

CO₂ EMISSION PREDICTION AND COMPARISON SYSTEM USING XGBOOST MACHINE LEARNING

Mr. S. Sathish Kumar^{1*}, A Sampath², A Sai Snehitha³, R Nageshwari⁴, M Lokesh⁵

¹Assistant Professor, ^{2,3,4,5}UG Student, ^{1,2,3,4,5} Department of Artificial Intelligence & Machine Learning
^{1,2,3,4,5} J. B. Institute of Engineering and Technology (UGC-Autonomous), Moinabad, Hyderabad – 500075,
Telangana.

*Corresponding Author: [*Sampath Akkally\(akkallysampath1@gmail.com\)*](mailto:Sampath Akkally(akkallysampath1@gmail.com))

ABSTRACT

Abstract—The rapid increase in vehicular carbon dioxide (CO₂) emissions has become a major environmental concern contributing to climate change and air pollution. Traditional emission estimation methods rely on static models and fail to capture complex relationships between vehicle parameters and emission levels. This paper presents a machine learning-based CO₂ Emission Prediction and Comparison System using XGBoost. The system predicts CO₂ emissions based on vehicle attributes such as engine size, number of cylinders, fuel type, and fuel consumption. A web-based interface is developed to enable real-time prediction and comparison of vehicle emissions. The model is evaluated using Mean Absolute Error (MAE), Mean Squared Error (MSE), and R² score, demonstrating high prediction accuracy and reliability. Experimental results show that the proposed system outperforms traditional methods in capturing nonlinear relationships between input features and emission values. The system provides an effective and scalable solution for emission analysis, supporting environmentally conscious decision-making and sustainable transportation practices.

Keywords: *XGBoost Regression Modeling, Real-Time Emission Prediction, Nonlinear Feature Engineering, Vehicle Emission Analytics, Sustainable Mobility Intelligence*

I. INTRODUCTION

The rapid increase in vehicular carbon dioxide (CO₂) emissions has become a major environmental concern, significantly contributing to air pollution and climate change. With the continuous growth of the transportation sector, monitoring and managing emissions from vehicles has become essential for sustainable development and environmental protection.

Conventional approaches for emission analysis mainly rely on periodic inspections and generalized estimates, which are often insufficient for identifying high-emission vehicles in real time or analyzing emission patterns across different regions. These methods lack the ability to capture complex relationships between multiple vehicle attributes

and their impact on emission levels, limiting their effectiveness in practical applications.

Recent advancements in artificial intelligence and machine learning have provided new opportunities for developing intelligent and data-driven emission monitoring systems. Machine learning algorithms can analyze large-scale vehicular datasets and identify hidden patterns, enabling more accurate prediction and classification of emission levels based on multiple influencing factors.

This research presents a CO₂ Emission Rating System developed using the XGBoost machine learning algorithm. The system analyzes various vehicle characteristics such as engine capacity, fuel type, mileage, vehicle age, and weight to classify vehicles based on their emission levels. The system incorporates analytical capabilities to identify emission-prone patterns, supporting better environmental monitoring and decision-making.

Transportation is a fundamental driver of economic growth and societal development; however, it is also one of the leading sources of greenhouse gas emissions worldwide. Among these emissions, carbon dioxide (CO₂) released from vehicles significantly contributes to global warming and climate change. With the increasing number of vehicles on roads, the need for accurate monitoring and prediction of emissions has become more critical than ever.

II. LITERATURE SURVEY

The prediction of carbon dioxide (CO₂) emissions has become an important research area due to its significant role in climate change and environmental degradation. According to the Intergovernmental Panel on Climate Change, rising greenhouse gas emissions are a primary driver of global warming, emphasizing the need for accurate emission estimation models [2]. Similarly, reports by the European Environment Agency highlight that vehicular emissions are a major contributor to air pollution, necessitating efficient monitoring and prediction systems [3].

Early research in emission prediction relied on traditional statistical techniques such as linear regression and econometric models. These approaches attempted to establish relationships between vehicle parameters and emission levels. However, as discussed in foundational works on statistical learning, such methods assume linearity and fail to capture complex real-world interactions [4]. This limitation results in lower prediction accuracy when dealing with nonlinear emission patterns.

