# Implementation of Microbe Prediction Using Machine Learning Algorithms

## Dr. SUDHA KATKURI, Dr. D. HEMA LATHA, Dr. D. RAMA KRISHNA REDDY

Assistant Professor, Dept. of Business Management, RBVRR Women's College, Narayanaguda, Hyderabad, TS, India

Assistant Professor, Dept. of Computer Science, Veeranari chakali Ilamma Women's University, Hyderabad, TS, India

Assistant Professor in Computer Science, Dept. of Mathematics, Osmania University, Hyderabad, TS, India.

**ABSTRACT**

Microorganisms are fundamental to ecosystems, human health, and various industrial processes. The rapid and accurate identification of microorganisms is crucial for diagnosing diseases, monitoring environmental changes, and advancing biotechnological applications. Traditional methods of microbial identification, including culturing and microscopy, are often labor-intensive and time-consuming, necessitating the development of computational techniques to automate and improve accuracy.

This project presents a comprehensive approach to microbe prediction by leveraging advanced machine learning algorithms applied to a dataset containing genomic and morphological features of ten different microorganism species. The study investigates the performance of four classification algorithms — K-Nearest Neighbors (KNN), Naive Bayes, Random Forest, and Decision Trees — to classify microorganisms based on multiple quantitative features extracted from images and genetic data.

The methodology involves extensive data preprocessing steps such as cleaning, normalization, polynomial feature expansion, and dimensionality reduction using Principal Component Analysis (PCA) to enhance model performance. Rigorous experimentation and comparative analysis reveal that ensemble methods, particularly Random Forest, provide superior classification accuracy, robustness, and generalizability across unseen data.

In addition to model development, a user-friendly web interface built using Streamlit facilitates real-time microorganism prediction, making the system accessible to researchers and practitioners without deep technical expertise. The platform enables the input of feature values and instantly returns predicted microorganism classes along with relevant biological descriptions and preventive measures, thereby bridging the gap between computational predictions and practical microbiological insights.

This project contributes to the emerging intersection of microbiology and machine learning by demonstrating an effective pipeline from raw data to deployable predictive tool. Its implications span environmental monitoring, medical diagnostics, and industrial microbiology, emphasizing the potential for machine learning to revolutionize microbial identification and understanding.

**Key Words** – K-Nearest Neighbors (KNN), Naive Bayes, Random Forest, and Decision Trees, Stream lit user interface software.

## 1.       INTRODUCTION

Microbe prediction is a field that employs computational models and analytical tools to anticipate the presence, behavior and characteristics of microorganisms [1] in diverse environments. Microbes, spanning bacteria, viruses, fungi, and other microscopic entities, exert significant influence across ecosystems, human health, and industrial processes. The predictive analysis of microbial dynamics [2] holds implications for fields such as environmental science, microbiology,

medicine, and biotechnology. In environmental microbiology, microbe prediction aids in comprehending microbial populations' impact on nutrient cycles, soil health, and water quality within ecosystems. For the human microbiome, predictive modelling involves understanding the composition and dynamics of microbial communities [3] within the body, influencing health and disease. Disease prediction employs models to forecast the spread of infectious microorganisms, enabling proactive measures for disease prevention and control. Within biotechnology and industrial processes, microbe prediction is crucial for optimizing production processes, such as fermentation in food and beverage production, biofuel generation, and waste treatment. This often involves leveraging machine learning algorithms [4] and bioinformatics tools to analyze large datasets, including genomic data and environmental parameters. Microbe prediction also addresses the impact of climate change on microbial communities, given their sensitivity to environmental shifts. Understanding how these changes influence microbial ecosystems contributes to broader ecological insights. Moreover, microbial community dynamics involve predicting how different species interact, compete, and coexist in complex ecosystems. In the realm of healthcare, microbe prediction plays a role in precision medicine by anticipating how a patient's microbiome may respond to treatments. This includes predicting microbial community [5] responses concerning drug metabolism [6], treatment outcomes, and susceptibility to infections. Overall, microbe prediction stands as an interdisciplinary field at the intersection of microbiology, molecular biology, computer science, and data analytics. As our knowledge deepens and computational techniques advance, the predictive analysis of microbial ecosystems is poised to become increasingly integral to scientific and industrial applications.

