# Investigation of Machine Learning Models for Foodborne Disease Classification

**Wuletawu Iyasu and Degif Teka***

Corresponding Author's Email: degiftk@gmail.com

*Department of Computer science, Institute of Technology, Hawassa University, Hawassa, Ethiopia.*

*Abstract*: Foodborne disease is a disease that has a high prevalence in low- and middle-income countries around the world. There are many people affected by foodborne disease in Ethiopia, due to various causes. There are high rate of infections; the control of most foodborne diseases in Ethiopia is low due to a lack of knowledge, medication and support to healthcare professionals for better diagnoses. Machine learning applications in the healthcare and biomedical domain are popular for the early detection of diseases and help to make a better diagnosis. Machine learning models can learn from past data, identify patterns and make decisions with a minimal human intervention. Though there are studies that apply machine learning for medical diagnosis and other fields there is a lack of study conducted to classify the foodborne diseases which are common in Ethiopia. This study focuses on foodborne diseases by selecting some of the prevalent foodborne illnesses in Ethiopia including Typhoid fever, Giardiasis, and Amoebiasis in consultation with medical experts. To achieve the objective of the study the researcher used an experimental research design. Data is collected from two Hospitals. After preprocessing the collected data, the researcher trained the developed model using state-of- art machine learning algorithms including Decision Tree, Random Forest, XGBoost and Stacking ensemble learning method. Based on the experiment conducted, the Stacking ensemble learning method model outperforms the others with an accuracy of 98.06% (98.1%), followed by Random Forest, XGBoost, and Decision Tree with accuracy of 97.5%, 96.9%, and 96.5% respectively. The result obtained by the study depicts that, the Stacking ensemble learning model is suitable for diseases classification.

*Keywords:* Disease Classification, Disease Prevalence, Machine Learning Model, Stacking ensemble learning method, foodborne Disease.

## 1.     Introduction

Access to safe and nutritious food is a fundamental human right and essential for well-being (Hossain, 2018). However, foodborne diseases, caused by consuming contaminated food or water, pose a significant threat to public health globally (Mahon, 1998). They are usually associated with contaminated foods and pathogens contained in foods. According to a World Health Organization (WHO) report, 600 million people worldwide suffered from diseases caused by eating contaminated

food every year, of whom 420000 die (WHO, 2008). In Ethiopia, food safety is a major concern due to lack of infrastructure and basic pre-requisites for food safety such as clean water and environment, washing facilities, compounded by limited implementation of food safety regulations, and a lack of incentives for producers to improve food safety (Gazu, 2023). Frequent foodborne illness outbreaks both domestically and internationally pose a major threat to social stability and public health, making them a global public health and food safety concern.

Machine Learning (ML) is a subfield of artificial intelligence that can learn from experience, identify patterns, and make decisions with minimal human intervention using its algorithms. Classification is one of the machine learning tasks that is used to classify different patterns based on past data. These applications in the healthcare domain have led to early detection of disease and helped to make a better diagnosis (Alanazi, 2022).

Several studies are conducted on disease classification using machine learning based on symptoms and image analysis. Among them, the classification of heart disease, kidney disease, stroke, breast cancer, and pneumonia disease are most commonly investigated in the literature using machine learning algorithms (Mohammad, 2021). But there is a lack of research conducted around the world on foodborne disease classification in order to detect these diseases early before spreading and identify the specific disease from other diseases which shares symptoms in common. From the reviewed literature there is an attempt to classify foodborne diseases such as Salmonella (which can manifest as Typhoid fever), Norovirus (known to cause acute gastroenteritis with vomiting and diarrhea), E. coli (causing a range of symptoms from mild to severe cramps and diarrhea, sometimes with blood), and Vibrio parahaemolyticus (a bacterial infection acquired through consuming raw or undercooked seafood). However, the most common foodborne diseases in Ethiopia are not included in the previous study (Mohammad, 2021). Thus, there is a need for a classifier model that can classify the most common foodborne disease that occur in Ethiopia in order to analyze the prevalence of a disease in a specific area and ease the detection of specific diseases early to make a better decision while diagnosing. Therefore, the proposed study mainly focused foodborne diseases with high prevalence in the study area, particularly Typhoid fever, Giardiasis, and Amoebiasis, which are also common throughout Ethiopia. It is essential to recognize that these prevalence rates can fluctuate across different regions within Ethiopia. The varied cultural practices across Ethiopia likely contribute to differing prevalence of foodborne diseases in various regions (Teferi, 2020). The proposed study significantly helps medical doctors to make a proper diagnosis of specific foodborne diseases at less time.

