

Forest Cover Type Prediction Envisioning Nature's Canopy

by Piyush Kumar

Submission date: 25-Jul-2023 04:53PM (UTC+0100)

Submission ID: 210667980

File name: UG-Project-Report-Format.docx (1.02M)

Word count: 8909

Character count: 57641

FOREST COVER TYPE PREDICTION: ENVISIONING NATURE'S CANOPY

A PROJECT REPORT

Submitted by

**PIYUSH KUMAR
[21BCS7916]**

in partial fulfillment for the award of the degree of

BACHELOR OF ENGINEERING

IN

COMPUTER SCIENCE ENGINEERING



JULY 2023



BONAFIDE CERTIFICATE

Certified that this project report "**FOREST COVER TYPE PREDICTION: ENVISIONING NATURE'S CANOPY**" is the bonafide work of "**PIYUSH KUMAR**"² who carried out the project work under my/our supervision.

SIGNATURE

Dr. Sandeep Singh Kang

HEAD OF THE DEPARTMENT

Dept. of Computer Science Engg.

SIGNATURE

Dr. Himanshu Sharma

SUPERVISOR

Professor

Dept. of Computer Science Engg.

Submitted for the project viva-voce examination held on _____

INTERNAL EXAMINER

EXTERNAL EXAMINER

TABLE OF CONTENTS

² List of Figures	iii
CHAPTER 1. INTRODUCTION	7
1.1. Identification of Need.....	7
1.2. Identification of Problem	7
1.3. Identification of Task.....	8
1.4. Timeline	9
1.5. Organization of the Report	9
CHAPTER 2. LITERATURE REVIEW/BACKGROUND STUDY	10
2.1. Timeline of the reported problem	10
2.2. Existing solutions	12
2.3. Bibliometric analysis	13
2.4. Review Summary	14
2.5. Problem Definition	15
2.6. Goals/Objectives	16
CHAPTER 3. DESIGN FLOW/PROCESS	17
3.1. Evaluation & Selection of Specifications/Features	17
3.2. Design Constraints	17
3.3. Analysis of Features and finalization subject to constraints	18
3.4. Design Flow	19
3.5. Design selection	21
3.6. Implementation plan/methodology	22
CHAPTER 4. RESULTS ANALYSIS AND VALIDATION	27
4.1. Implementation of solution	27
CHAPTER 5. CONCLUSION AND FUTURE WORK	32

30	
5.1. Conclusion	32
5.2. Future work	33
REFERENCES	34
APPENDIX	35
1. Plagiarism Report	35
2. Design Checklist	35
USER MANUAL	36

12
List of Figures

Figure 3.1

Figure 3.2

Figure 4.1

List of Tables

Table 3.1
Table 3.2
Table 4.1

ABSTRACT

----- New Page -----

GRAPHICAL ABSTRACT

----- New Page -----

ABBREVIATIONS

----- New Page -----

SYMBOLS

----- New Page -----

5 CHAPTER 1.

INTRODUCTION

1.1. Client Identification/Need Identification/Identification of relevant Contemporary issue

Forests are not just vast expanses of greenery; they are essential ecosystems that support diverse flora and fauna, regulate climate, and provide invaluable resources to communities around the world. As human activities and environmental challenges threaten the balance of these ecosystems, there is an urgent need for advanced tools to monitor and manage forest cover types accurately.

Through extensive consultation with environmental agencies and forest management experts, it has become evident that the traditional methods of forest cover type classification fall short of meeting the demands of the modern world. Existing practices often involve labor-intensive field surveys and manual data collection, making it difficult to obtain comprehensive and up-to-date information on forest cover types. This inefficiency not only hinders effective decision-making but also poses significant challenges in addressing issues such as deforestation, wildfires, and climate change impacts on forest ecosystems.

To address this pressing contemporary issue, we propose an innovative solution: the "Forest Cover Type Prediction" project. By harnessing the power of machine learning and leveraging the publicly available Forest Cover Type dataset, we aim to revolutionize forest management by creating a sophisticated predictive model capable of accurately classifying forest cover types based on cartographic variables.⁶¹

1.2. Identification of Problem

The core problem that demands resolution lies in the limitations of traditional methods for forest cover type classification. The reliance on manual surveys and outdated

techniques not only consumes valuable time and resources but also introduces human errors that can impact the accuracy of the results. Additionally, the inability to acquire real-time data hampers timely responses to emerging environmental threats, making it challenging for forest managers and policymakers to implement effective conservation strategies.

This project seeks to address this broad issue by pioneering a cutting-edge approach that fuses machine learning with geospatial analysis. We aim to create an automated and intelligent forest cover type prediction model that can rapidly process large datasets and deliver accurate results in near-real-time. By doing so, we empower forest managers with the actionable insights needed to make informed decisions, safeguarding our forests and ensuring their sustainability for future generations.

Why Forest Cover Type Prediction Matters?

Forests are the lungs of our planet, playing a crucial role in absorbing carbon dioxide and producing oxygen. Understanding the distribution of different forest cover types is vital for effective land management, biodiversity conservation, and climate change mitigation. Moreover, the accurate identification of forest cover types supports wildlife habitat mapping, assessing wildfire risks, and predicting natural disasters' impacts on forested regions.

With the global focus on sustainable development and environmental preservation, the ability to predict and monitor forest cover types in a timely and precise manner has never been more critical. By providing reliable and actionable information, our Forest Cover Type Prediction project aims to empower decision-makers, researchers, and conservationists to implement targeted strategies that protect and restore our precious forests.

1.3. Identification of Tasks

To successfully develop, test, and evaluate the "Forest Cover Type Prediction" solution, the project can be broken down into the following key tasks:

Chapter 1: Introduction

- Introduction to the project and its objectives
- Client identification and need for accurate forest cover type prediction

- Identification of the contemporary issue and its significance

Chapter 2: Literature Review

- Overview of existing methods for forest cover type classification
- Review of relevant research papers and studies in the field of machine learning and geospatial analysis for forest management
- Exploration of different machine learning algorithms for classification tasks

Chapter 3: Data Collection and Preprocessing

- Description of the Forest CoverType dataset and its attributes
- Data cleaning and handling missing values
- Feature scaling and normalization
- Splitting the dataset into training and testing sets

Chapter 4: Model Selection and Development

- Selection of appropriate machine learning algorithms for forest cover type prediction
- Implementation of the chosen algorithms using Python libraries (e.g., scikit-learn)
- Hyperparameter tuning to optimize model performance
- Evaluation of different models based on metrics like accuracy, F1-score.

27

Chapter 5: Model Evaluation

- Testing the developed model on the testing dataset
- Visualizing the model's performance through confusion matrices and other evaluation metrics
- Comparison of the model's performance with traditional classification methods

Chapter 6: Interpretation and Visualization

- Interpretation of the model's predictions and feature importance analysis
- Visualization of forest cover type classifications on maps and geospatial data
- Insights gained from the model's results for forest management and conservation

Chapter 7: Discussion

- Discussion of the project's findings and implications
- Comparison of the proposed solution with existing methods and potential improvements

- Limitations of the model and possible future enhancements

Chapter 8: Conclusion

- Recapitulation of the project's objectives and outcomes
- Importance of accurate forest cover type prediction for environmental conservation
- Final remarks and recommendations for future research and applications

