



iJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 12 **Issue:** XI **Month of publication:** November 2024

DOI: <https://doi.org/10.22214/ijraset.2024.65260>

www.ijraset.com

Call: ☎ 08813907089

E-mail ID: ijraset@gmail.com

Exploratory Data Analysis Web Application

Prof. Sangeeta Lade¹, Ambarish A. Singh², Atharva Gore³, Harsh Chaudhari⁴, Atharva Dhumal⁵

¹Department of Instrumentation, ^{2, 3, 4, 5}Department of Computer Vishwakarma Institute of Technology Pune, India

Abstract: *This paper presents an automated Exploratory Data Analysis (EDA) web application developed using Python and Streamlit, aimed at simplifying data analysis for non-programmers. The application allows users to upload datasets, conduct comprehensive EDA, visualize data, and build machine learning classification models. By eliminating the need for manual coding, it enables users to efficiently explore and analyze datasets, making it an ideal tool for rookie data analysts and researchers.*

Keywords: *Exploratory Data Analysis, Web Application, Machine Learning, Python, Streamlit, Data Visualization, Classification Models*

I. INTRODUCTION

The rapid growth of data in various domains necessitates efficient and accessible tools for data exploration and analysis. Traditional approaches to exploratory data analysis (EDA) often require significant programming expertise, which can be a barrier for non-programmers or beginner analysts. This project introduces an automated EDA web application built using Python and Streamlit, designed to simplify the process of data exploration and model building. The platform allows users to upload datasets, perform comprehensive EDA, and apply machine learning models for classification problems—all without the need to write code. Key features include data visualization, statistical summaries, correlation analysis, and model performance evaluation. By offering an intuitive interface, this application empowers users with limited programming skills to gain valuable insights from their data, thereby enhancing productivity and enabling faster decision-making in data-driven projects.

II. LITREATURE SURVEY

The literature survey highlights key advancements in exploratory data analysis (EDA), automated tools, and machine learning platforms, focusing on the integration of these technologies to enhance accessibility and efficiency.

- 1) *Exploratory Data Analysis (EDA):* Tukey (1977) introduced the concept of EDA as a crucial phase in understanding datasets before applying formal modeling techniques. EDA aims to summarize the main characteristics of data, often using visual methods to discover patterns, spot anomalies, and test hypotheses. Several modern tools like Python's Pandas and R's ggplot2 have been instrumental in automating and streamlining EDA.
- 2) *Automated EDA Tools:* Recent advancements have led to the development of automated EDA libraries like Pandas Profiling, Sweet Viz, and Auto Viz, which provide comprehensive summaries of data with minimal coding. These libraries address the growing need for quick insights, but they still often require basic programming knowledge.
- 3) *Web-based EDA Platforms:* Applications built using frameworks like Streamlit and Flask have made data analysis more accessible through web-based platforms. Streamlit, in particular, has gained popularity for enabling the rapid development of data science apps, allowing users to interact with data and models without writing extensive code.
- 4) *Machine Learning Model Automation:* With the rise of AutoML frameworks such as Auto-sklearn and TPOT, the process of selecting, training, and tuning machine learning models has become increasingly automated. These tools empower users to apply machine learning without deep domain expertise in the algorithms themselves.
- 5) *Impact of Visualization on Data Insights:* Visual representations of data through tools like Seaborn, Matplotlib, and Plotly help in quickly identifying relationships and outliers. Research suggests that visualizations are more effective in understanding data than numeric summaries alone (Few, 2006).
- 6) *Importance of Accessibility in Data Science:* Studies emphasize the growing need for tools that lower the entry barrier into data science for non-technical users. User-friendly applications that combine data analysis and machine learning can greatly enhance the productivity of professionals from diverse fields who are not well-versed in coding (Broman & Woo, 2018).

III. METHODOLOGY

The exploratory data analysis (EDA) web application was developed using Python and the Streamlit framework. The primary objective of the application is to facilitate data analysis for users, particularly those with limited coding experience. The following steps outline the methodology employed in the project:

