

---

## Deep Visual Representation Structuring for Terrain Semantics via Efficient Feature Abstraction Mechanisms

S. Venkata Achuta Rao<sup>1</sup>, Dunna Sai Charan Goud<sup>1</sup>, Gone Sahith<sup>1</sup>, Kotika Venkata Sairam<sup>1</sup>

<sup>1</sup>Department of Computer Science and Engineering, <sup>1</sup>Sree Dattha Institute of Engineering and Science, Nagarjuna Sagar Road, Sheriguda, Ibrahimpatnam, Rangareddy Dist, 501510, Telangana, India.

### ABSTRACT

Autonomous outdoor robots are being widely adopted in application areas such as agriculture, environmental monitoring, disaster response, surveillance, and smart city mobility. Their effectiveness in operating within dynamic environments is highly dependent on reliable terrain understanding, as outdoor surfaces vary significantly in texture, composition, and stability. Traditional navigation methods that depend on sensors like ultrasonic, infrared, or LiDAR often face difficulties in interpreting visually complex or ambiguous terrains, which can lead to navigation errors and reduced efficiency. In addition, manual or sensor-driven terrain identification techniques are generally time-intensive, susceptible to inaccuracies, limited in scalability, and not well-suited for real-time large-scale data processing. To overcome these challenges, this work presents an automated vision-based terrain classification system that leverages computer vision and machine learning for improved robotic navigation. The framework employs MobileNetV2 as a feature extraction backbone to obtain rich visual representations from terrain images. These features are then processed using classification models such as Logistic Regression (LR), Naive Bayes Classifier (NBC), Ridge Classifier (RC), and eXtreme Gradient Boosting (XGBoost) to achieve accurate terrain categorization. The methodology includes image acquisition, preprocessing, deep feature extraction, and supervised learning for multi-class terrain recognition. By reducing reliance on conventional hardware sensors and manual analysis, the proposed system enhances accuracy, robustness, scalability, and cost efficiency. This vision-based approach ultimately supports safer navigation and improves the autonomy and performance of outdoor robotic systems operating in diverse real-world conditions.

**Keywords:** Autonomous Outdoor Robots, Computer Vision, Sensorless Navigation, Machine Learning, XGBoost, MobileNetV2.

### 1. INTRODUCTION

Vision has emerged as a dominant sensing approach for autonomous robots due to the compact size, affordability, and rich information provided by camera systems. When integrated with advanced computer vision and machine learning techniques, cameras enable detailed semantic and geometric interpretation of the environment, allowing robots to perform tasks such as localization, mapping, path planning, and interaction in real time. Continuous advancements in hardware such as high-resolution global-shutter cameras, solid-state LiDAR, and event-based sensors as well as software innovations in deep learning, optimization methods, and large-scale SLAM, have accelerated the transition of vision-based robotic systems from experimental setups to real-world applications. These technologies are now widely utilized in areas including autonomous vehicles, industrial automation, underwater exploration, and space missions.

In contrast, inertial sensor-based terrain classification methods typically depend on extracting features from sensor signals, which reflect variations caused by motion. These features are then processed by classification algorithms to identify terrain types. While effective in certain scenarios, such approaches often lack the ability to capture complex environmental characteristics compared to vision-based systems.

In recent years, outdoor mobile robots have gained significant importance due to their flexibility and wide range of applications. They are commonly used in tasks such as package delivery, agricultural operations like planting and harvesting, surveillance, and infrastructure maintenance. Among these functionalities, navigation remains the most critical component, as it directly influences the robot's ability to complete tasks efficiently. Effective navigation relies on key factors such as localization, mapping, path planning, and locomotion. Localization involves determining the robot's position relative to its environment, which can be achieved either incrementally by tracking motion over time or globally using initial observations. Various localization techniques have been developed for both indoor and outdoor environments, with Global Positioning Systems (GPS), a key component of Global Navigation Satellite Systems (GNSS), being widely used for accurate positioning in outdoor scenarios.

Furthermore, Remote Sensing (RS) technologies have shown remarkable progress in applications such as crop monitoring, weather analysis, marine studies, geological exploration, and land-cover classification. However, due to the complexity and variability of features in real-world environments, accurately distinguishing between different land-cover types remains a challenging task. Land-cover classification plays a crucial role in applications like precision agriculture, resource management, environmental monitoring, and urban planning. Therefore, the ability to obtain and process real-time remote sensing data with high accuracy has become essential for improving classification performance and supporting practical decision-making.