To make a classification in machine learning, it needs datasets to learn from the data and then classify the expected outcomes from unseen data. So, we collected the data on the most common foodborne diseases from the Hospitals. Then, the collected dataset is preprocessed in such a way that it

is appropriate for the selected machine learning algorithms. Then, the models were trained and their performance is evaluated by unseen data. Finally, a suitable algorithm is selected to develop a model based on the experiment conducted. To preprocess the dataset and develop the model, we used python programming language.

There are many machine learning algorithms such as Logistic Regression (LR), Naïve Bayes (NB), K-Nearest Neighbors (KNN), Decision Tree (DT), Support Vector Machines (SVM), Random Forest (RF) which are discussed in the literature (Mohammad, 2021). For the proposed model development, the researcher used the state of art machine learning algorithms such as Decision Tree, Random Forest and eXtreme Gradient Boosting (XGBoost) because those algorithms are promising algorithms so as to have a model with a good performance. In addition, the algorithm like Random Forest is used to handle data imbalance (Agarwal, 2020). Furthermore, the researcher implemented a stacking ensemble method which leverages the strengths of individual models to potentially achieve superior classification performance (Brownlee, 2021).

The following are the major contributions of the proposed research:

The research successfully developed a machine learning model that utilizes various state-of-the-art algorithms, including Decision Tree, Random Forest, XGBoost, and a stacking ensemble method. This combination enhances the predictive performance and accuracy of foodborne disease classification in the Ethiopian context.

The study fills a critical gap in the existing literature by focusing on the classification of foodborne diseases prevalent in Ethiopia, an area that has been underexplored. By addressing local public health concerns, the research contributes valuable insights for healthcare professionals in diagnosing foodborne illnesses more effectively.

The research collects and preprocesses a considerable dataset of foodborne disease cases from hospitals. This comprehensive dataset not only aids in the model development but also serves as a valuable resource for future studies and analyses in the field of foodborne disease management.

By providing a machine learning model that healthcare professionals can utilize, the research improves the diagnostic process for foodborne diseases. This can lead to faster and more accurate identification of illnesses, ultimately contributing to better patient outcomes and potentially reducing the burden on healthcare facilities (Wang et al., 2020).

The application of advanced machine learning techniques, particularly the stacking ensemble method, demonstrates the potential for improved classification accuracy over traditional models. The results indicate that this approach can significantly outperform individual algorithms, providing a framework that could be applied to other classification problems beyond foodborne diseases.

The study provides a detailed framework for data preprocessing, model development, and evaluation, offering a guide for future research in this domain. The findings support policymakers in designing better monitoring and control strategies for foodborne diseases, contributing to improved public health outcomes.

## 2.        Literature Review

This section presents review of related literature that helps get understanding of the problem from the wider literature and identify the gap. It helps also learn how to design and develop models from literature recommendations. Oguntimilehin et al (2013) developed a classifier model that can identify Typhoid fever diseases using Rough Set machine learning techniques which is a valuable data analysis and knowledge discovery paradigm that plays a significant role within the broader field of Machine Learning (ML). The researchers worked with a dataset comprising 150 instances. The Rough Set algorithm achieved a training set accuracy of 95% and a testing set accuracy of 96% based on the experiment results. The model was evaluated using a confusion matrix to determine the number of correct and incorrect classifications per class while identifying Typhoid fever disease in unseen data. Although the problem domain is relevant to the one proposed, the researchers did not include diseases that share symptoms with Typhoid fever.

