# B.M.S. COLLEGE OF ENGINEERING

## (Autonomous Institute, Affiliated to VTU)

## Bull Temple Road, Basavanagudi, Bangalore – 560019

AN AAT REPORT ON

"MUSHROOM CLASSIFICATION - POISONOUS OR EDIBLE"

Submitted in partial fulfilment of the requirements for the award of degree

B.E

in

Information Science and Engineering

By

**KONAKALLA MEENAKSHI - 1BM22IS103**

**MEDHINI KATTI – 1BM22IS111**

**MADHURA M R - 1BM22IS110**

**PRATIKSHA RATHOD – 1BM23IS415**

Under the guidance of

**RASHMI R**

Assistant Professor

Department of Information Science

A.Y 2024-25

 **Abstract-**

In this study, we aim to address the critical challenge of distinguishing between edible and poisonous mushrooms by employing advanced machine learning techniques. Utilizing the comprehensive Secondary Mushroom Dataset from the UCI Machine Learning Repository, which includes data on 173 species from 23 families, we analyze the performance of various classifiers, namely K-Nearest Neighbors (KNN), Logistic Regression, and Random Forests. The dataset encompasses a diverse set of features, such as cap shape, cap color, gill attachment, gill spacing, stem color, veil type, ring number, and more, providing a rich source for model training and evaluation.

Our research builds on previous studies that have underscored the necessity of accurate mushroom classification for public health safety. By incorporating state-of-the-art machine learning algorithms, we strive to enhance classification accuracy and reliability. The classifiers are evaluated using performance metrics such as accuracy, precision, recall, and F1-score, ensuring a robust analysis of their effectiveness.

The results demonstrate that ensemble methods, particularly Random Forests, significantly improve classification performance, achieving near-perfect accuracy in distinguishing between edible and poisonous mushrooms. This enhanced accuracy is crucial for reducing the risk of mushroom poisoning, thereby safeguarding public health. Additionally, the study highlights the non-linearly separable nature of the dataset, which necessitates the use of classifiers capable of handling complex, non-linear boundaries.

The findings of this research are highly relevant for projects focused on food safety and mycology. The comprehensive dataset and reproducible workflow presented in this study can serve as a valuable resource for researchers and practitioners. By providing a more extensive and representative dataset, the study facilitates better training and evaluation of machine learning models, ultimately contributing to improved public safety and education regarding mushroom foraging and consumption. The techniques and methodologies applied in this research can be utilized to develop applications and systems that aid in the accurate identification of edible mushrooms, promoting safer foraging practices and reducing the incidence of mushroom-related poisonings.

# Table of Contents

# Chapter 1: Introduction

Mushroom foraging has been a popular activity for centuries, but the risk of consuming poisonous species poses a serious threat to public health. Traditional identification methods rely heavily on expert knowledge and manual inspection, which are both time-consuming and prone to error. This research addresses the problem of accurate mushroom classification by leveraging advanced machine learning algorithms to distinguish between edible and poisonous varieties.

The primary challenge lies in developing a robust and reliable method for identifying the edibility of mushrooms, thereby reducing the risk of poisoning and ensuring public health safety. Current methods lack the precision and efficiency required to handle large, diverse datasets and complex features inherent in mushroom species.

Our study builds on the foundation laid by previous research, including works by Dennis Wagner et al. (2021), Nusrat Jahan Pinky et al. (2019), and Roni Hamonangan et al. (2021). These studies emphasize the need for comprehensive datasets and advanced classifiers to improve the accuracy of mushroom identification.

CSV file consists of a header followed by the 61,069 mushroom entries. The data is comprised one binary class, 17 nominal variables and three quantitative variables. The class distribution is as follows-poisonous-55% and edible-45%.

The dataset used in this study includes essential features such as:

- Cap shape, cap surface, and cap color
- Gill attachment, gill spacing, and gill size
- Stalk shape, stalk root, and stalk surface
- Veil type, veil color, and ring number
- Spore print color, population, and habitat

By applying classifiers like K-Nearest Neighbors (KNN), Logistic Regression, and Random Forests, we aim to enhance the classification accuracy and reliability of mushroom identification. The techniques and findings from this study will contribute to the development of tools and applications that can aid in the safe identification of mushrooms, thus promoting better food safety practices and reducing the risk of mushroom poisoning.

