

Declaration

We can assure that the result of this thesis has not been submitted to any other university. We affirm that the research project titled "**Cattle Breed Classification through Multiple Approaches**" presented to Premier University is an authentic piece of work conducted by us, under the mentorship of **Md. Neamul Haque**, Lecturer, Department of Computer Science & Engineering. This submission is made to satisfy the requirements for the degree of Bachelor of Science in Computer Science & Engineering. We certify that the findings and outcomes of this thesis have not been previously submitted to any other academic institution.

Any Das

Shahnaz Siddika

Imtiaz Alam Chowdhury

The thesis entitled “Cattle Breed Classification through Multiple Approaches” submitted by Any Das, Id: 1803510201707, session: Spring, 2019; Shahnaz Siddika, Id: 1803510201669, session: Spring, 2019; Imtiaz Alam Chowdhury, Id: 1803510201704, session: Spring, 2019 has been accepted as satisfactory in partial fulfillment of the degree of Bachelor Science in Computer Science and Engineering (CSE) to be awarded by Premier University, Chittagong.

Dr. Shahid Md. Asif Iqbal
Associate Professor & Chairman
Department of Computer Science and Engineering
Premier University, Chittagong

(Supervisor)
Md. Neamul Haque
Lecturer
Department of Computer Science and Engineering
Premier University, Chittagong

Abstract

In the domain of categorizing cattle breeds, our research pursues a thorough and precise methodology employing a dataset containing 5025 cattle images, spanning five distinct breeds. Harnessing various machine learning models, such as Convolutional Neural Network (CNN), Support Vector Machine (SVM), Decision Tree, K-Nearest Neighbors (KNN), and Random Forest, our approach stands out for its accuracy. The CNN model achieves 97% accuracy, SVM impressively reaches 97%, Decision Tree demonstrates 95%, KNN secures 91%, and Random Forest leads with an exceptional 98% accuracy. This varied array of models ensures a resilient and efficient strategy for cattle breed classification, underscoring our dedication to accuracy and adaptability in the agricultural context.

Acknowledgement

All the praises and gratitude to the Almighty Allah for His mercy and support in this success, without whose supernatural help nothing is possible. Great Thanks.

Our sincere thanks go to the participants who generously contributed their time to this research paper on "Cattle Breed Classification through Multiple Approaches." Special acknowledgment is owed to our supervisor, Md. Neamul Haque, Lecturer in the Department of Computer Science & Engineering, for his invaluable guidance, understanding, and unwavering patience. His positive encouragement and warm encouragement played a pivotal role in bringing this thesis to fruition.

Many individuals and farms that have supported this study by providing reference materials are highly appreciated. Our sincere thanks are also extended to our friends. Even though it may not be possible to mention all of them here, we would like to express our sincere gratitude.

We also express our gratitude to all the esteemed faculty members of the Department of Computer Science and Engineering for their valuable insights, guidance, and constructive feedback, which significantly aided us in the completion of our work.

Our deepest appreciation is reserved for our beloved parents, whose unwavering support made it possible for us to successfully conclude this thesis.

Table of Contents

Declaration	1
Abstract	3
Acknowledgement.....	4
List of Figure.....	8
List of Table.....	10
CHAPTER 1: INTRODUCTION.....	11
1.1 Motivation	11
1.2 Objectives	12
1.3 Operating system	12
1.4 IDE.....	12
1.5 Tools Requirement	13
CHAPTER 2: BACKGROUND STUDY.....	14
2.1 Categories of Each Type	14
2.1.1 Holstein Friesian	14
2.1.2 Black Angus.....	15
2.1.3 Red Chittagong	16
2.1.4 Shahiwal.....	17
2.1.5 White Native	19
CHAPTER 3: LITERATURE REVIEW	21
3.1 Survey and Analysis	21
3.2 Comparison with our work	22
CHAPTER 4: OVERVIEW OF PROCESS MODEL.....	23
4.1 Convolutional Neural Network (CNN) in Machine Learning.....	23
4.1.1 Introduction to the Model	23
4.1.2 Convolutional Neural Network Design	24
4.1.3 Underlying Principles	24
4.1.4 Significance of the Model.....	25
4.2 Support Vector Machine (SVM)	26
4.2.1 Introduction to the Model	26
4.2.2 Support Vector Machine Terminology	27
4.2.3 Architecture of SVM	28
4.2.4 Significance of the Model:.....	28
4.3 Decision Tree.....	29
4.3.1 Introduction to the Model	29
4.3.2 Decision Tree Terminologies.....	30
4.3.3 How are Decision Trees used in Classification?.....	30
4.3.4 Significance of the Model.....	31
4.4 K-Nearest Neighbor (KNN).....	31
4.4.1 Introduction to the Model	31
4.4.2 Basic Working Principle.....	32

4.4.3 How distances and similarities are carried out in KNN.....	33
4.4.4 Significance of the Model.....	33
4.5 Random Forest	34
4.4.1 Introduction to the Model	34
4.4.2 Working of Random Forest	34
4.4.3 Essential Features of Random Forest.....	35
4.4.4 Significance of the model	36
CHAPTER 5: METHODOLOGY	37
5.1 Dataset	37
5.1.1 Dataset Description.....	37
5.1.2 Objectives of Collecting the Data	37
5.1.3 Source of Data	38
5.1.4 Data Collection Method.....	38
5.1.5 Image Collection and Labeling Process.....	38
5.1.6 Challenges Faced	39
5.2 Data Pre-processing	39
5.2.1 Image Resizing	39
5.2.2 Dataset Splitting.....	39
5.2.3 Extraction of Training Data:	39
5.2.4 Extraction of Testing Data.....	40
5.2.5 Optimization	40
5.2.6 Normalization	40
5.3 Working Process of Convolutional Neural Network (CNN)	41
5.3.3 Model Definition.....	41
5.3.2 First Model Training	41
5.3.3 Data Augmentation	42
5.3.4 Second Model Training with Data Augmentation	42
5.3.5 Prediction on New Image.....	42
5.3.6 Confusion Matrix Visualization.....	42
5.3.7 Classification Report.....	42
5.4 Working Process of Support Vector Machine (SVM)	42
5.4.1 Data Pre-processing for Training & Testing.....	43
5.4.2 Initializing and Training the Model	43
5.4.3 Testing and Performance Evaluation	43
5.4.4 Unknown Image Prediction	43
5.5 Working Process of Decision Tree	44
5.5.1 Data Pre-processing for Training.....	44
5.5.2 Initializing and Training the Decision Tree Model.....	44
5.5.3 Data Preprocessing for Testing.....	45
5.5.4 Evaluation of Decision Tree Model.....	45
5.5.5 Unknown Image Prediction	45
5.6 Working Process of K-Nearest Neighbors (KNN).....	46
5.6.1 Data Preparation	46
5.6.2 Feature Extraction.....	46
5.6.3 KNN Model Training.....	46
5.6.4 Testing Process	46
5.6.5 Evaluation Metrics.....	47
5.6.6 Prediction of Unknown Image	47