Radhika et al (2020) developed a classifier model for disease classification based on symptoms using a Decision Tree and Electronic Health Record Analysis. The model was trained to classify 40 symptom-based diseases, including Allergy, GERD, AIDS, Diabetes, Tuberculosis, Asthma, Dengue, Heart attack, Malaria, Typhoid Fever, and more. The Decision Tree initially classifies the disease with the highest confidence level for each of the classified diseases. The author then developed a predictive model to classify symptom-based diseases using the Decision Tree model. The study included foodborne disease like Typhoid Fever which is relevant to the proposed study but the study did not focus on most common foodborne diseases that are considered here in this paper. Better accuracy might also be obtained if the ensemble machine learning algorithms like RF, XGBOOST and Stacking ensemble learning were used  because these algorithms can potentially increase the accuracy of the model by boosting the gradient and taking the average of the weak learners performance from the decision tree, and predicting the outcomes (Kavzoglu and Teke, 2022).

A study by Wang et al. (2020) used machine learning to develop a model to classify foodborne diseases pathogens. The model used the following features to classify foodborne diseases: location, time of illness, age of patient, type of food consumed, and certain symptoms. The study divided 50,216 samples into training and test sets at a ratio of 7:3. The size of the training set was 35,151 samples, and the size of the test set was 15,065 samples. The study analyzed foodborne disease case data and used

machine learning methods such as decision trees, random forests, K-Nearest Neighbors, stochastic gradient descent, and extremely randomized trees to classify foodborne disease pathogens. The gradient boost decision tree model obtained the highest accuracy, with accuracy approaching 69% in identifying 4 pathogens: Salmonella, Norovirus, Escherichia coli, and Vibrio parahaemolyticus (Wang, 2020).

The study Leo et al (2019) aimed to develop a classifier model that can predict the occurrence of cholera disease by using various machine learning algorithms such as K-Nearest Neighbors (K-NN), Decision Tree, XGBoost, ExtraTree, AdaBoost, Random Forest, and Linear Discriminant Analysis (LDA). The researcher conducted a comparative analysis to identify the best machine learning algorithm that can classify cholera disease. After comparing the algorithms, the researcher selected XGBoost as the best cholera predictor based on the experiment results. Although this study is relevant to the proposed study, it only focused on cholera disease classification. One of the strengths of this study is that it used ensemble-based algorithms like Random Forest and XGBoost to classify cholera disease.

The study by Grampurohi et al (2020) aimed to use machine learning algorithms to analyze and classify 41 diseases, with the goal of assisting physicians in diagnosing diseases at an early stage. The paper presents a comprehensive comparative study of three algorithms, namely Decision Tree, Random Forest, and Naïve Bayes, and their performance on a medical record. Each algorithm yielded an accuracy of 95%. The researcher attempted to classify multiple diseases using machine learning algorithms including foodborne diseases like Typhoid Fever which is relevant to the proposed study but the study did not focus on most common foodborne diseases.

This study by Bhuiyan et al (2023) represents a pioneering effort to develop a typhoid fever prediction model that can anticipate outcomes prior to clinical trials. The study utilized a dataset containing 1746 entries with 29 attributes. Both machine learning (ML) and deep learning (DL) approaches were employed and a total of ten algorithms are evaluated. Among these, the XGBoost classifier emerges as the most promising, achieving an impressive accuracy of 97.87%. This study is relevant to proposed study but it did not include diseases that share symptoms with Typhoid fever.

A study by Kumar et al. (2023) used artificial intelligence (AI) to develop a model to predict the pathogens in foodborne diseases. They used data from the National Foodborne Disease Surveillance Reporting System on Kaggle, which has 12 attributes. They divided the dataset into training (75%) and testing (25%) datasets. They then compared and contrasted the accuracy of several popular machine learning algorithms for predicting the pathogens in foodborne diseases, including decision trees, random forests, k-nearest neighbors, stochastic gradient descent, and extremely randomized trees. They also developed an ensemble model that combined all of these algorithms. The ensemble model outperformed all of the other classifiers, achieving an average accuracy of 97.26%. The research study is relevant to the proposed research but it only focused on foodborne diseases which are common in India.