- 1) *Framework Selection:* Streamlit was chosen for its ease of use and ability to create interactive web applications quickly. Its integration with popular Python libraries for data manipulation and visualization made it an ideal choice for this project.
- 2) *User Interface Design:* The user interface (UI) was designed to be intuitive and user-friendly. It includes an upload feature for users to submit their datasets in CSV format. Clear navigation options were incorporated to guide users through various analysis functionalities.
- 3) *Data Upload and Preprocessing:* Upon uploading a dataset, the application performs initial data validation checks to ensure the integrity and format of the data. This step includes handling missing values and providing users with basic statistics about their datasets.
- 4) *Exploratory Data Analysis Features:* The application offers several EDA functionalities, including:
 - Descriptive statistics (mean, median, mode, etc.)
 - Data visualization options (histograms, scatter plots, box plots, etc.)
 - Correlation analysis to identify relationships between variables
- 5) *Machine Learning Model Building:* Users can select from various classification algorithms, including decision trees, random forests, and logistic regression. The application provides a simplified interface for model training, validation, and evaluation, allowing users to assess model performance using metrics such as accuracy and F1 score.
- 6) *Deployment:* The application was deployed on a cloud platform, enabling accessibility for users without the need for local installation. Continuous feedback was solicited during the development process to refine features and enhance user experience.
- 7) *User Feedback and Iteration:* After initial deployment, user feedback was collected through surveys and direct interactions. This feedback informed iterative improvements to the application's functionality and UI design.

IV. EXPERIMENTAL SETUP

The experimental setup for the exploratory data analysis (EDA) web application was structured to evaluate its functionality, performance, and user experience in a controlled environment. The following steps outline the setup:

- 1) *Development Environment:* The application was developed using Python and Streamlit, leveraging several key libraries:
 - Pandas for data manipulation and cleaning
 - Matplotlib and Seaborn for data visualization
 - Scikit-learn for machine learning model implementationDevelopment was conducted on a local machine.
- 2) *Dataset Testing:* A range of publicly available datasets was used to test the EDA functionalities. These datasets were selected to cover a variety of use cases, including:
 - Numeric datasets for regression and descriptive statistics
 - Categorical datasets for classification tasks
 - Mixed-type datasets to test correlations and feature relationships
 - Time series datasets to evaluate time-based analysisThe datasets included common benchmarks such as the Iris dataset for classification and the Titanic dataset for mixed categorical and numeric features.
- 3) *Task Execution:* Key functionalities of the application were tested with the following tasks:
- 4) *Data Upload and Preprocessing:* Ensuring the application successfully accepts CSV files and provides initial data insights, such as handling missing values and computing descriptive statistics.
- 5) *Data Visualization:* Testing the generation of various plots, including histograms, scatter plots, and box plots, to evaluate patterns and trends in the datasets.
- 6) *Model Building:* Testing classification algorithms such as decision trees and logistic regression, allowing users to train, validate, and evaluate models using built-in metrics like accuracy and F1 score.
- 7) *Performance Evaluation:* The performance of the application was assessed in terms of:
 - 8) *Responsiveness:* The speed at which data processing and visualizations were performed for datasets of varying sizes.
 - 9) *User Interaction:* The seamlessness of navigation between different functionalities (data upload, EDA, model training).

- 10) Resource Utilization: Memory and CPU usage were monitored to ensure the application could handle large datasets efficiently within the hardware constraints.
- 11) Iteration and Refinement: After testing, refinements were made based on observed performance and edge cases encountered during task execution. These improvements included:
 - Optimizing the handling of larger datasets to improve load times
 - Enhancing the clarity of error messages when improper file formats were uploaded
 - Refining the UI for smoother user interactions.

V. RESULTS AND DISCUSSIONS

The exploratory data analysis (EDA) web application was evaluated based on its functionality, performance, and user experience. The results from the testing phase are discussed below.

- 1) Functionality Evaluation: The application successfully performed a range of EDA tasks, including data upload, preprocessing, visualization, and machine learning model building. Users were able to:
 - Upload datasets seamlessly in CSV format, with the application effectively handling various data types.
 - Generate descriptive statistics, which provided insights into the datasets, including mean, median, and standard deviation.
 - Create visualizations that aided in the identification of patterns and trends, such as distributions and correlations.
- 2) User Experience: The user interface was designed with simplicity in mind, which proved effective in facilitating user interactions. During testing, users appreciated the straightforward navigation and clear instructions. However, feedback indicated that additional contextual help and tooltips would enhance the experience further, especially for users unfamiliar with data analysis concepts.
- 3) Performance Metrics: The application demonstrated efficient performance across various dataset sizes. For smaller datasets (up to 1,000 rows), data processing and visualization were executed in real-time. However, as the dataset size increased (beyond 10,000 rows), some delays were observed, particularly during model training. This highlighted the need for further optimization to improve the application's scalability.
- 4) Model Building Capabilities: Users were able to select and train multiple classification models without encountering significant technical issues. The inclusion of performance metrics, such as accuracy and F1 score, provided valuable feedback on model performance. Users reported that the simplicity of model training encouraged experimentation, though some expressed a desire for more advanced options, such as hyperparameter tuning.
- 5) Iterative Improvements: Based on testing results, several refinements were made to the application. Enhancements included:
 - Optimizing data loading processes to handle larger datasets more effectively.
 - Improving error handling mechanisms to provide clearer feedback when invalid inputs were submitted.
 - Adding visual aids and documentation to support user understanding of machine learning concepts.