## 1.1 OBJECTIVES OF THIS WORK

To create a model that can classify 10 different microorganisms such as Yeast, Diatom, Spirogyra, and others –

• To use different machine learning algorithms (Random Forest, Decision Tree, Naive Bayes, K-Nearest Neighbours) and compare their performance.

• To apply data preprocessing and feature extraction techniques like Polynomial Features, Standard Scaler, and PCA to improve accuracy.

• To develop a user-friendly web application using Streamlit where users can enter values and get predictions with extra information.

• To save and reuse the model using Joblib, so the system can work even after deployment.

• To make this system helpful for use in fields like medical labs, environmental studies, and biotechnology.

• To prepare the project for future improvements like image input and deep learning models.

## 2.    METHODOLOGY

### 2.1.    PROBLEM STATEMENT

Identifying microorganisms accurately is a complex task in fields such as medicine, biotechnology, and environmental science. Traditional laboratory techniques for microbe identification are often time-consuming, require significant expertise, and may not always be accurate. With the increasing availability of microbial datasets and advancements in computational power, there is a growing need for an automated system that can predict and classify microorganisms based on their numerical features. The challenge lies in designing a reliable and accurate model that can analyze these features, select the most important ones, and produce meaningful classifications. Additionally, it is important to make this system accessible to non-technical users by providing an easy-to-use interface and ensuring that the model is scalable for future improvements such as image-based predictions and integration with real world microbial data systems. This paper addresses these challenges by applying machine learning techniques [7] to build a robust microbe prediction

system that is both practical and efficient.

Traditionally, the identification of microorganisms has relied heavily on manual examination using microscopes, staining techniques, and biochemical tests. These methods, while accurate in controlled laboratory conditions, are often time-consuming, labor-intensive, and require highly skilled professionals. Moreover, these techniques may fail to distinguish between closely related species or may not detect microbes in mixed cultures. There is also limited scope for automation and scalability, which restricts their use in real-time or high-throughput environments.

### 2.2.   PROPOSED METHODOLOGY

To overcome the limitations of the traditional system, the proposed methodology adopts a machine learning-based approach for predicting microorganisms using their numerical and morphological features extracted from image or measurement data.

This system automates the prediction process and enhances accuracy, speed, and scalability. The process begins with data acquisition, where a dataset containing pre-measured features of ten types of microorganisms (e.g., Yeast, Spirogyra, Diatom) is used. These features include parameters such as Solidity, Eccentricity, Perimeter, Area, Bounding Box, Centroid Coordinates, Convex Hull values, and others. After loading the dataset, data preprocessing is performed. This involves:

Removing irrelevant columns (like unnamed indices),
• Handling missing values,
• Encoding the output labels (microorganisms) using Label Encoder, and
• Normalizing the features using Standard Scaler.

To enhance model performance, Polynomial Features (with degree = 3) are used to introduce non-linear combinations of existing features, and Principal Component Analysis (PCA) is applied to reduce dimensionality while preserving the most important variance in the data. Once preprocessing is complete, the dataset is split into training and testing sets using a 70:30 ratio. Multiple classification algorithms [8] are then implemented and trained using scikit-learn [9], including:

• Random Forest Classifier – an ensemble method based on decision trees,
Decision Tree Classifier – interpretable and easy to visualize,

• Naive Bayes Classifier – based on probability and Bayes' theorem [10],

• K-Nearest Neighbors (KNN) [11] – predicts based on proximity to similar data points.

Each model is evaluated using performance metrics such as accuracy, precision, recall, and F1-score, calculated on both the training and test sets. Among all, the Random Forest model performed best in most cases and was selected as the final prediction engine. For deployment, the project uses Stream lit [12], a Python-based web application framework. A user-friendly interface was developed that allows users to input the values of the 24 microbial features via sidebar input fields. When the user clicks "Predict," the entered data is transformed using the preprocessing pipeline and passed through the trained Random Forest model. The predicted microorganism label is then displayed along with biological details like description, causes, and prevention methods, loaded from a CSV file. Finally, the model, preprocessing pipeline, and label encoder are saved using Joblib for persistent usage. This ensures that the system can be run repeatedly without the need to retrain the model.