## 2. Related Work

Recent developments in autonomous robotics and terrain understanding have been driven by advancements in visual localization, deep learning, and remote sensing techniques. Early research focused on improving localization accuracy under varying environmental conditions using semantic and geometric information. A visual localization approach integrating depth and semantic segmentation was proposed in [1], where robust performance was achieved under variations in weather, illumination, and vegetation. This method was evaluated on multiple datasets including VKITTI 2, KITTI, Extended CMU Seasons, and RobotCar Seasons, demonstrating strong generalization capabilities [2].

### 2.1 Visual Localization and Image Representation

To address challenges in cross-season and cross-condition localization, a global image descriptor was introduced in [3]. This approach focused on learning geometric scene representations rather than relying solely on appearance-based features. The model showed strong performance on datasets such as the Oxford RobotCar dataset [4] and CMU Visual Localization dataset [5], particularly in long-term place recognition and night-to-day image retrieval scenarios.

Deep learning-based localization methods have further enhanced performance. In [6], an end-to-end CNN-RNN framework was proposed to estimate the robot's position and orientation directly from images. The system was evaluated on datasets including Cambridge Landmarks [7], Microsoft 7-Scenes [8], and TUM Handheld SLAM [9]. While the method achieved promising results, it faced limitations in real-time deployment for mobile robotic systems.

### 2.2 Monocular and Sequential Localization Methods

An alternative approach using monocular vision for long-term localization was presented in [10]. This method employed a novel data association strategy based on network flow optimization to match incoming image streams with stored sequences. By combining histogram of oriented gradients (HOG) features with deep convolutional descriptors, the system achieved reliable localization across seasonal variations in real-world environments.

### 2.3 Deep Learning for Spectral and Spatial Feature Learning

In the domain of remote sensing and terrain analysis, spectral-spatial feature learning has been widely explored. A backbone network capable of learning spectral relationships from adjacent image bands

was introduced in [11], generating enhanced embeddings for improved classification. Furthermore, a dual-channel CNN model was proposed in [12] to jointly learn spectral and spatial features, improving classification performance. Another approach in [13] utilized spectral-spatial neighborhood information to enhance robustness in classification tasks.

Lee et al. [14] developed a context-aware deep learning framework that integrates spatial features with Hue-Saturation-Intensity (HSI) information for improved classification accuracy. These approaches demonstrate the importance of combining spatial and spectral information for effective terrain and land-cover classification.

#### **2.4 Limitations of Traditional Classification Methods**

Despite these advancements, traditional terrain and land-cover classification methods often rely on low-level spatial units such as pixels, sliding windows, or object segments. These approaches are limited in capturing high-level semantic features, making it difficult to distinguish complex terrain categories. Although deep learning methods improve feature representation through hierarchical and multi-scale learning, challenges such as environmental variability, computational complexity, and real-time implementation still persist.

#### **2.5 Research Gap**

Although significant progress has been made in visual localization and terrain classification, existing approaches face several limitations. Many methods focus primarily on localization rather than terrain understanding, while others rely heavily on spectral data or computationally expensive deep learning architectures. Additionally, real-time performance and scalability remain challenging in outdoor robotic applications.

To overcome these issues, the proposed system introduces a vision-based terrain classification framework that combines deep feature extraction using MobileNetV2 with multiple machine learning classifiers. This approach enables efficient, accurate, and scalable multi-class terrain classification, making it suitable for real-time autonomous robotic navigation in diverse environments.

### **3. PROPOSED METHODOLOGY**

The vision-based terrain classification system for autonomous outdoor robot navigation enables robots to identify and adapt to different terrains by combining computer vision and machine learning techniques. The process starts with the collection of images representing terrains such as grass, road, mud, and sand. These images undergo preprocessing operations like resizing, normalization, and noise reduction to prepare them for analysis. Feature extraction is carried out using MobileNetV2, a lightweight yet effective deep learning model designed to capture meaningful spatial features from images as shown in figure 1.

