

# **CSV DATA INSIGHTS**

Divyabarathi P ,Vaishnavi Sriraman

Bannari Amman Institute of Technology divyabarathip@bitsathy.ac.in\_, svaishnavi2806@gmail.com

#### ABSTRACT

To create an all-inclusive and intuitive CSV file preparation and analysis tool. Through the uploading of CSV files and the retrieval of vital data, this technology enables users to receive profound insights into their datasets. The program automatically offers information upon file upload, including the quantity of columns, the detection of absent and duplicate rows, and a thorough breakdown of the data types linked to each column. Additionally, our research uses univariate analysis to provide customers with useful statistical insights into the distribution and properties of their data for both numerical and categorical columns. With the help of this functionality, users may make more informed decisions by having a deeper grasp of the distribution patterns in their data. Our application provides full insights when working with either numerical variables that require statistical measurements or categorical variables that require frequency analysis. We increase user engagement by using Plotly's power to present the results through an interactive dashboard. Users can easily examine patterns, trends, and anomalies in their datasets thanks to this visual representation. To further enhance the utility of our tool, we have implemented a user interface (UI) that facilitates data preprocessing. Users can conveniently remove unwanted columns and perform data type conversions directly through the interface. This streamlines the data preparation process, saving time and effort for individuals working with diverse datasets. Our project addresses the need for a comprehensive and accessible tool for users dealing with CSV files, providing a holistic solution for file analysis, visualization, and preprocessing. The

combination of informative insights and interactive features ensures that users can efficiently manage and optimize their datasets with ease.

#### **CHAPTER - I**

#### INTRODUCTION

The effective administration and investigation of datasets has become critical in the age of expanding data use. Our project aims to address this demand by offering a complete solution for people working with CSV files. The primary goal of the project is to enable users to automatically extract insightful information from their datasets when they submit files.

After uploading a CSV file, users receive a comprehensive summary that includes the basic attributes of their data, such as the number of columns, the presence of missing and duplicate rows, and an indepth examination of the data types linked to each column. In addition, the project goes beyond traditional file information retrieval by integrating univariate analysis for both numerical and categorical columns, providing users with a more profound comprehension of the distribution and properties of their data.

Additionally, the project uses Plotly to provide an interactive dashboard that improves user engagement and the comprehension of data trends. With the help of this dashboard, users may examine trends within

their datasets in an easy manner by converting raw data into dynamic, aesthetically pleasing displays. This capability is enhanced with a user interface (UI) that facilitates smooth data preparation. The interface makes it simple for users to carry out operations like deleting particular columns and converting data types, which expedites the data preparation process.

Our initiative seeks to close the gap between raw data and actionable insights in an environment where effective data analysis and preprocessing are essential to well-informed decision-making. Through the integration of interactive data visualization, userfriendly preprocessing, robust univariate analysis, and automated information extraction, our tool aims to be a useful resource for people in a variety of fields, easing the challenges related to managing and exploring CSV files.

#### **APPLICATIONS**:

Our project's application extends to various domains where effective data exploration and preprocessing are pivotal. Below are some practical applications highlighting the versatility and usefulness of our tool:

#### **Business Analytics:**

Business analysts can leverage the tool to quickly assess the structure and quality of large datasets, facilitating efficient decision-making. Univariate analysis and interactive dashboards enable them to derive actionable insights from categorical and numerical data.

#### Data Science and Research:

Data scientists and researchers can benefit from the automated information extraction and univariate analysis functionalities. The tool aids in the initial exploration of datasets, saving time and providing a solid foundation for further in-depth analysis.

# CHAPTER 2

# LITERATURE REVIEW

[1] Shengjia Cao, YunhanZeng, Shangru Yang, Songlin Cao have introduced data visualization and its technologies, using python's matplotlib and pyecharts libraries to understand data visualization. Demonstrated examples for the correct use of chars, visual processing of data in various fields for the effective visualization

[2] R. Toasa, M. Maximiano, C. Reis and D. Guevara studied the available techniques of data visualization, implemented a generic and dynamic dashboard based of real time information to interact with the data, based on initial set of hints, charts, tables and reports and conducted literature review of data visualization and its existing dashboard platforms

[3] G. Singh, A. Kumar, J. Singh and J. Kaur, visualized the dataset of the covid-pandemic to provide beneficial information for finding possible solution using Power BI, analyzed the data, represented the data in stacked column charts, tables and maps and explained how effective the three ways are to understand the impact of covid to the world

[4] T. Liang, S. Lu and Q. Liu, proposed a strategy to construct a data visualization tool or system based on big data analysis technology which combines big data analysis technology, relying on the cloud platform, to show the various varieties of data in a variety of ways, which are convenient for users to analyze and study the insights behind the data. This results the enhanced user experience and intuitiveness of data information eventually improved the efficiency of data analysis

[5] Stančin., et al. considers more than 20 libraries and separate them into six groups: core libraries, data preparation, data visualization, machine learning, deep learning and big data. The authors recommends the libraries such as pandas for data preparation; Matplotlib, seaborn or Plotly for data visualization; scikit-learn for machine leraning; TensorFlow, Keras and PyTorch for deep learning; and Hadoop Streaming and PySpark for big data.