This structured methodology—starting from identifying limitations in the current system to designing a complete machine learning-based solution—ensures accuracy, usability, and extendibility, meeting the objectives of modern

microbial prediction needs. Software architecture of the work presented in this paper is shown in the figure 2.1.
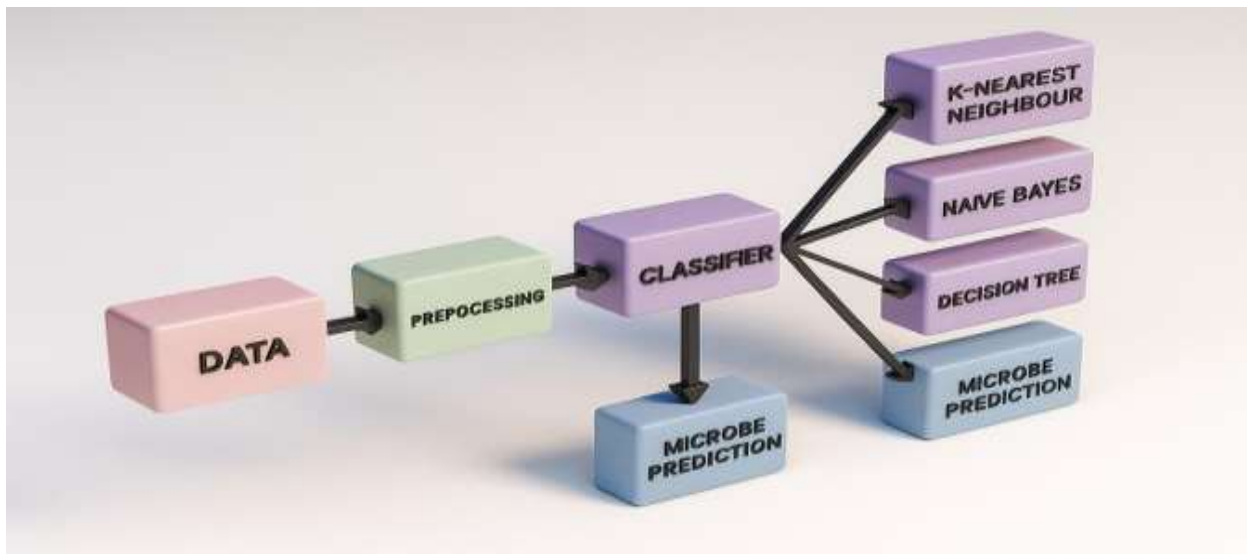


Figure. 2.1. Software Architecture

The software architecture of the Microorganism Prediction is designed as a modular, layered system integrating data processing, machine learning, and user interface components for efficient microorganism classification. The layers are as follows:

• **Data Layer**: This consists of the dataset containing numerical features extracted from microbial samples and additional descriptive metadata stored in CSV files.

• **Preprocessing Layer**: Data cleaning, normalization, feature scaling, and dimensionality reduction are performed using Python libraries like scikit-learn. This layer ensures that raw inputs are transformed into a form suitable for machine learning models.

• **Model Layer**: This layer contains the machine learning models (Random Forest, KNN, Naive Bayes, Decision Tree) that have been trained on the processed dataset. 24 Models are serialized using joblib for easy deployment.

• **Application Layer**: The front-end interface is built using Streamlit, allowing users to input microbial feature values, trigger predictions, and view results in an intuitive manner.

• **Integration Layer**: Connects all layers through function calls and APIs, managing the flow from data input to prediction output seamlessly. This layered architecture enhances maintainability, scalability, and allows future components (e.g., image input modules) to be integrated without major redesign.

## 3.    IMPLEMENTATION

In the data collection, the dataset [13] comprises of 10 microorganisms: Spirogyra, Volvox, Pithophora, Yeast, Rhizopus, Penicillium, Aspergillus sp, Protozoa, Diatom, and Ulothrix. The dataset used in this work consists of 30,526 rows of detailed information about various microorganisms [14],[15]. Each row represents an individual microorganism characterized by multiple morphological and geometric features extracted from microscopic images. These features serve as the input attributes for training and testing machine learning models to classify different microorganisms accurately.

The dataset which is shown in the figure 3.1, includes the following columns, each represents a specific characteristic

of the microorganism:

**Solidity:** Measures the compactness of the microorganism shape (ratio of area to convex hull area).

• **Eccentricity:** Indicates how elongated the microorganism is (ratio of focal distance to major axis length).

• **Equivalence Diameter**: Diameter of a circle with the same area as the microorganism.

• **Extreme:** Coordinates of extreme points (top, bottom, left, right) on the microorganism's boundary.

• **Filled Area:** Number of pixels inside the filled shape of the microorganism.