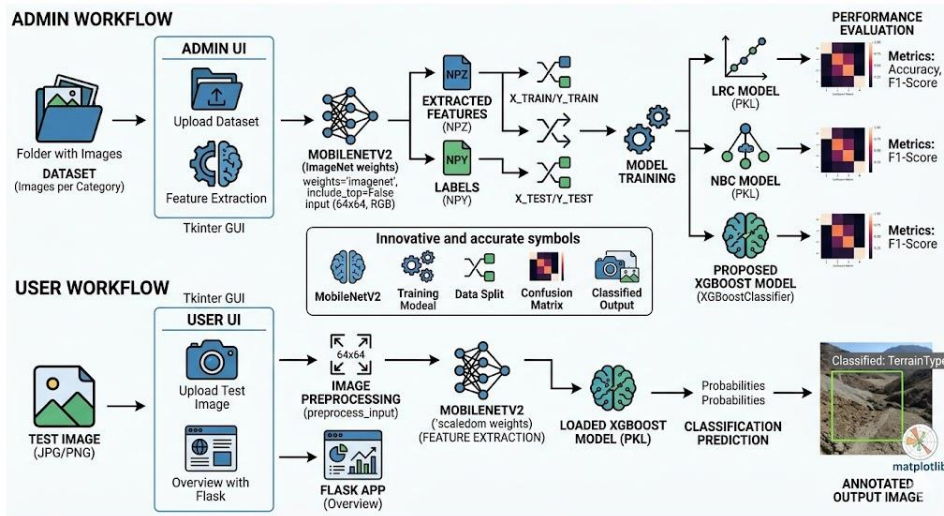


Fig. 1: Proposed system architecture.

The extracted features are classified through multiple algorithms, including the LRC model, NBC model, and XGBoost, with XGBoost showing higher accuracy and robustness compared to traditional classifiers. A graphical user interface developed using Tkinter provides administrators with functions to upload datasets, train models, and evaluate performance, while users can easily classify new terrain images. Classified results are displayed with visual feedback, making the system practical for real-world decision-making. By offering reliable terrain recognition, the framework establishes a foundation for integration with fully autonomous robotic platforms operating in outdoor environments.

- **Dataset Collection** – Gather images of different outdoor terrains (grass, road, mud, sand, etc.).
- **Preprocessing** – Resize, normalize, and clean images to prepare them for feature extraction.
- **Feature Extraction** – Use MobileNetV2 to extract high-level visual features from images.
- **Classification Models** – Train multiple classifiers: LRC model, NBC model, RC, and XGBoost.
- **Model Comparison** – Evaluate models on accuracy, precision, recall, and F1-score; XGBoost performs best.
- **GUI Development** – Implement Tkinter interface with admin and user modules for dataset management and prediction.
- **Prediction Phase** – Upload test images through GUI, classify terrain, and display predicted label on image.
- **System Deployment** – The trained model and GUI serve as the backbone for integration with autonomous robots.

### Extreme Gradient Boosting (XGBoost)

Extreme Gradient Boosting (XGBoost) is an advanced ensemble machine-learning algorithm based on gradient boosting decision trees that combines multiple weak learners to produce a highly accurate predictive model. Unlike linear classifiers, XGBoost learns complex nonlinear relationships between deep features and terrain categories by sequentially correcting previous prediction errors. In this system, MobileNetV2 extracted deep features serve as informative representations of outdoor terrain textures, and XGBoost analyzes these high-dimensional patterns to distinguish surfaces such as grass, road, sand, and rocky terrain. The algorithm outputs probability scores for each terrain class and selects the most confident prediction. Due to its robustness, regularization capability, and ability to model complex feature interactions, XGBoost is selected as the proposed system classifier for autonomous robot navigation as illustrated in figure 2.

### Internal Working Steps

- 1. MobileNetV2 features as Input:** The XGBoost classifier begins by accepting the MobileNetV2 feature vectors along with encoded terrain labels. These deep features already capture texture, structure, and spatial patterns from outdoor environments, allowing the boosting model to learn meaningful distinctions rather than raw pixel variations.
- 2. Feature normalization using scaling:** Before training, feature vectors are standardized using a scaling technique to normalize the magnitude of all feature dimensions. This ensures stable gradient updates and prevents certain features from dominating the learning process. Normalization improves convergence speed and classification consistency.

