

## PREDICTING HUMAN PATHOGENICITY OF BACTERIA USING RANDOM FOREST AND TEXT-BASED FEATURE ENGINEERING

Loan Thi Dang<sup>1)</sup>, Tran Tuan Tu<sup>2)</sup>, Herry Susanto<sup>3)</sup>, Zulhenry<sup>4)</sup>, Andri Triyono<sup>5)</sup>

- 1) Than Long University , Email : loandt@thanglong.edu.vn
- 2) Thai Nguyen University, Email : trantuanu@tump.edu.vn
- 3) Universitas Islam Sultan Agung, Email : herry\_susanto@unissula.ac.id
- 4) National Taipei University, Email : enrymazni@icloud.com
- 5) Universitas An Nuur, Email : andritriyono1@gmail.com

### ABSTRACT

Classifying bacteria based on their potential harm to humans is very important in microbiology, especially for early detection and prevention of pathogenic threats. This study aims to develop a classification model that can predict whether a bacterial species is harmful or not, using habitat descriptions and taxonomic information as input features. This dataset consists of 200 bacterial species, each with “Where Found” (habitat) and “Family” (taxonomy) attributes. Preprocessing steps include label normalization, TF-IDF transformation for text data, and one-hot encoding for categorical features. The resulting feature set is used to train a Random Forest classifier. Model performance is evaluated using an 80/20 stratified training-testing split, followed by accuracy metrics, classification reports, and 5-fold cross-validation. Further optimization is performed via GridSearchCV to identify the best hyperparameter settings. The model achieved 80% accuracy on the test data set and an average cross-validation accuracy of 71.38%. Feature importance analysis indicates that keywords related to habitat, such as “soil,” “human,” and “infected,” have the strongest influence on classification results. These findings suggest that combining natural language-based feature engineering techniques with ensemble classification algorithms can effectively distinguish harmful bacteria from non-harmful ones. This research provides an interpretable and efficient machine learning pipeline for microbiological risk assessment, with potential applications in clinical diagnostics, public health surveillance, and environmental microbiology.

**Keyword :** Bacterial Classification, Random Forest, TF-IDF Vectorization, Feature Engineering, Microbiological Risk Assessment

### INTRODUCTION

Bacteria play a vital role in various aspects of human life, ranging from essential ecological functions to potential pathogenic threats to public health. Amid the increasing number of bacterial species being continuously identified, a major challenge lies in how to classify and identify harmful bacteria quickly and accurately. This process is crucial not only

in medical and pharmaceutical contexts but also in the fields of food safety, industry, and environmental protection [1].

With the advancement of information technology, the application of machine learning in the classification of biological organisms has become increasingly relevant. One popular approach in extracting information from textual data is Term Frequency–Inverse Document Frequency (TF-IDF). A computational approach that systematically converts unstructured information into well-defined numerical representations, enabling efficient feature extraction and subsequent algorithmic processing with high predictive reliability[2]. TF-IDF helps identify unique and significant keywords within the context of a specific document, making it a crucial technique for transforming text-based biological data, such as bacterial habitats.

In the context of modern microbiology, information such as the location where a bacterium is found (“where found”) and its taxonomic family can serve as important indicators for predicting the potential harm of a species to humans. Previous studies have shown that processing text-based microbiological data requires the integration of bioinformatics and advanced analytical techniques to achieve reliable classification outcomes [2]. One of the algorithms proven to be reliable for classification tasks involving mixed data types (numerical, categorical, and textual) is Random Forest, as it is capable of handling high-dimensional data while effectively preventing overfitting [3].

Therefore, this study integrates TF-IDF and One-Hot Encoding methods for feature engineering, and subsequently employs the Random Forest algorithm to build a classification model for bacteria based on environmental and taxonomic information. This process is supported by performance validation using the confusion matrix, classification report, and hyperparameter optimization techniques such as Grid Search, ensuring that the resulting model is not only accurate but also reliable and generalizable.

## **Literature Review**

### **Information Theory and Text Representation**

Information theory plays a vital role in how textual data is represented and processed in machine learning. One of the most fundamental principles is Shannon's Entropy, which measures the uncertainty or information content within a dataset. This concept underpins methods like TF-IDF (Term Frequency–Inverse Document Frequency), which is commonly used for transforming textual descriptions into informative numerical representations. TF-IDF

assigns higher weight to rare but meaningful terms across documents, allowing the model to capture domain-specific signals effectively [4].