[6] McKinney., et al. described that Pandas library provide labelled and structure based data for grouping and aggregation of data. This paper focused on how to contain multiple tables in each to other.

[7] . Mitrpanont., et al. Python panda's library use to data analysis as well as weka. Weka also work on data manipulation and data analysis. In this Panda's Data Frame function set any time of data like json format data set in tabular format using function. It's store in structure format. Python is support Data science using any type of extension file call in tabular or structure format data and search top and bottom data. Panda's support Matplotlib library it's support graphics and 3-D animation. It's use files data display in graph format

[8] Abhinav Nagpal, Goldie Gabrani. The rapid growth of computers has led to a revival of data science and analytics, a branch of computer science. Python, a low-learning language, has emerged as a suitable programming solution for developing data science applications. Its ever-evolving libraries make it an ideal choice for data analytics. This paper discusses the features and characteristics of Python, its reasons for its rapid growth, and its position at the forefront of data science applications, research, and development.

[9] Shailendra Chaudekar, Python is used in this study's data analysis. Investigating the programming language. The basic steps of data analysis, including cleaning, transforming, and modeling data, are briefly covered in this article with a focus on exploratory data analysis. Analyzing data from an existing dataset and drawing conclusions Demonstrated some graphical data analysis from the dataset using multiple Python tools and techniques. [10] Shengjia Cao, YunhanZeng, Shangru Yang, Songlin Cao, Data visualization uses appropriate graphical representations, such as histograms, bar charts, column charts, and other statistical charts, to communicate data and picture information more effectively and efficiently. With its abundance of third-party libraries, open-source communities, and constantly improved documentation for data visualization, Python is a great option. This article demonstrates the realization of data visualization using the Matplotlib and Pyecharts libraries for Python.

[11] Chanin Nantasenamat, Avratanu Biswas, J.M. Nápoles-Duarte, Mitchell I. Parker, Roland L. Dunbrack Jr explores the usage of Streamlit for the development of software and tools in the field of bioinformatics.

# CHAPTER 3 OBJECTIVES AND METHODOLOGY

# 3.1 Aim:

To develop an advanced CSV file analysis and preprocessing tool that empowers users to efficiently explore, understand, and manipulate their datasets. By leveraging cutting-edge technologies and intuitive design, we strive to simplify the complexities associated with data management, providing users with a comprehensive solution for seamless file analysis and preparation.

# 3.2 Objectives:

- Develop a user-friendly web-based tool to automate the analysis and preprocessing of CSV files.
- Provide users with comprehensive information about the uploaded CSV file,

 USREM
 International Journal of Scientific Research in Engineering and Management (IJSREM)

 Volume: 08 Issue: 03 | March - 2024
 SJIF Rating: 8.176
 ISSN: 2582-3930

including the number of columns, missed rows, duplicate rows, and data types of columns.

- Perform univariate analysis of categorical and numerical columns to identify patterns and distributions within the data.
- Implement an interactive dashboard using Plotly to visualize the analysis results and provide users with insights into their data.
- Design a user interface (UI) with intuitive controls to allow users to preprocess the data, such as removing columns and converting data types.
- Ensure the scalability and efficiency of the tool to handle large CSV files and datasets efficiently.
- Validate the tool's effectiveness through user testing and feedback to identify areas for improvement.

# 3.3 Languages and Packages:

#### 3.3.1 Python:

Python's popularity can be attributed to a number of factors, such as its object-oriented design, modularity, portability, testing, and self-documentation capabilities; additionally, its presence of a Numeric library facilitates the efficient handling and storing of massive amounts of numerical data.

Additionally, it offers high-level data structures including lists, classes, exceptions, dynamic type modules, associative arrays, dynamic binding, and intelligent memory management. It can be expanded by adding external libraries, and its little kernel can be used. A vast library of standard extensions, written in Python and other languages like C or C++, are included in its distribution. These extensions cover a variety of tasks including manipulating strings, testing and profiling tools, debugging, web-related utilities, operating system services, and regular expressions akin to Perl. New modules can be created to expand the language.