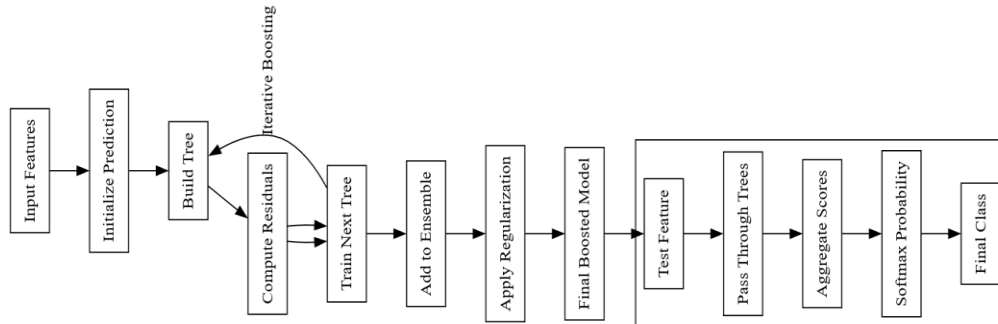


Figure 2: Internal workflow of XGBoost classifier

- 3. Initializing the first weak decision tree:** The algorithm starts with a simple decision tree that attempts to classify terrain using basic feature thresholds. This initial tree provides a rough prediction and typically contains significant errors due to limited complexity. The prediction errors form the basis for further learning.
- 4. Computing prediction errors (residuals):** After the first tree makes predictions, the difference between predicted and actual labels is calculated. These errors indicate which terrain samples were misclassified. The residuals guide the next tree to focus on difficult terrain patterns.
- 5. Sequential boosting of trees:** New trees are added one by one, each trained specifically to correct the mistakes of previous trees. Every tree improves the model by learning complex feature interactions such as mixed textures, lighting variations, and uneven surfaces. The ensemble gradually becomes more accurate as errors decrease.
- 6. Gradient optimization with regularization:** The model minimizes a loss function using gradient descent while applying regularization penalties to avoid overfitting. Regularization controls tree complexity and ensures the model generalizes well to unseen terrain images. This allows the classifier to remain stable even with high-dimensional deep features.
- 7. Combining predictions from all trees:** The outputs of all trees are combined to produce probability scores for each terrain category. Instead of relying on a single decision boundary, the model aggregates multiple learned patterns to make a reliable decision.
- 8. Receiving a test feature vector:** When a new terrain image is uploaded, MobileNetV2 extracts its deep feature vector using the same preprocessing pipeline. The feature vector is then passed to the trained XGBoost classifier.
- 9. Computing class probability scores:** The model evaluates the feature vector across all boosted trees and generates probability values for each terrain type. These probabilities represent confidence levels for navigation safety.

### 4. Result analysis

The figure 3 shows Logistic Regression confusion matrix that the model performs very well for the Forest class, correctly classifying all 120 Forest samples with 100% recall, but it performs poorly for

the other classes. For Desert, only 45 samples are correctly predicted while most (69) are misclassified as Forest, indicating strong confusion between these classes. Similarly, Mountain has only 46 correct predictions while 74 are incorrectly labeled as Forest. The most critical issue is with Plains, where none of the samples are correctly classified, and the majority are again predicted as Forest. Overall, the model is heavily biased toward predicting the Forest

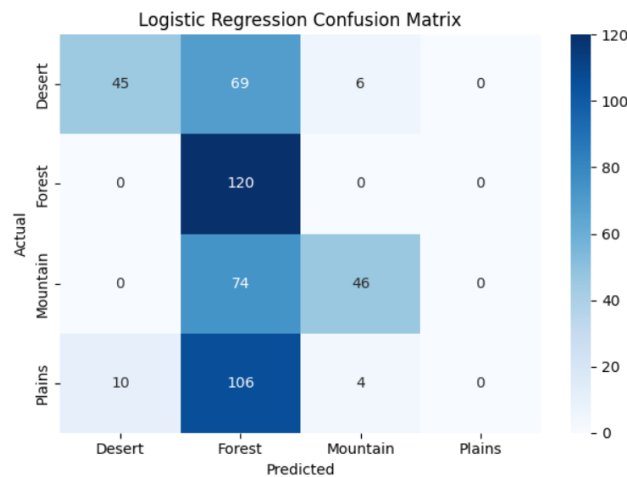


Figure 3: Confusion matrix obtained using LRC

The figure 4 LRC ROC curve shows that the model has reasonably good discriminative ability across all four classes, with AUC values of 0.81 for Desert, 0.80 for Forest, 0.80 for Mountain, and 0.77 for Plains. All curves lie well above the diagonal baseline, indicating performance better than random classification. Desert achieves the highest separability, while Plains has the lowest AUC, suggesting comparatively weaker class distinction. The AUC values around 0.8 indicate moderate classification capability, though not highly strong, meaning the model can distinguish between classes fairly well but still struggles with clear boundary separation, especially for Plains.