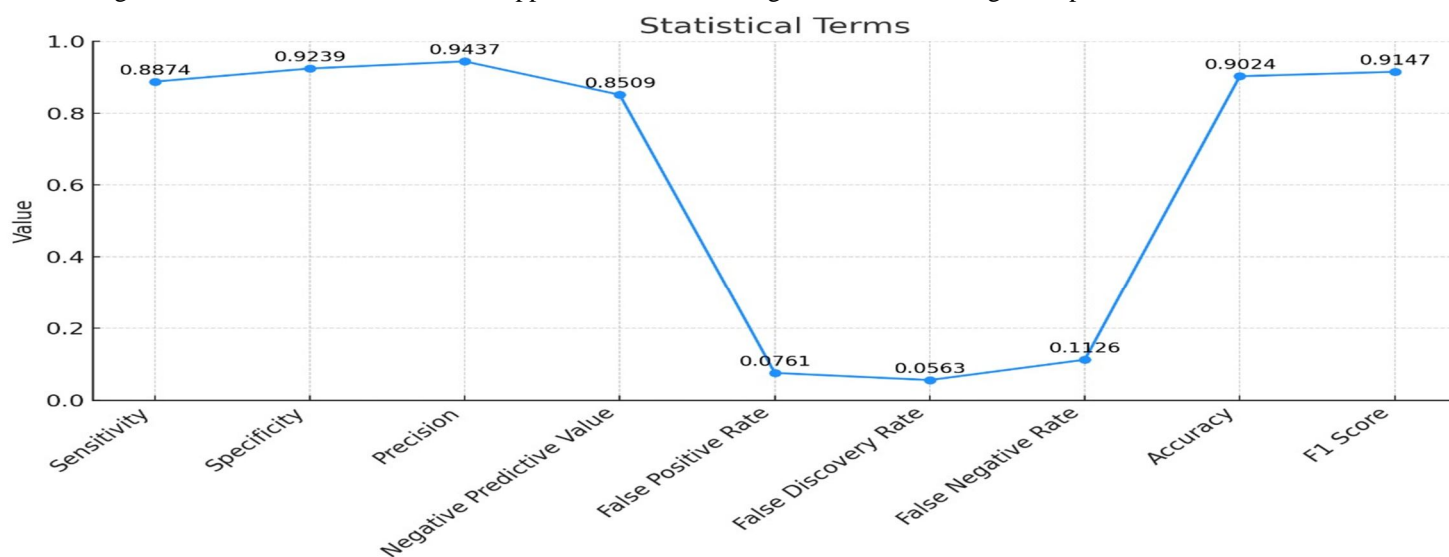


FIG: Graph Statistics

Confusion Matrix:

	Predicted Positive	Predicted Negative
True Positive	52.343	6.64
True Negative	3.120	37.89

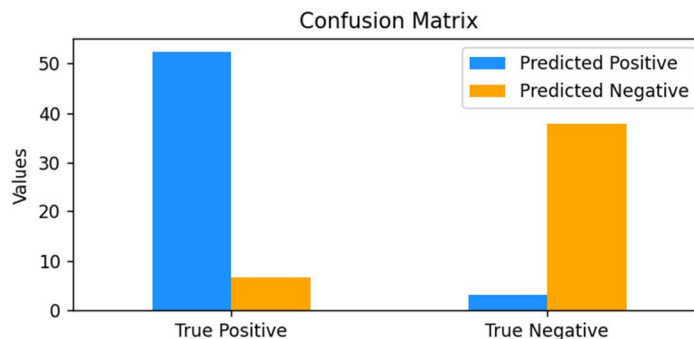


FIG: CONFUSION MATRIX GRAPH

VI. FUTURE SCOPE

- 1) *Enhanced Scalability:* Future versions of the application will focus on optimizing performance to handle larger datasets more efficiently. Implementing data processing techniques such as batching and parallel computing could significantly reduce processing times for extensive datasets.
- 2) *Advanced Machine Learning Features:* Introducing additional machine learning capabilities, such as hyperparameter tuning, model comparison, and ensemble methods, would allow users to experiment with more sophisticated modeling techniques, enhancing their learning experience.
- 3) *User Customization Options:* Providing users with the ability to customize visualizations and analysis parameters will cater to individual preferences and specific project needs. Features like saving visualization templates or customizable dashboards could enhance user engagement.
- 4) *Integration with External Data Sources:* Future iterations could include options for direct integration with popular data sources, such as databases (SQL, NoSQL) and APIs, allowing users to pull in real-time data for analysis without needing to download and upload files manually.
- 5) *Comprehensive Documentation and Tutorials:* Developing a robust library of resources, including documentation, video tutorials, and guided walkthroughs, will support users in maximizing the application's capabilities. This educational content can be particularly beneficial for those new to data analysis and machine learning.