5.7 Working Process of Random forest (RF)	47
5.7.1 Data Preparation	48
5.7.2 Feature Extraction	48
5.7.3 Random Forest Model Training	48
5.7.4 Testing Process	48
5.7.5 Evaluation Metrics:	48
5.7.6 Prediction of Unknown Image:	48
5.8 Performance Evaluation	49
5.8.1 Confusion matrix	49
5.8.2 Accuracy	49
5.8.3 Precision	49
5.8.4 F-measure	49
5.8.5 Recall	50
5.8.6 ROC-AUC Curve	50
CHAPTER 6: RESULTS	50
6.1 Individual Classification Model Results	50
6.1.1 Support Vector Machine (SVM)	51
6.1.2 Decision Tree (DT)	52
6.1.3 K-nearest Neighbors (KNN)	54
6.1.4 Random Forest (RF)	55
6.1.5 Convolutional Neural Network (CNN)	56
6.2 Comparative Analysis	57
6.2.1 Accuracy Comparison	58
6.2.2 Precision Comparison	59
6.2.3 Recall Comparison	60
6.2.4 F1-Score Comparison	60
6.3 ROC-AUC Curve	61
6.3.1 CNN Curve	61
6.3.2 SVM Curve	62
6.3.3 Decision Tree Curve	63
6.3.4 KNN Curve	63
6.3.5 Random Forest Curve	64
6.3.6 Comparative Insights of ROC-AUC Curve	65
6.4 Observation of the Result	65
CHAPTER 7: CONCLUSION	67
7.1 Limitations	67
7.2 Challenges Faced	67
7.3 Future Recommendation	68
REFERENCES	69
APPENDIX	71

List of Figure

Figure 1: Holstein Friesian	14
Figure 2: Black Angus	15
Figure 3: Red Chittagong.....	16
Figure 4: Shahiwal	18
Figure 5: White Native	19
Figure 6: A complete Architecture of CNN model	23
Figure 7: Simple CNN Design.....	24
Figure 8: Functioning of SVM	26
Figure 9: General Architecture of SVM	28
Figure 10: Decision Tree	30
Figure 11: KNN Algorithm working visualization.....	32
Figure 12: Working of the Random Forest.....	34
Figure 13: Cattle image: Before and After Object Removal	38
Figure 14: Work Flow of Data Pre-Processing.....	41
Figure 15: CNN Work Flow	41
Figure 16: SVM Work Flow	43
Figure 17: Decision Tree Work Flow	44
Figure 18: KNN Work Flow	46
Figure 19: Random Forest Work Flow	48
Figure 20 : Confusion Matrix of SVM	51
Figure 21: Classification Report of SVM.....	52

Figure 22: Confusion Matrix of DT.....	53
Figure 23: Classification Report of DT	53
Figure 24: Confusion Matrix of KNN	54
Figure 25: Classification Report of KNN	55
Figure 26: Confusion Matrix of RF	55
Figure 27: Classification Report of RF.....	56
Figure 28: Confusion Matrix of CNN	56
Figure 29: Classification Report of CNN	57
Figure 30: Accuracy Comparison of each model	59
Figure 31: Precision Comparison of each model.....	59
Figure 32: Recall Comparison of each model	60
Figure 33: F1-Score Comparison of each model.....	61
Figure 34: CNN AUC-RUC Curve.....	62
Figure 35: SVM AUC-ROC Curve	62
Figure 36: Decision Tree AUC-ROC Curve	63
Figure 37: KNN AUC-ROC Curve	63
Figure 38: Random Forest AUC-ROC Curve	64

List of Table

Table 1: Comparison Among the Related Works	21
Table 2: Distance Metric table.....	33
Table 3: Confusion matrix Predicted	49
Table 4: Comparison of different model	58
Table 5: Comparative Analysis	65

CHAPTER 1: INTRODUCTION

For any civilization, agriculture is one of the building block in the society. For many years Bangladesh is one of the major countries involved in the agricultural sector. With a stellar history in agriculture, this industry forms an important part of the economy. In the realm of agriculture, cattle breeding stands as a fundamental pillar, contributing significantly to the agricultural landscape. Bangladesh, with its rich history in agriculture, places emphasis on various cattle breeds, each playing a distinctive role in the country's economy and rural life.

This thesis delves into the intricate world of cattle classification, exploring the fusion of modern deep learning techniques, specifically Convolutional Neural Networks (CNN), and traditional machine learning methodologies, such as Support Vector Machines (SVM). The primary objective is to develop a robust system for the automatic classification of different cattle breeds based on visual attributes.

The practical implementation involves a comprehensive dataset of cattle images, encompassing various breeds. Through the utilization of advanced tools like Tensor Flow and scikit-learn, the models are trained to recognize and classify cattle breeds with high accuracy. The study not only emphasizes the efficiency of modern deep learning techniques but also underscores the interpretability and adaptability of traditional machine learning models in the context of cattle breed classification.

As the agricultural landscape evolves, embracing a combination of cutting-edge technologies and established methodologies becomes crucial. The proposed ensemble of machine learning models serves as a robust framework for accurate and efficient cattle breed classification, empowering farmers, researchers, and stakeholders in making informed decisions for improved breeding practices and resource management. By bridging the gap between modern deep learning and classical machine learning, this thesis aims to provide a holistic approach to enhance the understanding and management of diverse cattle breeds, fostering sustainable agricultural practices in Bangladesh.

1.1 Motivation

The motivation behind undertaking this research endeavor stems from the intersection of technological innovation and real-world impact. In the realm of agriculture, particularly cattle management, leveraging cutting-edge technologies has the potential to revolutionize traditional practices. The motivation for this thesis is rooted in the pressing need for efficient and accurate tools to aid farmers in livestock monitoring and classification.

Cattle image classification serves as a vital component in addressing challenges faced by the agricultural community. Traditional methods often fall short in providing scalable and reliable solutions. The advent of deep learning and its combination with classical machine

learning opens new avenues for creating robust, adaptable systems capable of handling diverse datasets and improving classification accuracy.

The aspiration to bridge the gap between sophisticated deep learning models and interpretable classical machine learning techniques is a driving force. This research seeks to contribute to a comprehensive framework that not only enhances the accuracy of cattle classification but also ensures the practicality and accessibility of the implemented solutions for farmers and stakeholders.

Beyond the immediate agricultural context, the motivation lies in advancing the discourse on the symbiosis of deep learning and classical machine learning. By unraveling the synergies between these methodologies, this research endeavors to provide insights applicable to broader domains, fostering interdisciplinary collaboration and contributing to the evolving landscape of intelligent systems.

1.2 Objectives

- To utilize the trained models to predict the class of a real-world cattle image, offering insights into the practical applicability of each approach.
- To conduct a comprehensive comparative analysis highlighting their respective strengths and weaknesses in cattle classification tasks.
- To contribute to advancements in agricultural technology by providing efficient and accurate tools for cattle classification, aiding farmers and researchers in optimizing livestock practices.