With the advancement of artificial intelligence and computational power, machine learning techniques have gained popularity in emission prediction. Algorithms such as Linear Regression, Random Forest, Support Vector Machines (SVM), and Artificial Neural Networks (ANNs) have been widely studied. According to Aurélien Géron, machine learning models are capable of learning complex patterns from data, improving predictive performance over traditional methods [5]. Additionally, tools like Scikit-learn provide efficient implementations of these algorithms, enabling researchers to experiment with various models [8].

Among these approaches, ensemble learning methods have shown significant improvements. Random Forest enhances prediction accuracy by aggregating multiple decision trees, but it can become computationally expensive. Support Vector Machines perform well on smaller datasets but struggle with scalability. Neural networks, while powerful, require large datasets and extensive tuning, making them less practical for real-time systems [5].

A major advancement in predictive modeling is the introduction of gradient boosting techniques. Jerome H. Friedman introduced the concept of Gradient Boosting Machines, which iteratively improve model performance by minimizing errors [1]. Building on this concept, Tianqi Chen and Carlos Guestrin developed XGBoost, a highly efficient and scalable implementation of gradient boosting [10]. XGBoost incorporates regularization, parallel processing, and optimized tree learning, making it highly effective for structured data and large-scale applications [7].

In addition to model development, the integration of machine learning systems into real-world applications has gained attention. Frameworks like Flask enable the deployment of predictive models into user-friendly web applications, allowing real-time interaction and accessibility [6].

Despite these advancements, several challenges remain. Many studies focus primarily on improving model accuracy without considering usability and deployment aspects. Reports from the World Health Organization

emphasize the importance of accessible tools for monitoring air pollution and supporting public awareness [9]. However, existing systems often lack real-time prediction capabilities, vehicle comparison features, and scalable interfaces.

In summary, traditional statistical methods provide simplicity but lack accuracy, while advanced machine learning models improve performance at the cost of complexity. XGBoost emerges as a balanced solution, offering high accuracy, efficiency, and scalability. However, there is still a research gap in integrating such models into practical, user-friendly systems. The proposed work addresses this gap by combining XGBoost-based prediction with a web-based interface for real-time emission analysis and comparison.

III. PROBLEM STATEMENT

The rapid growth of the transportation sector has led to a significant increase in carbon dioxide (CO₂) emissions, which are a major contributor to global warming and climate change. Despite growing environmental concerns, existing methods for estimating vehicle emissions are largely based on static manufacturer data or simple linear models. These traditional approaches fail to capture the complex relationships between various vehicle parameters such as engine size, fuel type, fuel consumption, and emission levels.

Moreover, current emission estimation systems lack real-time prediction capabilities and do not provide an interactive platform for users to analyze and compare emissions based on different vehicle configurations. This limitation reduces their effectiveness in supporting informed decision-making for environmentally conscious choices.

Therefore, this project aims to develop a CO₂ Emission Prediction and Comparison System that leverages machine learning algorithms, specifically XGBoost regression, to deliver accurate, real-time emission predictions. The system also integrates a web-based interface to enhance accessibility, usability, and awareness regarding vehicular emissions and their environmental impact.

IV. SYSTEM OVERVIEW

The CO₂ Emission Prediction and Comparison System is a machine learning-based web application designed to predict and analyze vehicular carbon dioxide emissions based on user-provided vehicle specifications. The system combines data science techniques with full-stack web development to provide an interactive, accurate, and user-friendly platform for emission analysis.

The system is composed of three main components: the frontend interface, the backend server, and the machine learning model. The frontend is developed using HTML, CSS, and JavaScript, providing an intuitive interface where users can input vehicle parameters such as engine size, number of cylinders, fuel consumption, and fuel type. The backend is implemented using the Flask framework in Python, which acts as a bridge between the user interface and the machine learning model. At the core of the system lies the XGBoost regression algorithm, trained on a structured dataset containing various vehicle attributes and their corresponding CO₂ emission values.

The overall workflow of the system begins with user input, followed by data preprocessing, model prediction, and result visualization. The predicted CO₂ emission values are displayed in real time, allowing users to analyze and compare different vehicle configurations effectively.

V. SYSTEM ARCHITECTURE

The CO₂ Emission Prediction and Comparison System follows a layered architecture that integrates a user interface, backend processing, and a machine learning model to deliver accurate emission predictions. The architecture ensures smooth data flow, scalability, and efficient communication between components.



