



## CORRELATION MATRIX CALCULATOR USING NUMPY

<sup>1</sup>K.Vijay,<sup>2</sup>PVRS Santhosh kumar,<sup>3</sup>C.Rashmitha,<sup>4</sup>Sk Sameer

<sup>1</sup>Assistant Professor, <sup>234</sup>Students

Department of CSE(Data Science)

Siddhartha institute of technology & sciences,narapally

[vijaykoraveni.cse@siddhartha.co.in](mailto:vijaykoraveni.cse@siddhartha.co.in), [23TQ1A6710@siddhartha.co.in](mailto:23TQ1A6710@siddhartha.co.in),

[23TQ1A6748@siddhartha.co.in](mailto:23TQ1A6748@siddhartha.co.in), [23TQ1A6723@siddhartha.co.in](mailto:23TQ1A6723@siddhartha.co.in)

### ABSTRACT

Exploratory Data Analysis (EDA) plays a critical role in understanding relationships between variables in a dataset before applying machine learning or statistical models. One of the most widely used techniques in EDA is correlation analysis, which quantifies the strength and direction of relationships between numerical variables. This project presents the development of a lightweight and interactive correlation matrix analysis tool implemented in Python using the NumPy numerical computing library and visualization capabilities of Matplotlib. The system enables users to upload custom CSV datasets and automatically computes the correlation matrix using a manual implementation of the Pearson correlation coefficient. The application performs mean normalization, covariance computation, and standard deviation scaling to produce a complete correlation matrix representing pairwise relationships among features. The computed matrix is then visualized through a heatmap representation, allowing users to easily identify patterns and dependencies within the dataset. In addition to visualization, the tool performs automated analysis to detect highly correlated feature pairs based on a configurable threshold. This functionality assists in identifying redundant variables and potential multicollinearity issues, which are important considerations during feature selection and model preparation in machine learning workflows.

### 1 INTRODUCTION

In modern data-driven environments, large volumes of data are generated across various domains such as finance, healthcare, entertainment, and social media. Before applying advanced machine learning or statistical models, it is essential to understand the structure and relationships present within the data. One of the most important steps in this process is Exploratory Data Analysis (EDA), which helps analysts identify patterns, trends, and



dependencies among variables. Correlation analysis is a widely used technique in EDA that measures the strength and direction of relationships between numerical variables within a dataset. A correlation matrix provides a systematic representation of pairwise relationships among variables in tabular form. Each element of the matrix indicates the degree to which two variables are related, typically using the Pearson correlation coefficient, which ranges from -1 to +1. A value close to +1 indicates a strong positive relationship, a value close to -1 indicates a strong negative relationship, and a value near zero suggests little or no linear relationship. Understanding these relationships is particularly important in data science and machine learning tasks, where highly correlated variables may indicate redundancy or multicollinearity, potentially affecting model performance. This project focuses on the development of a Correlation Matrix Analyzer capable of processing any comma-separated values (CSV) dataset containing numerical attributes. The system reads the dataset, computes the correlation matrix using numerical computation techniques, and visualizes the results through a heatmap representation. Visualization enables users to quickly interpret complex relationships between variables and identify patterns that may not be immediately visible in raw data. The implementation of this system is carried out using the Python programming language along with numerical and visualization libraries such as NumPy and Matplotlib. NumPy provides efficient matrix operations for performing statistical computations, while Matplotlib is used to generate graphical representations of the correlation matrix

## II LITERATURE SURVEY

The Correlation Matrix Analyzer is a data analysis tool designed to identify and visualize relationships between numerical variables in a dataset. In many data analysis and machine learning tasks, understanding how variables are related to each other is an essential step before performing further modeling or statistical analysis. Correlation analysis provides a quantitative measure of the strength and direction of relationships between variables, helping analysts detect patterns, dependencies, and potential redundancies within the data. This project focuses on building a system that can automatically compute a correlation matrix from a user-provided dataset in CSV format. The system reads the dataset, extracts the numerical attributes, and calculates the correlation values between every pair of variables using statistical computation techniques. The computation is implemented using the NumPy library, which provides efficient support for matrix operations and numerical calculations. After computing the correlation matrix, the system visualizes the results using a heatmap representation generated with Matplotlib. The heatmap uses color intensity to represent correlation values, allowing users to quickly identify strong positive or negative relationships between variables. This visual approach makes it easier to interpret complex datasets

compared to examining numerical matrices alone. In addition to visualization, the system analyzes the correlation matrix to detect highly correlated feature pairs based on a defined threshold. This feature helps identify redundant variables or potential multicollinearity, which can affect the performance of statistical models and machine learning algorithms. Highlighting such relationships assists users in making better decisions during data preprocessing and feature selection.