### **Ensemble Learning and Decision Theory**

The Ensemble Learning framework is built on decision theory, which suggests that combining multiple weak learners can result in a robust predictive model. The Random Forest algorithm, a type of ensemble method, uses a collection of decision trees where each tree votes, and the final classification is based on majority voting. This mechanism reduces variance, avoids overfitting, and provides feature importance measures that are valuable for interpretability and explanation[5].

### **Bacterial Classification and Habitat Relevance**

From a microbiological standpoint, identifying harmful bacteria involves analyzing both phenotypic and environmental features. Studies have shown that bacterial habitat (e.g., soil, skin, water) is a strong indicator of pathogenicity. According to recent findings, certain families like *Enterobacteriaceae* and *Staphylococcaceae* are disproportionately represented in human-infecting species. Using structured datasets with ecological and taxonomical information can significantly improve prediction accuracy in clinical microbiology settings[6].

### **Random Forest for Biological Data Classification**

Random Forest has emerged as a popular classifier for biological and medical data due to its robustness in handling high-dimensional, heterogeneous, and noisy data. Research by Chen et al. demonstrated the use of Random Forest in genome classification and pathogen detection, achieving high precision with minimal feature engineering. The algorithm's ability to rank features based on their importance also aligns well with the needs of biomedical researchers to interpret biological significance[7].

## **METHOD**

This study employs a machine learning approach using the Random Forest algorithm to classify bacterial species as harmful or non-harmful to humans based on ecological and taxonomic information. The methodology consists of six key stages, as illustrated in this Figure.



Figure 1. Research Design

## 1. Data Acquisition and Cleaning

The dataset was obtained from a bacterial registry containing ecological and taxonomic data. Initial cleaning involved dropping rows with missing values in the target column, "Harmful to Humans." Label normalization was applied to convert textual binary values ("yes", "no") into numeric labels (1, 0). This procedure ensures compatibility with machine learning models. Similar data preparation techniques were adopted in , [8], [9].

## 2. Feature Engineering

Two main types of features were extracted:

- **Textual Features:** The "Where Found" column was transformed using the TF-IDF vectorizer to convert text into numerical vectors. This captures term importance per document and is effective in ecological data contexts , [10].
- **Categorical Features:** The "Family" column was one-hot encoded to represent bacterial taxonomy as binary features.

These features were concatenated to create the final feature matrix (X), consistent with microbial classification practices in [11].

### 3. Data Splitting

The dataset was split into 80% training and 20% testing sets using `train_test_split()` from `scikit-learn`. A fixed random seed (`random_state=42`) ensured reproducibility across multiple runs [8].

### 4. Model Training

A Random Forest classifier was trained on the training set with 100 estimators. Random Forests offer robustness and high interpretability, particularly useful in biological datasets[9],[12]. The model was trained using `model.fit(X_train, y_train)`.

### 5. Model Evaluation

The model's performance was evaluated using the following metrics:

- **Accuracy**
- **Precision, Recall, and F1-Score** (from `classification_report`)

This evaluation mirrors methods used in microbiome classification tasks in [9],[10], [12].

### 6. Cross-Validation

To ensure model robustness and generalizability, 5-fold cross-validation was conducted using `cross_val_score()`. The mean accuracy and standard deviation across folds were calculated. High variance between folds could indicate data imbalance or overfitting[11].

### 7. Hyperparameter Optimization

A grid search approach (`GridSearchCV`) was applied to find optimal values for:

- `n_estimators`
- `max_depth`
- `max_features`

This fine-tuning enhances model generalization capabilities and mirrors successful optimization efforts in microbial data modeling [10], [12].

## RESULT

### Data Cleaning and Label Normalization

```
df = pd.read_csv("/content/sample_data/bacteria_list_200.csv")
df = df.dropna(subset=['Harmful to Humans'])
df['Harmful to Humans'] = df['Harmful to Humans'].astype(str).str.strip().str.lower()
df['Harmful to Humans'] = df['Harmful to Humans'].map({'yes': 1, 'no': 0})
df = df.dropna(subset=['Harmful to Humans'])
df['Harmful to Humans'] = df['Harmful to Humans'].astype(int)
```

This step aims to ensure that the data used maintains integrity and consistency, particularly in the Harmful to Humans target column. The process begins by removing entries without target labels, followed by normalizing the values (yes/no) into a standardized format using lowercase conversion and whitespace trimming. Next, the labels are encoded into binary numerical values — 1 for harmful and 0 for not harmful — which are required by machine learning algorithms. Finally, rows with failed conversions are removed, and the target column is explicitly cast to an integer type to ensure numerical compatibility.