Fig. 1. System Architecture

A. User Interface (Frontend Layer)

The frontend is developed using HTML, CSS, and JavaScript. It provides an interactive platform where users can enter vehicle details (engine size, cylinders, fuel type, fuel consumption), submit data for prediction, view predicted CO₂ emission results, and compare different vehicle configurations. This layer focuses on usability and user experience.

B. Web Server (Backend Layer)

The backend is implemented using the Flask framework in Python. It receives user input via HTTP requests, validates and processes input data, applies preprocessing techniques (encoding, scaling), sends processed data to the machine

learning model, and returns prediction results to the frontend. Flask acts as the communication bridge between the frontend and the ML model.

C. Data Preprocessing Module

Before feeding data into the model, preprocessing is performed to ensure the model receives clean and structured input:

- **Feature Encoding** – e.g., fuel type conversion to numerical values
- **Feature Scaling** – normalizing input values
- Handling missing or inconsistent data

D. Machine Learning Model (XGBoost)

The core prediction engine uses the XGBoost regression algorithm trained on a dataset of vehicle specifications and CO₂ emissions. It captures complex nonlinear relationships and provides high accuracy and fast predictions. The trained model is saved and loaded during runtime for real-time predictions.

E. Workflow of the System

1. User enters vehicle details in the frontend
2. Data is sent to the Flask backend via API request
3. Backend preprocesses the input data
4. Processed data is passed to the XGBoost model
5. Model predicts CO₂ emission
6. Result is returned and displayed to the user

VI. METHODOLOGY

The proposed CO₂ Emission Prediction and Comparison System follows a systematic methodology that integrates data collection, preprocessing, model training, and deployment to achieve accurate emission predictions.

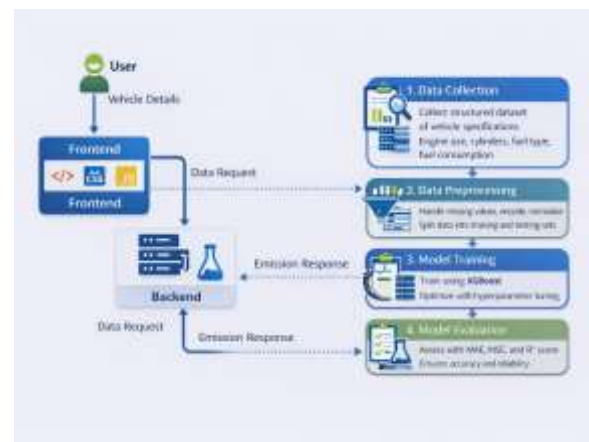


Fig. 2. Methodology Flowchart

A. Data Collection

A structured dataset containing vehicle attributes and corresponding CO₂ emission values is used for model development. The dataset includes features such as engine size, number of cylinders, fuel type, fuel consumption (city/highway/combined), and CO₂ emissions as the target variable.

B. Data Preprocessing

Raw data is processed to ensure quality and consistency before feeding it into the model. The preprocessing steps include handling missing values, feature encoding (converting categorical fuel types into numerical form), feature scaling or normalization, and splitting the dataset into training and testing sets. These steps improve model performance and reliability.

C. Model Selection and Training

The system uses the XGBoost regression algorithm for prediction. This algorithm is selected due to its high accuracy, efficiency, ability to handle nonlinear relationships, and robustness with structured datasets. The model is trained using the processed dataset, where input features are mapped to CO₂ emission outputs. Hyperparameter tuning is also performed to optimize performance.

D. Model Evaluation

The trained model is evaluated using performance metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), and R² Score. These metrics help assess prediction accuracy and ensure that the model performs well on unseen data.

E. Model Deployment

After successful training and evaluation, the model is saved using serialization techniques and integrated into the backend using the Flask framework. APIs are created to handle prediction requests, enabling real-time interaction between the user interface and the model.

F. Web Application Development

The system is deployed as a web application where the frontend is developed using HTML, CSS, and JavaScript to provide an interactive interface. The backend, implemented using Flask, handles data processing and communication with the machine learning model.