The literature showed that previous research primarily concentrated on classifying generic and chronic diseases, such as heart disease, diabetes, stroke, lung cancer, breast cancer, and malaria. Researchers have employed various machine learning algorithms like SVM, KNN, NB, DT, and RF for disease classification. However, there is a notable gap: no existing machine learning model specifically classifies common foodborne diseases in Ethiopia. In contrast, a study conducted in China developed a machine learning model that classifies foodborne diseases common in that country. Unfortunately, previous research on foodborne disease classification did not include conditions like giardiasis, Amoebiasis both of which are prevalent in Ethiopia especially in the case study considered. This lack of comprehensive classification support poses challenges for healthcare professionals who encounter foodborne diseases with similar symptoms. Consequently, the burden of foodborne diseases in Ethiopia remains high, constituting a significant public health issue. To address this gap, we focused on classifying common foodborne diseases in Ethiopia using machine learning algorithms. The study collected disease data from the specific research area and analyzed the prevalence of foodborne diseases within that region.

## 3. Experimental datasets & data preprocessing

### 3.1 Dataset Description

The data for the proposed study was obtained from two selected hospitals within the Wolaita Zone, southern Ethiopia. The Total number of foodborne diseases data collected from hospitals for model development is 3100 confirmed cases with 24 attributes based on the sampling technique applied. When collecting the data the sample frame that forms the population like age group, gender, and year interval that the sample taken was considered. Table 1 below shows the samples gathered for each foodborne disease along with the sample frames, and attributes of each disease (symptoms) as well as, the type of data collected. In order to avoid class imbalance, the researcher determined the year range during data collection. For instance, the confirmed cases of Typhoid fever in one year (2012 E.C) is greater than the confirmed cases of Giardiasis in three years (2012-2015 E.C) as shown in Table 1.

*Table 1: Dataset Description*

| Foodborne diseases | No. of sample taken | Age group | Gender | Year interval where the data Taken | Confirmed symptoms(attributes) | Data type |
|---|---|---|---|---|---|---|
| Amoebiasis | 1060 | All | Both | 2012-2015 E.C | Diarrhea, Vomiting, Abdominal Cramp, dysentery, Weakness. | Categorical |

| | | | | | | |
|---|---|---|---|---|---|---|
| Giardiasis | 834 | All | Both | 2012-2015 E.C | Diarrhea, Fatigue, Abdominal discomfort, Weight loss, Nausea, | Categorical |
| Typhoid Fever | 1206 | All | Both | 2014 E.C | Fever, chills, headache diarrhea, Abdominal flank, Muscle or joint pain, cough, Loss of appetite. | Categorical |

### 3.2. Data preprocessing

In the research, data preprocessing involved several critical steps to prepare the collected dataset of foodborne disease cases for model development. Initially, the data was meticulously cleaned by identifying and handling missing values using the mode method which is a preferable method for categorical numeric data imputation, thereby ensuring integrity and completeness (Dubey, 2024).

The preprocessing also included standardization and normalization of features, with specific attention to the 'age' column, which was normalized using the min-max technique to scale values between 0 and 1. Categorical variables, particularly those representing symptoms, were converted into a binary format where '1' indicated the presence and '0' the absence of symptoms. To do so, one hot encoder is used to transform. It is a powerful technique in machine learning for transforming nominal categorical data into a numeric format. It preserves independence by creating new binary variables for each unique category, ensuring no artificial ordinal relationship is imposed. This approach avoids misinterpretation, as label encoding assigns integer values based on alphabetical order, which can introduce misleading rankings (Bobbitt, 2022). Furthermore, feature selection was conducted using random forest feature importance analysis to retain only the most relevant attributes, thereby enhancing model performance while reducing complexity (Dubey, 2024).