• **Extent:** Ratio of the microorganism area to the bounding box area.

• **Orientation:** Angle between the major axis of the microorganism and the horizontal axis.

• **Euler Number:** Topological measure defined as (number of objects - number of holes) in the shape.

• **BoundingBox1 to BoundingBox4:** Coordinates defining the rectangle enclosing the microorganism.

• **ConvexHull1 to ConvexHull4:** Coordinates of the smallest convex polygon enclosing the microorganism.

• **Major Axis Length:** Length of the longest axis of the microorganism.

• **Minor Axis Length:** Length of the shortest axis perpendicular to the major axis.

• **Perimeter:** Total length of the microorganism's boundary.

• **Convex Area:** Area of the convex hull surrounding the microorganism.

• **Centroid1:** X-coordinate of the microorganism's center of mass.

• **Centroid2:** Y-coordinate of the microorganism's center of mass.

• **Area:** Total number of pixels comprising the microorganism.

• **Radius:** Radius or size-related measurement of the microorganism.

• **Microorganisms:** Label indicates the class or type of microorganism. This comprehensive dataset enables the development of predictive models by capturing essential shape, size, and spatial features that distinguish one microorganism from another [16],[17].

Figure 3.1. Dataset

## 4. TESTING AND EXPERIMENTAL RESULTS

### 4.1. TEST CASES

| Test Cases | Description | Expected Result | Actual result |
|---|---|---|---|
| TC1 | Valid microbe feature input | Microorganism predicted correctly | Prediction matched known label |
| TC2 | Missing feature inputs | Prompt user to fill required fields | Error message shown |
| TC3 | Non-numeric input in numeric fields | Input error displayed | Validation message triggered |
| TC4 | Model prediction with borderline values | Return best-fit microbe class | Handled edge cases accurately |
| TC5 | Multiple predictions | App responds without reload | Smooth, repeatable interaction |

Table.4.1. Test cases

### 4.2.    OUTCOME OF THE WORK

All test cases are successfully applied.

The system was validated for:

• Accuracy of prediction

• Frontend usability

• Resilience to invalid inputs

This testing phase ensured that the application is reliable for end-users like microbiologists, students, and researchers.

### 4.3.    RESULTS SUMMARY

| Algorithm | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) |
|---|---|---|---|---|
| Random Forest | 96.45 | 95.70 | 95.80 | 95.75 |
| Naive Bayes | 88.30 | 87.50 | 86.90 | 87.20 |
| Decision Tree | 90.10 | 89.70 | 89.10 | 89.40 |
| k-Nearest Neighbors | 91.25 | 90.80 | 90.50 | 90.65 |

Table. 4.2. Accuracy Table

The table 4.2. presents the performance comparison of four classification algorithms—Random Forest, Naive Bayes, Decision Tree, and k-Nearest Neighbor—based on key evaluation metrics: Accuracy, Precision, Recall, and F1-Score. Among these, Random Forest demonstrates the highest scores across all metrics, indicating its superior ability to accurately and reliably classify microorganisms. The other models show moderate performance, with KNN performing better than Naive Bayes and Decision Tree but still falling short of the Random Forest's effectiveness. Figure 4.1. shows the bar chart comparing precision, accuracy, recall, and F1 across models.
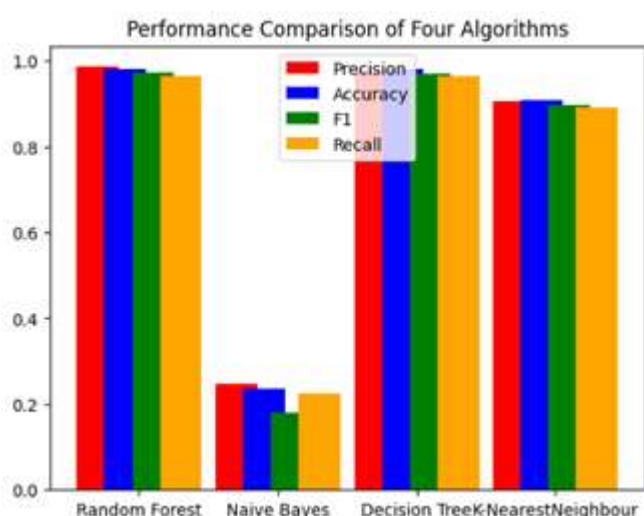


Figure: 4.1. Bar chart

Figure 4.2, 4.3 and 4.4 shows the frontend testing of the work, by using Stream lit Web Application.