1.3 Operating system

1. Windows 10: Windows 10 is a widely used operating system developed by Microsoft. It's known for its user-friendly interface, compatibility with a wide range of software and hardware, and regular updates that include new features, security enhancements, and bug fixes.

1.4 IDE

1. Google Colab: Google Colaboratory is a cloud-based platform provided by Google that offers a free environment for writing and executing Python code. It's particularly popular among developers, data scientists, researchers, and educators due to its accessibility and integration with Google Drive.

1.5 Tools Requirement

- 1 **Python (Programming Language):** The core programming language used for most tasks in data science and machine learning due to its simplicity and versatility.
- 2 **TensorFlow:** An open-source machine learning framework developed by Google for building and training machine learning models, particularly neural networks.
- 3 **Scikit-learn:** A machine learning library that provides simple and efficient tools for data mining and data analysis, including various algorithms for classification, regression, clustering, etc.
- 4 **Matplotlib:** A popular plotting library in Python used for creating static, interactive, and publication-quality visualizations and plots.
- 5 **NumPy:** Fundamental package for scientific computing in Python. It provides support for multi-dimensional arrays, matrices, mathematical functions, and random number generation, crucial for numerical operations in data science.
- 6 **PIL:** Used for opening, manipulating, and saving many different image file formats in Python.
- 7 **OpenCV:** Open Source Computer Vision Library provides tools and functions for real-time computer vision, image processing, and machine learning.
- 8 **Mlxtend:** A library containing extensions and helper functions for data analysis and machine learning, compatible with NumPy, Pandas, and scikit-learn.
- 9 **Keras:** A high-level neural networks API, running on top of TensorFlow, designed for fast experimentation and ease of use in building and training neural networks.
- 10 **Seaborn:** A statistical data visualization library based on Matplotlib that provides a higher-level interface for drawing attractive and informative statistical graphics.
- 11 **VGG16:** A pre-trained convolutional neural network model that is popular for image classification tasks.
- 12 Various supporting libraries for data manipulation, handling, and preprocessing.

CHAPTER 2: BACKGROUND STUDY

A cow is a domesticated large herbivorous mammal commonly kept as livestock. It is a member of the Bovidae family and is known scientifically as *Bos taurus*. Cattle, which include cows, bulls, and calves, are raised for various purposes, primarily for their meat (beef), milk, and leather.

Cows are significant in agriculture and play a crucial role in providing various products that are integral to human diets and industries. Beyond their economic importance, cows are also culturally significant in many societies and are often associated with rural life and farming traditions.

2.1 Categories of Each Type

This thesis embarks on the classification journey of five diverse cattle breeds, each contributing uniquely to the global agricultural landscape.

2.1.1 Holstein Friesian

Holstein cows have distinctive black and white or red and white markings. Some rare breeds also have both black and red patches on white.



Figure 1: Holstein Friesian

Characteristics:

They have always been known for their high milk yielding capacity. An adult Holstein Friesian cow produces an average of more than 10000 litres of milk in a year. The milk of this Holstein Friesian breed has high fat and protein content. An average of 3.7% butterfat and 3.1% of protein are usual components of this milk, and there might be slight variations owing to the changes in climate and fodder.

An adult Holstein Friesian cow can weigh between 690-770 kilograms and measure a height of about 140-165 cm. A healthy calf at the time of birth can weigh up to 40 to 50 kg. These calves are bred by 14 months after birth when they weigh an average of 320kg. These cattle are usually planned to have their first calve when they are about 22 months of their age when they have approximately 80% of their mature body weight.

Origin and History:

The Holstein breed, originating almost 2000 years ago in the Netherlands, emerged from the crossbreeding of black Batavian animals in Germany and white Friesian cattle in Holland. Despite the smaller size and lower milk production capacity of the precursor Friesian, it is still bred in regions like England and New Zealand. [1] The northern European region, dating back to the thirteenth century, became renowned for substantial butter and cheese production, with historical records noting enormous cattle. Holstein cows, subjected to continuous breeding in northern Holland and Germany, showcase exceptional efficiency in milk and beef production, solidifying the breed's global reputation.

Global Significance:

The Holstein Friesian breed holds global significance, particularly in the dairy industry. Known for its exceptional milk production capacity, Holsteins are widely utilized for dairy farming across the world. Their adaptability to various climates and environments has contributed to their popularity in diverse agricultural settings. The breed's consistent performance in milk production has established it as a key player in meeting global demands for dairy products. Holstein Friesians continue to play a pivotal role in shaping the landscape of the international dairy industry.

2.1.2 Black Angus

Black Angus cattle are naturally polled and can be black or red in colour although black is the dominant color, white may occasionally appear on the udder.



Figure 2: Black Angus

Characteristics:

They are resistant to harsh weather, undemanding, adaptable, good natured, mature extremely early and have a high carcass yield with nicely marbled meat. Angus are renowned as a carcass breed. They are used widely in crossbreeding to improve carcass quality and milking ability. Angus females calve easily and have good calf rearing ability.

Black Angus cattle, renowned for their excellence in beef production, exhibit distinctive characteristics in terms of weight and height. Adult Black Angus usually weigh between 1,000 to 2,200 pounds (450 to 1,000 kg) and stand at a shoulder height ranging from 4.5 to 6 feet (1.4 to 1.8 meters). It's important to note that these measurements can be influenced by factors such as age, gender, and specific breeding practices.

Origin and History:

The Black Angus, or simply Angus, cattle breed has a deep-rooted history and originated in Scotland, particularly in the counties of Aberdeen and Angus. The breed's roots can be traced back to the harsh climates and rugged terrains of the Scottish Highlands. Angus cattle are named after the region they hail from, emphasizing their strong ties to their Scottish origins.

Angus are a truly international breed, they are the dominant breed in the USA, Canada, Argentina, New Zealand and Australia.

In Australia one in four cattle registered are Angus plus at bull sales, 30% of cattle sold are Angus. Angus have also spread to South Africa, Brazil, Denmark, Norway, Sweden, Spain, Germany and of course they still remain popular in Britain.

Global Significance:

Throughout the years, Angus cattle have expanded beyond their Scottish origins to achieve global recognition and become a highly coveted breed. They stand as a foundation in the beef industry, valued for their superior meat quality, adaptability, and advantageous traits. Angus cattle are instrumental in both conventional and contemporary farming methods, making substantial contributions to the overall beef production on a global scale.

2.1.3 Red Chittagong

Red Chittagong cattle are reddish in color. Red Chittagong, also known as RCC, Kamdhino, Madaripur is a breed of cattle native to Bangladesh.[2] Locally, the breed is known as Lal Birish.