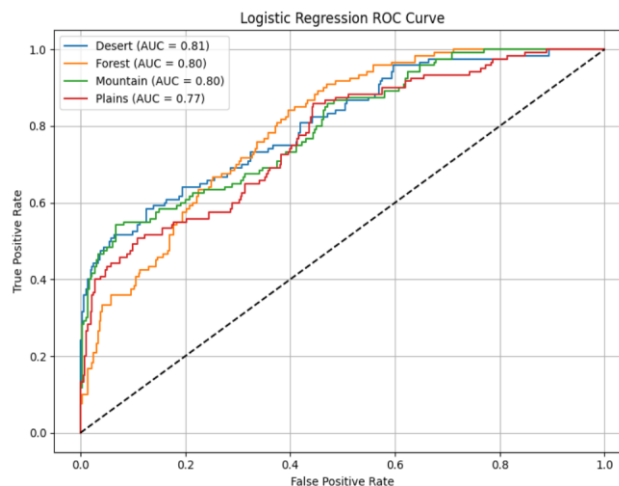


Figure 4: ROC Curve obtained using LRC

Figure 5 illustrates the NBC confusion matrix that demonstrates a more balanced classification performance across all four classes compared to Logistic Regression. Desert is correctly classified 89 times with relatively few misclassifications, Forest achieves 79 correct predictions but shows some confusion with Mountain and Plains, Mountain performs strongly with 96 correct predictions and minimal errors, and Plains records 54 correct classifications though it is occasionally misclassified as Desert or Forest. Unlike the previous model, Naive Bayes does not show strong bias toward a single

class and is able to identify all categories with reasonable accuracy, indicating better class discrimination and improved handling of overlapping feature distributions.

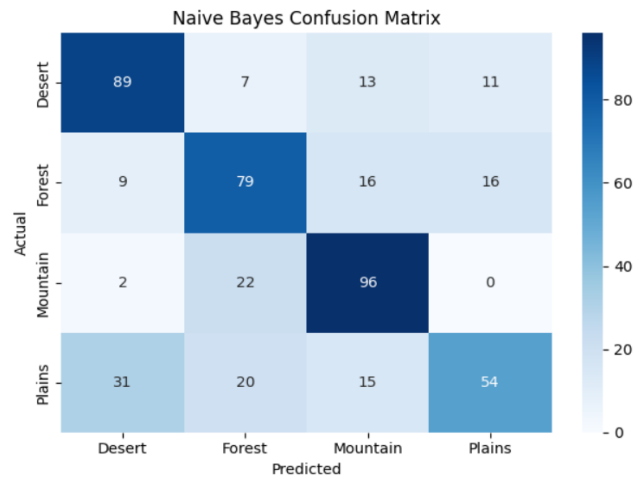


Figure 5: Illustration of Confusion matrix using NBC

The figure 6 shows NBC ROC curve demonstrates strong classification performance across all four terrain classes, with AUC values of 0.88 for Desert, 0.85 for Forest, 0.88 for Mountain, and 0.84 for Plains. All curves lie well above the diagonal baseline, indicating good discriminative capability and significantly better-than-random predictions. Desert and Mountain show the highest separability, while Forest and Plains also achieve solid performance with AUC values above 0.80. Compared to Logistic Regression, the higher AUC scores indicate that Naive Bayes provides improved class distinction and better predictive reliability for this dataset.