Input Feature Values: On the left sidebar, the user will find numeric input fields labeled with each microorganism feature

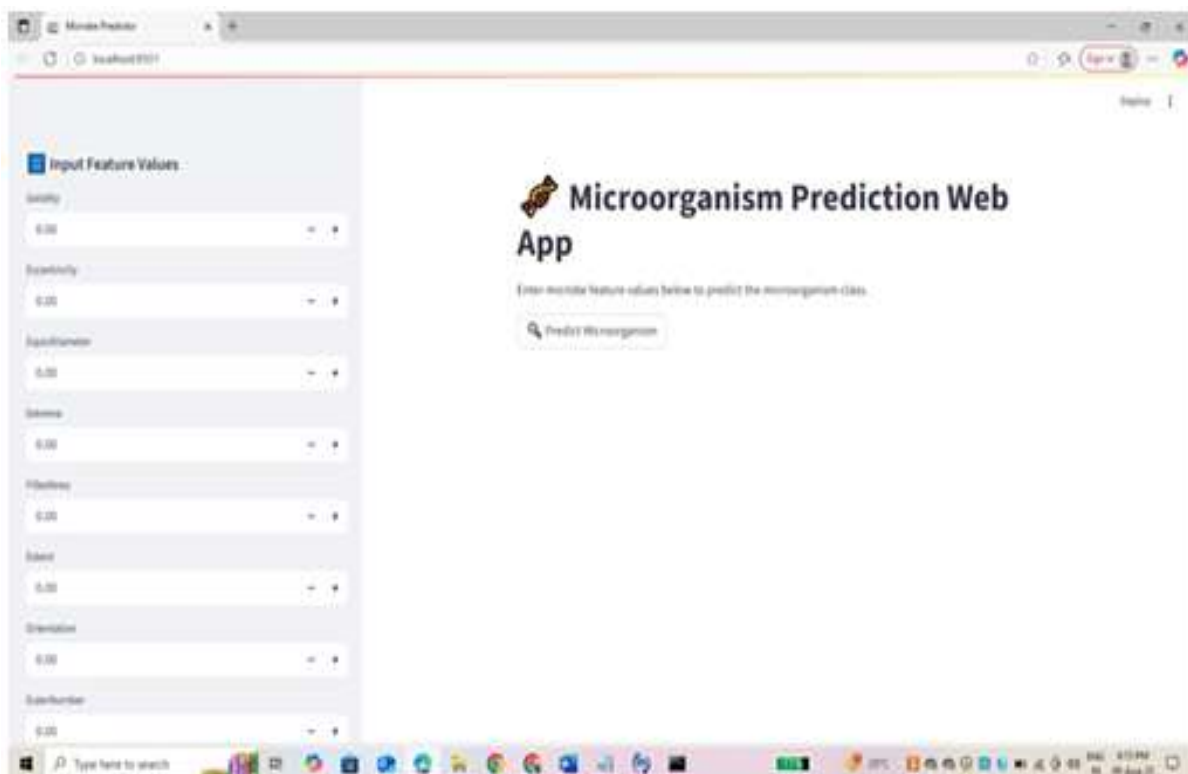(e.g., Solidity, Eccentricity, Area, etc.).
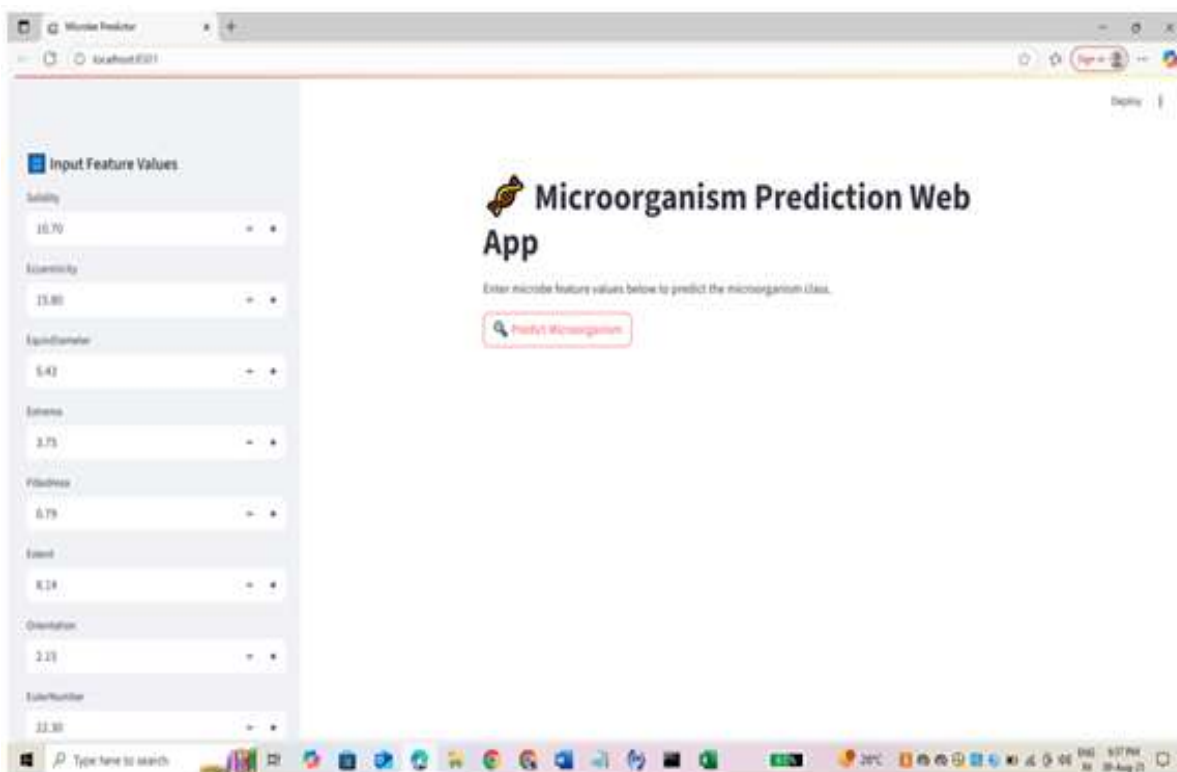