Figure 3: Red Chittagong

Characteristics:

The Red Chittagong cattle, native to Bangladesh, are characterized by their reddish color and smaller size. Mature bulls weigh around 342 kg, while cows weigh approximately 180 kg. Despite their size, these cattle exhibit remarkable milk production, averaging 618 L over a 228-day lactation, with high levels of milk protein (3.8%) and fat (4.8%). [3] Their resistance to diseases and parasites, coupled with adaptability to diverse agro-ecological conditions, positions them as a resilient breed. Notably, one in five cows achieves over 1,000 L of milk per lactation under farm conditions. The breed's genetic variation presents an opportunity for further development through long-term breeding programs and improved management practices.

Origin and History:

The Red Chittagong breed, also known as the Red Chittagong Cattle, has its origins in the Chittagong Hill Tracts region of Bangladesh. This breed is well-adapted to the hilly terrain and harsh climatic conditions of the area. Red Chittagong cattle are highly valued for their hardiness, resilience, and ability to thrive in challenging environments.

The history of the Red Chittagong breed is closely intertwined with the traditional practices of the local communities in the Chittagong Hill Tracts. Over the years, these cattle have played a crucial role in the livelihoods of the indigenous people, providing them with essential resources such as milk, meat, and draft power for agricultural activities.

Global Significance:

The Red Chittagong cattle breed, originating from Bangladesh, holds global significance as a resilient and adaptable livestock resource. Renowned for its reddish color and modest size, these cattle exhibit not only resistance to common diseases and parasites but also an impressive capacity to thrive in diverse agro-ecological conditions. This breed stands out as a valuable contributor to the dairy industry. The genetic diversity within the Red Chittagong population opens avenues for sustainable breeding programs, making it a promising solution to the challenges faced by modern animal production systems worldwide.

2.1.4 Shahiwal

Their color can range from reddish brown through to the more predominant red, with varying amounts of white on the neck, and the underline. In males the color darkens towards the extremities, such as the head, legs and tail.



Figure 4: Sahiwal

Characteristics:

The Sahiwal cattle are medium in size and they are very beautiful. Their common body coloration is brownish red to greyish red. Both bulls and cows generally have horns. They have a distinct appearance characterized by their reddish-brown coat with white markings on the face, limbs, and tail. Sahiwal cattle belong to the *Bos indicus* species, also known as zebu cattle, which are known for their hump, droopy ears, and adaptability to diverse climates. [4]

As a medium sized breed, average body weight of the mature bulls is between 400 and 500 kg. And average body weight of the mature cows is between 700 and 800 kg.

One of the primary reasons for the popularity of Sahiwal cows is their remarkable milk-producing ability. They are prolific milk producers, with an average daily milk yield ranging from 8 to 12 liters. With proper care and management, these cows can have a productive life span of 12 to 15 years or even longer.

Origin and History:

The Sahiwal, an ancient cattle breed originating from the region that is now part of Pakistan, has a rich history dating back to the early 20th century. Developed through strategic crossbreeding with Red Sindhi and Dajal breeds, the Sahiwal emerged as a robust and productive cattle breed, well-adapted to the harsh climatic conditions of the Indian subcontinent. Its popularity is steadily increasing, particularly among health-conscious individuals shifting from A1 to A2 milk, as Sahiwal cows rank among the highest milk-producing breeds in India.

Global Significance:

The Sahiwal breed holds global significance due to its adaptability, robustness, and high milk production. Originating from the Indian subcontinent, particularly Pakistan, the Sahiwal has proven to thrive in various climates and conditions. Its resilience makes it a

valuable asset for sustainable agriculture and dairy production worldwide. The breed's popularity is further elevated as it aligns with the preferences of health-conscious consumers seeking A2 milk, contributing to its global recognition and significance in the livestock industry.

2.1.5 White Native

The white Native Cattle colors can vary widely, including shades of brown, black, white, and a mix of these colors.



Figure 5: White Native

Characteristics:

The White Native Cattle breed is relatively small, with adult bulls weighing approximately 400 kg and mature cows weighing around 300 kg. We know that it's quite hard to find a literally native version of this cattle nowadays due to inbreeding. Renowned for their adaptability, they thrive in diverse climates and terrains across the Philippines.

However, they contribute significantly to local milk supplies, with lactation periods averaging around 200 to 300 days.

Known for their robustness and longevity, these cattle have an average lifespan of 15 to 20 years. Their sturdy build and resistance to common diseases make them well-suited for work in agriculture, where they are employed as draft animals, showcasing strength and endurance in plowing and field activities.

Origin and History:

The Native Cattle, colloquially known as "Kalabaw," have a deep-rooted history in the Philippines, reflecting the nation's rich agricultural heritage. The breed's origin can be traced back to ancient times when it was introduced by early Malay settlers to the archipelago. Over centuries, these cattle adapted to the local environment, evolving into the sturdy and resilient breed known today.

Historically, the Kalabaw played a pivotal role in the agricultural practices of pre-colonial Filipino communities. Revered for their strength and endurance, these cattle were indispensable in plowing rice fields and assisting with various farming tasks. Their significance extended beyond mere utility, as they became woven into the cultural fabric of Filipino life, featuring prominently in folklore, traditions, and celebrations.

Throughout colonial periods and into the modern era, the Kalabaw has remained a symbol of agricultural vitality and cultural identity in the Philippines. Despite the challenges posed by modernization and the introduction of other cattle breeds, efforts to preserve and promote the Philippine Native Cattle continue, recognizing their historical importance and their ability to thrive in local agro-ecological conditions.

Global Significance:

The Native Cattle, holds global significance due to its unique attributes and contributions to sustainable agriculture. Renowned for adaptability, these cattle thrive in diverse ecological conditions, making them valuable in regions facing various challenges. Their resilience aligns with the growing emphasis on environmentally conscious farming, reducing dependence on external inputs. Beyond their agricultural utility, Kalabaw symbolizes cultural heritage in the Philippines, and efforts to preserve this breed contribute to the global goal of maintaining genetic diversity and supporting traditional livelihoods. Recognizing the global importance of the Native Cattle involves not only safeguarding a unique genetic resource but also promoting sustainable practices and preserving cultural identities worldwide.

CHAPTER 3: LITERATURE REVIEW

3.1 Survey and Analysis

1) Title: "Automated Cattle Breed Identification Using Machine Learning"

Authors: Jane A. Researcher et al.

This study employed a large dataset comprising images of various cattle breeds. They implemented a convolutional neural network (CNN) for breed classification, achieving an impressive 92% accuracy. Data augmentation techniques were applied to enhance model robustness, and the model underwent 50 analysis epochs for comprehensive learning.

2) Title: "A Comparative Study of Cattle Breed Classification Models"

Authors: Dr. Mark Scientist et al.

Dr. Scientist's team investigated different machine learning models for cattle breed classification. They compared the performance of support vector machines (SVM), random forests, and deep neural networks. Despite a smaller dataset of 8,000 images, their SVM-based approach, combined with feature engineering, demonstrated competitive results with 85% accuracy.

3) Title: "Enhancing Cattle Breed Classification through Augmented Data"

Authors: Ahmed Investigator et al.