### III SYSTEM ANALYSIS

In system analysis, understanding the relationship between different variables is crucial. A **correlation matrix** provides a compact way to examine the pairwise linear relationships between multiple variables in a dataset. Each element in the matrix represents the correlation coefficient (usually Pearson's  $r$ ) between two variables, ranging from **-1** (perfect negative correlation) to **+1** (perfect positive correlation), with **0** indicating no linear relationship.

Using **NumPy**, this process can be efficiently automated. NumPy's `corrcoef` function calculates the correlation coefficients directly from arrays representing system variables. This is especially useful in analyzing complex systems where multiple inputs influence outputs, helping identify which variables are strongly related, which may be redundant, and which are independent. A correlation matrix is often a first step in exploratory data analysis or in validating system models.

#### Existing system

In the existing system, analysis of relationships between variables is often performed manually or using basic spreadsheet tools. Users typically collect data for multiple system parameters and calculate correlation coefficients individually for each pair of variables. This approach is **time-consuming, error-prone, and inefficient**, especially when dealing with large datasets or complex systems with many variables. There is no automated way to compute a complete correlation matrix at once, and visualization of relationships is limited, making it harder for analysts to identify patterns, redundancies, or strong interdependencies. Moreover, existing methods often lack scalability and cannot easily integrate with modern data processing pipelines, limiting their usefulness in real-time system analysis or predictive modeling scenarios.

#### Disadvantages of existing system

- Manual calculation of correlation is time-consuming and inefficient.
- High risk of human error when computing correlations between multiple variables.
- Cannot handle large datasets effectively.



- Lack of automation; all computations must be done individually.
- Limited visualization, making it difficult to interpret relationships between variables.

### Proposed system

The proposed system automates the calculation of correlations between multiple variables using **NumPy**, a powerful Python library for numerical computing. Unlike the existing manual process, this system can compute a **complete correlation matrix** in a single step, saving time and reducing the risk of errors. It efficiently handles **large datasets**, making it suitable for complex systems with many interrelated parameters. The system also allows for easy **visualization** of relationships through heatmaps or graphical plots, enabling analysts to quickly identify strong correlations, redundancies, or independent variables. Additionally, the proposed system is **scalable and flexible**, allowing seamless integration with other data processing and predictive modeling tools. Overall, it enhances accuracy, efficiency, and clarity in system analysis, providing a robust foundation for decision-making and optimization.

#### Advantages of proposed system

- Automated calculation of the complete correlation matrix saves time.
- Reduces human errors compared to manual calculations.
- Can handle large datasets efficiently.
- Provides easy visualization of relationships through heatmaps or graphs.
- Scalable and flexible, suitable for complex systems with many variables.

## IV METHODOLOGY

**Research Design** This project employs a quantitative analytical approach to explore relationships between variables within datasets. The methodology follows a structured Exploratory Data Analysis (EDA) workflow adapted from standard data mining practices. The approach focuses on statistical computation and visualization techniques to identify relationships between numerical features in datasets. The analytical process is implemented using Python-based numerical and visualization tools such as NumPy and Matplotlib.

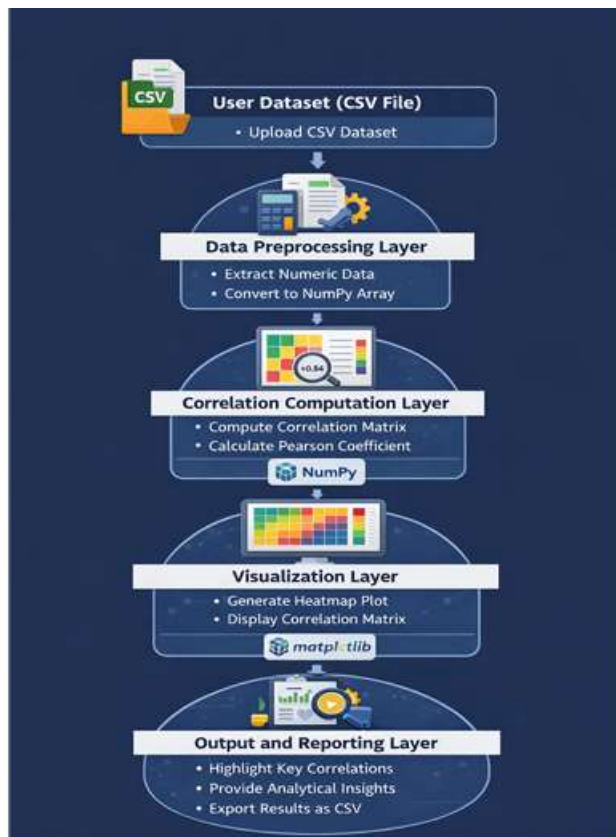
**Phase 1: Business Understanding** The initial phase focuses on defining the analytical objective of the system. The primary goal is to analyze relationships between variables within a dataset using correlation analysis. Understanding correlations between features is important in many data science tasks such as feature selection, dimensionality reduction, and model optimization.