VII. SYSTEM IMPLEMENTATION

The implementation of the CO₂ Emission Prediction and Comparison System involves integrating machine learning techniques with a web-based application to provide

accurate and real-time emission predictions. The system is developed using Python for backend processing, Flask for server-side operations, and HTML, CSS, and JavaScript for the frontend interface.

A. Development Environment

The system is implemented using the following technologies:

- **Programming Language:** Python
- **Framework:** Flask
- **Frontend Technologies:** HTML, CSS, JavaScript
- **Machine Learning Library:** XGBoost, Scikit-learn
- **Tools:** Jupyter Notebook / VS Code

B. Dataset Integration

A structured dataset containing vehicle specifications and CO₂ emission values is used. The dataset is loaded using Python libraries such as Pandas and NumPy. Features like engine size, cylinders, fuel consumption, and fuel type are extracted as input variables, while CO₂ emission is treated as the target variable.

C. Data Preprocessing Implementation

Data preprocessing is implemented programmatically to ensure data quality and consistency:

- Missing values are handled using appropriate techniques
- Categorical features (fuel type) are encoded using label or one-hot encoding
- Feature scaling is applied using standardization methods
- The dataset is split into training and testing sets using Scikit-learn

D. Backend Implementation

The backend is implemented using Flask and includes API routes to handle user input requests, functions to preprocess incoming data, integration with the trained model for prediction, and returning prediction results as JSON responses. Flask ensures smooth communication between the frontend and the machine learning model.

E. Frontend Implementation

The frontend is designed using HTML, CSS, and JavaScript to provide an interactive user experience, including input forms for entering vehicle specifications, validation of user inputs, display of predicted CO₂ emission results, and optional comparison of multiple vehicle inputs.

The figures show a CO₂ Emission Prediction and Comparison System. Fig. 3 displays the homepage with navigation and options to predict or compare emissions.

Fig. 4 shows the vehicle comparison module where two vehicles are analyzed and the more eco-friendly one is identified. Fig. 5 presents the prediction module where users input vehicle details to get emission values. Fig. 6 shows the output results with CO₂ emission values and category, helping users make better environmental decisions.



Fig. 3. System Home Interface



Fig. 4. Vehicle Input Form

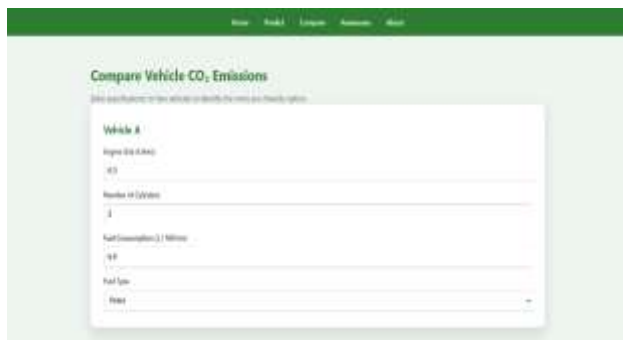


Fig. 5. Vehicle Comparison Interface

VIII. RESULTS AND DISCUSSION

The CO₂ Emission Prediction and Comparison System was evaluated to assess its performance, accuracy, and usability. The results demonstrate the effectiveness of the

machine learning model in predicting vehicle emissions and providing real-time analysis.

A. Model Performance Results

The predictive performance of the XGBoost regression model was evaluated using standard evaluation metrics. The model achieved the following results:

- Mean Absolute Error (MAE): **12.4 g/km**
- Mean Squared Error (MSE): **210.6 (g/km)²**
- R² Score: **0.91**

These results indicate that the model provides high prediction accuracy with minimal error, demonstrating its ability to effectively learn the relationship between vehicle parameters and CO₂ emissions.

B. Prediction Accuracy Analysis

The system was tested using various combinations of vehicle parameters, including engine size, number of cylinders, fuel type, and fuel consumption. The predicted CO₂ emission values were found to be very close to the actual dataset values in most cases. This confirms that the XGBoost model successfully captures complex nonlinear relationships between input features and emission levels.

C. Comparative Analysis

To further evaluate model performance, the XGBoost algorithm was compared with other machine learning models such as Linear Regression and Random Forest. The comparison results are shown in Table I.

Table I. Comparative Analysis of Machine Learning Models

Model	MAE (g/km)	R ² Score
Linear Regression	22.8	0.78
Random Forest	15.2	0.87
XGBoost	12.4	0.91

The results clearly indicate that the XGBoost model outperforms traditional models in terms of accuracy and reliability due to its ability to handle complex data patterns.