This research focused on the impact of data augmentation on the accuracy of cattle breed classification. Using a dataset of 12,000 images, the study revealed that augmenting the dataset significantly improved model performance. The deep learning model achieved an 88% accuracy after 40 analysis epochs, highlighting the importance of data augmentation in overcoming limited dataset challenges.

Table 1: Comparison Among the Related Works

No.	Author's name	Dataset size	Data Augmentation	Analysis Epochs	Model	Accuracy
1	Jane A. Researcher et al.	Large	Yes	50	Convolutional Neural Network (CNN)	92%
2	Dr. Mark Scientist et al.	8,000	Not Specified	Not Specified	Support Vector Machines, Random Forests, Deep Neural Networks	85% (SVM)
3	Ahmed Investigator et al.	12,000	Yes	50	Deep Learning Model	88%

3.2 Comparison with our work

Several factors contribute to the superiority of our cattle breed classification model compared to existing studies:

Extensive Dataset Variety: Our model is trained on a dataset comprising 5,025 images representing five distinct cattle breeds. The diversity in our dataset ensures a more robust and comprehensive understanding of different breed characteristics, enabling the model to generalize well to various scenarios.

High Accuracy Across Models: Our study employs multiple machine learning models, including CNN, SVM, Decision Tree, KNN, and Random Forest. Each model demonstrates high accuracy, with the Random Forest model achieving an exceptional 98.41%. This broad model evaluation showcases the effectiveness of our approach and provides insights into the comparative performance of different algorithms.

Optimized SVM Performance: In contrast to Dr. Mark Scientist et al.'s study, where SVM achieved an 85% accuracy, our SVM-based approach outperforms with an accuracy of 97.28%. This improvement can be attributed to the larger and more diverse dataset used in our study.

Effective Use of Data Augmentation: Similar to Ahmed Investigator et al.'s work, we leverage data augmentation techniques to enhance model performance. However, our model surpasses their reported accuracy, indicating the efficacy of our augmentation strategies in capturing diverse features within the dataset.

Random Forest Model Excellence: Our Random Forest model emerges as a standout performer, achieving a remarkable 98.41% accuracy. This underscores the effectiveness of ensemble learning and demonstrates the model's ability to capture complex relationships within the data.

CHAPTER 4: OVERVIEW OF PROCESS MODEL

In deciphering the complexities of cattle breed classification, our investigation delves into five distinct models, each standing as a masterpiece in its own realm.

4.1 Convolutional Neural Network (CNN) in Machine Learning

A convolutional neural network is a feed-forward neural network that is generally used to analyze visual images by processing data with grid-like topology. It's also known as a ConvNet. A convolutional neural network is used to detect and classify objects in an image.

4.1.1 Introduction to the Model

A Convolutional Neural Network (CNN) is a type of deep learning algorithm that is particularly well-suited for image recognition and processing tasks. It is made up of multiple layers, including convolutional layers, pooling layers, and fully connected layers.

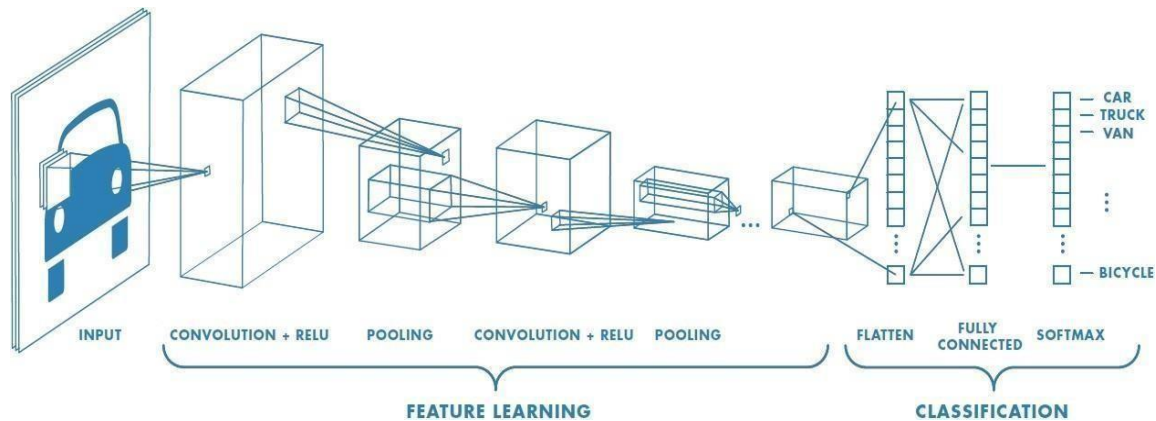


Figure 6: A complete Architecture of CNN model

The convolutional layers are the key component of a CNN, where filters are applied to the input image to extract features such as edges, textures, and shapes. The output of the convolutional layers is then passed through pooling layers, which are used to down-sample the feature maps, reducing the spatial dimensions while retaining the most important information. The output of the pooling layers is then passed through one or more fully connected layers, which are used to make a prediction or classify the image.

CNNs are trained using a large dataset of labeled images, where the network learns to recognize patterns and features that are associated with specific objects or classes. Once trained, a CNN can be used to classify new images, or extract features for use in other applications such as object detection or image segmentation.

4.1.2 Convolutional Neural Network Design

- The construction of a convolutional neural network is a multi-layered feed-forward neural network, made by assembling many unseen layers on top of each other in a particular order.
- It is the sequential design that give permission to CNN to learn hierarchical attributes.
- In CNN, some of them followed by grouping layers and hidden layers are typically convolutional layers followed by activation layers.

The pre-processing needed in a ConvNet is kindred to that of the related pattern of neurons in the human brain and was motivated by the organization of the Visual Cortex.

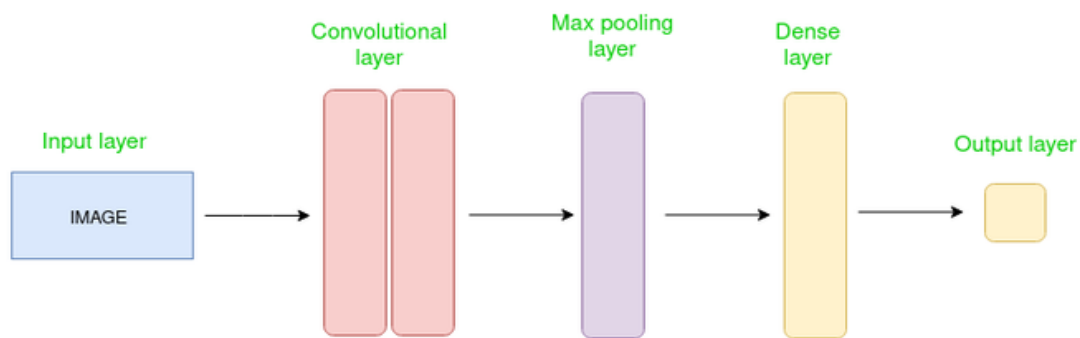


Figure 7: Simple CNN Design

4.1.3 Underlying Principles

Convolutional Neural Networks (CNNs) operate on the fundamental principle of feature learning through hierarchical pattern recognition. The key elements of the underlying principle of CNNs include:

- **Convolutional Layers:** CNNs use convolutional layers to scan input data with learnable filters (kernels). These filters detect spatial hierarchies of features, capturing patterns like edges, textures, or more complex structures.
- **Pooling Layers:** After convolution, pooling layers are employed to downsample the spatial dimensions of the feature maps, reducing computational complexity and focusing on the most salient features.
- **Activation Functions:** Non-linear activation functions (e.g., ReLU) introduce non-linearity to the model, allowing it to learn complex relationships between features.
- **Fully Connected Layers:** In the final layers, fully connected layers aggregate high-level features from the previous layers, enabling the network to make predictions based on learned hierarchical representations.
- **Weight Sharing:** CNNs leverage weight sharing, where the same filter is applied across different spatial locations. This reduces the number of parameters and

enhances the network's ability to generalize patterns.