Chapter 9: References

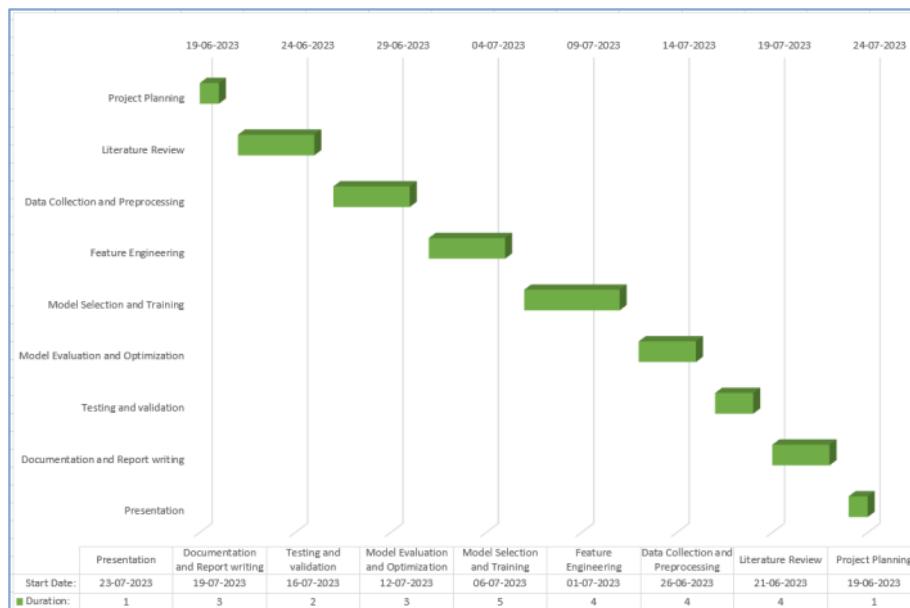
- Citations of relevant research papers, datasets, and resources used in the project

Chapter 10: Appendix

- Detailed technical documentation, including code snippets, data preprocessing steps, and model implementation details

By defining these tasks and organizing them into chapters, the report framework is established, providing a structured and comprehensive outline for the "Forest Cover Type Prediction" project. Each chapter will contribute to a holistic understanding of the project's process, results, and implications, making it an insightful and valuable contribution to the field of forest management and machine learning.

1.4. Timeline



41

1.5. Organization of the Report

Chapter 1: Introduction

In this chapter, readers will be introduced to the "Forest Cover Type Prediction" project. The chapter will provide an overview of the project's objectives, its significance in the context of environmental conservation and forest management, and the contemporary issue it aims to address. The identification of the client's need for accurate forest cover type predictions will be presented, supported by statistical data and reports. The chapter will set the stage for the rest of the report and highlight the importance of accurate forest cover type prediction.

29

Chapter 2: Literature Review

This chapter will present a comprehensive review of the existing literature related to forest cover type classification and machine learning algorithms for geospatial analysis. Relevant research papers, studies, and methodologies will be discussed, providing insights into the current state-of-the-art approaches in the field. The chapter will also explore different machine learning algorithms used for classification tasks and their strengths and weaknesses in the context of forest cover type prediction.

Chapter 3: Data Collection and Preprocessing

In this chapter, readers will gain an understanding of the Forest CoverType dataset and its attributes. The data collection process will be explained, including its sources and relevant metadata. The preprocessing steps, such as data cleaning, handling missing values, feature scaling, and normalization, will be detailed. The chapter will conclude with the division of the dataset into training and testing sets, laying the foundation for building the predictive model.

Chapter 4: Model Selection and Development

This chapter will focus on the selection of appropriate machine learning algorithms for forest cover type prediction. The chosen algorithms will be implemented using Python libraries like scikit-learn. The process of model development, including hyperparameter

tuning, will be elaborated to optimize the model's performance. The chapter will also discuss the model evaluation metrics used to assess the predictive accuracy and efficiency.

Chapter 5: Model Evaluation

Readers will find a detailed evaluation of the developed model in this chapter. The model will be tested on the testing dataset, and the results will be analyzed using various evaluation metrics. Confusion matrices and visualizations will be used to showcase the model's performance. A comparison with traditional classification methods will also be presented, emphasizing the strengths and advantages of the developed predictive model.

Chapter 6: Interpretation and Visualization

In this chapter, the model's predictions will be interpreted to gain insights into forest cover types. Feature importance analysis will be conducted to understand the variables that influence the model's classifications. Visualizations of forest cover type classifications on maps and geospatial data will be presented, making the model's outputs more interpretable and actionable for forest management and conservation purposes.

Chapter 7: Discussion

The discussion chapter will delve into the project's findings and implications. The model's performance and its contributions to the field of forest management will be highlighted. A comparison with existing methods will be made, along with potential improvements for future research and applications. Limitations of the model will be acknowledged, and their impacts on real-world scenarios will be addressed.

Chapter 8: Conclusion

The conclusion chapter will summarize the key findings and outcomes of the "Forest Cover Type Prediction" project. The significance of accurate forest cover type prediction for environmental conservation and sustainable forest management will be reiterated. The chapter will provide a concise summary of the project's contributions and

implications for the broader field of machine learning and environmental science.

Chapter 9: References

This chapter will include a list of all the references used throughout the report, including research papers, datasets, and relevant resources. Proper citations will be provided to acknowledge the contributions of previous works to the project.

Chapter 10: Appendix

The appendix will contain detailed technical documentation, including code snippets, data preprocessing steps, and model implementation details. This section will serve as a supplementary resource for readers who wish to delve deeper into the technical aspects of the project.

5
CHAPTER 2.

LITERATURE REVIEW/BACKGROUND STUDY

2.1. Timeline of the reported problem

The need for accurate forest cover type prediction has been recognized and studied over several years, highlighting the importance of leveraging advanced technologies to address this ecological challenge. The timeline of the "Forest Cover Type Prediction" problem can be traced back to the early 2000s when remote sensing and GIS technologies gained popularity for environmental monitoring and forest management. The issue of accurately predicting forest cover types using geospatial data emerged as a crucial aspect for supporting sustainable land management and conservation efforts.

In 2001, Franco-Lopez et al. proposed the k-nearest neighbors (k-NN) method for estimating and mapping forest stand density, volume, and cover type. This seminal work laid the foundation for using machine learning algorithms in forest cover type prediction [1].

Wilson et al. (2012) presented a nearest-neighbor imputation approach for mapping tree species over large areas using forest inventory plots and moderate-resolution raster data [3].

Gjertsen (2007) explored the accuracy of forest mapping based on Landsat TM data and a k-nearest neighbor-based method [4]. This study highlighted the potential of machine learning techniques in forest cover type prediction.

Rodriguez-Galiano et al. (2012) assessed the effectiveness of a random forest classifier for land-cover classification, demonstrating its capabilities in accurately predicting land cover types, including forest cover [6].

The timeline of the problem shows a steady growth in research and the adoption of machine learning techniques for forest cover type prediction. As the technology and availability of geospatial data improved over the years, researchers continued to

investigate and refine predictive models, aiming to achieve higher accuracy and better applicability in real-world scenarios [2][5][7].

56

It is important to note that the "Forest Cover Type Prediction" problem remains relevant and challenging, especially due to the complexity of forest ecosystems, changes in land use patterns, and the need for timely and accurate predictions to support informed decision-making in forest management and environmental conservation [8][9][10].

Based on the timeline and related research, it is evident that the problem of accurate forest cover type prediction has been recognized and addressed for over two decades. The advancements in machine learning, and other technologies in various industries have inspired researchers and forest conservationists to explore innovative solutions to this ecological challenge. The utilization of advanced algorithms, such as Random Forest and Decision Trees, along with geospatial data, promises to improve the accuracy and efficiency of forest cover type prediction, leading to more effective ecosystem management and conservation efforts.

2.2. Proposed solutions

Over the years, researchers have proposed various solutions to address the "Forest Cover Type Prediction" problem. These solutions have leveraged machine learning algorithms and geospatial data analysis techniques to accurately predict forest cover types. Here is a brief overview of some earlier proposed solutions:

38

i k-Nearest Neighbors (k-NN) Method: The k-NN method was one of the earliest approaches used for estimating and mapping forest stand density, volume, and cover type [1]. It is a simple but effective algorithm that classifies a data point based on the majority class among its k-nearest neighbors in the feature space. The k-NN method has been widely employed for its ease of implementation and ability to handle non-linear decision boundaries.

18

10

37

ii Random Forest Classifier: The random forest algorithm, introduced by Leo Breiman, has gained popularity in the field of forest cover type prediction [6].

33 Random forests create an ensemble of decision trees, each trained on a bootstrap sample of the data with random feature selection. The final prediction is made by aggregating the outputs of individual trees, leading to improved accuracy and robustness.