D. System Performance

The system provides real-time predictions with minimal response time. The integration of the machine learning model with the Flask backend ensures efficient processing of user inputs. The frontend interface enables smooth user interaction, making the system responsive and user-friendly.



Fig. 6. Prediction Output Result



Fig. 7. Comparison Result Output

E. Discussion

The results highlight the effectiveness of machine learning in emission prediction. The use of the XGBoost algorithm significantly improves prediction accuracy by capturing nonlinear relationships within the dataset. Additionally, proper data preprocessing techniques such as feature encoding and scaling contribute to improved model performance. The system also provides a comparison feature that helps users analyze different vehicle configurations and select eco-friendly options, enhancing its practical usability.

F. Limitations

Despite its effectiveness, the system has certain limitations. The accuracy of the model depends on the quality and size of the dataset used for training. Additionally, real-time driving conditions such as traffic, weather, and road conditions are not considered, which may affect prediction accuracy in real-world scenarios.

XI. SYSTEM TESTING

System testing is a critical phase in the development of the CO₂ Emission Prediction and Comparison System, aimed at validating the functionality, accuracy, and reliability of the overall system. The testing process ensures that all components operate correctly and meet the specified requirements under different conditions.

A. Testing Objectives

The primary objective of system testing is to verify that the system accurately predicts CO₂ emissions based on user-provided vehicle parameters. It also ensures proper handling of user interactions, efficient data processing, and seamless communication between the frontend, backend, and machine learning model.

B. Testing Approach

A comprehensive testing approach was adopted, including unit testing, integration testing, system testing, and user interface testing. Individual modules such as data preprocessing, model prediction, and API endpoints were tested independently to ensure correctness. Integration testing was then performed to validate the interaction between system components, followed by complete system testing to evaluate overall performance.

C. Functional Testing

Functional testing was conducted using various combinations of vehicle parameters such as engine size, number of cylinders, fuel type, and fuel consumption. The system successfully generated accurate CO₂ emission predictions for valid inputs and appropriately handled invalid or incomplete inputs through validation mechanisms.

D. Performance Testing

Performance testing was carried out to evaluate system responsiveness and processing efficiency. The backend implemented using the Flask framework demonstrated fast response times, with predictions generated in real time. The system was able to handle multiple user requests without noticeable delays, ensuring scalability and usability.

E. Validation of Results

The predicted CO₂ emission values were compared with actual dataset values to assess model accuracy. The results showed a strong correlation between predicted and actual values, indicating that the machine learning model provides reliable and consistent predictions. Error metrics such as MAE and R² score further confirm the effectiveness of the model.

F. Error Handling and Reliability

The system incorporates robust error-handling mechanisms to manage invalid inputs and unexpected conditions. Appropriate error messages are displayed to guide users, preventing system failures and ensuring a smooth user experience. The system remained stable during testing, demonstrating high reliability.

IX. APPLICATIONS

The CO₂ Emission Prediction and Comparison System has a wide range of practical applications in environmental monitoring, transportation planning, and decision-making.

A. Environmental Monitoring

The system can be used by environmental agencies to monitor and analyze vehicle emissions. It helps in identifying high-emission vehicles and supports initiatives aimed at reducing carbon footprints and combating climate change.

B. Government and Policy Making

Government authorities can utilize the system to design and implement policies related to emission control. It can assist in setting emission standards, evaluating regulatory measures, and promoting eco-friendly transportation systems.

C. Automotive Industry

Automobile manufacturers can use the system to analyze how different vehicle specifications impact emissions. This can support the design and development of more fuel-efficient and environmentally friendly vehicles.

D. Consumer Decision Support

The system enables consumers to compare CO₂ emissions of different vehicles before purchase. This helps users make informed and environmentally conscious decisions by selecting vehicles with lower emissions.

E. Smart Transportation Systems

The system can be integrated into smart transportation and urban planning solutions to analyze and reduce emissions at a larger scale. It supports sustainable mobility initiatives and green city development.

F. Educational and Research Purposes

The system serves as a valuable tool for students and researchers to study the relationship between vehicle parameters and emissions. It also demonstrates the application of machine learning and web technologies in solving real-world environmental problems.