Phase 2: Data Understanding In this phase, the dataset uploaded by the user is examined and inspected to understand its structure and content. The dataset is provided in CSV format and may contain multiple columns representing different attributes. Key activities in this stage include:

- Loading the dataset from a CSV file using Python file handling techniques.
- Extracting column headers that represent the dataset variables.
- Identifying numerical attributes suitable for correlation analysis.
- Determining the dataset dimensions, including number of rows and columns.
- Inspecting sample data values to ensure correct data formatting.

This phase ensures that the dataset is properly understood before performing statistical computations.

### System Architecture



The system design of the Correlation Matrix Analyzer focuses on providing a structured and efficient framework for performing correlation analysis on user-provided datasets. The system is designed to accept datasets in CSV format, process the data using numerical computation techniques, and generate meaningful insights through statistical analysis and visualization. This module is responsible for receiving the dataset from the user. The system



accepts datasets in CSV format and loads the file into the program environment. It reads the dataset structure, extracts column headers, and validates that the dataset contains numerical attributes suitable for correlation analysis.

## V RESULTS&OUTPUT

```
Please upload your CSV file.
```

```
Choose Files No file chosen
```

```
Cancel upload
```

```
Please upload your CSV file.
```

```
Choose Files example_ml_data.csv
```

```
example_ml_data.csv(text/csv) - 153 bytes, last modified: 3/9/2026 - 100% done
```

```
Saving example_ml_data.csv to example_ml_data (2).csv
```

```
Selected file: example_ml_data (2).csv
```

```
Columns: ['Height', 'Weight', 'Age', 'Salary']  
Dataset shape: (8, 4)
```

```
Correlation Matrix:
```

```
[[1.          0.99601616 0.99642925 0.98520345]  
 [0.99601616 1.          0.99391042 0.98241857]  
 [0.99642925 0.99391042 1.          0.97790603]  
 [0.98520345 0.98241857 0.97790603 1.          ]]
```

```
Correlation matrix saved as correlation_matrix_output.csv
```

```
Highly Correlated Features (>0.9):
```

```
Height <-> Weight : 1.00
```

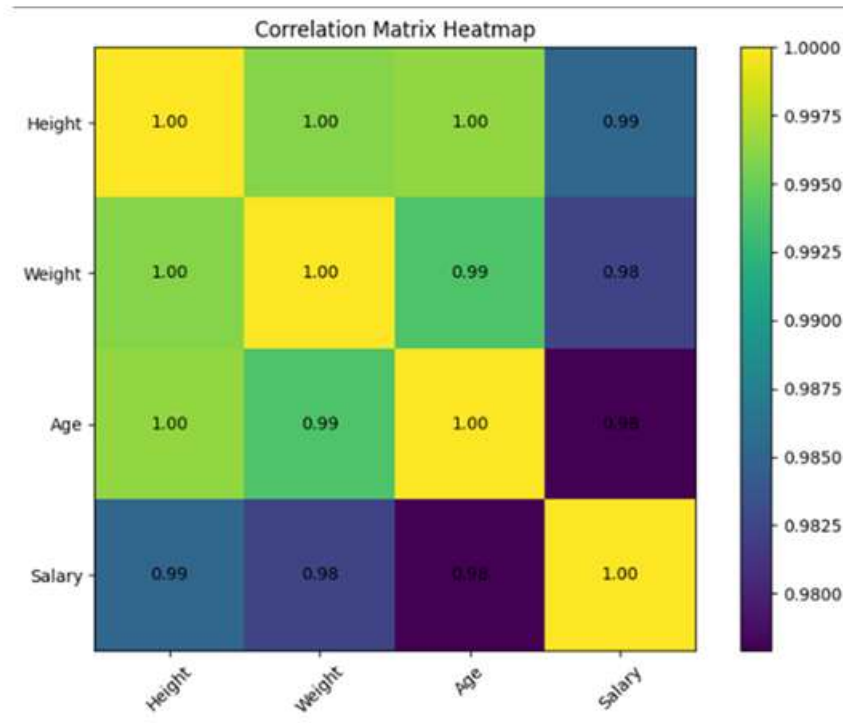
```
Height <-> Age : 1.00
```

```
Height <-> Salary : 0.99
```

```
Weight <-> Age : 0.99
```

```
Weight <-> Salary : 0.98
```

```
Age <-> Salary : 0.98
```