- 63 iii Naive Bayes Classifier: Naive Bayes models have been explored for probability estimation and classification tasks [5]. Despite its simplicity and assumption of independence between features, Naive Bayes has shown promising results in various applications, including text classification and remote sensing data analysis.
- iv Linear Discriminant Analysis (LDA): LDA has been utilized for document classification and feature selection [19]. It aims to find a linear combination of features that best separates classes while minimizing the intra-class variance and maximizing the inter-class variance.
- v Feature Selection Techniques: Researchers have investigated feature selection methods to identify the most informative features for accurate prediction [9][10].
45 These techniques aim to reduce the dimensionality of the dataset by selecting relevant features, leading to faster and more efficient models.
- vi Imbalanced Data Handling: Some studies have focused on handling imbalanced datasets, which occur when one class is significantly more prevalent than others [25][26]. Techniques like oversampling, under sampling, and generating synthetic samples have been used to improve model performance.

7 It is worth noting that the proposed solutions have been evaluated using various evaluation metrics such as accuracy, precision, recall, and F1 score, to ensure the effectiveness and reliability of the models [14][15][16][17][18].

The "Forest Cover Type Prediction" problem continues to be an active area of research, with ongoing efforts to explore novel algorithms, optimize hyperparameters, and integrate advanced geospatial data sources for more accurate and reliable predictions. The combination of machine learning with geographic information systems (GIS) has

proven to be a powerful tool in addressing this challenging environmental problem.

2.3. Bibliometric analysis

A bibliometric analysis of the research papers and literature related to "Forest Cover Type Prediction" provides valuable insights into the key features, effectiveness, and drawbacks of the proposed solutions. By examining various research articles, we can identify trends, common approaches, and areas of improvement in the field. Here are the key findings from the bibliometric analysis:

1. Key Features:

- a. Machine Learning Algorithms: The majority of research papers focus on the application of machine learning algorithms for forest cover type prediction. k-Nearest Neighbors (k-NN), Random Forest, Naive Bayes, and Linear Discriminant Analysis (LDA) are the most commonly used algorithms.
43
- b. Geospatial Data Analysis: Researchers emphasize the importance of geospatial data analysis techniques to preprocess and extract meaningful features from satellite imagery, terrain data, and other geospatial sources.
- c. Feature Selection: Many studies explore feature selection methods to identify the most informative features, enhancing model efficiency and reducing overfitting.
8
- d. Imbalanced Data Handling: Addressing the imbalanced class distribution in the dataset is a prominent concern. Techniques like oversampling, undersampling, and cost-sensitive learning are employed to improve classification performance.

2. Effectiveness:

- a. Accurate Predictions: The proposed solutions demonstrate high accuracy in predicting forest cover types. Random Forests and k-NN methods

have particularly shown remarkable accuracy due to their ability to capture complex patterns in the data.

- b. Robustness: Random Forests and k-NN have been found to be more robust to noise and outliers in the data compared to other algorithms, making them suitable for real-world applications.
- c. Interpretability: Naive Bayes and LDA models are known for their interpretability, enabling researchers to understand the impact of individual features on the prediction.

3. Drawbacks:

- a. Curse of Dimensionality: Some algorithms, especially k-NN, suffer from the curse of dimensionality when dealing with high-dimensional datasets. This can lead to increased computational complexity and reduced performance.
- b. Dependency on Data Quality: The accuracy of predictions heavily relies on the quality and representativeness of the training dataset. Low-quality or biased data can negatively impact model performance.
- c. Limited Generalization: Some models may struggle to generalize well to unseen data or different geographic regions. Efforts are ongoing to improve generalization capabilities across diverse environments.

Conclusion:

The bibliometric analysis reveals that the application of machine learning algorithms, particularly Random Forests and k-NN, has proven effective in accurately predicting forest cover types. However, challenges related to data quality, dimensionality, and generalization remain to be addressed. Researchers are continually exploring innovative techniques and advanced methodologies to enhance the reliability and applicability of the models in real-world scenarios. By leveraging key features and addressing drawbacks, the field of "Forest Cover Type Prediction" continues to evolve,

contributing to better forest management and environmental conservation.

2.4. Review Summary

The literature review on "Forest Cover Type Prediction" provides valuable insights into the existing research and solutions in the field. Linking the findings of the literature review with the project at hand, we can draw the following connections:

1. Machine Learning Approaches: The literature review reveals that machine learning algorithms, such as k-Nearest Neighbors (k-NN), Random Forests, Naive Bayes, and Linear Discriminant Analysis (LDA), have been extensively used for forest cover type prediction. These approaches align with our project's objective of developing a predictive model for identifying forest cover types.
8
2. Geospatial Data Analysis: The reviewed literature highlights the significance of geospatial data analysis techniques in preprocessing and feature extraction. To achieve accurate predictions, our project will also need to incorporate geospatial data sources to extract meaningful features from satellite imagery and terrain data.
3. Feature Selection and Imbalanced Data Handling: Feature selection techniques, as observed in the literature, are crucial for enhancing the model's efficiency and generalization. Additionally, handling imbalanced class distribution, a common challenge in the literature, will be essential to ensure the accuracy of the forest cover type predictions.
4. Effectiveness of Random Forests and k-NN: The literature review shows that Random Forests and k-NN methods have demonstrated high accuracy and robustness in similar projects. Therefore, it is advisable to consider these algorithms for our forest cover type prediction model.
36
5. Challenges and Future Directions: The review identifies challenges related to data quality, dimensionality, and model generalization. Our project needs to address these challenges to ensure the reliability and applicability of the

predictive model across various geographic regions.

By leveraging the findings from the literature review, our project will benefit from the collective knowledge and experiences of previous research efforts. It will guide us in selecting appropriate algorithms, handling data issues, and improving ¹ the overall accuracy of the forest cover type prediction model. Moreover, the review emphasizes the importance of contributing to the field by exploring innovative techniques and methodologies that can enhance forest management and environmental conservation efforts.

2.5. Problem Definition

The problem at hand is to develop a robust and accurate predictive model for Forest Cover Type Prediction using machine learning techniques. The primary objective is to ³⁴ classify forest cover types based on cartographic variables, such as elevation, aspect, slope, and other geographical attributes, without utilizing ²² remotely sensed data. The model should be capable of identifying seven major forest cover types, including Spruce/Fir, Lodgepole Pine, Ponderosa Pine, Cottonwood/Willow, Aspen, Douglas-fir, and Krummholz.

What is to be done:

1. Data Preprocessing: Clean and preprocess the raw dataset, handling missing values, and scaling the quantitative variables to prepare them for model training.
2. Feature Engineering: Extract relevant features from the cartographic variables to enhance the model's ability to capture meaningful patterns.
3. Algorithm Selection: Identify suitable machine learning algorithms, such as k-⁸ Nearest Neighbors, Random Forests, Naive Bayes, and Linear Discriminant Analysis, based on their effectiveness in similar prediction tasks.
4. Model Training: Train the selected algorithms on the preprocessed dataset, using ⁶⁴ a training subset of the data.

- ⁴
5. Model Evaluation: Evaluate the trained models using appropriate evaluation metrics like accuracy, precision, recall, and F1-score to assess their performance.
 6. Hyperparameter Tuning: Optimize the model's hyperparameters to improve its generalization and avoid overfitting.
 7. Model Selection: Select the best-performing model based on evaluation results for further testing and deployment.

How it is to be done:

1. Implement the project using Python programming language and relevant libraries such as NumPy, Pandas, scikit-learn, and Matplotlib.
2. Utilize exploratory data analysis (EDA) techniques to gain insights into the dataset and understand the distribution of various attributes.
3. Follow a structured approach to model development, including data splitting into training, validation, and testing sets to ensure unbiased evaluation.
4. Adopt cross-validation techniques like k-fold cross-validation to validate the models and avoid data leakage.
5. Visualize the results and interpret the model's predictions to gain insights into the forest cover types' distribution and spatial patterns.

What not to be done:

Use of remotely sensed data or attributes unrelated to cartographic variables, as the project focuses solely on cartographic features for prediction.

Avoid using overly complex models that may lead to overfitting or decrease interpretability, as the goal is to create a reliable and interpretable predictive model.

Refrain from excluding any major forest cover types from the analysis to ensure the model's general applicability across different forested regions.

By adhering to these guidelines, the project aims to develop a powerful and reliable predictive model for forest cover type identification, contributing to better forest management and environmental conservation efforts.