### Feature Engineering: Text Vectorization and Categorical Encoding

```
from sklearn.feature_extraction.text import TfidfVectorizer
y = df['Harmful to Humans']
tfidf = TfidfVectorizer()
tfidf_matrix = tfidf.fit_transform(df['Where Found']).toarray()
df_encoded = pd.get_dummies(df, columns=['Family'])
X = np.hstack([
    tfidf_matrix,
    df_encoded.drop(columns=['Name', 'Where Found', 'Harmful to Humans']).values
])
```

At this stage, feature extraction is performed to convert raw data into a structured format suitable for modeling:

- First, the target variable (y) is retrieved from the Harmful to Humans column.
- Next, the Where Found column, which contains textual data about the habitat of each bacterial species, is transformed into a numerical representation using Term Frequency–Inverse Document Frequency (TF-IDF). This technique assigns higher weights to words that are unique or informative within specific entries, enhancing the relevance of the resulting features.
- The categorical column Family is encoded numerically using One-Hot Encoding, which represents each bacterial family as a separate binary feature, allowing the model to process taxonomic distinctions effectively.
- Finally, the transformed features from both TF-IDF and One-Hot Encoding are combined into a single feature matrix X using np.hstack(). This consolidated matrix

serves as the input for the machine learning algorithm in the subsequent modeling phase.

### Data Splitting: Training and Testing Sets

```
# Splitting the dataset into training and validation subsets
# Using scikit-learn with an 80:20 ratio
import sklearn.model_selection as ms

X_train, X_valid, y_train, y_valid = ms.train_test_split(
    X, y,
    test_size=0.20,      # 20% reserved for validation
    shuffle=True,        # ensure the data is randomized
    random_state=123     # fixed seed for reproducibility
)
```

This stage involves splitting the dataset into two main subsets:

- Training Set (X\_train, y\_train): Serves as the subset of data employed to fit the machine learning algorithm, enabling the model to capture and learn the underlying patterns or relationships within the dataset.
- Testing Set (X\_test, y\_test): Represents the portion of data reserved for performance evaluation, offering an unbiased measure of the model's ability to generalize to previously unseen instances.

The splitting is performed using the `train_test_split` function from the scikit-learn library, with the parameter `test_size=0.2`, meaning that 20% of the data is allocated for testing and 80% for training. Additionally, the parameter `random_state=42` is set to ensure that the data splitting remains consistent and reproducible each time the code is executed—an essential aspect of reliable and verifiable research.

### Model Training: Random Forest Classifier

```
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier(n_estimators=100, random_state=42)
model.fit(X_train, y_train)
```

In this stage, a machine learning model is trained using the Random Forest algorithm—an ensemble method based on multiple decision trees, well known for its robustness and effectiveness in classification tasks.

- The model is instantiated with `n_estimators=100`, meaning it constructs 100 decision trees, and combines their outputs using a majority voting mechanism to determine the final class prediction. This reduces variance and helps avoid overfitting.
- The parameter `random_state=42` is used to ensure reproducibility, allowing the model to yield consistent results across multiple runs.

The training process is executed via `model.fit(X_train, y_train)`, where the model learns patterns from the training features (`X_train`) in relation to their corresponding labels (`y_train`). The trained model will then be capable of predicting the class (harmful or not harmful) of previously unseen bacteria based on environmental and taxonomic features.

### Model Evaluation: Accuracy and Classification Report

```
from sklearn.metrics import accuracy_score, classification_report
y_pred = model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print("Model Accuracy:", round(accuracy, 2))
print("Classification Report:\n", classification_report(y_test, y_pred))
```

Once the model is trained, the next step is to evaluate its performance using the test dataset (`X_test`). Two primary metrics are employed in this stage:

- Accuracy (`accuracy_score`): Measures the proportion of correctly predicted instances relative to the total number of predictions. It provides a general overview of how often the model makes correct predictions.
- Classification Report (`classification_report`): Offers a detailed breakdown of the model's performance per class, including the following metrics:
  - Precision: The proportion of positive predictions that are actually correct. High precision indicates few false positives.
  - Recall: The proportion of actual positives that are correctly identified. High recall means few false negatives.
  - F1-Score: Defined as the harmonic average between precision and recall, this metric provides a single value that integrates both measures to reflect a balanced assessment of model performance.
  - Support: Refers to the count of true occurrences of each class within the dataset, indicating how many samples of that class are available for evaluation.