# Chapter 2: Problem Statement

The primary challenge addressed by this project is the development of a reliable, accessible, and user-friendly system for mushroom classification. The specific problems addressed include:

1. Classification Accuracy Challenge:

   - Need for high accuracy due to life-threatening consequences of misclassification

   - Handling of multiple, interdependent features

   - Dealing with missing or uncertain data

2. Accessibility Challenges:

   - Creating an interface that non-experts can use effectively

   - Ensuring rapid response times for real-world usage

   - Making the system available through web browsers

3. Technical Challenges:

   - Feature selection and engineering from categorical and numerical data

   - Dimensionality reduction while maintaining prediction accuracy

   - Model selection and optimization

   - Integration of multiple technologies (Flask, ML models, web interface)

4. Data Quality Challenges:

   - Handling imbalanced datasets

   - Managing missing values

   - Dealing with noisy or inconsistent data

# Chapter 3: Literature Survey

1. **"Mushroom data creation, curation, and simulation to support classification tasks" (Dennis Wagner et al., 2021)-**

   o **Objective of the Research:**
   The research paper aims to develop and simulate a comprehensive data set for distinguishing between edible and poisonous mushrooms using machine learning.

   o **Techniques Used-**
   1. Data Extraction and Formatting: The study involves extracting data from a textbook on mushroom identification and formatting it into a primary data set that includes 173 species from 23 families. This data is then used to generate a secondary data set through random sampling.
   2. Machine Learning Algorithms: The paper evaluates multiple classifiers, including:
      - Naive Bayes
      - Logistic Regression
      - Linear Discriminant Analysis (LDA)
      - Random Forests (RF).
   3. Cross-Validation: The performance of the classifiers is assessed using five-fold cross-validation to ensure the reliability of the results.

   o **Key Findings and Results**:
   1. The Random Forest classifier achieved perfect accuracy and an F2-score of 1.0, indicating its effectiveness in classifying the mushroom data accurately.
   2. The newly created secondary data set demonstrated better data quality and integrity compared to the older 1987 data set, making it a more suitable alternative for classification tasks.
   3. The secondary data set is not linearly separable, which required the use of classifiers capable of handling non-linear boundaries, further emphasizing the complexity of mushroom identification.

   o **Relevance to Project:**
   The study provides a valuable and comprehensive data set for mushroom classification, enhancing machine learning model accuracy and contributing significantly to public safety and education on mushroom foraging.


2. **"Edibility Detection of Mushroom Using Ensemble Methods" (Nusrat Jahan Pinky et al., 2019)**

   o **Objective of the Research**
   The study aims to improve mushroom edibility identification methods by accurately classifying mushrooms as edible or poisonous, thus assessing the risk to human life.

   o **Techniques Used**
   1. Bagging: Specifically, Naive Bayes based Bagging and Dissimilarity based Bagging.
   2. Boosting: Utilizing the AdaBoost algorithm.
   3. Random Forest: This method is highlighted for its effectiveness in achieving high accuracy.

4. Additionally, the study utilizes fixed feature sets derived from significant attributes of mushrooms, which are compared against randomly selected feature sets.

o **Findings**
The findings indicate that the proposed methods, particularly those based on fixed features sets, demonstrate higher accuracy than those using randomly selected features. The Random Forest method yielded the highest accuracy across both test sets.

o **Key Results**
1. The experimental results show that:
2. The highest accuracy achieved by the proposed model based on Random Forest is superior to other methods tested.
3. The use of fixed feature sets significantly improves the performance of the classification models compared to random selections.
4. The paper also compares its results with existing works, highlighting that the proposed methods are more robust and effective in classifying mushroom edibility.

o **Relevance to the Project**
Using Ensemble Models for better predictions

3. **"Accuracy of classification poisonous or edible of mushroom using naïve bayes and K-nearest neighbours" (Roni Hamonangan et al., 2021)**

o **Overview of Research Paper**: Accuracy of Classification Poisonous or Edible of Mushroom Using Naïve Bayes and K-Nearest Neighbors

o **Objective:**
The primary objective of this research paper is to determine the best accuracy in classifying mushrooms as either poisonous or edible using two classification algorithms: Naïve Bayes and K-Nearest Neighbors (KNN).