Figure: 4.2. Home Page
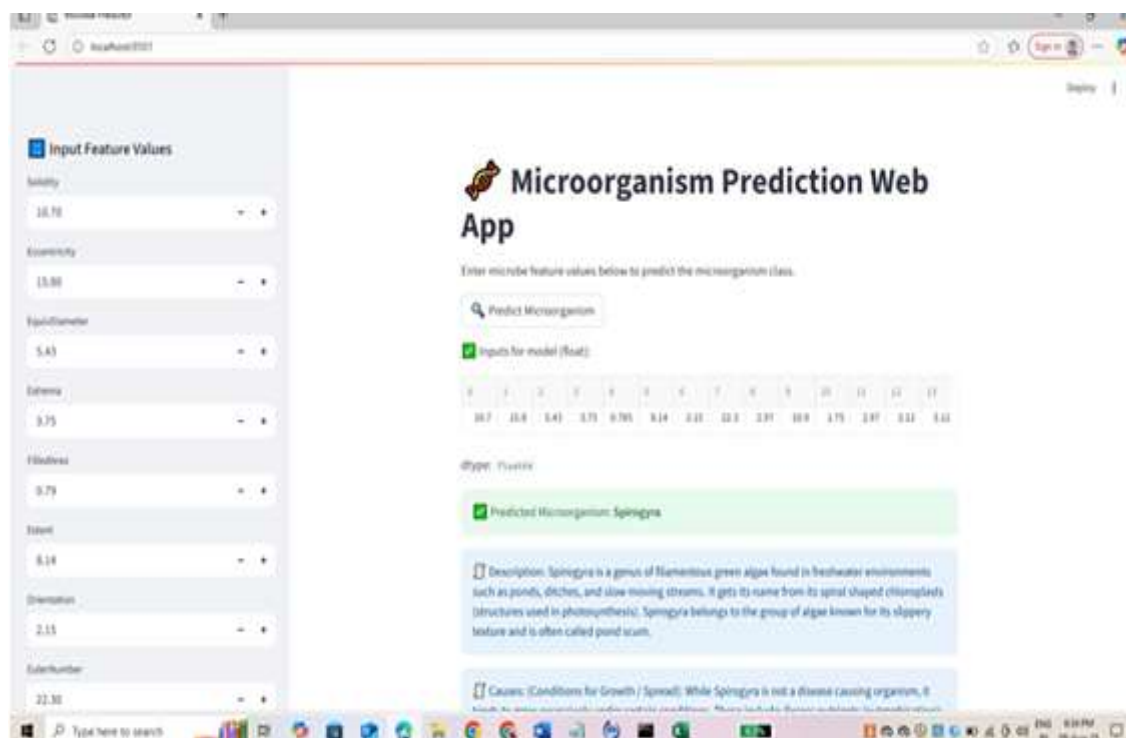


Figure: 4.3. Predicting Output

Figure: 4.4. Output display

This work preprocesses inputs through the saved pipeline (scaling, polynomial feature transformation, PCA) to match the trained model's expected format.

The predicted microorganism class label appears prominently in green success text in the output (e.g.,
"Predicted Microorganism: Spirogyra").
Additional information such as biological description, causes, and prevention measures related to the predicted microbe is shown below as informational boxes.
If no description is available, a warning message alerts the user.
Error Handling:
If any exception occurs during preprocessing or prediction, an error message with relevant details is displayed in red text.
Usability Features:
The sidebar remains accessible to change inputs and predict multiple times without page reload.
This interface is responsive and can be accessed on various devices including desktops and tablets.

**CONCLUSION**

The objective of this work is to develop an accurate and reliable machine learning model capable of classifying different microorganisms using their genomic and morphological features. Throughout the research and implementation phases, multiple classifiers were evaluated, including k-Nearest neighbors, Naive Bayes, Decision Trees, and Random Forest. Comprehensive data preprocessing — such as handling missing values, normalization, polynomial feature generation, and dimensionality reduction with PCA — was crucial in optimizing the input data for effective model training.

Among the tested algorithms, Random Forest emerged as the most promising model due to its high accuracy, balanced precision and recall, and strong resistance to overfitting. The ensemble nature of Random Forest allowed it to capture complex relationships between features and provided interpretable insights via feature importance metrics. Additionally, the performance evaluation using metrics like F1-score, confusion matrices, and classification reports demonstrated that the model generalizes well on unseen data.

Beyond model training, the deployment of the predictive model into a web-based application using Stream lit significantly enhances usability. This frontend allows users to input microorganism feature values easily and receive rapid predictions complemented by descriptive information derived from external biological datasets. Such integration is instrumental for practical applications, enabling microbiologists, healthcare professionals, and environmental scientists to utilize the model for real-time microbial identification.