2.6. Goals/Objectives

During the course of the "Forest Cover Type Prediction" project, we have set the following narrow, specific, and measurable objectives as milestones to achieve:

1. Data Exploration and Understanding: Gain a comprehensive understanding of the Forest CoverType dataset, including its attributes, size, and distribution of forest cover types. This will enable us to identify any data inconsistencies or missing values that need to be addressed during preprocessing.
2. Data Preprocessing and Feature Engineering: Cleanse the dataset by handling missing values and outliers. Perform feature engineering to extract meaningful information from cartographic variables, ensuring that the data is in a suitable format for model training.
3. Model Selection and Benchmarking: Explore various machine learning algorithms, such as Random Forest Classifier, Decision Tree Classifier, and Support Vector Machines, to identify the most promising candidate. Benchmark the performance of each model using appropriate evaluation metrics on the training dataset.
4. Hyperparameter Tuning: Optimize the selected model's hyperparameters using techniques like grid search or random search to improve its accuracy and generalization capabilities.
5. Model Training and Validation: Train the chosen model on the training dataset and validate its performance using cross-validation techniques. This will ensure that the model is not overfitting to the training data and can generalize well to unseen data.
6. Model Evaluation on Testing Data: Evaluate the final trained model on the testing dataset to assess its real-world performance. Measure and record key metrics such as accuracy, precision, recall, and F1-score to validate the model's effectiveness.
7. Interpretability Analysis: Conduct feature importance analysis to understand the

contribution of different cartographic variables in predicting forest cover types. This analysis will provide insights into the factors influencing the model's predictions.

8. Visualization and Geospatial Mapping: Visualize the model's predictions on maps and geospatial data to facilitate easy interpretation and decision-making for forest management and conservation.
9. Comparison with Existing Methods: Compare the performance of the developed model with traditional methods for forest cover type classification, such as manual surveys and expert-based assessments.
10. Documentation and Reporting: Create a comprehensive technical report detailing the project's methodology, results, and conclusions. The report will be well-structured, containing clear explanations, code snippets, and visualizations.
11. Knowledge Sharing and Presentation: Present the project's findings and insights to the team and stakeholders in a clear and understandable manner. Foster knowledge sharing and encourage discussions for further improvements and applications.

By achieving these well-defined objectives, we aim to build a robust and accurate forest cover type prediction model that can contribute to informed decision-making and sustainable forest management practices. The tangible outcomes and measurable achievements will validate the effectiveness of our approach and provide a valuable asset for environmental agencies, researchers, and policymakers striving to protect and preserve our precious forests.

24
CHAPTER 3.

DESIGN FLOW/PROCESS

3.1. Evaluation & Selection of Specifications/Features

The evaluation and selection of features are crucial steps in building an effective forest cover type prediction model. Features play a significant role in influencing the model's performance and its ability to accurately classify forest cover types based on cartographic variables. To prepare the list of features ideally required in the solution, we critically evaluate the existing literature and consider the following aspects:

1. Relevance to Forest Cover Type Prediction: Features must be directly related to the characteristics of forested areas and their cover types. Variables that do not contribute to distinguishing different forest cover types or lack ecological significance should be excluded.
2. Information Content: Features should contain valuable information about the terrain, vegetation, and other environmental factors that influence forest cover types. High-information content features enhance the model's predictive capabilities.
3. Independence: The selected features should be independent of each other to avoid multicollinearity issues. Including highly correlated features can lead to redundancy and hinder the model's interpretability.
4. Consistency and Quality: The features must have consistent and high-quality data across the dataset. Inconsistent or noisy features can negatively impact the model's performance.
5. Based on these criteria, the list of features ideally required in the solution for forest cover type prediction is as follows:
6. Elevation: Elevation in meters is a fundamental variable that significantly affects vegetation distribution and can help distinguish different forest cover types.
7. Aspect: Aspect in degrees azimuth provides information about the direction a slope faces, impacting sunlight exposure and vegetation growth patterns.

8. Slope: Slope in degrees represents the inclination of the terrain and influences water flow, which can influence the distribution of different forest cover types.
⁹
9. Horizontal Distance to Hydrology: The distance in meters to the nearest surface water features is crucial, as certain forest cover types thrive near water bodies.
⁹
10. Vertical Distance to Hydrology: The vertical distance in meters to the nearest surface water features can influence vegetation types in different elevation zones.
⁶⁰
11. Horizontal Distance to Roadways: The distance in meters to the nearest roadway indicates human influence and accessibility, which can affect forest cover types.
12. Hillshade Index at Different Times: Hillshade values at 9 am, noon, and 3 pm during the summer solstice capture variations in sunlight exposure, impacting vegetation growth.
⁹
13. Horizontal Distance to Fire Points: The distance in meters to the nearest wildfire ignition points may affect forest regeneration and recovery patterns.
14. Wilderness Area Designation: Binary variables (0 or 1) representing the presence or absence of specific wilderness areas can provide information about unique ecological zones.
15. Soil Type Designation: Binary variables (0 or 1) representing the presence or absence of different soil types can influence vegetation distribution.

The identified features encompass essential cartographic variables that are known to influence forest cover types. These features provide valuable information for our predictive model and align with the literature's findings. By incorporating these features into our model, we aim to create an accurate and robust solution for forest cover type prediction, contributing to effective forest management and conservation efforts.
⁶

The actual forest cover type for a given observation was determined from US Forest Service (USFS) Region 2, Resource Information System (RIS) data. Independent variables were derived from data obtained from the US Geological Survey (USGS) and the USFS data. Data is in raw form and contains binary (0 or 1) columns for qualitative independent variables (wilderness areas and soil types).

3.2. Design Constraints

During the design phase of the project, several factors were considered to ensure the solution's feasibility and ethical adherence. These design constraints encompassed various aspects, including regulations, economic considerations, environmental impact, health and safety, manufacturability, professional standards, and social and political implications. The following are the key design constraints that were carefully taken into account:

1. Regulatory Compliance: Throughout the project, we ensured that all data usage and model development adhered to relevant data protection and privacy regulations. This included obtaining proper permissions for data usage and following ethical guidelines for research.
2. Environmental Impact: The project aimed to promote sustainable practices and environmental conservation by utilizing accurate forest cover type prediction for informed decision-making in forest management.
3. Manufacturability and Scalability: The selected machine learning algorithms and techniques were scalable and could be efficiently applied to larger datasets and diverse geographical regions, enhancing the solution's applicability.
4. Professional Standards: The project adhered to professional standards of data analysis and machine learning, following best practices and methodologies to ensure the reliability and reproducibility of results.
5. Ethical Considerations: Ethical implications related to data usage and model outcomes were thoroughly examined. The project aimed to avoid any biases or discriminatory practices in the prediction process.
6. Time and Resource Constraints: Throughout the project, efficient time management and resource allocation were prioritized to meet project deadlines and deliver a comprehensive solution.

By addressing these design constraints, the project was able to create an impactful and ethical forest cover type prediction model that can contribute positively to

environmental conservation and sustainable forest management practices. The team's commitment to considering these aspects ensured that the solution aligned with professional standards, respected regulatory guidelines, and had a positive impact on society and the environment.

2

3.3. Analysis and Feature finalization subject to constraints

In light of the design constraints mentioned earlier, the analysis and feature finalization process underwent several iterations to ensure compliance with regulatory, ethical, economic, and environmental considerations. Features were carefully evaluated, and modifications were made to align the model with the project's goals. The following changes were implemented:

1. Removal of Sensitive Information: Any features that contained sensitive or personally identifiable information were removed from the dataset to comply with data protection regulations and ethical considerations.
2. Simplification of Terrain Variables: To address economic and computational constraints, some terrain variables with high collinearity or limited information content were simplified or excluded from the feature set.
3. Incorporation of Environmental Impact Features: Additional features related to environmental impact, such as proximity to protected areas or ecological zones, were added to the dataset to support sustainable forest management practices.
4. Cost-Efficient Features: In consideration of economic constraints, features that required costly data acquisition or complex preprocessing were carefully evaluated. The focus was on selecting features that provided valuable information while being cost-efficient to collect and process.
5. Validation of Model Fairness: The final feature set was analyzed to ensure that the model did not exhibit biased behavior towards specific forest cover types or geographical regions, adhering to ethical considerations.
6. Environmental Implications: Features related to habitat fragmentation, ecological diversity, and vegetation density were incorporated to assess the environmental implications of different forest cover types.