*Table 2: Feature importance of each attribute in dataset*

| No. | Feature's name | Score of feature |
|---|---|---|
| 1 | Weakness | 0.243319 |
| 2 | Chills | 0.196843 |
| 3 | Nausea | 0.120400 |
| 4 | early fatigue | 0.088068 |
| 5 | Age | 0.058681 |
| 6 | AC(abdominal cramp) | 0.044468 |

| | | |
|---|---|---|
| 7 | Headache | 0.044365 |
| 8 | Diarrhea | 0.039404 |
| 9 | Vomiting | 0.028162 |
| 10 | AD(abdominal discomfort) | 0.023372 |
| 11 | Fever | 0.022850 |
| 12 | Back Pain | 0.016247 |
| 13 | Joint pain | 0.013212 |
| 14 | Gender | 0.011061 |
| 15 | Poor appetite | 0.007643 |
| 16 | Total body burn sensation | 0.006271 |
| 17 | Epigastria Pain | 0.006240 |
| 18 | Cough | 0.005984 |
| 19 | Sweaty | 0.005933 |
| 20 | AF(abdominal Flank) | 0.005789 |
| 21 | Flank pain | 0.005487 |
| 22 | Dysuria | 0.004345 |
| 23 | Belching | 0.001857 |

As, indicated in Table 2, the first 4 features namely weakness, chills, nausea and early fatigue in the table are the most important features to predict the target classes. Based on random forest feature importance analysis some of the features like flank pain, dysuria and belching rated less score and they are subjected to be dropped. . However, according to domain expert consultation, all features have their own shares while making foodborne disease diagnosis. So that, there is no need of dropping features. Thus, the researcher used all features for foodborne disease classification.

### 3.3 Proposed framework for foodborne disease classification

As shown in figure 1, It begins with the systematic collection and preprocessing of data, which includes cleaning, normalization, and feature selection to create a high-quality dataset. The framework utilizes a combination of machine learning algorithms, including Decision Tree, Random Forest, and eXtreme Gradient Boosting (XGBoost), alongside a stacking ensemble approach that harnesses the strengths of each model to improve classification performance. To optimize the models further, the researcher applied hyper parameter tuning through a grid search method, where various parameter settings were explored to identify the most effective combinations that yield the best performance metrics such as accuracy, precision, recall, and F1 score.
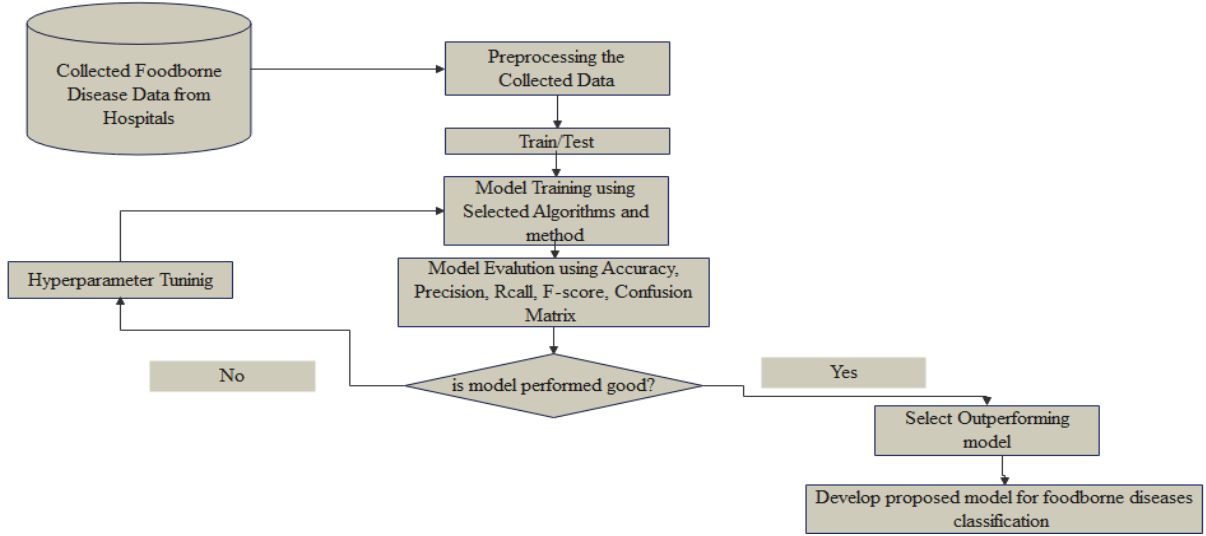
*Figure 1: Proposed research design of foodborne disease classification*

