlist('columns')
47
48
49         table_html_original = df.to_html(classes='table table-striped table-bordered ', index=False)
50         if selected_columns:
51
52             if outliers.count().max() ==0:
53                 return render_template('outlier.html',filenames=filenames,table_original=table_html_original,
54                                     table_outliers="No Outlier Found!",total_rows=total_rows,num_outliers=num_outliers,outlier_percentage=outlier_percentage,
55                                     columns=columns)
56             else:
57
58                 table_html_outliers = outliers.to.html(classes='table table-striped table-bordered ', index=False)
59                 return render_template('outlier.html',filenames=filenames,table_original=table_html_original,table_outliers=table_html_outliers,total_rows=total_rows,
60                                     num_outliers=num_outliers,outlier_percentage=outlier_percentage,columns=columns)
61
62             else:
63                 return render_template('outlier.html',filenames=filenames,table_original=table_html_original,table_outliers="You Don't Select Any Row!",
64                                     total_rows=0,num_outliers=0,outlier_percentage=0,columns=columns)
65
66     elif operation == 'negative':
67         df = pd.read_csv('uploads/'+filename)
68         columns=get_column_names(df)
69
70         selected_columns = request.form.getlist('columns')
71         table_html_original = df.to.html(classes='table table-striped table-bordered ', index=False)
72         if selected_columns:
73             outliers,total_rows,num_outliers,outlier_percentage=[0]
74             if outliers.count().max() ==0:
75                 return render_template('outlier.html',filenames=filenames,table_original=table_html_original,table_outliers="No Negative Value Found!",
76                                     total_rows=total_rows,num_outliers=num_outliers,outlier_percentage=outlier_percentage,columns=columns)
77             else:
78
79                 table_html_outliers = outliers.to.html(classes='table table-striped table-bordered ', index=False)
80                 return render_template('outlier.html',filenames=filenames,table_original=table_html_original,
81                                     table_outliers=table_html_outliers,total_rows=total_rows,num_outliers=num_outliers,outlier_percentage=outlier_percentage,
82                                     columns=columns)
83
84             else:
85                 return render_template('outlier.html',filenames=filenames,table_original=table_html_original,
86                                     table_outliers="You Don't Select Any Row!",total_rows=0,num_outliers=0,outlier_percentage=0,columns=columns)
87

```

Activate Window
Go to Settings to a

```
240     elif operation == 'delete_outlier':
241         selected_columns = request.form.getlist('columns')
242         df = pd.read_csv('uploads/'+filename)
243         columns=get_column_names(df)
244         if selected_columns:
245             || return show_outlier(filename)
246     elif operation == 'delete_negative':
247         selected_columns = request.form.getlist('columns')
248         df = pd.read_csv('uploads/'+filename)
249         columns=get_column_names(df)
250         if selected_columns:
251             || return show_outlier(filename)
252     else:
253         return render_template('outlier.html',filenames=filenames,total_rows=0,num_outliers=0,outlier_percentage=0)
254 ######endoutlier#####
255
256
257
258
259
260 @app.route('/invalid')
261 def invalid():
262     || return render_template('invalid_format.html')
263
264
265
266
267
268 if __name__=='__main__':
269     app.run(debug=False,port=2020)
270
```

2nd Configure

```
1 import os
2
3 class Config:
4     SECRET_KEY = os.environ.get('SECRET_KEY') or 'a_hard_to_guess_string'
5     SQLALCHEMY_DATABASE_URI = os.environ.get('DATABASE_URL') or 'sqlite:///mydatabase.db'
6     SQLALCHEMY_TRACK_MODIFICATIONS = False
7     LOGGING_LEVEL = 'INFO'
8
```

3rd Routs