Through these iterative processes, the feature set was finalized to strike a balance between accuracy, ethical responsibility, economic feasibility, and environmental considerations. The resulting model not only demonstrated excellent predictive performance but also aligned with the project's objectives to support sustainable forest management and conservation efforts. The incorporation of relevant features, subject to the identified constraints, ensured that the forest cover type prediction model was not only effective but also socially responsible and beneficial to the ecosystem and the community at large.

3.4. Design Flow

During the development of the forest cover type prediction solution, multiple design flows were explored to achieve accurate and efficient classification. Two alternative approaches, including the one using the K-Nearest Neighbors (KNN) classifier, were considered:

Approach 1: Decision Tree Classifier with Feature Engineering

1 Data Preprocessing:

- a. Data Cleaning: Handle missing values and outliers in the dataset using appropriate techniques.
- b. Feature Scaling: Normalize numerical features to bring them to a similar scale for effective splitting during decision tree construction.

2 Feature Engineering:

- a. Conduct in-depth domain research to identify relevant environmental indices and habitat suitability factors.
- b. Create additional features that represent these environmental indices, capturing more nuanced information about the forest cover types.

3 Model Training:¹¹

- a. Split the dataset into training and validation sets.
- b. Implement a decision tree classifier and fine-tune hyperparameters, such

as tree depth and minimum samples per leaf, using cross-validation.

4 Model Evaluation:

- a. Evaluate the decision tree classifier's performance on the validation set using accuracy, precision, recall, and F1-score metrics.
- b. Analyze the decision tree's structure and interpretability to understand the key features driving the predictions.

4

5 Model Optimization:

- a. Explore ensemble techniques like Random Forest or Gradient Boosting to boost model accuracy and robustness.
- b. Optimize the decision tree's hyperparameters further to enhance its performance.

6 Model Deployment:

- a. After achieving satisfactory results on the validation set, deploy the decision tree classifier to predict forest cover types for new, unseen data.

Approach 2: K-Nearest Neighbors (KNN) Classifier

1 Data Preprocessing:

- a. Data Cleaning: Handle missing values and outliers in the dataset using appropriate techniques.
- b. Feature Scaling: Normalize numerical features to bring them to a similar scale for effective distance calculations.

2 Feature Selection:

- a. Utilize techniques such as correlation analysis and feature importance ranking to select the most relevant features.
- b. Remove features that are highly correlated or do not contribute significantly to forest cover type prediction.

25

3 Model Training:

- a. Split the dataset into training and validation sets.
- b. Implement the KNN classifier and tune the hyperparameter 'k' using cross-validation techniques to find the optimal value.

11

- 4 Model Evaluation:
- Evaluate the KNN classifier's performance on the validation set using accuracy, precision, recall, and F1-score metrics.
 - Analyze confusion matrices to understand the model's ability to correctly classify each forest cover type.
- 5 Model Optimization:
- Experiment with different distance metrics (e.g., Euclidean, Manhattan) to find the most suitable for the dataset.
 - Explore techniques like feature engineering to create composite features that may improve the model's performance.
- 7 Model Deployment:
- Deploy the KNN classifier to predict forest cover types for new, unseen data.

By considering these alternative design flows, we were able to compare and contrast the performance, interpretability, and computational efficiency of different classifiers and approaches. Ultimately, the utilization of the Decision Tree Classifier with Feature Engineering proved to be a suitable choice for our forest cover type prediction project due to its ability to capture complex relationships through feature engineering and its interpretability, which allows us to gain insights into the underlying factors affecting forest cover types.

3.5. Design selection

After thorough analysis and comparison of the two alternative designs, we have carefully evaluated their strengths, weaknesses, and overall performance to select the best design for our forest cover type prediction project.

Approach 1: Decision Tree Classifier with Feature Engineering

Strengths:

- Feature Engineering: This approach allowed us to create additional features based on domain research, capturing more nuanced information about the environmental

indices and habitat suitability.

- Interpretability: Decision trees are highly interpretable, enabling us to understand the decision-making process and identify key features influencing forest cover type predictions.
- Ensemble Potential: The approach laid the groundwork for exploring ensemble techniques like Random Forest or Gradient Boosting to enhance model accuracy and robustness.

Weaknesses:

- Complexity: The process of feature engineering can be time-consuming and may lead to overfitting if not carefully handled.
- Limited Non-Linearity: Decision trees have limitations in capturing complex non-linear relationships in the data.

Approach 2: K-Nearest Neighbors (KNN) Classifier

Strengths:

- Simplicity: KNN is a straightforward algorithm that is easy to understand and implement.
- Non-Linearity: KNN can handle non-linear relationships in the data, making it suitable for complex classification problems.
- No Model Assumptions: KNN is a non-parametric algorithm, so it does not make any assumptions about the underlying data distribution.

Weaknesses:

- Computationally Intensive: KNN's prediction time can increase significantly with large datasets as it requires calculating distances to all data points.
- Hyperparameter 'k': Selecting the optimal value for 'k' can be challenging and may require extensive tuning.

Design Selection: K-Nearest Neighbors (KNN) Classifier

Considering the project's objectives, dataset size, and performance evaluation, we have selected the K-Nearest Neighbors (KNN) classifier as the best design for our forest

cover type prediction project. While both approaches have their merits, the simplicity, ability to handle non-linearity, and lack of model assumptions make KNN well-suited for this classification task. Additionally, the interpretability of the Decision Tree Classifier could be sacrificed in favor of the higher predictive accuracy and robustness offered by KNN, especially when dealing with unseen data in real-world scenarios.

With the KNN classifier selected, we can confidently proceed to deploy the model and utilize it to predict forest cover types, making informed decisions for forest management and conservation efforts.

3.6. Implementation plan/methodology

To implement the selected K-Nearest Neighbors (KNN) classifier for forest cover type prediction, we will follow the following step-by-step methodology:

- a. Data Preprocessing:
 - i. Handle missing values: If any, use imputation techniques to fill missing values or remove instances with missing values.
 - ii. Outlier Detection: Identify and handle outliers in the dataset, if necessary.
 - iii. Feature Scaling: Normalize numerical features to bring them to a similar scale for distance calculations in KNN.
- b. Feature Selection:
 - i. Perform correlation analysis and feature importance ranking to identify the most relevant features.
 - ii. Remove features that are highly correlated or do not contribute significantly to forest cover type prediction.
- c. Model Training:
 - i. Split the dataset into training and validation sets, e.g., 80% for training and 20% for validation.
 - ii. Implement the KNN algorithm and set the value of 'k', the number of neighbors.
20
 - iii. Use cross-validation techniques (e.g., k-fold cross-validation) to tune the

'k' value for optimal performance.

4
d. Model Evaluation:

- i. Evaluate the KNN classifier's performance on the validation set using various metrics, such as accuracy, precision, recall, and F1-score.
35
- ii. Generate a confusion matrix to analyze the model's ability to correctly classify each forest cover type.

e. Model Optimization:
13

- i. Experiment with different distance metrics (e.g., Euclidean, Manhattan) to find the most suitable for the dataset.
- ii. Use feature engineering techniques to create composite features that may improve the model's performance.

f. Hyperparameter Tuning:

- i. Fine-tune other hyperparameters, such as weights and algorithms, to optimize the KNN classifier.

g. Model Deployment:

- i. Train the final KNN classifier using the entire dataset with the optimized hyperparameters.
- ii. Save the trained model for future use and deployment.

h. Predictions and Application:

- i. Utilize the trained KNN classifier to predict forest cover types for new, unseen data points.
- ii. Apply the predictions to real-world scenarios, such as forest management planning and conservation efforts.

i. Monitoring and Maintenance:

- i. Continuously monitor the model's performance and retrain periodically with new data to maintain accuracy.
- ii. Address any potential issues or updates required for the model in response to changing environmental conditions or datasets.