Designing a stack ensemble classifier involves creating a hierarchical model that leverages the strengths of multiple base learners to improve predictive accuracy in foodborne disease classification. In this research, the base learners selected for the ensemble included Decision Tree, Random Forest, and eXtreme Gradient Boosting (XGBoost), each offering unique advantages in handling complex data relationships and enhancing model interpretability. The ensemble classifier was structured such that each base learner was trained independently on the preprocessed dataset, generating initial predictions for the target disease categories. These predictions were then fed into a meta-learner, specifically a logistic regression model, which was trained to combine the outputs from the base learners effectively. This stacking approach allowed the meta-learner to learn how best to weight the predictions of the base models, thereby improving overall classification performance. Hyperparameter tuning was also applied to optimize both the base learners and the meta-learner, ensuring that the ensemble was fine-tuned for maximum accuracy. This robust design not only capitalizes on the individual strengths of the algorithms but also minimizes their weaknesses, resulting in a comprehensive and effective classification framework for foodborne diseases.

## 4.        Results and Discussions

The results of the study on foodborne disease classification demonstrated the superiority of the stacking ensemble model over individual machine learning algorithms, achieving remarkable performance metrics. Upon evaluation using a 70/30 train-test split, the stacking model attained a training accuracy of 99.40% and a test accuracy of 97.53%, while in the 80/20 split it recorded 99.35%

training accuracy and 98.06% test accuracy. Additionally, the stacking ensemble showed impressive precision, recall, and F1 scores across multiple classes, with precision values reaching up to 0.99 for amoebiasis and 0.97 for typhoid fever, indicating its effectiveness in correctly identifying disease cases. The hyperparameter tuning of the base learners Decision Tree, Random Forest, and XGBoost contributed significantly to these results, allowing for optimized performance that leveraged their complementary strengths. Furthermore, the research addressed the critical questions regarding the importance of specific attributes in classifying foodborne diseases and identified the optimal machine learning model, with the stacking ensemble emerging as the most suitable due to its enhanced accuracy, robustness, and generalizability in real-world applications, thereby filling a gap in the existing literature on foodborne disease classification methods.

*Table 3: Comparisons of all models in terms of precision, recall, and f1_score*

| Model | Train test split | Precision | Recall | F1_score |
|---|---|---|---|---|
| DT | 70/30 | 0.9667 | 0.9667 | 0.9666 |
| | 80/20 | 0.9647 | 0.9645 | 0.9645 |
| RF | 70/30 | 0.9712 | 0.9710 | 0.9708 |
| | 80/20 | 0.9761 | 0.9758 | 0.9757 |
| XGBOOST | 70/30 | 0.9724 | 0.9720 | 0.9719 |
| | 80/20 | 0.9696 | 0.9694 | 0.9693 |
| Stacking | 70/30 | 0.9755 | 0.9753 | 0.9752 |
| | 80/20 | 0.9808 | 0.9806 | 0.9806 |