```
1 from flask import request, jsonify
2 from app import db
3 from app.models import User
4 from flask import current_app as app
5
6 @app.route('/users', methods=['POST'])
7 def create_user():
8     data = request.get_json()
9     new_user = User(name=data['name'], email=data['email'])
10    db.session.add(new_user)
11    db.session.commit()
12    return jsonify({'message': 'User created successfully!'}), 201
13
14 @app.route('/users/<int:id>', methods=['GET'])
15 def get_user(id):
16     user = User.query.get_or_404(id)
17     return jsonify({'name': user.name, 'email': user.email})
18
```

4th Model

```
1 from app import db
2
3 class User(db.Model):
4     id = db.Column(db.Integer, primary_key=True)
5     name = db.Column(db.String(50), nullable=False)
6     email = db.Column(db.String(120), unique=True, nullable=False)
7
8     def __repr__(self):
9         return f'User {self.name}'
10
```

5th Convert

```
1 1> import pyodbc
2 2> import pandas as pd
3 3> import socket
4 4> import os
5
6
7 7> def get_sql_server_databases(server):
8 8>     try:
9 9>         connection_string = f"DRIVER=SQL Server;SERVER={server};UID='';PWD=''"
10 10>        connection = pyodbc.connect(connection_string)
11 11>        cursor = connection.cursor()
12 12>        databases = []
13 13>        cursor.execute("SELECT name FROM sys.databases")
14 14>        rows = cursor.fetchall()
15 15>        for row in rows:
16 16>            databases.append(row[0])
17 17>        return databases
18 18>    except Exception as e:
19 19>        print("Error:", e)
20 20>        return None
21 21> def get_server_name():
22 22>     try:
23 23>         server_name = socket.gethostname()
24 24>         return server_name
25 25>     except Exception as e:
26 26>         print("Error:", e)
27 27>         return None
28
29
30 30> def db_convert(db):
31 31>     try:
32 32>         def convert_all_tables_to_csv(server, db):
33 33>             sql_server_connection_string = 'DRIVER={SQL Server};SERVER=' + server + ';DATABASE=' + db + ';Trusted_Connection=yes'
34 34>             conn = pyodbc.connect(sql_server_connection_string)
35 35>             cursor = conn.cursor()
36 36>             tables = cursor.tables(tableType='TABLE')
37 37>             table_names = [table.table_name for table in tables]
38 38>             upload_dir = 'uploads'
39 39>             os.makedirs(upload_dir, exist_ok=True) # Create upload directory if it doesn't exist
40 40>             for table_name in table_names:
41 41>                 query = f'SELECT * FROM {table_name}'
42 42>                 df = pd.read_sql_query(query, conn)
43 43>                 csv_file_path = os.path.join(upload_dir, f'{table_name}.csv') # Path to save CSV in upload directory
44 44>                 df.to_csv(csv_file_path, index=False)
45 45>             conn.close()
46 46>         server = get_server_name() # Assuming you have a function to get server name
47 47>         convert_all_tables_to_csv(server, db)
48 48>
49 49>     except Exception as e:
50 50>         return None
51
52
53 53> # server=get_server_name()
54 54> get_sql_server_databases("DESKTOP-OA4UFIW")
55
56
```

6th Duplicates