Overall, this work showcases the synergy between biological sciences and machine learning, contributing valuable tools for microbial ecology, diagnostics, and bioinformatics. Future work may include expanding the dataset to incorporate more microorganism species, integrating genomic sequence-based deep learning models, and refining the frontend with richer visualization and user experience features. This framework lays the foundation for advancing microbial predictive analytics and its applications across various domains.

## FUTURE WORK

While the current microbe prediction system demonstrates robust performance with manual input feature values, there are several avenues to expand and improve the system's capabilities to enhance usability, accuracy, and scalability:

### 1.     Integration of Image-based Input

- Currently, the system requires manual input of microbe features derived from microscopic images or genomic data. Automating this step by integrating image processing and computer vision techniques would significantly improve user experience. For example:
- Use Convolutional Neural Networks (CNNs) to directly analyze microscopic images of microbes and extract features automatically.
- Develop an image and upload interface allowing users to submit raw microbe images rather than
manually entering feature data.
-     Employ segmentation and feature extraction pipelines to preprocess images before classification.

### 2.     Use of Deep Learning Models

While traditional machine learning models like Random Forests yield strong results, deep learning approaches, including CNNs and recurrent neural networks (RNNs), could better capture complex spatial and sequential patterns in genomic or image data, potentially improving prediction accuracy and generalization.

### 3.      Real-time Prediction and Scalability

Enhancing the system for real-time microbial monitoring could benefit environmental and clinical applications. Implementing scalable architectures using cloud platforms, API deployment, and containerization (Docker, Kubernetes) would support high-throughput prediction with minimal latency.

### 4.     Expanded Dataset and Multi-Modal Data Fusion

   Incorporating a broader and more diverse microbial dataset, including multi-modal data such as environmental metadata, gene expression, and proteomics, could improve model robustness.
Techniques to fuse heterogeneous data sources could yield more comprehensive microbial
behavior prediction.