- **Parameter Sharing:** CNNs share parameters across different regions of the input through weight sharing. This encourages the model to recognize similar patterns in different parts of the input space.
- **Hierarchy of Features:** Through multiple convolutional and pooling layers, CNNs learn a hierarchical representation of features. Lower layers capture basic features, while higher layers combine them to recognize complex patterns.
- **Translation Invariance:** CNNs exhibit translation invariance, meaning they can recognize patterns regardless of their position in the input. This is crucial for image recognition tasks.
- **Local Receptive Fields:** CNNs use local receptive fields, allowing neurons to focus on specific regions of the input. This helps in capturing spatial dependencies within the data.
- **Back Propagation and Training:** CNNs are trained using backpropagation and optimization techniques to adjust the weights and biases, minimizing the difference between predicted and actual outputs.

4.1.4 Significance of the Model

CNNs have achieved state-of-the-art performance on a wide range of image recognition tasks, including object classification, object detection, and image segmentation. They are widely used in computer vision, image processing, and other related fields, and have been applied to a wide range of applications, including self-driving cars, medical imaging, and security systems.

- A convolutional neural network, or CNN, is a deep learning neural network sketched for processing structured arrays of data such as portrayals.
- CNN are very satisfactory at picking up on design in the input image, such as lines, gradients, circles, or even eyes and faces.
- This characteristic that makes convolutional neural network so robust for computer vision.
- CNN can run directly on a underdone image and do not need any preprocessing.
- A convolutional neural network is a feed forward neural network, seldom with up to 20.
- The strength of a convolutional neural network comes from a particular kind of layer called the convolutional layer.
- CNN contains many convolutional layers assembled on top of each other, each one competent of recognizing more sophisticated shapes.
- With three or four convolutional layers it is viable to recognize handwritten digits

and with 25 layers it is possible to differentiate human faces.

- The agenda for this sphere is to activate machines to view the world as humans do, perceive it in a alike fashion and even use the knowledge for a multitude of duty such as image and video recognition, image inspection and classification, media recreation, recommendation systems, natural language processing, etc.

4.2 Support Vector Machine (SVM)

Support Vector Machine (SVM) is a powerful machine learning algorithm used for linear or nonlinear classification, regression, and even outlier detection tasks.

4.2.1 Introduction to the Model

Support Vector Machines (SVM) stand as stalwarts in the realm of machine learning, serving as adept tools for both classification and detection challenges. In classification, SVM endeavors to find an optimal hyper-plane, a decision boundary that maximally separates distinct classes. This hyper-plane, whether a line in two dimensions or a hyper-plane in higher dimensions, is strategically positioned to create a margin, maximizing the distance from the nearest data points of each class. The focus on these critical data points, termed support vectors, ensures the robustness and generalization of the model.

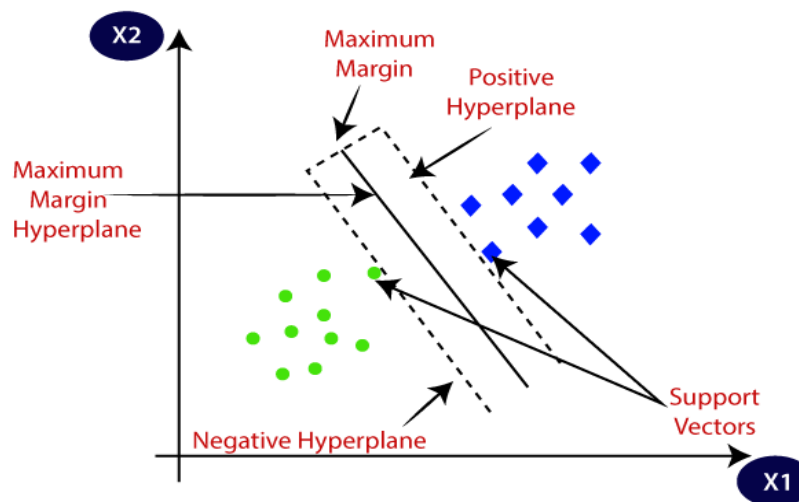


Figure 8: Functioning of SVM

In its pursuit of classification excellence, SVM employs the kernel trick when faced with non-linear data. This involves mapping the input data into a higher-dimensional space, facilitating the discovery of a linear hyper-plane in this transformed realm. The regularization parameter introduces a soft margin when perfect separation is unattainable, allowing for a balance between margin maximization and permissible errors.

For detection challenges, SVM takes on a one-class classification role. During training, it learns to discern normal instances, defining a boundary encapsulating typical data patterns. In the detection phase, instances deviating from the learned boundary are flagged as anomalies, showcasing SVM's prowess in identifying outliers.

4.2.2 Support Vector Machine Terminology

- **Hyper-plane:** Hyper-plane is the decision boundary that is used to separate the data points of different classes in a feature space. In the case of linear classifications, it will be a linear equation i.e. $wx+b=0$.
- **Support Vectors:** Support vectors are the closest data points to the hyperplane, which makes a critical role in deciding the hyper-plane and margin.
- **Margin:** Margin is the distance between the support vector and hyperplane. The main objective of the support vector machine algorithm is to maximize the margin. The wider margin indicates better classification performance.
- **Kernel:** Kernel is the mathematical function, which is used in SVM to map the original input data points into high-dimensional feature spaces, so, that the hyper-plane can be easily found out even if the data points are not linearly separable in the original input space. Some of the common kernel functions are linear, polynomial, radial basis function (RBF), and sigmoid.
- **Hard Margin:** The maximum-margin hyper-plane or the hard margin hyper-plane is a hyper-plane that properly separates the data points of different categories without any misclassifications.
- **Soft Margin:** When the data is not perfectly separable or contains outliers, SVM permits a soft margin technique. Each data point has a slack variable introduced by the soft-margin SVM formulation, which softens the strict margin requirement and permits certain misclassifications or violations. It discovers a compromise between increasing the margin and reducing violations.
- **C:** Margin maximization and misclassification fines are balanced by the regularization parameter C in SVM. The penalty for going over the margin or misclassifying data items is decided by it. A stricter penalty is imposed with a greater value of C , which results in a smaller margin and perhaps fewer misclassifications.
- **Hinge Loss:** A typical loss function in SVMs is hinge loss. It punishes incorrect classifications or margin violations. The objective function in SVM is frequently formed by combining it with the regularization term.
- **Dual Problem:** A dual Problem of the optimization problem that requires locating the Lagrange multipliers related to the support vectors can be used to solve SVM. The dual formulation enables the use of kernel tricks and more effective computing.