The feature importance analysis presented in Section 3.2 indicated that attributes such as Weakness, Chills, Nausea, and Early fatigue were among the most predictive variables for foodborne disease classification. However, based on consultations with medical experts, it was decided not to exclude features with lower statistical importance, since every symptom has potential diagnostic value in real clinical settings. Consequently, all features were retained in the modeling process. The strong performance of ensemble methods, particularly Random Forest and Stacking, suggests that these models were able to automatically assign higher weights to the most informative features while still incorporating less dominant attributes. This approach ensured that the models remained clinically relevant and aligned with expert reasoning, while at the same time achieving robust and accurate predictions.
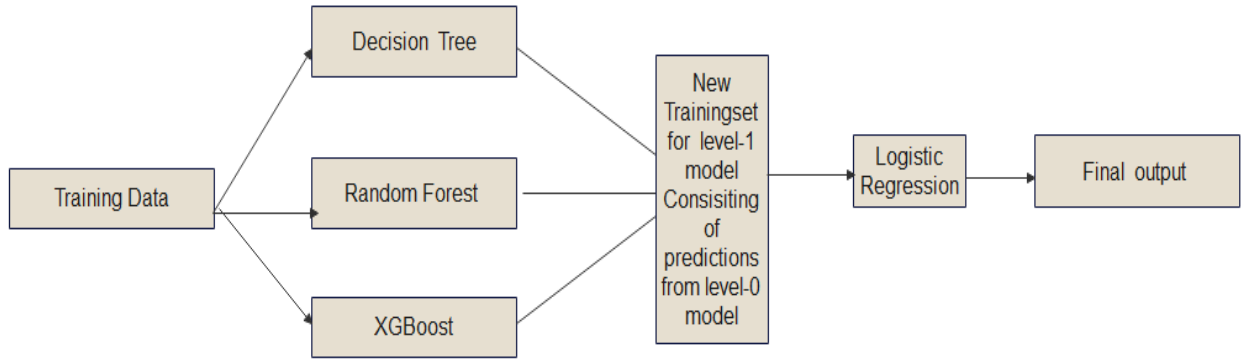
*Figure 2: Stacked ensemble learning classifier.*

Decision Tree has performed less with an accuracy 96.5%. This might be associated with its sensitivity to small variations. RF has performed better than DT with accuracy 97.5%. RF mitigates the weaknesses of a single DT by building multiple trees on random subsets of data and features. Stacking Ensemble has performed better than all with accuracy of around 98.1%. Stacking harnesses the power of the diverse learners by leveraging the individual outputs. The base models such as DT, RF and XGB, may capture different strengths or weaknesses. The meta-learner learns the optimal way to weigh and combine these predictions, effectively correcting individual errors and creating a more powerful, generalizable model.

## 5.      Conclusion and Recommendations

For the proposed study, the data was collected from two Hospitals, Gesuba primary hospital (GPH) and Wolaita Sodo university comprehensive specialized hospital (WSUCSH). The confirmed patient cases of foodborne diseases were collected based on purposive sampling and expert recommendations. A total of 3100 sample data was collected from both GPH and WSUCSH by recording the confirmed cases of foodborne diseases from 2012 to 2015 E.C. Experimental research design is applied in order to select a suitable machine learning model for the proposed study based on the experiments conducted. To prepare and analyze the collected data, different python libraries and modules such as Pandas, NumPy and Sklearn were used. After preprocessing the collected data, different state of art machine learning such as Decision Tree, Random Forest, XGBoost, and stacking ensemble learning model were applied to train a model and their performance was evaluated using various evaluation metrics such as confusion matrix, accuracy, precision, recall, and f1_score. To train a model, 70/30 and 80/20 splits were applied and the 80/20 split results shows better performance. In the 80/20 split, the Stacking ensemble model outperforms others with an accuracy of 98.06 (98.1%) followed by Random forest with an accuracy of 97.5%, XGBoost with 96.9%, and Decision Tree with 96.5%.

We planned to cover the foodborne diseases which are most prevalent in Ethiopia but the researcher couldn't cover all of them due to budget limits and lack of data on some foodborne diseases like Campylobacter Infections, rabies, anthrax, brucellosis, leptospirosis, echinococcosis, Shiga Toxin-Producing Escherichia coli (STEC). Thus, the researcher covered three foodborne diseases namely Amoebiasis, Giardiasis, and Typhoid Fever classification using machine learning and stacking ensemble learning model. Moreover, by retaining all features identified during data preprocessing, the study respected domain expert advice that each symptom contributes to diagnostic decision-making. The superior performance of ensemble methods, especially the stacking model, highlights their ability to integrate both highly ranked features and those of lower statistical importance, resulting in a balanced and clinically meaningful classifier. In the future, the researcher would like to recommend classification of foodborne diseases prevalent in Ethiopia by including those which are not covered here.

**Declaration of Conflicting Interests**

The author declare that they have no conflicts of interest.