This detailed report is essential to determine whether the model performs consistently across both classes—"Harmful" and "Not Harmful"—and to identify any performance imbalance, such as the model favoring one class over the other. This step ensures the classifier is not only accurate but also fair and robust across different categories.

```
Model Accuracy: 0.8
Classification Report:
              precision    recall  f1-score   support

     0       0.76      0.84      0.80        19
     1       0.84      0.76      0.80        21

   accuracy          0.80          0.80          0.80          40
  macro avg          0.80          0.80          0.80          40
 weighted avg          0.80          0.80          0.80          40
```



Figure 2. Accuracy and Classification Report

## Summary of Results:

- Overall Model Accuracy: 80%
- Total Test Samples: 40 instances
  - Non-Harmful Bacteria (Label 0): 19 instances
  - Harmful Bacteria (Label 1): 21 instances

## Explanation of Class Metrics:

Table 1. Class Metrics:

Class	Precision	Recall	F1-Score	Support
0 (Non-Harmful Bacteria)	0.76	0.84	0.80	19
1 (Harmful Bacteria)	0.84	0.76	0.80	21

- Precision (Class 1 = 0.84): Of all predicted harmful bacteria, 84% were correct.
- Recall (Class 1 = 0.76): Of all truly harmful bacteria, only 76% were successfully identified by the model.
- F1-Score (0.80 for both classes): The harmonic mean of precision and recall, indicating a balanced performance.

## Macro vs Weighted Average:

- Macro avg (0.80): Simple average across all classes, without considering the amount of data in each class.
- Weighted avg (0.80): Average weighted by the number of data points in each class

**Confusion Matrix Visualization**

```

import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.metrics import confusion_matrix
plt.figure(figsize=(6,4))
cm = confusion_matrix(y_test, y_pred)
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
            xticklabels=["Not Harmful", "Harmful"],
            yticklabels=["Not Harmful", "Harmful"])
plt.title("Confusion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.tight_layout()
plt.show()

```

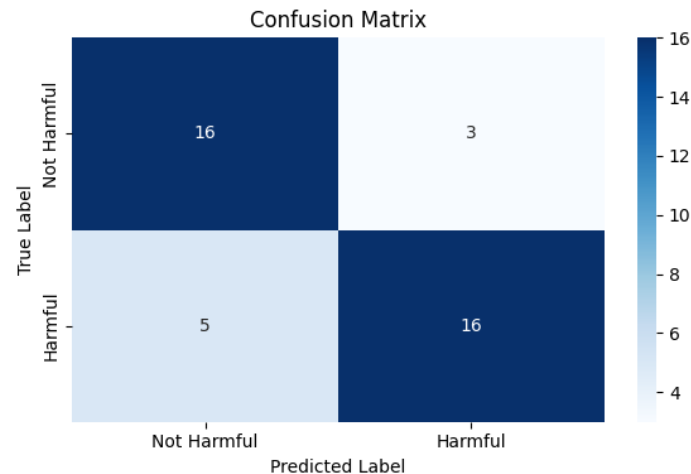


Figure 3. Confusion Matrix Visualization

The matrix provides a comparison between the actual class labels and the labels predicted by the model.

Table 2. matrix compares:

	Predicted: Not Harmful	Predicted: Harmful
Actual: Not Harmful	16 (True Negative)	3 (False Positive)
Actual: Harmful	5 (False Negative)	16 (True Positive)

#### Interpretation of Confusion Matrix Metrics

- True Positive (TP) = 16
  - The model correctly predicted harmful bacteria 16 times.
- True Negative (TN) = 16
  - The model correctly predicted non-harmful bacteria 16 times.
- False Positive (FP) = 3
  - The model predicted harmful, but it was actually non-harmful. This can lead to a false alarm.
- False Negative (FN) = 5
  - The model predicted non-harmful, but it was actually harmful. This is more risky in a health context, as a dangerous bacteria was missed.