Features are:

1. High readability and little coding

2. Mobility and Adaptability

3. Independent of Platforms

4. Equilibrium High-Level and Low-Level Programming

Python trails only Ruby in the number of lines to code. Thus, for many ML and AI-based applications, this is the primary reason why Python is chosen by many enterprises.

3.3.2. Pandas:

An efficient way to work with data structures is with the Python panda. A superior platform for statistical computation and data analytics is offered by the Pandas library. Additionally, the Pandas library has SQL tools for data processing, including the ability to integrate data using different joins (inner, right, and left). The Panel data analytics package, which seeks to offer comparable capabilities and has incorporated several features including automated data alignment and hierarchical indexing, is where the name "Pandas" originated.

Pandas library is a contemporary object-oriented highlevel programming library that includes a significant number of add-on packages.

Pandas has an integrated library called Numpy, which is used for numerical data.

Numpy uses the arraydata type for operations like indexing, sorting, and reshaping. Numpy's array data type allows for homogeneous data. Pandas supports Matplotlib for graphical data presentation and file saving in several formats, including excel, csv, and JSON.

Pandas supports three types of data structures: series, data frames, and panels. The series data structure consists of a one-dimensional array. It supports homogeneous data types. Data Frames are two-dimensional arrays that may hold a variety of data types and be sized and modifiable. The panel has a three-dimensional array.



#### 3.3.3. Plotly:

Plotly is an open-source Python graphing package that allows you to create stunning and interactive infographics. It's an excellent tool for identifying trends in a dataset before diving into machine learning models.

Features:

- Visualizations are interactive, unlike Seaborn and Matplotlib.
- The high-level Express API makes it simple to create complex visuals.
- Plotly Dash is a framework for hosting visualizations and machine learning projects.
- Can produce HTML code for your visualizations, which you can then integrate on the website.

# 3.3.4 Matplotlib Vs Plotly

Matplotlib and Plotly are Python libraries for data visualization. Matplotlib is a popular toolkit that is ideal for making static visualizations, but Plotly is a more advanced tool that is better suited to producing complicated charts more effectively. Matplotlib is more specific in specifying each plot element, making it an excellent place for novice Python users to begin, but Plotly is well-suited for building interactive plots that can be seen in a web browser.

3.3.5 Plotly vs Other Visualization libraries Matplotlib and Plotly are Python libraries for data visualization. Matplotlib is a popular toolkit that is ideal for making static visualizations, but Plotly is a more advanced tool that is better suited to producing complicated charts more effectively. Matplotlib is more specific in specifying each plot element, making it an excellent place for novice Python users to begin, but Plotly is well-suited for building interactive plots that can be seen in a web browser.

#### 3.3.6 Streamlit

Streamlit is the fastest method to create data applications. It is an open-source Python toolkit that enables you to create web apps for sharing analytical results, creating sophisticated interactive experiences, and iterating over new machine learning models. Furthermore, designing and deploying Streamlit apps is extremely rapid and versatile, decreasing application development time from days to hours.

# 3.4 Methodology:

3.4.1 Requirement Gathering:

- Conducted research to understand user needs and requirements for the tool.
- Document the specific features and functionalities desired by users to guide the development process.

3.4.2 Design and Architecture:

- Design the overall architecture of the tool, including frontend UI components and backend processing modules.
- Choose appropriate technologies and frameworks for implementing the UI (e.g., Streamlit) and backend logic (e.g., Python).
- Plan the data flow and interactions between different components of the system.

3.4.3 Implementation:

- Develop the frontend UI using Streamlit, incorporating features for file upload, data analysis, and visualization.
- Implement backend processing logic in Python, including functions for analyzing CSV files, performing univariate analysis, and generating interactive dashboards using Plotly.
- Integrate the frontend and backend components to ensure seamless communication and data exchange.

3.4.4 Testing and Validation:

- Conduct thorough testing of the tool to ensure accuracy, reliability, and robustness.
- Perform unit tests for individual components as well as integration tests to validate the overall functionality.



#### 3.4.5 Optimization and Performance Tuning:

- Optimize the tool for performance and scalability, especially when handling large datasets.
- Implement caching mechanisms and data processing optimizations to improve efficiency and responsiveness.

#### 3.4.6 Documentation:

- Prepare comprehensive documentation for the tool, including user manuals, installation guides, and developer documentation.
- Document the design decisions, implementation details, and usage instructions to facilitate understanding and usage of the tool.