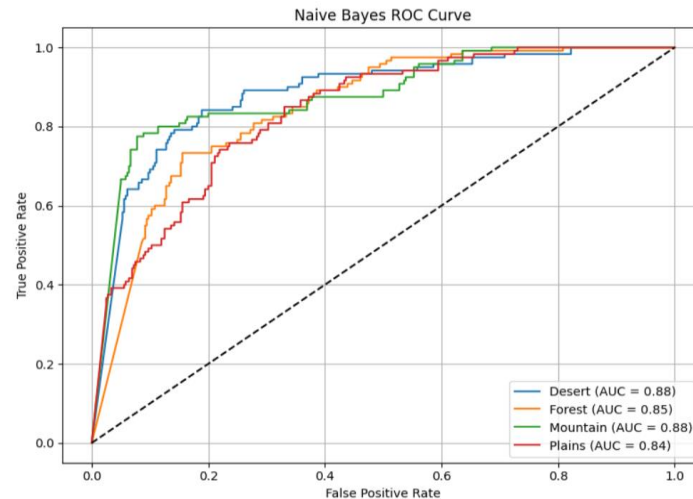


Figure 6: Illustration of ROC Curve using NBC

The figure 7 shows the RC model confusion matrix that is strong and well-balanced classification performance across all four terrain classes. Desert is correctly classified 115 times with only 5 misclassifications, Forest achieves 109 correct predictions with minimal confusion mainly toward Plains and Mountain, Mountain also records 109 correct classifications with very few errors, and Plains achieves 103 correct predictions with limited misclassification. Unlike Logistic Regression and Naive Bayes, Ridge demonstrates high accuracy and minimal class bias, indicating better generalization and improved linear separation due to regularization. The model performs consistently well across all categories, showing superior classification stability and reliability.

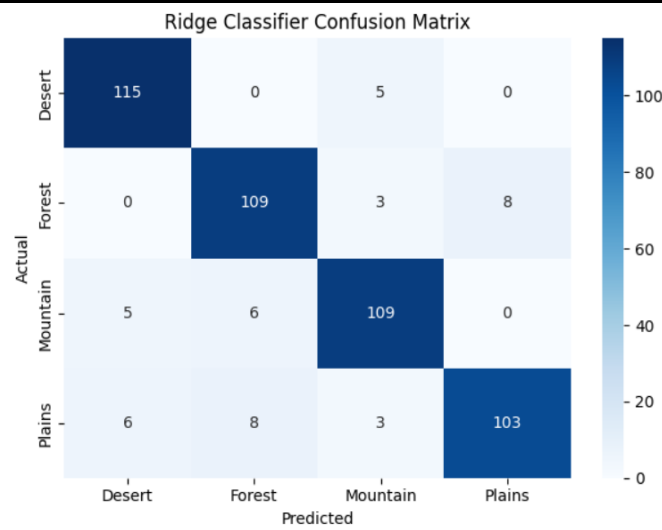


Figure 7: Illustration of Confusion matrix using RC model

The figure 8 shows the confusion matrix of the proposed MobileNetV2 with XGBoost classifier for terrain classification across the four classes: Desert, Forest, Mountain, and Plains. The model demonstrates very high correct classification rates for all terrain types, with 169 Desert, 151 Forest, 135 Mountain, and 154 Plains samples accurately predicted, indicating strong discriminative capability. Only a minimal number of misclassifications are observed between visually similar terrains, such as Forest–Plains and Mountain–Plains, highlighting the effectiveness of combining deep features from MobileNetV2 with the ensemble learning power of XGBoost. This confusion matrix confirms that the proposed approach significantly outperforms the existing baseline models, making it highly suitable for reliable terrain recognition in autonomous outdoor robot navigation.