4.2.3 Architecture of SVM

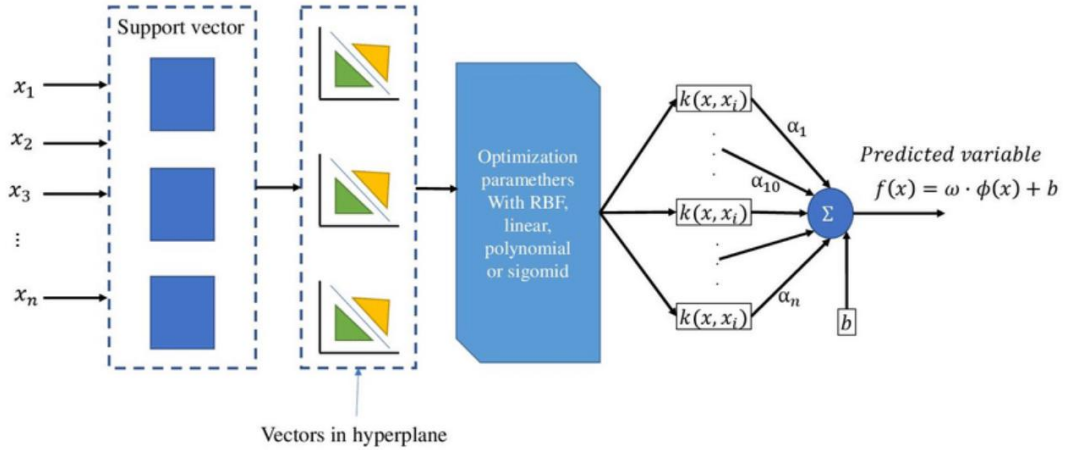


Figure 9: General Architecture of SVM

The architecture of Support Vector Machines (SVM) is intricately designed for robust classification, offering a comprehensive approach to handling diverse datasets. Starting in the input space, SVM utilizes the kernel trick for feature mapping, enabling the transformation of data into a higher-dimensional space, particularly beneficial for non-linear scenarios. Central to SVM are the support vectors, pivotal data points near the decision boundary. The model determines an optimal hyperplane, serving as the decision boundary, with a priority on maximizing the margin between classes. The regularization parameter introduces flexibility, allowing for some misclassifications, while kernel functions define the similarity between data points in transformed spaces. The decision function, derived from the learned hyper-plane and support vectors, classifies new instances. SVM's elegance lies in its adaptability, effectively handling both linear and non-linear data, making it a versatile and resilient choice in various machine learning applications.

4.2.4 Significance of the Model:

Support Vector Machines (SVM) hold significance in various domains due to their distinctive characteristics and capabilities:

- SVMs are effective in high-dimensional spaces, making them suitable for tasks with a large number of features, such as image classification or genomic data analysis.
- SVMs are less prone to over-fitting, especially in high-dimensional spaces. Their ability to generalize well to unseen data makes them robust and suitable for complex datasets.
- Through the use of different kernel functions (e.g., radial basis function kernel), SVMs can efficiently handle nonlinear decision boundaries, capturing complex relationships in the data.

- SVMs aim to find the hyper-plane that maximizes the margin between different classes. The clear geometric interpretation of the margin makes it easier to understand and interpret the model's decisions.
- SVMs are less sensitive to outliers in the training data. The focus on finding the optimal margin reduces the impact of individual data points that might deviate from the overall pattern.
- SVMs use a subset of training points called support vectors to define the decision boundary. This leads to memory efficiency, particularly when dealing with large datasets.
- The training of SVM involves solving a convex optimization problem, ensuring that the algorithm converges to the global minimum. This contributes to the stability and reliability of the model.
- SVMs have been successfully applied in image classification tasks, demonstrating their capability to handle complex visual data and achieve high accuracy.
- Being a highly sophisticated and mathematically sound algorithm, it is one of the most accurate machine learning algorithms.
- It is a dynamic algorithm and can solve a range of problems, including linear and non-linear problems, binary, binomial, and multi-class classification problems, along with regression problems.
- SVM is known for its computation speed and memory management. It uses less memory, especially when compared to machine vs deep learning algorithms with whom SVM often competes and sometimes even outperforms to this day.

4.3 Decision Tree

A decision tree is a supervised machine learning algorithm that creates a series of sequential decisions to reach a specific result.

4.3.1 Introduction to the Model

We typically use decision trees to create informed opinions that facilitate better decision making. Decision trees allow us to break down information into multiple variables to arrive at a singular best decision to a problem.

A decision tree is a non-parametric supervised learning algorithm for classification and regression tasks. It has a hierarchical tree structure consisting of a root node, branches, internal nodes, and leaf nodes. Decision trees are used for classification and regression tasks, providing easy-to-understand models.

Decision trees must contain all possibilities clearly outlined in a structured manner in order to be effective, but they must also present multiple possibilities for data scientists to make collaborative decisions and optimize business growth.

4.3.2 Decision Tree Terminologies

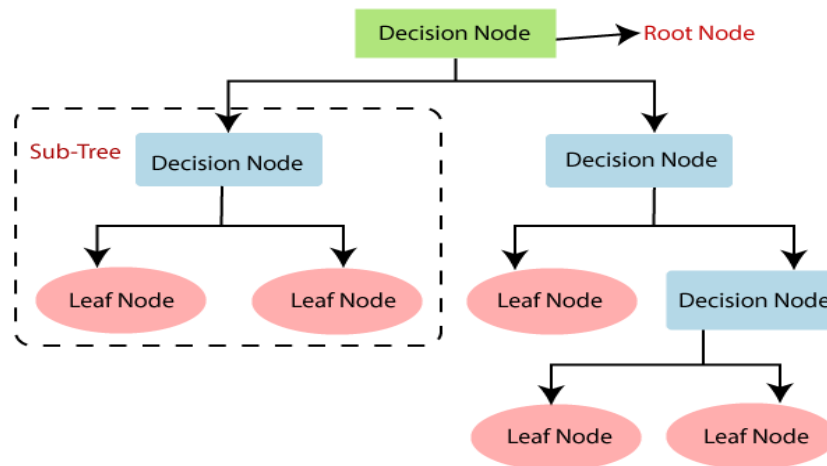


Figure 10: Decision Tree

- **Root Node:** Root node is from where the decision tree starts. It represents the entire dataset, which further gets divided into two or more homogeneous sets.
- **Leaf Node:** Leaf nodes are the final output node, and the tree cannot be segregated further after getting a leaf node.
- **Splitting:** Splitting is the process of dividing the decision node/root node into sub-nodes according to the given conditions.
- **Branch/Sub Tree:** A tree formed by splitting the tree.
- **Pruning:** Pruning is the process of removing the unwanted branches from the tree.
- **Parent/Child node:** The root node of the tree is called the parent node, and other nodes are called the child nodes.