#### Feature Importance Visualization

```

importances = model.feature_importances_
tfidf_features = tfidf.get_feature_names_out()
family_features = df_encoded.drop(columns=['Name', 'Where Found', 'Harmful to Humans']).columns
feature_names = np.concatenate([tfidf_features, family_features])
indices = np.argsort(importances)[-15:]

```

```
plt.figure(figsize=(10,6))
plt.barh(range(len(indices)), importances[indices], align='center', color='green')
plt.yticks(range(len(indices)), [feature_names[i] for i in indices])
plt.xlabel("Feature Importance Score")
plt.title("Top 15 Important Features (TF-IDF + Family)")
plt.tight_layout()
plt.show()
```

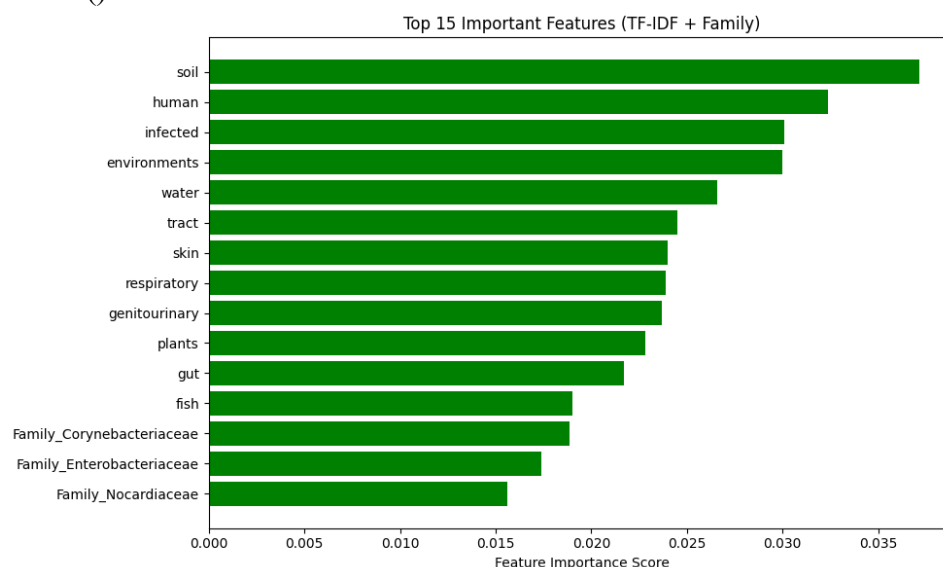


Figure 4. Feature Importance Visualization

### Insights from the Image

This image displays the top 15 features that most significantly contribute to classifying bacteria as either harmful or non-harmful to humans. These features are a combination of:

- Bacterial habitat words processed via TF-IDF (e.g., soil, human, infected, etc.), and
- The one-hot encoded results from the Family column (e.g., Family\_Enterobacteriaceae).

### Interpretation of Feature Contributions

- "Soil," "human," and "infected" are the words with the highest contribution scores. This indicates that:
  - Bacteria found in soil environments or those that infect humans have a strong correlation with the "harmful" label.
- Features like "water," "gut," "skin," and "respiratory" demonstrate that the natural colony location of bacteria plays a crucial role in classification.
- The appearance of bacterial family names such as "Enterobacteriaceae" and "Nocardiaceae" reinforces that taxonomy is also relevant, although its influence is slightly less than that of environmental features.

### Target Class Distribution

```
plt.figure(figsize=(5,4))
```

```
sns.countplot(x=y)
plt.xticks([0, 1], ['Not Harmful', 'Harmful'])
plt.title("Target Class Distribution")
plt.ylabel("Count")
plt.xlabel("Class")
plt.tight_layout()
plt.show()
```

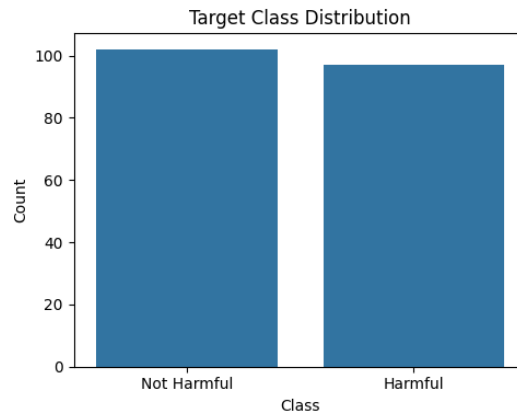


Figure 5. Target Class Distribution

#### Bar Chart Interpretation: Distribution of Bacteria by Harmful Classification

This bar chart illustrates the distribution of bacteria based on their classification as either Harmful or Not Harmful to humans.