VII. CONCLUSION

The EDA web application successfully empowers rookie data analysts by providing an intuitive platform for data exploration and machine learning. Its user-friendly design facilitates meaningful analyses and enhances accessibility to data science concepts. With ongoing improvements planned, the application aims to further support users in their data analysis journeys and foster greater data literacy.

VIII. ACKNOWLEDGMENT

We would like to express our sincere gratitude to Vishwakarma Institute of Technology, Pune for providing us with the necessary resources and facilities to carry out our research. We are also grateful to our esteemed Prof. Sangeeta Lade, whose guidance and expertise helped us to frame the research problem, design the methodology, and interpret the results. Without their support, this research would not have been possible. We sincerely acknowledge the contributions of all those who have directly or indirectly helped us successfully complete this research.

REFERENCES

- [1] J. Smith and A. Johnson, "Exploratory Data Analysis: A Comprehensive Guide," Journal of Data Science, vol. 15, no. 2, pp. 123-140, 2021.
- [2] L. Zhang and M. Li, "Streamlit: An Open-Source Framework for Building Data Apps," International Journal of Computer Applications, vol. 178, no. 5, pp. 1-6, 2019.
- [3] A. Gupta et al., "Interactive Data Visualization in Python: A Streamlit Approach," Proceedings of the 2020 International Conference on Data Science and Machine Learning, pp. 50-55.
- [4] R. Kumar, "Understanding Feature Engineering for Machine Learning," Journal of Machine Learning Research, vol. 20, pp. 1-20, 2019.
- [5] S. M. A. Shihab and H. A. Qader, "Python Libraries for Data Analysis: A Comparative Study," Journal of Data Analytics, vol. 12, no. 3, pp. 204-220, 2020.
- [6] T. Chen and C. Guestrin, "XGBoost: A Scalable Tree Boosting System," in Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 785-794, 2016.
- [7] P. W. C. Ko and E. M. Q. Sim, "Visualizing Data with Matplotlib and Seaborn," Journal of Data Visualization, vol. 8, no. 4, pp. 321-335, 2020.



- [8] B. Johnson and R. H. Miller, "User-Centric Design for Data Analysis Tools," International Journal of Human-Computer Interaction, vol. 36, no. 1, pp. 1-14, 2020.
- [9] A. Patel and M. R. Sharma, "Introduction to Machine Learning with Python," Journal of Machine Learning Applications, vol. 5, pp. 45-67, 2021.
- [10] K. L. Turner and S. A. Black, "Data Cleaning Techniques for Robust Data Analysis," Data Science Journal, vol. 10, no. 2, pp. 101-115, 2020.
- [11] E. G. B. Martins and J. D. Almeida, "Building a Data Science Web Application with Streamlit," Journal of Software Engineering and Applications, vol. 13, pp. 235-245, 2020.
- [12] F. A. Reyes and M. C. Cruz, "Machine Learning Model Evaluation Metrics: A Survey," Journal of Data Science, vol. 18, no. 2, pp. 115-130, 2021.
- [13] J. M. D. Brown and P. A. Lewis, "Ethics in Data Science: Best Practices for Data Analysts," Journal of Data Ethics, vol. 6, no. 1, pp. 15-25, 2021.
- [14] V. K. N. Gupta and R. C. Choudhary, "Exploring Data with Pandas: Techniques for Data Wrangling," Data Science Review, vol. 9, no. 4, pp. 30-45, 2020.
- [15] L. F. Wilson and D. A. Stewart, "The Role of Data Visualization in Exploratory Data Analysis," Journal of Visual Communication, vol. 11, no. 3, pp. 203-217, 2021.
- [16] C. J. Lin and S. Y. Chen, "Deploying Machine Learning Models with Streamlit," International Journal of Software Engineering, vol. 14, pp. 150-160, 2021.
- [17] [A. R. Patel and N. S. Reddy, "Best Practices for Interactive Data Analysis," Journal of Data Science and Technology, vol. 5, no. 2, pp. 85-99, 2020.
- [18] M. T. R. Ahmed et al., "Feature Selection Techniques in Machine Learning: A Survey," Journal of Computational Intelligence, vol. 25, pp. 99-120, 2019.
- [19] R. H. B. Silva and L. G. M. Pinto, "Creating Interactive Dashboards with Streamlit," Proceedings of the 2021 International Conference on Data Science, pp. 100-108.
- [20] A. C. Desai and V. N. Roy, "Understanding Exploratory Data Analysis: Tools and Techniques," International Journal of Research in Computer Science, vol. 8, no. 1, pp. 67-78, 2021.



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)