#### 3.4.7 Deployment and Maintenance:

- Deploy the tool to a production environment, ensuring compatibility and stability across different platforms.
- Establish procedures for ongoing maintenance, including bug fixes, updates, and enhancements based on user feedback and evolving requirements.
- Provide user support and training to ensure successful adoption and usage of the tool.

By following this methodology, the project aims to deliver a high-quality and user-friendly tool for automating the analysis and preprocessing of CSV files, enabling users to gain valuable insights into their data and streamline their data analysis workflows effectively.

#### CHAPTER - 4

#### **PROPOSED WORK MODULES**

4.1 User Interface (UI) Development Module:

- Design and develop a visually appealing and user-friendly UI using modern web frameworks like Streamlit
- Create an intuitive interface with clear instructions for uploading CSV files and displaying analysis results.
- Implement interactive elements such as buttons, sliders, and dropdown menus for preprocessing options.
- Ensure responsiveness and compatibility across different devices and screen sizes for optimal user experience.

#### 4.2 CSV File Analysis Module:

- Develop algorithms to parse and analyze the uploaded CSV file efficiently.
- Implement logic to calculate the number of columns, identify missing rows, and detect duplicate rows.
- Utilize data manipulation libraries such as Pandas to extract and analyze data types of each column.
- Perform univariate analysis on categorical and numerical columns to compute summary statistics, frequencies, and distributions.
- 4.3 Plotly Dashboard Creation Module:
  - Design and build an interactive dashboard using Plotly, a powerful visualization library in Python.
  - Develop charts, graphs, and plots to visualize the results of the CSV file analysis in an intuitive and informative manner.
  - Customize the dashboard layout and styling to enhance readability and user engagement.
  - Implement dynamic updates to the dashboard based on user interactions and preprocessing actions.



SJIF Rating: 8.176

#### 4.4 Data Preprocessing Module:

- Implement preprocessing functionalities to manipulate the data based on user preferences and requirements.
- Develop logic to remove columns, handle missing values (e.g., imputation or deletion), and convert data types (e.g., numeric to categorical).
- Ensure data integrity and consistency throughout the preprocessing steps to avoid errors or inconsistencies in the analysis results.
- Provide feedback to the user on the impact of preprocessing actions on the dataset through the UI.

4.5 Backend Processing and Integration Module:

- Develop backend logic to handle file upload, data analysis, and preprocessing tasks using Python.
- Integrate the backend with the frontend UI to enable seamless communication and data exchange.
- Implement error handling and validation mechanisms to handle edge cases and unexpected inputs gracefully.
- Optimize processing algorithms for performance and scalability, especially when dealing with large CSV files or datasets.

4.6 Testing and Quality Assurance Module:

- Conduct thorough testing of each module to ensure functionality, reliability, and robustness.
- Perform unit tests for individual functions and integration tests for the entire system to validate end-to-end functionality.
- Utilize testing frameworks and tools like pytest to automate test execution and identify potential bugs or issues.
- Solicit feedback from beta testers and stakeholders to identify usability issues and areas for improvement.

- 4.7 Documentation and Deployment Module:
  - Prepare comprehensive documentation for the project, including user manuals, installation guides, and developer documentation.
  - Document the design decisions, implementation details, and usage instructions for each module to facilitate understanding and future maintenance.
  - Deploy the application to a production environment, ensuring proper configuration and setup for optimal performance and security.
  - Provide ongoing support and maintenance, including bug fixes, updates, and enhancements based on user feedback and evolving requirements.

By organizing the project into these proposed work modules, the development process can be managed effectively, with clear objectives and responsibilities for each component. This approach ensures systematic implementation and thorough testing of the CSV file analysis and preprocessing tool, resulting in a highquality and user-friendly application that meets the needs of its intended users.

# CHAPTER -5

#### **RESULTS AND DISCUSSIONS**

Upon completion of the project, the developed CSV file analysis and preprocessing tool demonstrated robust functionality and effectiveness in fulfilling its objectives. The results obtained from the tool can be summarized as follows:

5.1 Comprehensive Information Extraction:

The tool successfully extracted essential information from the uploaded CSV file, including the number of columns, missed rows, and duplicate rows. Users were provided with detailed insights into the structure and quality of their data, enabling them to make informed decisions regarding data preprocessing and analysis.



#### 5.2 Data Type Identification:

The tool accurately determined the data types of each column in the CSV file, distinguishing between categorical and numerical data. This information was crucial for performing subsequent univariate analysis and preprocessing operations, ensuring the integrity and consistency of the data throughout the analysis process.