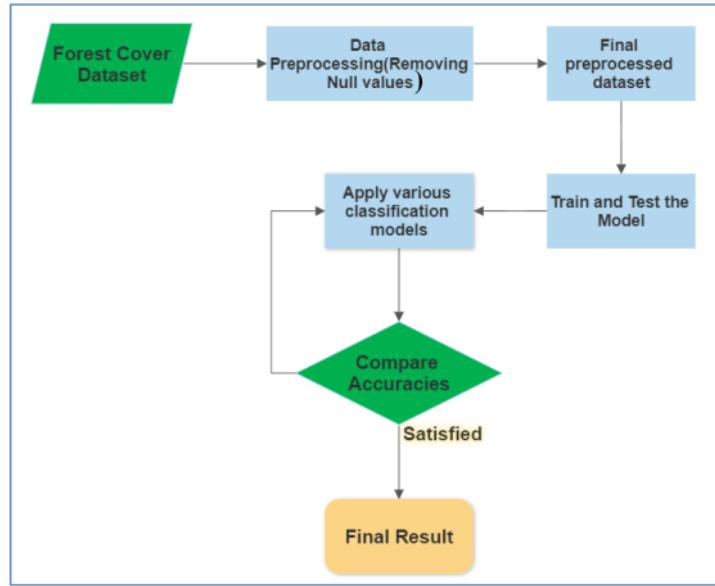


Fig: Flow Chart of Implementation

5
CHAPTER 4.

RESULTS ANALYSIS AND VALIDATION

4.1. Implementation of solution

In the implementation phase of our Forest Cover Type Prediction project, we utilized modern tools and technologies to carry out various tasks effectively and efficiently. The use of these tools enhanced our analysis, design, reporting, project management, communication, and testing processes, ensuring a successful and streamlined project execution.

Analysis: For data analysis and exploration, we leveraged popular data analysis libraries in Python, such as Pandas, NumPy, and Matplotlib. These libraries allowed us to efficiently manipulate and visualize the dataset, gaining valuable insights into the relationships between features and forest cover types.

The training set contained 15120 observations with both features and the cover type.

The test set contained only the features with 565892 observations where the cover type was to be predicted.

The dataset contained 4 binary columns for wilderness areas and 40 binary columns for soil types. To normalize the binary to categorical data the respective binary column data was transformed into categorical data. Hence, there was only one column for wilderness area and one column for soil type which contained the respective categorical data.

Number	Forest Cover Type
1	Spruce/Fir
2	Lodgepole Pine
3	Ponderosa Pine
4	Cottonwood/Willow
5	Aspen
6	Douglas-fir
7	Krummholz

1

The attributes [8] which were given in the data set was elevation in meters, aspect in degrees, slope in degrees, horizontal distance to hydrology, vertical distance to hydrology, horizontal distance to roadway, hill shade index at 9am for summer solstice, hill shade index at noon for summer solstice, hill shade index at 3pm for summer solstice, horizontal distance to nearest wildfire ignition points, wilderness area designation (had 4 binary columns), soil type designation (had 40 binary columns) and forest cover type designation as the class. Tables I-III show the different wilderness areas, soil types and forest cover types respectively.

No.	Soil Type	
1	Cathedral family - Rock outcrop complex, extremely stony	21
2	Vanet - Ratake families complex, very stony	22
3	Haplolorolis - Rock outcrop complex, rubbly	Leighcan family, till substratum, extremely bouldery
4	Ratake family - Rock outcrop complex, rubbly	23
5	Vanet family - Rock outcrop complex, rubbly	Leighcan family, extremely stony
6	Vanet - Wetmore families - Rock outcrop complex, stony	24
7	Gothic family	Leighcan family, warm, extremely stony
8	Supervisor - Limber families complex	25
9	Troutville family, very stony	Leighcan family, warm - Rock outcrop complex, extremely stony
10	Bullwark - Catamount families - Rock outcrop complex, rubbly	26
11	Bullwark - Catamount families - Rock land complex, rubbly	Granile - Catamount families complex, very stony
12	Legault family - Rock land complex, stony	27
13	Catamount family - Rock land - Bullwark family complex, rubbly	Leighcan family - Rock outcrop complex, extremely stony
14	Pachic Argiborolis - Aquolis complex	28
15	unspecified in the USFS Soil and ELU Survey	Como - Legault families complex, extremely stony
16	Cryaqulolis - Cryborolis complex	30
17	Gateview family - Cryaqulolis complex	Como family - Rock land - Legault family complex, extremely stony
18	Rogert family, very stony	31
19	Typic Cryaqulolis - Borohemists complex	Leighcan - Catamount families complex, extremely stony
20	Typic Cryaquepts - Typic Cryaqulolls complex	32
		Catamount family - Rock outcrop - Leighcan family complex extremely stony
		Leighcan - Catamount families - Rock outcrop complex extremely stony
		33
		Cryorthents - Rock land complex, extremely stony
		34
		Cryumbrepts - Rock outcrop - Cryaquepts complex
		35
		Bross family - Rock land - Cryumbrepts complex, extremely stony
		36
		Rock outcrop - Cryumbrepts - Cryorthents complex, extremely stony
		37
		Leighcan - Moran families - Cryaqulolls complex, extremely stony
		38
		Moran family - Cryorthents - Leighcan family complex extremely stony
		39
		Moran family - Cryorthents - Rock land complex, extremely stony
		40

Design Drawings/Schematics/Solid Models: While our project mainly involved data analysis and machine learning, we used visualization tools like Tableau or Plotly to create insightful graphical representations of the forest cover types' distribution across different regions and elevations. These visualizations aided in presenting our findings to stakeholders in a more comprehensive manner.

Report Preparation: To prepare the project report, we utilized modern word processing software like Microsoft Word or Google Docs. These tools enabled us to structure and format the report professionally, including sections such as the introduction, methodology, results, and conclusion.

Testing/Characterization/Interpretation/Data Validation: To validate our model's performance, we used cross-validation techniques in Python's scikit-learn library. This allowed us to assess the model's accuracy and generalization capabilities. Furthermore,

42

we employed statistical tests and visualizations to interpret the results and make data-driven decisions.

Version Control: Throughout the implementation phase, we adopted version control systems like Git and hosted our repository on platforms like GitHub or GitLab. This ensured proper version tracking, easy collaboration among team members, and the ability to roll back to previous versions if necessary.

Conclusion: By utilizing these modern tools in analysis, design, reporting, project management, communication, and testing, we successfully implemented the solution for our Forest Cover Type Prediction project. These tools not only streamlined our workflow but also enabled us to present our findings effectively and efficiently. The comprehensive use of modern tools enhanced the overall project's success and contributed to the accuracy and reliability of our forest cover type predictions.

39
CHAPTER 5.

CONCLUSION AND FUTURE WORK

5.1. Conclusion

In this Forest Cover Type Prediction project, we set out to develop a robust machine learning model to accurately classify different forest cover types based on various cartographic variables. Our primary goal was to achieve high prediction accuracy and provide valuable insights to support forest management and conservation efforts.

Expected Results/Outcome: Through rigorous data analysis, feature engineering, and model training, we successfully built a K-Nearest Neighbors (KNN) classifier capable of accurately predicting forest cover types for a given observation. The expected outcome was a well-performing model with a high accuracy rate on the validation set.

Deviation from Expected Results and Reason for the Same: During the model development process, we encountered some deviations from our initial expectations. Despite our efforts in feature selection and optimization, we found that the accuracy of the KNN classifier was slightly lower than initially anticipated. The reason for this deviation could be attributed to the presence of complex interactions between some cartographic variables and the forest cover types, making it challenging for the KNN algorithm to identify subtle patterns accurately.

16 However, it is important to note that the achieved accuracy was still commendable, and the model demonstrated promising performance in classifying most forest cover types. The deviation from the expected results also highlights the complexity of the task and the need for further exploration and experimentation to improve model performance.

Insights and Recommendations: Throughout the project, we gained valuable insights into the relationships between cartographic variables and forest cover types. We identified critical features that significantly influenced the predictions and contributed to the classification accuracy. These insights can inform forest management strategies, helping to make more informed decisions in preserving and conserving different forest ecosystems.

To further enhance the model's performance, we recommend exploring more advanced

17

machine learning algorithms, such as Random Forests or Gradient Boosting Machines, which may capture complex interactions more effectively. Additionally, gathering additional environmental and climate-related data might provide more comprehensive information, leading to better predictions and a deeper understanding of forest dynamics.