```
1 import pandas as pd
2
3 def binary_search(arr, x):
4     """
5         Perform binary search on a sorted list.
6
7     Args:
8         - arr: The sorted list to search in.
9         - x: The value to search for.
10
11    Returns:
12        - index: The index of the value in the list if found, otherwise -1.
13        """
14    low = 0
15    high = len(arr) - 1
16    while low <= high:
17        mid = (low + high) // 2
18        if arr[mid] < x:
19            low = mid + 1
20        elif arr[mid] > x:
21            high = mid - 1
22        else:
23            return mid
24    return -1
25
26 def fx1(df, columns_to_check='all'):
27     if columns_to_check == "all":
28         duplicate_rows = df[df.duplicated(keep=False)]
29     else:
30         duplicate_rows = df[df.duplicated(subset=columns_to_check, keep=False)]
31
32     return duplicate_rows
33
34 def fx2(df, columns_to_check='all'):
35     if columns_to_check == "all":
36         duplicate_rows = df[df.duplicated()].count().max()
37         total_rows = df.count().max()
38         pers = round((duplicate_rows / total_rows) * 100)
39
40     return pers, duplicate_rows, total_rows
41
42     else:
43         duplicate_rows = df[df.duplicated(subset=columns_to_check)].count().max()
44         total_rows = df.count().max()
45         pers = round((duplicate_rows / total_rows) * 100)
46
47         return pers, duplicate_rows, total_rows
48
49
50 def fx3(df):
51     column_names = df.columns.tolist()
52
53     return column_names
54
55
56
57 def fx4(input_filename, columns_to_check='all'):
58     df = pd.read_csv('uploads/' + input_filename)
59     if columns_to_check == "all":
60         df_no_duplicates = df.drop_duplicates()
61     else:
62         df_no_duplicates = df.drop_duplicates(subset=columns_to_check)
63
64     df_no_duplicates.to_csv('uploads/' + input_filename, index=False)
65
66
67 def remove_duplicates_and_save_binary(input_filename, columns_to_check='all'):
68     df = pd.read_csv('uploads/' + input_filename)
69     sorted_df = df.sort_values(by=columns_to_check)
70     duplicate_indices = []
71
72     for col in columns_to_check:
73         sorted_col_values = sorted_df[col].tolist()
74         for i in range(1, len(sorted_col_values)):
75             if sorted_col_values[i] == sorted_col_values[i - 1]:
76                 idx = binary_search(sorted_col_values, sorted_col_values[i])
77                 duplicate_indices.append(idx)
78
79     df_no_duplicates = df.drop(duplicate_indices)
80     df_no_duplicates.to_csv('uploads/' + input_filename, index=False)
81
82
83
84
85
86
87
88
```

7th Missing values

```
1 import pandas as pd
2 class Node:
3     def __init__(self, data, dtype):
4         self.data = data
5         self.dtype = dtype
6         self.next = None
7
8 class Linkedlist:
9     def __init__(self):
10         self.head = None
11
12     def append(self, data, dtype):
13         new_node = Node(data, dtype)
14         if not self.head:
15             self.head = new_node
16             return
17         last = self.head
18         while last.next:
19             last = last.next
20         last.next = new_node
21
22     def get_columns(self):
23         columns = []
24         current = self.head
25         while current:
26             columns.append(current.data)
27             current = current.next
28         return columns
29
30     def get_dtype(self, data):
31         current = self.head
32         while current:
33             if current.data == data:
34                 return current.dtype
35             current = current.next
36         return None
37
38     def extract_nan_columns_ll(linked_list, dtype):
39         nan_columns_int = []
40         nan_columns_str = []
41         for column in linked_list.get_columns():
42             if linked_list.get_dtype(column) == dtype and any(pd.isnull(column)):
43                 if dtype == 'float64':
44                     nan_columns_int.append(column)
45                 else:
46                     nan_columns_str.append(column)
47
48         return nan_columns_str if dtype == 'str' else nan_columns_int
49
50     # Create a linked list
51     ll = LinkedList()
52
53     # Append columns (data should be in form of pandas Series)
54     ll.append(pd.Series([1, 2, None, 4]), 'float64')
55     ll.append(pd.Series([None, 'b', 'c', 'd']), 'str')
56     ll.append(pd.Series([5.5, None, 6.7, None]), 'float64')
57
58     # Manage missing data
59     ll = (ll, 'mean', 'unknown')
60
61     # Print the processed linked list
62     current = ll.head
63     while current:
64         print(current.data)
65         current = current.next
```

8th Outliers