**References**

Hossain, K., D. Raheem, and S. Cormier. (2018). Food Security: A Basic Need for Humans, in Food Security Governance in the Arctic-Barents Region, Cham: Springer International Publishing, 2018, pp. 5–14. doi: 10.1007/978-3-319-75756-8_2.

Mahon, C.R. (1998). Foodborne illness: is the public at risk? Clinical laboratory science, journal of the American Society for Medical Technology, vol. 11 5, pp. 291–297.

WHO World Health Organization. (2008). Foodborne disease outbreaks: Guidelines for investigation and control. 20 Avenue Appia, 1211 Geneva 27, Switzerland.

Gazu, L. (2023). Foodborne disease hazards and burden in Ethiopia: A systematic literature review, 1990–2019, in Frontiers in Sustainable Food Systems. Available: https://api.semanticscholar.org/CorpusID:256938077

Alanazi, A. (2022). Using machine learning for healthcare challenges and opportunities," *Informatics in Medicine Unlocked*, vol. 30, p. 100924. doi: https://doi.org/10.1016/j.imu.2022.100924.

S. Mohammad, S. (2012). Multi Disease Detection and Predictions Based On Machine Learning.

Teferi, S. (2020). A Review on Food hygiene Knowledge, Practice and Food Safety in Ethiopia, Food Science and Quality Management. Available: https://api.semanticscholar.org/CorpusID:229037211

Agarwal, A. (2021). The five Most Useful Techniques to Handle Imbalanced Datasets - KDnuggets,"

Brownlee, J. (2021). Stacking Ensemble Machine Learning With Python.

H. Wang, W. Cui, Y. Guo, Y. Du, and Y. Zhou. (2020). Prediction of Foodborne Diseases Pathogens: A Machine Learning Approach (Preprint)," *JMIR Medical Informatics*, vol. 9, doi: 10.2196/24924.

Oguntimilehin, A. and Adetunmbi, A.O. and Abiola, O.B. (2013). A Machine Learning Approach to Clinical Diagnosis of Typhoid Fever. Available: http://eprints.abuad.edu.ng/id/eprint/73

Radhika S. (2020). Symptoms based disease prediction using decision tree and electronic health record analysis, European Journal of Molecular \& Clinical Medicine.

Kavzoglu, T. and A. Teke. (2022). Predictive Performances of Ensemble Machine Learning Algorithms in Landslide Susceptibility Mapping Using Random Forest, Extreme Gradient Boosting (XGBoost) and Natural Gradient Boosting (NGBoost). Arabian Journal for Science and Engineering, vol. 47, no. 6, pp. 7367–7385. doi: 10.1007/s13369-022-06560-8.

Leo, E. Luhanga, K. Michael. (2019). Machine Learning Model for Imbalanced Cholera Dataset in Tanzania," The Scientific World Journal, 12. doi: https://doi.org/10.1155/2019/9397578.

Grampurohit, S., C. Sagarnal. (2020). Disease Prediction using Machine Learning Algorithms, p. pp 1-7.

Bhuiyan, A., S. S. Rad, F. T. Johora, A. Islam, M. I. Hossain, and A. A. Khan. (2023). Prediction of Typhoid Using Machine Learning and ANN Prior to Clinical Test. International Conference on Computer Communication and Informatics (ICCCI), p 1–7.doi: 10.1109/ICCCI56745.2023.10128226.

Kumar, K., I. Kaur, and S. Mishra. (2023). Foodborne Disease Symptoms, Diagnostics, and Predictions Using Artificial Intelligence-Based Learning Approaches. A Systematic Review," Archives of Computational Methods in Engineering. doi: 10.1007/s11831-023-09991-0.

Zach Bobbitt. (2022). Label Encoding vs. One Hot Encoding: What's the Difference? Accessed: Apr. 03, Available: https://www.statology.org/label-encoding-vs-one-hot-encoding/

Akash Dubey A. (2024). Feature Selection Using Random forest| by Akash Dubey/Towards Data Science. Available: https://towardsdatascience.com/feature-selection-using-random-forest-26d7b747597f.