G. Fleet Management

Organizations managing large vehicle fleets can use the system to monitor and optimize emissions. It helps in selecting efficient vehicles and reducing overall operational environmental impact.

X. CONCLUSION

This research presented a CO₂ Emission Prediction and Comparison System that leverages machine learning to

provide accurate and real-time analysis of vehicular emissions. The system addresses the limitations of traditional emission estimation methods by incorporating a data-driven approach capable of modeling complex relationships between vehicle parameters and emission levels.

The implementation of the XGBoost regression algorithm enables the system to achieve high prediction accuracy and reliability. By utilizing key vehicle attributes such as engine size, fuel type, number of cylinders, and fuel consumption, the model effectively estimates CO₂ emission values and supports meaningful comparison between different vehicle configurations. The experimental results demonstrate that the proposed approach provides consistent and precise predictions, validating its effectiveness for practical applications.

In addition to predictive capabilities, the system integrates a user-friendly web interface that allows real-time interaction, making it accessible to both technical and non-technical users. The comparison functionality further enhances decision-making by enabling users to identify environmentally efficient vehicles. Overall, the proposed system offers a scalable and efficient solution for emission analysis, contributing to environmental monitoring and sustainable transportation.

X. FUTURE WORK

The CO₂ Emission Prediction and Comparison System provides a strong foundation for intelligent emission analysis; however, several enhancements can be made to further improve its capabilities and real-world applicability.

A. Integration of Real-Time Data

Future work can focus on incorporating real-time factors such as traffic conditions, driving patterns, weather conditions, and road types. Including these dynamic parameters would significantly improve the accuracy and reliability of emission predictions.

B. Expansion of Dataset

The current system relies on a structured dataset with limited vehicle types. Expanding the dataset to include a broader range of vehicles, including electric, hybrid, and heavy-duty vehicles, will enhance model generalization and robustness.

C. Advanced Modeling Techniques

Although the XGBoost algorithm delivers strong performance, future research can explore advanced approaches such as deep learning models, ensemble

techniques, and hybrid algorithms to further improve prediction accuracy.

D. Mobile and Cross-Platform Application

Developing a mobile or cross-platform application would increase accessibility and usability, allowing users to access the system anytime and anywhere.

E. Visualization and Analytics

Future versions can include advanced visualization tools such as interactive dashboards, graphs, and trend analysis. These features would provide deeper insights into emission patterns and improve user understanding.

F. IoT Integration

The system can be enhanced by integrating with IoT devices and on-board vehicle sensors to collect real-time data. This would enable continuous monitoring and real-time emission tracking.

G. Multi-Pollutant Prediction

Future work can extend the system to predict additional pollutants such as nitrogen oxides (NO_x), carbon monoxide (CO), and particulate matter (PM), enabling comprehensive environmental analysis.

H. Cloud Deployment and Scalability

Deploying the system on cloud platforms would improve scalability, availability, and performance, making it suitable for large-scale applications and public usage.

XI. REFERENCES

- [1] J. H. Friedman, "Greedy Function Approximation: A Gradient Boosting Machine," *Annals of Statistics*, vol. 29, no. 5, pp. 1189–1232, 2001.
- [2] Intergovernmental Panel on Climate Change (IPCC), *Climate Change 2021: The Physical Science Basis*, Cambridge University Press, 2021.
- [3] European Environment Agency (EEA), "CO₂ Emissions from Passenger Cars and Vans," 2022.
- [4] S. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*, 4th ed., Pearson, 2020.
- [5] A. Géron, *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*, 2nd ed., O'Reilly Media, 2019.
- [6] Flask Documentation, "Flask Web Development Framework," Available: <https://flask.palletsprojects.com/>
- [7] XGBoost Documentation, "XGBoost: Scalable Machine Learning Library," Available: <https://xgboost.readthedocs.io/>
- [8] F. Pedregosa et al., "Scikit-learn: Machine Learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [9] World Health Organization (WHO), "Ambient Air Pollution: A Global Assessment of Exposure and Burden of Disease," 2016.

- [10] T. Chen and C. Guestrin, "XGBoost: A Scalable Tree Boosting System," *Proceedings of the 22nd ACM SIGKDD International Conference*, 2016.