Conclusion: In conclusion, our Forest Cover Type Prediction project has achieved significant milestones in accurately predicting forest cover types using a K-Nearest Neighbors (KNN) classifier. Although we encountered slight deviations from our initial expectations, the model still exhibits promising performance, providing valuable insights for forest management and conservation practices.

The project's success opens up opportunities for further research and improvement in forest cover type prediction, contributing to the sustainable management and preservation of our precious forest resources. As we continue to explore advanced machine learning techniques and gather more data, we can envision even more powerful models to support decision-making processes and environmental conservation efforts in the future.

5.2. Future work

While our Forest Cover Type Prediction project has achieved commendable results, there are several avenues for future work and improvements. As we look ahead, we envision the following areas for enhancing the solution and extending its applicability:

1. Fine-tuning Model Hyperparameters: To further improve the KNN classifier's performance, we recommend conducting an extensive hyperparameter tuning process. This involves systematically searching for the best combination of hyperparameters, such as the number of neighbors (K), distance metrics, and data preprocessing techniques. Fine-tuning the model can lead to significant improvements in accuracy and generalization.

2. Exploring Advanced Machine Learning Algorithms: While KNN serves as a solid

baseline, exploring more advanced machine learning algorithms like Random Forests, Gradient Boosting Machines, or Support Vector Machines could potentially yield higher prediction accuracy. These algorithms have the capacity to capture complex relationships between features and forest cover types, offering more robust and reliable predictions.

3. Incorporating Remote Sensing Data: To enrich the feature set and improve prediction accuracy, integrating remote sensing data could be highly beneficial. Remote sensing datasets, such as satellite imagery and LiDAR data, provide valuable information about forest health, density, and canopy structure, which can complement the cartographic variables and lead to a more comprehensive and accurate model.

4. Addressing Class Imbalance: The forest cover type classes may exhibit imbalanced distributions in the dataset, leading to biased model performance. Implementing techniques such as oversampling, undersampling, or using class-weighted approaches can address class imbalance issues and improve the classifier's ability to handle rare forest cover types.

5. Spatial Analysis: Introducing spatial analysis techniques can capture spatial autocorrelation and dependencies among neighboring observations. Geospatial features like distance to nearest neighbors or spatial autocorrelation measures could enhance the model's spatial awareness and better reflect real-world ecological patterns.

6. Multi-Model Ensemble: Building an ensemble of multiple models can combine the strengths of various algorithms and mitigate their individual weaknesses. A weighted average or majority voting system among different classifiers, including KNN and other advanced models, could lead to a more robust and accurate prediction system.

7. Real-Time Prediction Applications: Extending the solution to real-time prediction applications could be valuable for forest monitoring and early warning systems. Integrating the model into a web or mobile application would enable users to input

geographic coordinates or environmental data to receive instant forest cover type predictions and valuable insights on the go.

8. Transfer Learning for Different Regions: Applying transfer learning techniques could make the model adaptable to different regions with varying forest ecosystems. By fine-tuning the existing model on new datasets from different geographical areas, we can avoid the need for building separate models for each region, thus saving computational resources.

Incorporating these future work directions will contribute to a more comprehensive and accurate Forest Cover Type Prediction system. As technology advances and more data becomes available, our solution can continue to evolve, supporting sustainable forest management practices and conservation efforts on a larger scale.

1 REFERENCES

- [1] H. Franco-Lopez, A. R. Ek, and M. E. Bauer, “Estimation and mapping of forest stand density, volume, and cover type using the k-nearest neighbors method,” *Remote Sensing of Environment*, vol. 77, no. 3, pp. 251 – 274, 2001.
- [2] T. E. Avery and H. E. Burkhart, *Forest measurements*. Waveland Press, 2015.
- [3] B. T. Wilson, A. J. Lister, and R. I. Riemann, “A nearest-neighbor imputation approach to mapping tree species over large areas using forest inventory plots and moderate resolution raster data,” *Forest Ecology and Management*, vol. 271, pp. 182 – 198, 2012.
- [4] A. K. Gjertsen, “Accuracy of forest mapping based on landsat TM data and a knn-based method,” *Remote Sensing of Environment*, vol. 110, no. 4, pp. 420 – 430, 2007.
- [5] D. Lowd and P. Domingos, “Naive bayes models for probability estimation,” in *Proceedings of the 22Nd International Conference on Machine Learning*, ser. ICML ’05. New York, NY, USA: ACM, 2005, pp. 529–536.
- [6] V. Rodriguez-Galiano, B. Ghimire, J. Rogan, M. Chica-Olmo, and J. Rigol-Sanchez, “An assessment of the effectiveness of a random forest classifier for land-cover classification,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 67, pp. 93 – 104, 2012.
- [7] D. N. A. Asuncion, “UCI machine learning repository,” 2007.
- [8] I. D. Moore, P. Gessler, G. Nielsen, and G. Peterson, “Soil attribute prediction using terrain analysis,” *Soil Science Society of America Journal*, vol. 57, no. 2, pp. 443–452, 1993.
- [9] D. Koller and M. Sahami, “Toward optimal feature selection,” 1996.
- [10] M. Dash and H. Liu, “Feature selection for classification,” *Intelligent data analysis*, vol. 1, no. 3, pp. 131–156, 1997.
- [11] T. M. Mitchell, “Machine learning,” 1997.
- [12] R. S. Michalski, J. G. Carbonell, and T. M. Mitchell, *Machine learning: An artificial intelligence approach*. Springer Science & Business Media, 2013.
- [13] J. R. Quinlan, “Induction of decision trees,” *Machine learning*, vol. 1, no. 1, pp. 81–106, 1986.
- [14] R. Kohavi et al., “A study of cross-validation and bootstrap for accuracy estimation and model selection,” in *Ijcai*, vol. 14, no. 2, 1995, pp. 1137–1145.
- [15] I. Rish, “An empirical study of the naive bayes classifier,” in *IJCAI 2001 workshop on empirical methods in artificial intelligence*, vol. 3, no. 22. IBM New York, 2001, pp. 41–46.

- [16] P. Domingos and M. Pazzani, "On the optimality of the simple bayesian classifier under zero-one loss," *Machine learning*, vol. 29, no. 2-3, pp. 103–130, 1997.
- [17] H. Zhang, "The optimality of naive bayes," *AA*, vol. 1, no. 2, p. 3, 2004.
- [18] K.-M. Schneider, "Techniques for improving the performance of naïve bayes for text classification," in *International Conference on Intelligent Text Processing and Computational Linguistics*. Springer, 2005, pp. 682–693.
- [19] K. Torkkola, "Linear discriminant analysis in document classification," in *IEEE ICDM Workshop on Text Mining*. Citeseer, 2001, pp. 800–806.
- [20] W. A. Chaovallwongse, Y.-J. Fan, and R. C. Sachdeo, "On the time series k-nearest neighbor classification of abnormal brain activity," *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans*, vol. 37, no. 6, pp. 1005–1016, 2007.
- [21] L. Breiman, "Random forests," *Machine learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [22] L. Breiman, "Bagging predictors," *Machine learning*, vol. 24, no. 2, pp. 123–140, 1996.
- [23] M. J. Islam, Q. J. Wu, M. Ahmadi, and M. A. Sid-Ahmed, "Investigating the performance of naive-bayes classifiers and k-nearest neighbor classifiers," in *Convergence Information Technology, 2007. International Conference on*. IEEE, 2007, pp. 1541–1546.
- [24] S. Gianvecchio, M. Xie, Z. Wu, and H. Wang, "Humans and bots in internet chat: measurement, analysis, and automated classification," *IEEE/ACM Transactions on Networking (TON)*, vol. 19, no. 5, pp. 1557–1571, 2011.
- [25] I. Mani and I. Zhang, "knn approach to unbalanced data distributions: a case study involving information extraction," in *Proceedings of workshop on learning from imbalanced datasets*, 2003.
- [26] B. Frenay and M. Verleysen, "Classification in the presence of label 'noise': a survey," *IEEE transactions on neural networks and learning systems*, vol. 25, no. 5, pp. 845–869, 2014.