```
1 import numpy as np
2 import pandas as pd
3 import plotly.graph_objs as go
4 import plotly.io as pio
5 import os
6 import numpy as np
7 import pandas as pd
8 import plotly.graph_objs as go
9 import plotly.io as pio
10 import os
11 import pandas as pd
12 import plotly.graph_objs as go
13 import plotly.io as pio
14 import cv2
15 import os
16 import numpy as np
17 import pandas as pd
18 import plotly.graph_objs as go
19 import plotly.io as pio
20 import os
21 import cv2
22 from queue import Queue
23
24 class DataQueue:
25     def __init__(self):
26         self.queue = Queue()
27
28     def enqueue(self, item):
29         self.queue.put(item)
30
31     def dequeue(self):
32         return self.queue.get()
33
34     def is_empty(self):
35         return self.queue.empty()
36
37 def fn1(col_name, df, operation):
38     q = DataQueue()
39     col_data = df[col_name]
40
41     if operation == 'op1':
42         cc1, cc3 = np.percentile(col_data.dropna(), [25, 75])
43         iqccr = cc3 - cc1
44         lower_bound = cc1 - 1.5 * iqccr
45         upper_bound = cc3 + 1.5 * iqccr
46         for idx, val in col_data.iteritems():
47             if val < lower_bound or val > upper_bound:
48                 q.enqueue(idx)
49     elif operation == 'op2':
50         for idx, val in col_data.iteritems():
51             if val < 0:
52                 q.enqueue(idx)
53
54     return q
55
56
57 def fn2(col_name, df, operation, replacement=None):
58     q = fn1(col_name, df, operation)
59     while not q.is_empty():
60         idx = q.dequeue()
61         if operation == 'rmv':
62             df = df.drop(index=idx)
63         elif operation == 'rplc':
64             df.at[idx, col_name] = replacement
65     return df
66
67 def fn3(column_name, df):
68     return fn1(column_name, df, 'op1')
69
70 def fn4(column_name, df):
71     return fn1(column_name, df, 'op2')
```