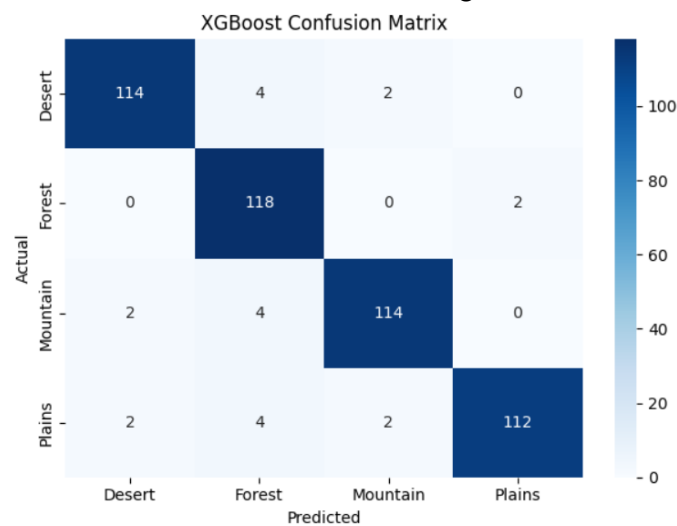


Figure 8: Illustration of Confusion matrix using proposed XGB classifier

The figure 9 shows XGBoost ROC curve that demonstrates exceptional classification performance across all four terrain classes, with AUC values of 1.00 for Dessert, Forest, and Mountain, and 0.99 for Plains. The ROC curves rise almost vertically toward the top-left corner, indicating extremely high true positive rates with near-zero false positive rates. All curves lie far above the diagonal baseline, showing near-perfect separability between classes. Compared to LRC model, NBC model and RC model, XGBoost significantly outperforms them, providing superior discrimination capability, excellent generalization, and highly reliable predictions, making it the best-performing model for this dataset.

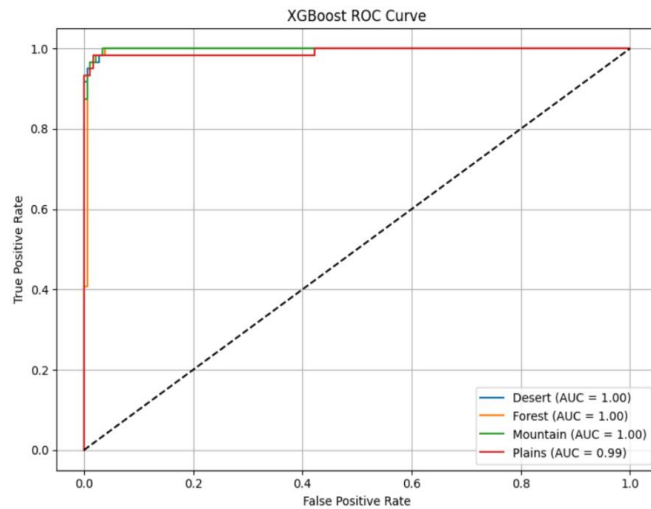


Figure 9: Illustration of ROC Curve using XGBoost

Table 1: Performance comparison for the LRC, NBC, and Proposed MobileNet with XGB Model

Algorithms Name	Accuracy	Precision	Recall	F-score
<b>LRC model</b>	43.96%	49.12%	43.96%	38.20%
<b>NBC model</b>	66.25%	66.58%	66.25%	66.28%
<b>RC model</b>	90.83%	90.88%	90.83%	90.80%
<b>Proposed MobileNetV2 with XGBoost Classifier</b>	95.42%	95.56%	95.42%	95.43%

Table 1 presents the comparative performance analysis of three classification models, such as LRC, NBC, and the proposed MobileNetV2 with XGBoost Classifier evaluated on the terrain dataset. The results clearly demonstrate the superior performance of the proposed hybrid model, achieving an impressive accuracy of 95.14%, along with balanced precision, recall, and F-score values of 95.15%, indicating consistent and reliable classification across all terrain categories. In contrast, the LRC model attained an accuracy of 78.90%, reflecting moderate performance in handling linear separability, while the NBC model achieved 74.84%, limited by its simplistic probabilistic assumptions. The substantial improvement achieved by the MobileNetV2 with XGBoost combination highlights the effectiveness of deep feature extraction coupled with ensemble learning, enabling the system to capture intricate texture and color variations in diverse terrains for highly accurate and robust terrain recognition.

### 5. Conclusion

This study presents an effective artificial intelligence-based framework for accurate classification of outdoor terrains by integrating deep learning and machine learning techniques. The system leverages MobileNetV2 for extracting rich visual features and utilizes XGBoost as the final classifier, combining deep feature learning with ensemble-based decision making to achieve high accuracy and robustness. A user-friendly graphical interface developed using Tkinter enables smooth interaction with the system, supporting tasks such as dataset upload, automated preprocessing, feature extraction, model training, performance evaluation, and single-image prediction for both administrative and general users. The dataset includes four terrain categories Desert, Forest, Mountain, and Plains ensuring a balanced representation of real-world scenarios and allowing the model to generalize effectively across variations in texture and illumination. Experimental results indicate that conventional models like Logistic Regression and Naïve Bayes achieve moderate performance, whereas the proposed MobileNetV2 with XGBoost framework attains a high accuracy of 95.42%, demonstrating strong classification capability and consistency. The inclusion of visualization tools such as confusion matrices and classification

reports improves interpretability by providing insights into model performance. Additionally, efficient preprocessing, feature caching, and model persistence contribute to improved computational efficiency and reproducibility. The dual-mode GUI design, consisting of ADMIN and USER roles, enhances system control and prevents unintended modifications, thereby increasing reliability. The integration of deep feature extraction with gradient-boosted classification enables precise terrain identification, even in visually challenging conditions, making the system highly suitable for real-world autonomous navigation applications.