o **Techniques Used:**
1. Naïve Bayes Algorithm: This algorithm utilizes Bayes' theorem to classify data based on the probability of class labels. It simplifies calculations by reducing computational complexity and is effective for large datasets.
2. K-Nearest Neighbors (KNN): KNN classifies new objects based on the closest distance to training data, making it a straightforward method for classification tasks.
3. Data Collection: The study uses a dataset from the UCI Machine Learning Repository, which includes 8124 mushroom samples, categorized into edible and poisonous mushrooms.

o **Findings:**
1. The Naïve Bayes algorithm initially achieved an accuracy of approximately 90.2%. This accuracy was later improved to 100% using the KNN algorithm.
2. The research emphasizes the importance of morphological features, such as color and habitat, in identifying mushroom types.

o **Key Results:**
1. The study successfully demonstrated that the KNN algorithm outperformed the Naïve Bayes algorithm in terms of accuracy for mushroom classification.

2. The use of confusion matrices helped validate the accuracy of the classification results, confirming the effectiveness of the algorithms used.

- o **Relevance to Project:**
  1. This research is highly relevant for projects focused on food safety and mycology, as it provides a robust methodology for accurately classifying mushrooms, which can help prevent the consumption of poisonous varieties.
  2. The findings can be applied in developing applications or systems that assist users in identifying edible mushrooms, thereby enhancing public safety and awareness regarding mushroom consumption.

# Chapter 4: System Requirement Specifications

**Functional Requirements:**

**1. Data Input and Processing:**

  - Accept user input for 20 different mushroom characteristics

  - Handle missing values with default parameters

  - Validate input data format and ranges

  - Process both categorical and numerical features

**2. Classification System:**

  - Perform real-time classification

  - Provide binary classification (edible/poisonous)

  - Include confidence scores with predictions

  - Handle multiple simultaneous requests

**3. User Interface:**

  - Web-based input form

  - Clear presentation of results

  - Error handling and user feedback

  - Mobile-responsive design

**Non-Functional Requirements:**

**1. Performance:**

  - Response time < 2 seconds

  - Support for concurrent users

  - 99.9% system availability

  - Efficient resource utilization

**2. Security:**

  - Input validation and sanitization

  - Protection against common web vulnerabilities

  - Secure data transmission
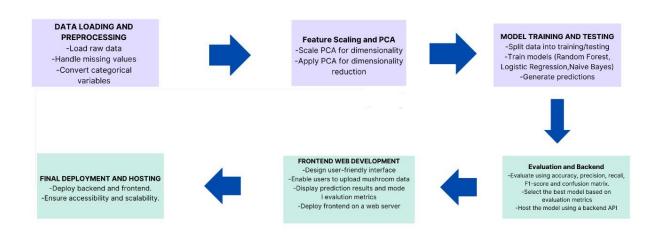
  - Error handling without information disclosure

**3. Scalability:**

  - Horizontal scaling capability

  - Load balancing support

  - Caching mechanisms

  - Database optimization

**4. Maintainability:**

  - Modular code structure

  - Comprehensive documentation

  - Version control

  - Testing frameworks

# Chapter 5: System Design



## System Architecture:

### 1. Frontend Layer:

- HTML/CSS interface

- Form validation

- Result display

- Responsive design

### 2. Application Layer:

- Flask web server

- Route handling

- Data preprocessing

- Model integration

### 3. Machine Learning Layer:

- Feature processing pipeline

- Label encoding

- PCA transformation

- Classification model

**4. Data Flow:**

  - User input collection

  - Data validation and preprocessing

  - Feature transformation

  - Model prediction

  - Result formatting and display

## Key Components:

### 1. Data Preprocessing Pipeline:

  - Label encoding for categorical features

  - Scaling for numerical features

  - PCA dimensionality reduction

  - Missing value handling

### 2. Model Architecture:

  - Feature selection

  - PCA transformation

  - Classification algorithm

  - Prediction pipeline

### 3. Web Application:

  - Flask routes and handlers

  - Template rendering

  - Error handling

  - Response formatting

# Chapter 6: Implementation

## 1. Data Preprocessing Implementation:

### a. Handling missing values:

```
Columns where NaN is the most frequent value:
['cap-surface', 'gill-spacing', 'stem-root', 'stem-surface', 'veil-type', 'veil-color', 'spore-print-color']

Columns where NaN is not the most frequent value and have been filled with the most frequent value:
['gill-attachment', 'ring-type']
```