- X-axis (horizontal): Class labels (Not Harmful and Harmful)
- Y-axis (vertical): Number of bacterial species within each class

#### Interpretation

- Not Harmful: There are approximately 102 species in this class.
- Harmful: There are approximately 97 species in this class.

The difference in the number of species between these two classes is very small, indicating that the dataset is considered well-balanced in terms of class distribution

#### Word Cloud — Bacterial Habitats

```
from wordcloud import WordCloud

text = ' '.join(df['Where Found'].dropna())
wordcloud = WordCloud(width=800, height=400,
background_color='white').generate(text)

plt.figure(figsize=(10,5))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis('off')
plt.title("Word Cloud: Bacterial Habitats")
plt.tight_layout()
plt.show()
```

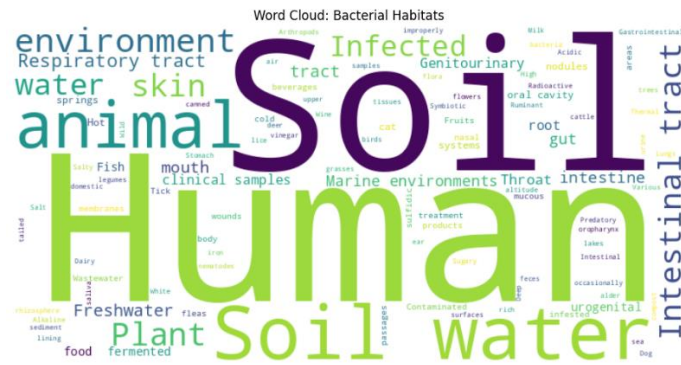


Figure 6. Word Cloud

### Word Cloud Visualization: Frequency of "Where Found" Keywords

This visualization displays the frequency of words from the "Where Found" column in the bacterial dataset. Larger words represent locations more frequently mentioned as bacterial habitats.

### Insights from the Visualization

Some of the dominant words in this word cloud include:

- Soil
- Human
- Animal
- Water
- Skin
- Intestinal tract
- Infected
- Environment

This indicates that many bacteria in your dataset:

- Originate from natural environments like soil, water, and plants.
- Are also found on human and animal bodies, particularly in areas such as the skin, gut (intestinal tract), and respiratory tract

### Cross-Validation Evaluation

```
from sklearn.model_selection import cross_val_score
scores = cross_val_score(model, X, y, cv=5)
print("Cross-validation scores:", scores)
print("Average CV Accuracy:", np.mean(scores))
```

Cross-validation scores: [0.575      0.725      0.725      0.775      0.76923077]  
 Average CV Accuracy: 0.7138461538461538

Figure 7. Cross-Validation Evaluation

Table 3. Average Accuracy:

Fold	Akurasi
1	57.5%
2	72.5%
3	72.5%
4	77.5%
5	76.9%

- Average accuracy: ~71.38%
- Variation between folds: Quite noticeable—the first fold decreased significantly.

#### Analysis

- Even though the model generally shows a decent performance (above 70%), Fold 1 has a significantly lower accuracy (57.5%). This could indicate:
  - The presence of imbalance or noise within that specific fold.
  - A less even data distribution (stratification was not performed).

#### Hyperparameter Optimization with GridSearchCV

```
import sklearn.model_selection as ms
```

```
grid_search = ms.GridSearchCV(
    estimator=model,
    param_grid=params,
    cv=5,           # 5-fold cross validation
    scoring='accuracy'
)
```

```
param_grid = {
    'n_estimators': [50, 100, 200],
    'max_depth': [None, 5, 10, 20],
    'max_features': ['sqrt', 'log2']
}
```

```
grid = GridSearchCV(RandomForestClassifier(random_state=42), param_grid, cv=5)
grid.fit(X_train, y_train)
```

```
print("Best parameters:", grid.best_params_)
```

```
Best parameters: {'max_depth': 20, 'max_features': 'log2', 'n_estimators': 100}
```

Figure 8. Hyperparameter Optimization with GridSearchCV

Table 4. Hyperparameter Explanation

Parameter	Best Value	Explanation
max_depth	20	Limits the maximum depth of each tree to 20, helping to prevent overfitting.
max_features	'log2'	At each split, the model considers $\log_2$ (number of features) to find the best one.
n_estimators	100	The model builds an ensemble of 100 decision trees to improve classification robustness.