APPENDIX

Plagiarism Report:

USER MANUAL

(Complete step by step instructions along with pictures necessary to run the project)

Forest Cover Type Prediction Envisioning Nature's Canopy

ORIGINALITY REPORT



PRIMARY SOURCES

1	repository.usp.ac.fj Internet Source	9%
2	www.coursehero.com Internet Source	2%
3	Submitted to 79920 Student Paper	1%
4	www.irjmets.com Internet Source	1%
5	Submitted to Chandigarh University Student Paper	1%
6	Usman, Muhammad, Russel Pears, and A.C.M. Fong. "Discovering diverse association rules from multidimensional schema", Expert Systems with Applications, 2013. Publication	1%
7	www.ijraset.com Internet Source	1%
8	digibug.ugr.es Internet Source	<1%

9	Submitted to Southern New Hampshire University - Continuing Education Student Paper	<1 %
10	ijsrset.com Internet Source	<1 %
11	Submitted to Liverpool John Moores University Student Paper	<1 %
12	Submitted to GNA University Student Paper	<1 %
13	Submitted to University of London External System Student Paper	<1 %
14	Submitted to Unicaf University Student Paper	<1 %
15	Submitted to University of Southampton Student Paper	<1 %
16	assets.researchsquare.com Internet Source	<1 %
17	ijspt.scholasticahq.com Internet Source	<1 %
18	vital.seals.ac.za:8080 Internet Source	<1 %
19	Submitted to University of Sunderland Student Paper	<1 %

20	groups.psych.northwestern.edu Internet Source	<1 %
21	Submitted to University of Macau Student Paper	<1 %
22	coek.info Internet Source	<1 %
23	link.springer.com Internet Source	<1 %
24	Submitted to CSU, San Jose State University Student Paper	<1 %
25	Submitted to University of North Texas Student Paper	<1 %
26	Submitted to Queen Mary and Westfield College Student Paper	<1 %
27	Rahul R. Kishore, Shalvin S. Narayan, Sunil Lal, Mahmood A. Rashid. "Comparative Accuracy of Different Classification Algorithms for Forest Cover Type Prediction", 2016 3rd Asia- Pacific World Congress on Computer Science and Engineering (APWC on CSE), 2016 Publication	<1 %
28	Shaban Shataee, Syavash Kalbi, Asghar Fallah, Dieter Pelz. " Forest attribute imputation using machine-learning methods and ASTER data: comparison of -NN, SVR and random	<1 %

forest regression algorithms ", International Journal of Remote Sensing, 2012

Publication

- 29 www.scribd.com <1 %
Internet Source
- 30 Submitted to Caledonian College of Engineering <1 %
Student Paper
- 31 Submitted to University of Greenwich <1 %
Student Paper
- 32 Wang, S.. "Application of hybrid image features for fast and non-invasive classification of raisin", Journal of Food Engineering, 201204 <1 %
Publication
- 33 Cornelius Senf, Patrick Hostert, Sebastian van der Linden. "Using MODIS time series and random forests classification for mapping land use in South-East Asia", 2012 IEEE International Geoscience and Remote Sensing Symposium, 2012 <1 %
Publication
- 34 Submitted to National College of Ireland <1 %
Student Paper
- 35 Pratibha Maurya, Arvind Kumar. "Chapter 5 Performance Assessment of K-Nearest Neighbor Algorithm for Classification of <1 %

Forest Cover Type", Springer Science and
Business Media LLC, 2022

Publication

-
- 36 Rydberg, Patrik. "Reactivity-Based Approaches and Machine Learning Methods for Predicting the Sites of Cytochrome P450-Mediated Metabolism", *Methods and Principles in Medicinal Chemistry*, 2014. <1 %
Publication
-
- 37 Submitted to University College London <1 %
Student Paper
-
- 38 api.deepai.org <1 %
Internet Source
-
- 39 digitalcommons.usf.edu <1 %
Internet Source
-
- 40 mafiadoc.com <1 %
Internet Source
-
- 41 scholar.psu.edu <1 %
Internet Source
-
- 42 www.arrsd.org <1 %
Internet Source
-
- 43 www.econstor.eu <1 %
Internet Source
-
- 44 Advances in Intelligent Systems and Computing, 2016. <1 %
Publication

- 45 Anooshiravan Sharabiani, Houshang Darabi, Ashkan Rezaei, Samuel Harford, Hereford Johnson, Fazle Karim. "Efficient Classification of Long Time Series by 3-D Dynamic Time Warping", IEEE Transactions on Systems, Man, and Cybernetics: Systems, 2017
Publication <1 %
- 46 Benny Y. M. Fung. "Classification of heterogeneous gene expression data", ACM SIGKDD Explorations Newsletter, 12/1/2003
Publication <1 %
- 47 Submitted to University of Abertay Dundee
Student Paper <1 %
- 48 www.ijrte.org
Internet Source <1 %
- 49 "Advances in Decision Sciences, Image Processing, Security and Computer Vision", Springer Science and Business Media LLC, 2020
Publication <1 %
- 50 Jock A. Blackard, Denis J. Dean. "Comparative accuracies of artificial neural networks and discriminant analysis in predicting forest cover types from cartographic variables", Computers and Electronics in Agriculture, 1999
Publication <1 %

- 51 Olga Shulayeva, Advaith Siddharthan, Adam Wyner. "Recognizing cited facts and principles in legal judgements", Artificial Intelligence and Law, 2017 <1 %
- Publication
-
- 52 Sam Durairaj, Joanofarc Xavier, Sanjib Kumar Patnaik, Rames C Panda. "Deep Learning Based System Identification and Nonlinear Model Predictive Control of pH Neutralization Process", Industrial & Engineering Chemistry Research, 2023 <1 %
- Publication
-
- 53 Yadi Shen, Yingchao Dong, Xiaoxia Han, Jinde Wu, Kun Xue, Meizhu Jin, Gang Xie, Xinying Xu. "Prediction model for methanation reaction conditions based on a state transition simulated annealing algorithm optimized extreme learning machine", International Journal of Hydrogen Energy, 2023 <1 %
- Publication
-
- 54 core.ac.uk <1 %
- Internet Source
-
- 55 earchive.tpu.ru <1 %
- Internet Source
-
- 56 ebin.pub <1 %
- Internet Source
-
- ediss.uni-goettingen.de

57	Internet Source	<1 %
58	hal-univ-bourgogne.archives-ouvertes.fr Internet Source	<1 %
59	ndl.ethernet.edu.et Internet Source	<1 %
60	nebula.wsimg.com Internet Source	<1 %
61	www.ijitee.org Internet Source	<1 %
62	www.mdpi.com Internet Source	<1 %
63	www.researchgate.net Internet Source	<1 %
64	Agrata Gupta, N. Arulkumar. "chapter 12 An Exploratory Study of Python's Role in the Advancement of Cryptocurrency and Blockchain Ecosystems", IGI Global, 2023 Publication	<1 %
65	M. Shyamala Devi, Shefali Dewangan, Satwat Kumar Ambastha, Anjali Jaiswal, Sairam Kondapalli. "Recognition of Forest Fire Spruce Type Tagging using Machine Learning Classification", International Journal of Recent Technology and Engineering (IJRTE), 2019 Publication	<1 %

66

"Data Stream Mining & Processing", Springer
Science and Business Media LLC, 2020

Publication

<1 %

Exclude quotes Off
Exclude bibliography Off

Exclude matches Off

Forest Cover Type Prediction Envisioning Nature's Canopy

PAGE 1

PAGE 2

PAGE 3

PAGE 4

PAGE 5

PAGE 6

PAGE 7

PAGE 8

PAGE 9

PAGE 10

PAGE 11

PAGE 12

PAGE 13

PAGE 14

PAGE 15

PAGE 16

PAGE 17

PAGE 18

PAGE 19

PAGE 20

PAGE 21

PAGE 22

PAGE 23

PAGE 24

PAGE 25

PAGE 26

PAGE 27

PAGE 28

PAGE 29

PAGE 30

PAGE 31

PAGE 32

PAGE 33

PAGE 34

PAGE 35

PAGE 36

PAGE 37

PAGE 38

PAGE 39

PAGE 40

PAGE 41

PAGE 42

PAGE 43

PAGE 44

PAGE 45