#### 5.3 Univariate Analysis:

The tool conducted thorough univariate analysis on both categorical and numerical columns, generating summary statistics, frequency distributions, and visualizations to aid in data exploration. Users were able to gain insights into the distribution and characteristics of their data, facilitating further analysis and decision-making.

#### 5.4 Interactive Dashboard:

The interactive dashboard created using Plotly provided users with a dynamic and intuitive interface for visualizing analysis results. Users could interact with various charts and plots to explore different aspects of their data and identify patterns or outliers effectively. The dashboard enhanced user engagement and facilitated data-driven decision-making.

#### 5.5 Preprocessing Capabilities:

The tool's preprocessing functionalities, including column removal and data type conversion, enabled users to clean and prepare their data for analysis efficiently. Users could customize preprocessing options based on their specific requirements, streamlining their data preparation workflow and improving overall productivity.

#### 5.6 Discussion:

The results obtained from the developed CSV file analysis and preprocessing tool demonstrate its effectiveness in providing users with valuable insights into their data and facilitating data preparation for further analysis. The following discussion highlights the key findings and implications of the project:

#### 5.7 Enhanced Data Understanding:

By providing users with comprehensive information about their CSV files and conducting univariate analysis, the tool enhances users' understanding of their data. Users can identify data quality issues such as missing values or duplicate entries, allowing them to take corrective actions to improve data integrity.

#### 5.8 Streamlined Data Preparation:

The tool's preprocessing capabilities contribute to streamlining the data preparation process, reducing manual effort and potential errors. Users can easily preprocess their data using the interactive UI, saving time and resources while ensuring data consistency and accuracy.

5.9 Facilitated Data Exploration:

The interactive dashboard created using Plotly empowers users to explore their data visually and interactively. Users can identify trends, patterns, and outliers more effectively, leading to deeper insights and better decision-making in data analysis projects.

5.10 User-Centric Design:

The user interface of the tool is designed with the end user in mind, prioritizing usability, intuitiveness, and accessibility. User feedback and usability testing were incorporated throughout the development process to ensure that the tool meets the needs and preferences of its intended users.

5.11 Scalability and Adaptability:

The tool's architecture and design allow for scalability and adaptability to different use cases and datasets. It can handle large CSV files and datasets with ease, making it suitable for a wide range of applications and industries.

In conclusion, the developed CSV file analysis and preprocessing tool represents a significant contribution to the field of data analysis and data preparation. Its robust functionality, user-friendly design, and interactive features make it a valuable asset for researchers, analysts, and data scientists seeking to extract insights from CSV data efficiently and effectively. Future enhancements and refinements to the tool could include additional analysis



capabilities, support for other file formats, and integration with external data sources for a more comprehensive data analysis experience.

#### **CHAPTER - 6**

# CONCLUSIONS AND SUGGESTIONS FOR FUTURE WORK

In conclusion, the CSV file analysis tool developed with Streamlit for the user interface and Python for the backend presents a valuable solution for users seeking efficient data analysis and preprocessing capabilities. The project successfully addresses the need for automating the analysis of CSV files, providing users with essential information such as the number of columns, missed and duplicate rows, and data types of columns. Additionally, the inclusion of univariate analysis for both categorical and numerical columns, along with interactive dashboards created using Plotly, enhances the tool's utility by enabling users to gain insights and visualize their data effectively.

#### Suggestion for future work

- Sign up and login pages should be integrated
- Maintaining the user details and dataset used by the respective users, so that the preprocessing process can't be repeated
- Making the project flexible for large volume of dataset
- Implementing data modeling, so that a user can upload more than one dataset and perform data analysis

#### REFERENCES

[1] Dr Ossama Embarak, Embarak, and Karkal. Data analysis and visualization using python. Springer, 2018.

[2] Michel Jambu. Exploratory and multivariate data analysis. Elsevier, 1991.

[3] Matthieu Komorowski, Dominic C Marshall, Justin D Salciccioli, and Yves Crutain. Exploratory data analysis. Secondary analysis of electronic health records, 2016.

[4] https://stackoverflow.com

[5] https://github.com

[6] KD Nuggets poll result https://www.kdnuggets.com/

[7] Guido Van Rossum et al. Python programming language. In USENIX annual technical conference, 2007.

[8]https://towardsdatascience.com/a-guide-to-pandasandmatplotlib-for-data-exploration-56fad95f951c

- [9] https://python.plainenglish.io/
- [10] https://streamlit.io/
- [11] https://docs.kanaries.net