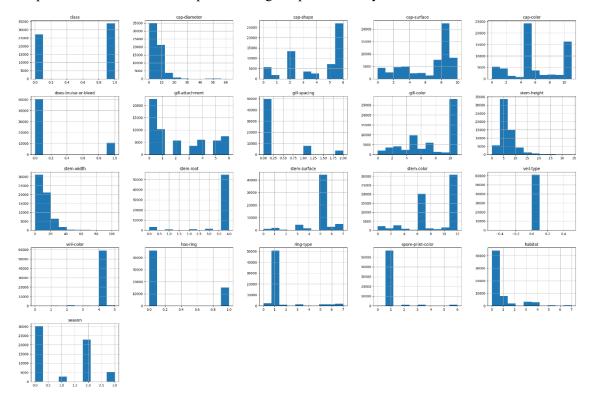The columns are split into those where NaN is most frequent and those where it is not-

### b. Encoding-Label encoding implementation-

categorical_columns = [

    'cap-shape', 'cap-surface', 'cap-color', 'does-bruise-or-bleed',

    'gill-attachment', 'gill-spacing', 'gill-color', 'stem-root',

    'stem-surface', 'stem-color', 'veil-type', 'veil-color',

    'has-ring', 'ring-type', 'spore-print-color', 'habitat', 'season'
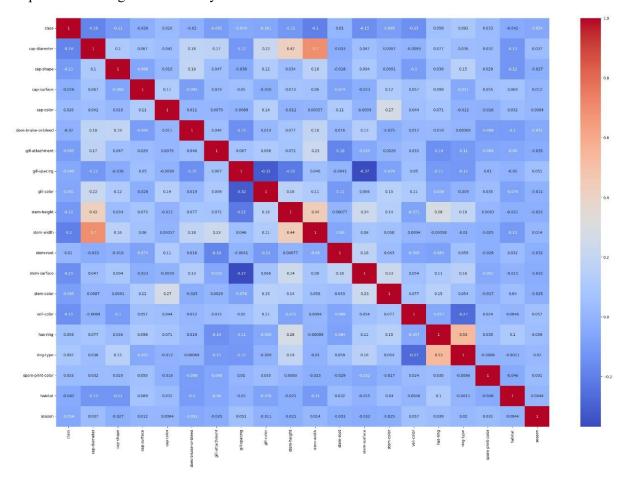
]

Columns like gill-attachment and habitat have high cardinality hence one-hot encoder is not preferred. Ordinal encoding is used in scenarios where there is a particular order for example-Low, medium and high. Since characteristics of a mushroom species like color of a cap or the shape do not involve any order label encoding is preferred.

### c. Data Visualization and Analysis-

To understand the nature of the data and to apply further preprocessing steps, histograms denoting the frequencies of the features are plotted using matplotlib library.

Heat map for understanding the correlation between the target column and the dependent features is implemented using seaborn library.



- **Diagonal Values (1.0):**
  All diagonal entries are **1.0** because each feature is **perfectly correlated** with itself.
- **Weak Correlations:**
  Most off-diagonal correlations are **close to 0** (blue shades).
  This indicates **weak linear relationships** between features, implying that features are **mostly independent**.
- **Notable Correlations:**
  Some features like:
  - **cap-diameter** and **stem-width** (~0.7).
  - **has-rings** and **ring-type** (~0.53).
  - **stem-height** and **stem-width** (~0.44).

  These show **moderate correlations**, which may contribute to **multicollinearity** in linear models.

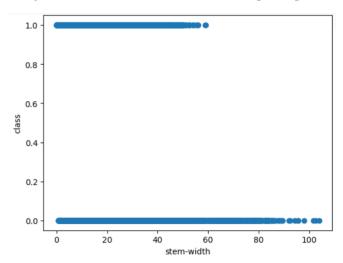- **Correlation with Target (class):**
  Features like:
  - **cap-diameter (-0.18)**
  - **gill-spacing (-0.49)**
  - **stem-height (-0.12)**
  - **stem-width (0.7)**

Show some degree of correlation with the **class** (target). These could be **important predictors** for classification.

**d. Polynomial Features-**

Polynomial features **transform existing features** into **higher-order terms** to **capture nonlinear relationships**. The features in the dataset do not have a linear correlation with the class feature. This means changes in the features do not predict changes in the class (edibility) in a straightforward linear manner.