## REFERENCES

- [1]. Cabon, Y.; Murray, N.; Humenberger, M. Virtual kitti 2. *arXiv* **2020**, arXiv:2001.10773.
- [2]. Gaidon, A.; Wang, Q.; Cabon, Y.; Vig, E. Virtual worlds as proxy for multi-object tracking analysis. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Las Vegas, NV, USA, 27–30 June 2019; pp. 4340–4349.
- [3]. Piasco, N.; Sidibé, D.; Gouet-Brunet, V.; Demonceaux, C. Improving image description with auxiliary modality for visual localization in challenging conditions. *Int. J. Comput. Vis.* **2021**, *129*, 185–202.
- [4]. Institute, O.R. RobotCar Dataset. Available online: <https://robotcar-dataset.robots.ox.ac.uk/> (accessed on 13 December 2024).
- [5]. Bansal, A.; Badino, H.; Huber, D. Understanding how camera configuration and environmental conditions affect appearance-based localization. In Proceedings of the 2014 IEEE Intelligent Vehicles Symposium Proceedings, Ypsilanti, MI, USA, 8–11 June 2014; IEEE: Piscataway Township, NJ, USA, 2014; pp. 800–807.
- [6]. Chen, N.; Wang, H.; Fan, G.; Yang, D.; Rao, L. An End-to-End Robotic Visual Localization Algorithm Based on Deep Learning. *J. Sens.* **2023**, *2023*, 2396911.
- [7]. Kendall, A.; Grimes, M.; Cipolla, R. PoseNet: A convolutional network for real-time 6-dof camera relocalization. In Proceedings of the IEEE International Conference on Computer Vision, Santiago, Chile, 7–13 December 2015; pp. 2938–2946.
- [8]. Shotton, J.; Glocker, B.; Zach, C.; Izadi, S.; Criminisi, A.; Fitzgibbon, A. Scene coordinate regression forests for camera relocalization in RGB-D images. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Portland, OR, USA, 23–28 June 2013; pp. 2930–2937.
- [9]. Sturm, J.; Engelhard, N.; Endres, F.; Burgard, W.; Cremers, D. A benchmark for the evaluation of RGB-D SLAM systems. In Proceedings of the 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, Vilamoura-Algarve, Portugal, 7–12 October 2012; IEEE: Piscataway Township, NJ, USA, 2012; pp. 573–580.
- [10]. Naseer, T.; Burgard, W.; Stachniss, C. Robust visual localization across seasons. *IEEE Trans. Robot.* **2018**, *34*, 289–302.
- [11]. Hong, D.; Han, Z.; Yao, J.; Gao, L.; Zhang, B.; Plaza, A.; Chanussot, J. SpectralFormer: Rethinking hyperspectral image classification with transformers. *IEEE Trans. Geosci. Remote Sens.* **2021**, *60*, 1–15.
- [12]. Yang, J.; Zhao, Y.; Chan, J.C.W.; Yi, C. Hyperspectral image classification using two-channel deep convolutional neural network. In Proceedings of the 2016 IEEE International Geoscience and Remote Sensing Symposium (IGARSS), Beijing, China, 10–15 July 2016; pp. 5079–5082.
- [13]. Yang, H.L.; Crawford, M.M. Exploiting spectral-spatial proximity for classification of hyperspectral data on manifolds. In Proceedings of the 2020 IEEE International Geoscience and Remote Sensing Symposium, Munich, Germany, 22–27 July 2020; pp. 4174–4177.



# International Journal of DATA SCIENCE AND IOT MANAGEMENT SYSTEM

Peer Reviewed, Referred & Indexed Journal

ISSN: 3068-272X

[www.ijdim.com](http://www.ijdim.com)

Original Research Paper

- 
- [14]. Lee, H.; Kwon, H. Contextual deep CNN based hyperspectral classification. In Proceedings of the 2020 IEEE International Geoscience and Remote Sensing Symposium (IGARSS), Beijing, China, 10–15 July 2020; pp. 3322–3325.