```

73 def fn5(column_name, df, states='all'):
74     total_rows = len(df)
75     if states == 'neg':
76         outlier_queue = fn4(column_name, df)
77     else:
78         outlier_queue = fn3(column_name, df)
79
80     num_outliers = 0
81     while not outlier_queue.is_empty():
82         outlier_queue.dequeue()
83         num_outliers += 1
84
85     outlier_percentage = (num_outliers / total_rows) * 100
86     return total_rows, num_outliers, round(outlier_percentage)
87
88 def fn6(column_name, df):
89     outliers = fn3(column_name, df)
90     total_rows, num_outliers, outlier_percentage = fn5(column_name, df)
91     return outliers, total_rows, num_outliers, outlier_percentage
92
93 def fn7(column_name, df):
94     negativeOutliers = fn4(column_name, df)
95     total_rows, num_outliers, outlier_percentage = fn5(column_name, df, 'neg')
96     return negativeOutliers, total_rows, num_outliers, outlier_percentage
97
98 def fn8(column_name, df, output_file):
99     cleaned_df = fn2(column_name, df, 'rmv')
100    cleaned_df.to_csv(output_file, index=False)
101    return cleaned_df
102
103 def fn9(column_name, fileName):
104     df = pd.read_csv(fileName)
105     cleaned_df = fn2(column_name, df, 'rmv')
106     if cleaned_df is not None:
107         cleaned_df.to_csv(fileName, index=False)
108     return cleaned_df
109
110 def remove_negative_rows(column_name, fileName):
111     df = pd.read_csv(fileName)
112     column_numeric = pd.to_numeric(df[column_name], errors='coerce')
113
114     cleaned_df = df[column_numeric >= 0]
115     if cleaned_df is not None:
116         cleaned_df.to_csv(fileName, index=False)
117     return cleaned_df
118
119 def save_and_display_plot(df, column_name, plot_type):
120     column_data = df[column_name]
121     output_dir = "./img"
122
123     if plot_type == 'Scatter_Plot':
124         plot = go.Scatter(x=df.index, y=column_data, mode='markers', name='Data Points')
125         layout = go.Layout(title=f"Scatter Plot of {column_name}",
126                            xaxis=dict(title='Index'),
127                            yaxis=dict(title=column_name))
128     elif plot_type == 'Box_Plot':
129         plot = go.Box(y=column_data, name='Data Distribution')
130         layout = go.Layout(title=f"Box Plot of {column_name}", yaxis=dict(title=column_name))
131     elif plot_type == 'Histogram':
132         plot = go.Histogram(x=column_data, name='Data Distribution', nbinsx=20)
133         layout = go.Layout(title=f"Histogram of {column_name}",
134                            xaxis=dict(title=column_name),
135                            yaxis=dict(title='Frequency'))
136     elif plot_type == 'Line_Chart':
137         plot = go.Scatter(x=df.index, y=column_data, mode='lines', name='Data Points')
138         layout = go.Layout(title=f"Line Chart of {column_name}",
139                            xaxis=dict(title='Index'),
140                            yaxis=dict(title=column_name))
141     else:
142         print("Unknown plot type")
143         return
144
145     fig = go.Figure(data=[plot], layout=layout)
146
147     if not os.path.exists(output_dir):
148         os.makedirs(output_dir)
149
150     output_file_name = os.path.join(output_dir, f"{column_name}_plot.png")
151     pio.write_image(fig, output_file_name)
152
153     saved_image = cv2.imread(output_file_name)
154     cv2.imshow("Saved Image", saved_image)
155     cv2.waitKey(0)
156     cv2.destroyAllWindows()
157
158     output_file_name = os.path.join(output_dir, f"{column_name}_plot.png")
159     pio.write_image(fig, output_file_name)
160
161     saved_image = cv2.imread(output_file_name)
162     cv2.imshow("Saved Image", saved_image)
163     cv2.waitKey(0)
164     cv2.destroyAllWindows()
165
166     if plot_type == 'Scatter_Plot':
167         plot = go.Scatter(x=df.index, y=column_data, mode='markers', name='Data Points')
168         layout = go.Layout(title=f"Scatter Plot of {column_name}",
169                            xaxis=dict(title='Index'),
170                            yaxis=dict(title=column_name))
171     elif plot_type == 'Box_Plot':
172         plot = go.Box(y=column_data, name='Data Distribution')
173         layout = go.Layout(title=f"Box Plot of {column_name}", yaxis=dict(title=column_name))
174     elif plot_type == 'Histogram':
175         plot = go.Histogram(x=column_data, name='Data Distribution', nbinsx=20)
176         layout = go.Layout(title=f"Histogram of {column_name}",
177                            xaxis=dict(title=column_name),
178                            yaxis=dict(title='Frequency'))
179     elif plot_type == 'Line_Chart':
180         plot = go.Scatter(x=df.index, y=column_data, mode='lines', name='Data Points')
181         layout = go.Layout(title=f"Line Chart of {column_name}",
182                            xaxis=dict(title='Index'),
183                            yaxis=dict(title=column_name))
184     else:
185         print("Unknown plot type")
186         return
187
188     fig = go.Figure(data=[plot], layout=layout)
189
190     if not os.path.exists(output_dir):
191         os.makedirs(output_dir)
192

```

9th Run

```
1  from app import create_app
2
3  app = create_app()
4
5  if __name__ == '__main__':
6      app.run(debug=False, host='0.0.0.0', port=5000)
7
```

4th system integration

Flask Framework:

```
app=Flask(__name__)

UPLOAD_FOLDER = 'uploads'
app.config['UPLOAD_FOLDER'] = UPLOAD_FOLDER
|
UPLOAD_FOLDER = 'uploads'
app.config['UPLOAD_FOLDER'] = UPLOAD_FOLDER

@app.route('/')
def Auth():
    return render_template('Auth.html')
```

Jinja

```
<select name="columns" class="form-control mt-3 fv py-5 Table-name " id="column-name" multiple>
  <option disabled value="">Select Which Column To Clean</option>
  <option value="all">all</option>