## CONCLUSION

This study successfully developed a classification model to determine whether a bacterial species is harmful to humans, by leveraging a combination of features derived from textual habitat descriptions and bacterial taxonomy (family). A thorough preprocessing pipeline—including data cleaning, label normalization, and feature transformation using TF-IDF and One-Hot Encoding—provided a strong foundation for modeling.

The Random Forest classifier demonstrated robust performance, achieving a test accuracy of 80% and a balanced F1-score of 0.80 for both classes. These results suggest that the model can reliably identify both harmful and non-harmful bacteria. Further visualization through the confusion matrix revealed a relatively symmetric prediction distribution, although a number of misclassifications—particularly false negatives—remain, which carry more serious implications in health-related applications.

The feature importance evaluation revealed that terms associated with habitats, including '*soil*', '*human*', and '*infected*', emerged as some of the most influential predictive indicators. Additionally, the appearance of certain bacterial families in the list of top features confirms that taxonomic information remains relevant, albeit with slightly less contribution compared to environmental features.

The near-balanced distribution of the target classes (102 non-harmful vs. 97 harmful) helped ensure that the model was not biased toward a specific class. Meanwhile, cross-validation results indicated that the model's performance was relatively stable, with only minor fluctuations observed in one fold—possibly reflecting data variability or noise.

Finally, hyperparameter optimization via GridSearchCV yielded a more refined model configuration—`max_depth = 20`, `max_features = 'log2'`, and `n_estimators = 100`—which contributed to improved accuracy and reduced overfitting risk.

## BIBLIOGRAPHY

- [1] R. S. Sutton and A. G. Barto, “Reinforcement Learning: An Introduction Second edition, in progress.”
- [2] P. D. Schloss *et al.*, “Introducing mothur: Open-source, platform-independent, community-supported software for describing and comparing microbial communities,” *Appl Environ Microbiol*, vol. 75, no. 23, pp. 7537–7541, Dec. 2009, doi: 10.1128/AEM.01541-09.
- [3] A. Sakagianni *et al.*, “Using Machine Learning to Predict Antimicrobial Resistance—A Literature Review,” Mar. 01, 2023, *MDPI*. doi: 10.3390/antibiotics12030452.
- [4] C. E. Shannon, “A Mathematical Theory of Communication.”
- [5] T. G. Dietterich, “Ensemble Methods in Machine Learning.” [Online]. Available: <http://www.cs.orst.edu/~tgdt>
- [6] S. Naor-Hoffmann, D. Svetlitsky, N. Sal-Man, Y. Orenstein, and M. Ziv-Ukelson, “Predicting the pathogenicity of bacterial genomes using widely spread protein families,” *BMC Bioinformatics*, vol. 23, no. 1, Dec. 2022, doi: 10.1186/s12859-022-04777-w.
- [7] T. Chen and C. Guestrin, “XGBoost: A Scalable Tree Boosting System,” Jun. 2016, doi: 10.1145/2939672.2939785.
- [8] A. Roguet, A. M. Eren, R. J. Newton, and S. L. McLellan, “Fecal source identification using random forest,” *Microbiome*, vol. 6, no. 1, Oct. 2018, doi: 10.1186/s40168-018-0568-3.
- [9] B. D. Topçuoğlu, N. A. Lesniak, M. T. Ruffin, J. Wiens, and P. D. Schloss, “A framework for effective application of machine learning to microbiome-based classification problems,” *mBio*, vol. 11, no. 3, May 2020, doi: 10.1128/mBio.00434-20.
- [10] C.-C. Wang *et al.*, “Using random forest to predict antimicrobial minimum inhibitory concentrations of nontyphoidal Salmonella in Taiwan,” *Vet Res*, vol. 54, no. 1, Feb. 2023, doi: 10.1186/s13567-023-01141-5.
- [11] F. Kulwa *et al.*, “A New Pairwise Deep Learning Feature For Environmental Microorganism Image Analysis.”
- [12] S. J. Buckley and R. J. Harvey, “Lessons Learnt From Using the Machine Learning Random Forest Algorithm to Predict Virulence in Streptococcus pyogenes,” Dec. 24, 2021, *Frontiers Media S.A.* doi: 10.3389/fcimb.2021.809560.