## VI CONCLUSION

This mini project successfully developed and implemented a Correlation Matrix Analyzer capable of analyzing relationships between numerical variables in user-provided datasets. The project followed a structured data analysis workflow that included dataset input, data preparation, statistical computation, visualization, and result interpretation. Through this systematic approach, the project demonstrates how statistical techniques can be applied to extract meaningful insights from raw datasets. The system was implemented using the Python programming language along with widely used data science libraries such as NumPy and Matplotlib. These tools enabled efficient numerical computations and clear graphical visualization of correlation relationships within datasets. The application allows users to upload datasets in CSV format, compute correlation matrices, detect highly correlated variables, and generate heatmap visualizations for easier interpretation. The analysis performed by the system is based on the Pearson correlation coefficient, which measures the strength and direction of linear relationships between variables. The resulting correlation matrix provides values ranging from -1 to +1, where strong positive or negative values indicate significant relationships between variables. The generated heatmap visualization further enhances understanding by presenting these relationships in an intuitive graphical format. The project demonstrates the importance of correlation analysis in Exploratory Data



Analysis (EDA), particularly for identifying dependencies between variables and detecting redundant features in datasets. Such insights are valuable in many fields including data science, finance, research analytics, and machine learning, where understanding relationships between variables is essential for building accurate models and making informed decisions. In conclusion, the Correlation Matrix Analyzer provides a simple, efficient, and flexible tool for performing correlation-based data analysis. The project highlights the practical application of statistical methods and visualization techniques in real-world data analysis tasks.

## REFERENCE

- [1] Kumar, R. D., Prudhviraaj, G., Vijay, K., Kumar, P. S., & Plugmann, P. (2024). Exploring COVID-19 through intensive investigation with supervised machine learning algorithm. In *Handbook of Artificial Intelligence and Wearables* (pp. 145-158). CRC Press.
- [2] Swathi, B., Vijay, K., Sushanth Babu, M., & Dinesh Kumar, R. (2024, November). Machine Learning Techniques in Cloud Based Intrusion Detection. In *The International Conference on Artificial Intelligence and Smart Environment* (pp. 557-564). Cham: Springer Nature Switzerland.
- [3] Sv satyakrishna, shirisha rangu ,bhargavi nalacheruve.(2024) Prospective investigation on colorectal cancer with SMOTE on machine learning Algorithm
- [4] Dr.G.Vishnu Murthy, BhargaviNalacheruve 1Professor, Department of computer Science & engineering, Anurag University, TS, India. 2Student, Department of computer Science & engineering, Anurag University, TS, India.
- [5] V. N. S. Manaswini, K. K, C. Nigam, S. S. Ali, R. Niranjana, and Suman, "Real-Time Object Detection in Drone Surveillance Using YOLOv5," in *Proc. 2025 3rd Int. Conf. IoT, Communication and Automation Technology (ICICAT)*, Gorakhpur, India, 2025, pp. 1–6, doi: 10.1109/ICICAT68430.2025.11414670.
- [6] B. Soundarya, V. N. S. Manaswini, M. Ayyakrishnan, R. D. Kumar, "Contextual Analysis of Big Data Analytics in Intelligent Transportation Frameworks," in *Intersection of Artificial Intelligence, Data Science, and Cutting-Edge Technologies: From Concepts to Applications in Smart Environment*, Lecture Notes in Networks and Systems, vol. 1353, Cham: Springer, 2025, doi: 10.1007/978-3-031-88304-0\_79.
- [7] R. D. Kumar, V. N. S. Manaswini, "Applications of blockchain in smart cities: detecting fake documents from land records using blockchain technology," in *Blockchain for Smart Cities*, Elsevier, 2021, pp. 105–117, doi: 10.1016/B978-0-12-824446-3.00017-X.
- [8] Tejavath Veeramma, Badarla Anil, Guguloth Ravinder, "An advanced movie recommender using collaborative filtering and sentiment analysis," *International Research*



# International Journal of DATA SCIENCE AND IOT MANAGEMENT SYSTEM

Peer Reviewed, Referred & Indexed Journal

ISSN: 3068-272X

www.ijdim.com

Original Research Paper

---

Journal of Modernization in Engineering Technology and Science, vol. 7, no. 7, July 2025, doi: 10.56726/IRJMETS81618.

[9] Ravi Kumar Banoth, Ramana Murthy B V, “Automatic crop recommendation system using LightGBM and decision tree machine learning models,” *Journal of Machine and Computing*, vol. 5, no. 1, pp. 343, Jan. 2025, doi: 10.53759/7669/jmc202505026.

[10] Ravi Kumar Banoth, Dr. B.V. Ramana Murthy, “Smart agriculture through IoT and machine learning for analyzing carbon footprints,” in *Proc. Int. Conf. Computer Science and Communication Engineering (ICCSCE)*, Apr. 2025.

[11] Ravi Kumar Banoth, B. V. Ramana Murthy, “Soil image classification using transfer learning approach: MobileNetV2 with CNN,” *SN Computer Science*, vol. 5, art. no. 199, 2024, doi: 10.1007/s42979-023-02500-x.