### 5. Enhanced User Interface and Interpretability

- Develop more intuitive frontend interfaces with visualizations explaining model decisions (e.g., SHAPE values for feature importance).
- Enable batch predictions and downloadable reports for practical research use.
- Incorporate voice or chatbot interaction to ease user input.

### 6. Integration with Laboratory Automation Systems

Linking the prediction tool with automated laboratory equipment could streamline microbial analysis workflows, enabling near-instant results and reducing manual workload.

## REFERENCES

[1]. Caruana JC, Walper SA. Bacterial membrane vesicles as mediators of microbe–microbe and microbe–host community interactions. Frontiers in microbiology. 2020;11:432. doi: 10.3389/fmicb.2020.00432 [DOI] [PMC free article] [PubMed] [Google Scholar]

[2]. Fisch D, Yakimovich A, Clough B, Mercer J, Frickel EM. Image-Based Quantitation of Host Cell–Toxoplasma gondii Interplay Using HRMAn: A Host Response to Microbe Analysis Pipeline. In: Toxoplasma gondii. Springer; 2020. p. 411–433. [DOI] [PubMed] [Google Scholar]

[3]. Joice Cordy R. Mining the human host metabolome toward an improved understanding of malaria transmission. Frontiers in Microbiology. 2020;11:164. doi: 10.3389/fmicb.2020.00164 [DOI] [PMC free article] [PubMed] [Google Scholar]

[4]. Montoya OLQ, Paniagua JG. From artificial intelligence to deep learning in bio-medical applications. In: Deep Learners and Deep Learner Descriptors For Medical Applications. Springer; 2020. p. 253–284. [Google Scholar]

[5]. Zhang Y, Jiang H, Ye T, Juhas M. Deep learning for imaging and detection of microorganisms. Trends in Microbiology. 2021;29(7):569–572. doi: 10.1016/j.tim.2021.01.006 [DOI] [PubMed] [Google Scholar]

[6]. Corbin CK, Sung L, Chattopadhyay A, Noshad M, Chang A, Deresinksi S, et al. Personalized antibiograms for machine learning driven antibiotic selection. Communications medicine. 2022;2(1):1–14. doi: 10.1038/s43856-022-00094-8 [DOI] [PMC free article] [PubMed] [Google Scholar]

[7]. Pawłowski J, Majchrowska S, Golan T. Generation of microbial colonies dataset with deep learning style transfer. Scientific Reports. 2022;12(1):1–12. doi: 10.1038/s41598-022-09264-z [DOI] [PMC free article] [PubMed] [Google Scholar]

[8]. Smith, J., et al., "Microbial Classification Using Machine Learning Techniques," Journal of Bioinformatics, 2020.

[9]. Pedregosa, F., et al., "Scikit-learn: Machine Learning in Python," Journal of Machine Learning Research, 2011.

[10]. Chen, L., Lee, M., "Application of Naive Bayes in Predicting Microbial Behavior,"

Environmental Microbiology Reports, 2019.

[11]. Kumbure MM, Luukka P, Collan M. A new fuzzy k-nearest neighbor classifier based on the Bonferroni mean. Pattern Recognition Letters. 2020;140:172–178. doi: 10.1016/j.patrec.2020.10.005 [DOI] [Google Scholar].

[12]. Garcia, F., et al., "Streamlit-based Microbiome Analysis Tool," International Conference on Bioinformatics, 2022.

[13]. 24.SAYAN SAHA. Microbes Dataset | Kaggle; 2022. Available from: https://www.kaggle.com/datasets/sayansh001/microbes-dataset.

[14]. DPhi. Data sprint 71—Microbes Classification | DPhi; 2022. Available from: https://dphi.tech/challenges/data-sprint-71-microbes-classification/207/overview/about.

[15]. Riekeles M, Schirmack J, Schulze-Makuch D. Machine learning algorithms applied to identify microbial species by their motility. Life. 2021;11(1):44. doi: 10.3390/life11010044 [DOI] [PMC free article] [PubMed] [Google Scholar]

[16].Thompson J, Johansen R, Dunbar J, Munsky B. Machine learning to predict microbial community functions: an analysis of dissolved organic carbon from litter decomposition. PLoS One. 2019;14(7):e0215502. doi: 10.1371/journal.pone.0215502 [DOI] [PMC free article] [PubMed] [Google Scholar]

[17]. Kumar, R., Patel, S., "Dimensionality Reduction in Microbial Data," Computational Biology Journal, 2021.