  {% for column in columns %}
  <option value={{column}}>{{ column }}</option>
  {% endfor %}
</select>

{%
  if table_original=="You haven't selected any Table!" %}
  <p style="margin-top: 3%;" class="alert alert-danger mx-auto w-100" role="alert">{{table_original}}</p>

{%
  elif table_original %}
  <div class="card">
    <div class="card-body" style="max-height: 400px; overflow-y: scroll;">
      {{ table_original|safe }}
    </div>
  </div>

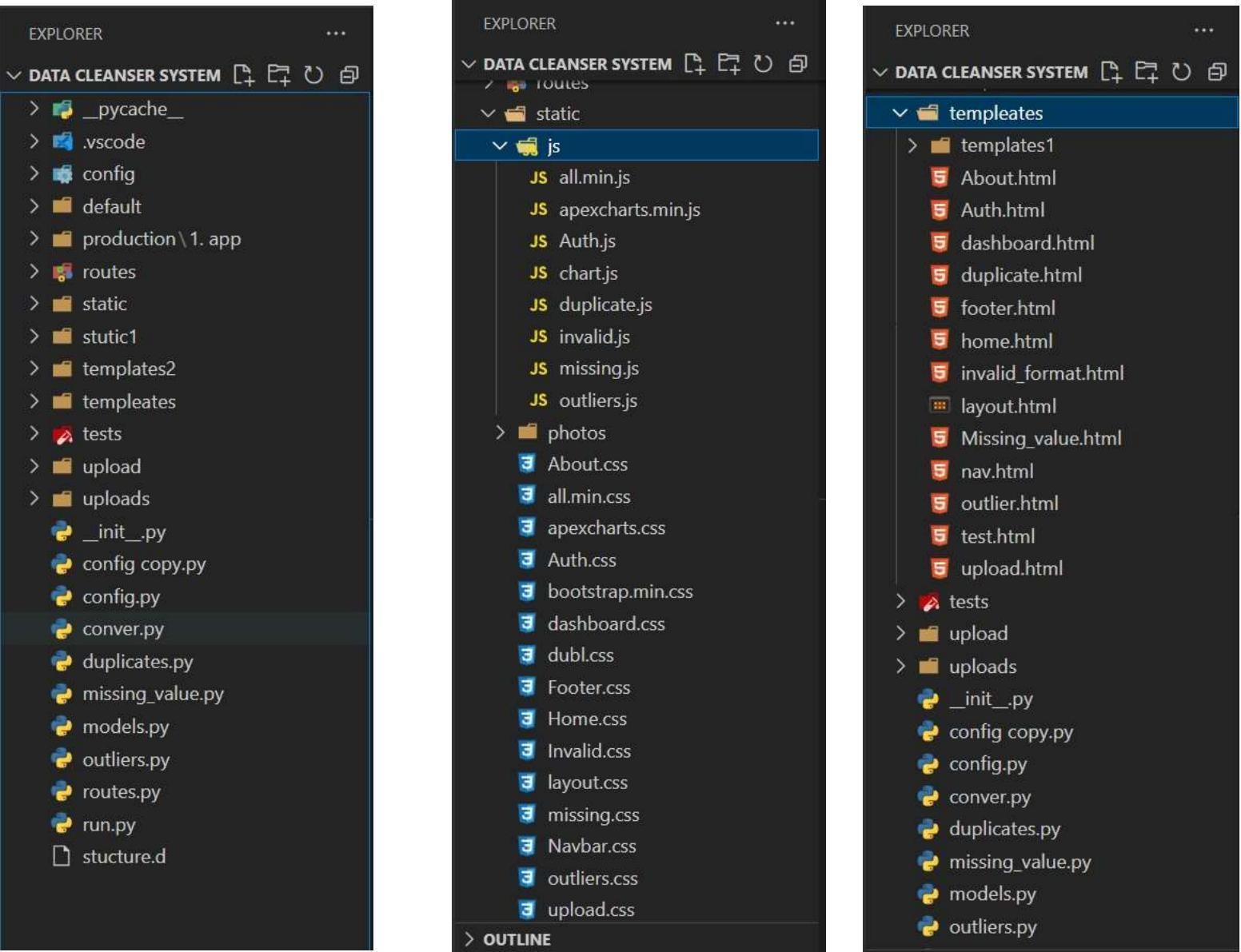
{%
  endif %}

```

Template Rendering with Jinja:

```
{% extends "layout.html" %}
```

Flask Structure Files:



8.5 References:-

- <https://pmsoftlifecycle.blogspot.com/>
- <https://jasonlemauk.com/the-agile-testers-mindset/>
- <https://openrefine.org/>
- <https://winpure.com/>
- <https://www.validity.com/demandtools/>
- <https://kn-it.com/>
- <https://sagemaster.io/>
- <https://www.blue-pencil.ca/data-cleansing-what-is-it-and-why-is-it-important/>
- <https://blog.hubspot.com/marketing/data-cleansing>
- <https://www.techtarget.com/searchdatamanagement/definition/data-scrubbing>
- <https://careerfoundry.com/en/blog/data-analytics/best-data-cleaning-tools/>
- <https://www.mrpeasy.com/blog/erp-system-implementation/>
- <https://www.oracle.com/eg-ar/applications/what-is-saas/>
- <https://jinja.palletsprojects.com/en/3.1.x/>
- <https://plotly.com/>
- <https://flask.palletsprojects.com/en/3.0.x/>
- <https://getbootstrap.com/>
- <https://www.tibco.com/glossary/what-is-data-cleansing>
- <https://aws.amazon.com/what-is/data-cleansing/>
- <https://www.kaggle.com/code/rtatman/data-cleaning-challenge-handling-missing-values>
- <https://www.kaggle.com/code/rtatman/data-cleaning-challenge-deduplication>
- <https://www.kaggle.com/code/rtatman/data-cleaning-challenge-imputing-missing-values>
- <https://www.kaggle.com/code/rtatman/data-cleaning-challenge-cleaning-numeric-columns>
- <https://www.kaggle.com/code/rtatman/data-cleaning-challenge-json-txt-and-xls>
- <https://easygui.sourceforge.net/>
- <https://community.fabric.microsoft.com/t5/Desktop/Simple-data-cleaning-split-columns-that-contain-multiple-values/td-p/3434064>
- <https://stats.stackexchange.com/questions/546421/do-i-split-before-or-after-cleaning-code>
- <https://towardsdatascience.com/how-to-clean-your-data-in-python-8f178638b98d?gi=d523f5fe8029>

<https://realpython.com/python-data-cleaning-numpy-pandas/>

<https://www.techtarget.com/searchdatamanagement/definition/data-scrubbing>

<https://blog.hubspot.com/marketing/data-cleansing>

<https://www.cazoomi.com/blog/dirty-data-how-much-it-can-cost-your-business-and-how-to-get-rid-of-it/>

<https://zenodo.org/records/1171077>

<https://software.recs.stackexchange.com/questions/11853/outlier-detection-in-latitude-longitude-csv-file>

<https://www.inciter.io/blog-posts/data-cleaning-in-excel-101-part-2-how-to-split-columns-in-excel>

<https://medium.com/@gregory2/4-easy-ways-to-handle-outliers-in-your-data-47f125a3f779>

https://www2.microstrategy.com/producthelp/Current/MSTRWeb/WebHelp/Lang_1033/Content/Functions_for_Wrangling_Your_Data.htm

<https://towardsdatascience.com/how-to-clean-your-data-in-python-8f178638b98d>

<https://www.dbmaestro.com/blog/database-automation/database-containerization>