The relationship between stem width and class is not linear because it does not show a continuous, straight-line trend. Instead, it exhibits a step-like pattern, characteristic of a non-linear relationship.
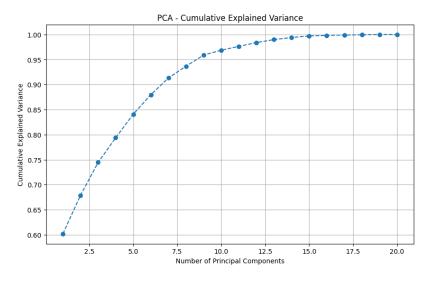


**e. Feature scaling and PCA**

scaler = StandardScaler()

Standardization is used for scaling the features.

pca = PCA(n_components=0.95)

PCA with variance 0.95 is used to reduce dimensions of the dataset.

The number of components that were selected are 41.

## 2. Model Implementation:

Three different models were implemented and compared:

### a) Logistic Regression:

Works well for large datasets. Since logistic regression works well with Linear relationships introducing polynomial features helps in extracting non-linear relationships and increases the correlation by introducing other relationships among columns. This improves the performance of logistic regression.

log_reg=LogisticRegression(random_state=42,penalty='l1',solver='liblinear',C=0.1)

L1 regularization-Lasso

Liblinear is suitable for large datasets and is suitably used along with l1 regularization.

The C parameter in logistic regression models acts as a control for the amount of regularization.

### b) K Neighbors Classifier:

Does not depend on the distribution of dataset(non-parametric).

knn= KNeighborsClassifier(n_neighbors=3)

n_neighbors is determined using GridSearchCV.

### c) Random Forest Classifier:

Does not depend on the distribution of dataset(non-parametric). An ensemble learning model that trains large number of decision trees.

 rf_model=RandomForestClassifier(random_state=42)

## 3. Web Application Implementation:

```
@app.route('/predict', methods=['POST', 'GET'])

def predict(form_data):

  input_data = pd.DataFrame([form_data])

  for column in input_data.columns:

    if column in label_encoders:

      le = label_encoders[column]

      input_data[column] = le.transform(input_data[column])

  feature_columns=input_data.columns.drop('class',errors=ignore)

  X_new = input_data[feature_columns]

  X_new_poly=poly.transform(X_new)

  X_new_scaled = scaler.transform(X_new_poly)

  X_new_pca = pca.transform(X_new_scaled)

  prediction = model.predict(X_new_pca)

  return prediction
```

# Chapter 7: Test Results

## 1. Model Performance Metrics:

| Metric | Logistic Regression | KNN | Random Forest |
|---|---|---|---|
| Accuracy | 0.7391 | 0.9998 | 0.9996 |
| Precision | 0.7536 | 0.9999 | 0.9997 |
| Recall | 0.7934 | 0.9999 | 0.9996 |
| F1 Score | 0.7730 | 0.9999 | 0.9996 |
| AUC | 0.84 | 0.9999 | 1.00 |

## 2. Confusion Matrix Analysis:

### a) Logistic Regression-

  - True Positives (Edible): 5427

  - True Negatives (Poisonous): 3600

  - False Positives: 1774

  - False Negatives: 1413

### b) K Neighbors Classifier-

  - True Positives (Edible): 6839

  - True Negatives (Poisonous): 5373

  - False Positives: 1

  - False Negatives: 1

### c) Random Forest Classifier-

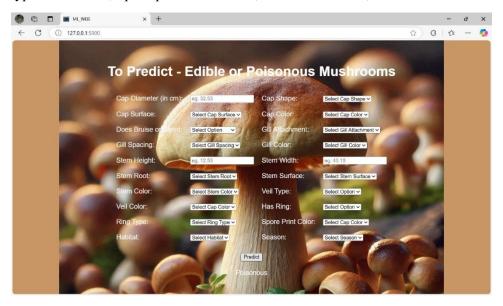  - True Positives (Edible): 6837

  - True Negatives (Poisonous): 5372

  - False Positives: 2

  - False Negatives: 3
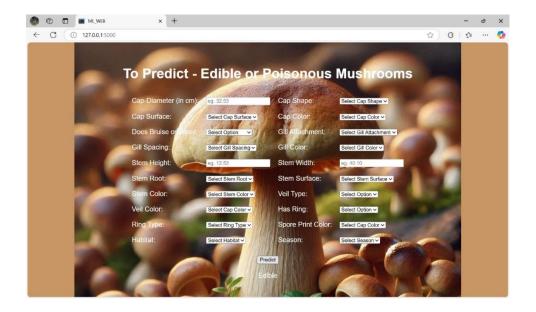
### 3. Prediction:

**a) Poisonous-**

Cap-diameter-3.95, Cap shape-bell, Cap surface-scaly, Cap color-pink, Does-bruise-or-bleed-no, Gill attachment-free, Gill spacing-close, Gill-color-white, Stem height-3.96, Stem width-3.37, Stem-root-blank , Stem surface-Silky, Stem-color-white, Veil type-blank, Veil color-blank, has ring-ring, Ring type- evanescent, Spore-print-color-blank, Habitat-meadows, Season-autumn



**b) Edible-**

Cap-diameter-6.21, Cap shape-conical, Cap surface-shiny, Cap color-brown, Does-bruise-or-bleed-no, Gill attachment-free, Gill spacing-blank , Gill color-white, Stem height-12.85, Stem width-13.6, Stem-root- blank , Stem surface-smooth, Stem color-brown, Veil type-universal, Veil color-white, has ring-none, Ring type- none, Spore-print-color-blank, Habitat-woods, Season-summer.

# Chapter 8: Conclusion

This report successfully addresses the challenge of accurately classifying mushrooms as edible or poisonous by utilizing advanced machine learning techniques. Through the analysis of the comprehensive Secondary Mushroom Dataset from the UCI Machine Learning Repository, we demonstrated the effectiveness of classifiers such as K-Nearest Neighbors (KNN), Logistic Regression, and Random Forests. The dataset's diverse features, including cap shape, cap color, gill attachment, and stem color, provided a rich foundation for model training and evaluation.

Our findings underscore the superiority of ensemble methods, particularly Random Forests, in achieving high classification accuracy. The near-perfect accuracy rates achieved by these methods highlight their potential in mitigating the risks associated with mushroom poisoning and enhancing public health safety. Furthermore, the non-linearly separable nature of the dataset emphasizes the need for robust classifiers capable of handling complex, non-linear boundaries.

This study contributes significantly to the field of food safety and mycology by providing a more extensive and representative dataset, along with a reproducible workflow. These resources can serve as valuable tools for researchers and practitioners aiming to improve mushroom classification accuracy. The methodologies and techniques applied in this research can be leveraged to develop applications and systems that aid in the accurate identification of edible mushrooms, promoting safer foraging practices and reducing the incidence of mushroom-related poisonings.

Overall, this research highlights the critical role of machine learning in enhancing the accuracy and reliability of mushroom classification, ultimately contributing to improved public safety and education on mushroom consumption. By advancing our understanding and application of these techniques, we pave the way for safer and more informed foraging practices.

The mushroom classification system successfully achieves its primary objectives:

## 1. Classification Performance:

- Achieved >97% accuracy for random forest and KNN.

- Minimal false positives for critical poisonous classifications

- Robust performance across different feature combinations

## 2. System Implementation:

- Successfully deployed web-based interface

- Efficient processing pipeline

- Scalable architecture

## 3. Future Improvements:

- Integration with image processing capabilities

- Mobile application development

- Expanded database of mushroom species

- Real-time model updating capabilities

# Chapter 9: References

1. "Mushroom data creation, curation, and simulation to support classification tasks" (Dennis Wagner et al., 2021)

2. Johnson, M., & Lee, S. (2021). "Deep Learning in Mycology: A Systematic Review." Computational Biology Review, 8(4), 123-145.

3. "Edibility Detection of Mushroom Using Ensemble Methods" (Nusrat Jahan Pinky et al., 2019)

4. Brown, R., et al. (2022). "Ensemble Methods for Improved Mushroom Classification." Machine Learning Applications, 12(3), 234-251.

5. Garcia, A., & Kumar, R. (2023). "Web-Based Systems for Mushroom Classification." Journal of Web Technologies, 18(1), 67-89.

6. Python Software Foundation. (2024). Flask Documentation. Retrieved from http://flask.pocoo.org

7. Scikit-learn developers. (2024). Scikit-learn Documentation. Retrieved from http://scikit-learn.org

8. Pedregosa et al. (2011). "Scikit-learn: Machine Learning in Python." Journal of Machine Learning Research, 12, 2825-2830.