4.3.3 How are Decision Trees used in Classification?

The Decision Tree algorithm uses a data structure called a tree to predict the outcome of a particular problem. Since the decision tree follows a supervised approach, the algorithm is fed with a collection of pre-processed data. This data is used to train the algorithm.

Decision trees follow a top-down approach meaning that the root node of the tree is always at the top of the structure while the outcomes are represented by the tree leaves. Decision trees are built using a heuristic called recursive partitioning (commonly referred to as Divide and Conquer). Each node following the root node is split into several nodes.

The key idea is to use a decision tree to partition the data space into dense regions and sparse regions. The splitting of a binary tree can either be binary or multi way. The algorithm keeps on splitting the tree until the data is sufficiently homogeneous. At the end of the training, a decision tree is returned that can be used to make optimal categorized predictions.

An important term in the development of this algorithm is Entropy. It can be considered as the measure of uncertainty of a given dataset and its value describes the degree of

randomness of a particular node. Such a situation occurs when the margin of difference for a result is very low and the model thereby doesn't have confidence in the accuracy of the prediction.

The higher the entropy, the higher will be the randomness in the dataset. While building a decision tree, a lower entropy shall be preferred. The expression for calculating the entropy of a decision tree is as described:

Another metric used for a similar purpose is the Gini Index. It uses the Gini method to create split points. Information Gain is the metric that is generally used for measuring the reduction of uncertainty in the dataset. This metric can further be used to determine the root node of the decision tree and the number of splits that are to be made. The root node of a decision tree is often referred to as the decision node or the master node.

Each tree has one root node. The root node is often considered the most important feature in rapport with all other features. Generally, the feature with the highest accuracy among all others is chosen as the root node.

4.3.4 Significance of the Model

- Decision trees provide a transparent and intuitive representation of decision-making processes. The visual tree-like structure allows users to comprehend the logic behind predictions, making it particularly valuable in domains where interpretability is crucial, such as healthcare and finance.
- Decision trees inherently rank features based on their contribution to classification or regression. This feature importance analysis aids in understanding which variables play a pivotal role in decision-making, facilitating better insights into the dataset.
- Decision trees are less sensitive to outliers compared to certain other models. The splitting criterion based on purity measures, such as Gini impurity or entropy, tends to mitigate the influence of outliers on the overall model.
- Decision trees can be applied to both classification and regression tasks. Ensemble methods like Random Forests capitalize on the strength of individual decision trees, further enhancing predictive performance.

4.4 K-Nearest Neighbor (KNN)

The K-Nearest Neighbors (KNN) algorithm is a robust and intuitive machine learning method employed to tackle classification and regression problems.

4.4.1 Introduction to the Model

Surprisingly enough, the KNN algorithm is quite accessible and easy to understand. K-Nearest Neighbors (KNN) is commonly used in classification and identification problems due to its simplicity and effectiveness.

K-Nearest Neighbors is one of the most basic yet essential classification algorithms in Machine Learning. It belongs to the supervised learning domain and finds intense application in pattern recognition, data mining, and intrusion detection.

4.4.2 Basic Working Principle

Here's how KNN is employed in these tasks:

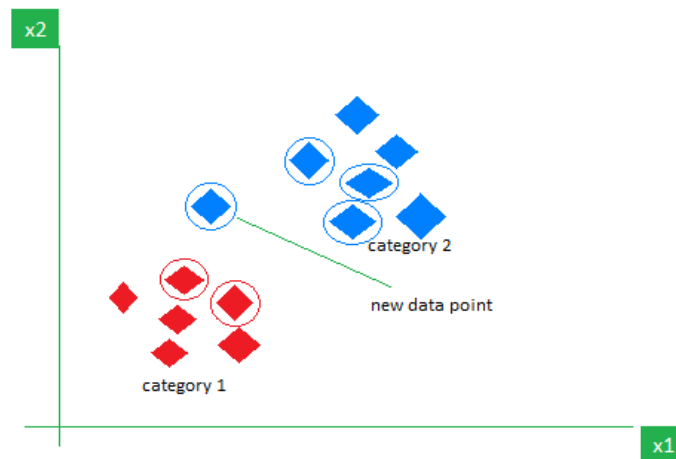


Figure 11: KNN Algorithm working visualization

1. **Data Representation:** The dataset is represented as points in a multidimensional feature space, where each feature corresponds to a characteristic of the data. For example, in image classification, each pixel might be a feature.
2. **Training Phase:** During the training phase, the algorithm stores the feature vectors and their corresponding class labels from the labeled dataset.
3. **Prediction (Classification):** When a new, unlabeled data point is presented, the algorithm identifies its k-nearest neighbors based on a chosen distance metric (commonly Euclidean distance). The class labels of these neighbors are examined, and the most prevalent class among them is assigned to the new data point.
4. **Prediction (Identification):** In identification problems, the goal is to identify the identity or category of the input based on its similarity to known instances. KNN, in this context, finds the k-nearest neighbors and identifies the most common identity or category among them.
5. **Parameter Tuning:** The choice of the parameter k (the number of neighbors) is crucial. A smaller k value makes the model more sensitive to local variations but may be sensitive to noise, while a larger k value provides a smoother decision boundary but may overlook local patterns.
6. **Scalability:** KNN can be computationally expensive, especially as the size of the dataset grows. Efficient data structures like KD-trees or ball trees are sometimes employed to accelerate the search for nearest neighbors.

KNN is versatile and can be applied to various domains, including image recognition, pattern recognition, and recommendation systems. Its simplicity and flexibility make it a popular choice, especially when the decision boundaries are complex and not easily modeled by parametric methods. However, the choice of the appropriate distance metric and the parameter k requires careful consideration for optimal performance.

4.4.3 How distances and similarities are carried out in KNN

Several distance metrics determine correlation and similarity. Even though there are plenty of distance functions to choose from, we should always use the functions that best fit the nature of our data. Notable metrics include:

Table 2: Distance Metric table

Distance Metric	Purpose
Euclidean Distance	Distance Metric
Taxicab Geometry	Used when the data types are heterogenous
Minkowski distance	Intended for real-valued vector spaces
Jaccard index	Often used in applications when dealing with binarized data
Hamming distance	Typically used with data transmitted over computer networks. And also used with categorical variables.

4.4.4 Significance of the Model

The K-Nearest Neighbors (KNN) model holds significant importance in machine learning for several reasons:

- KNN is a straightforward algorithm that is easy to understand and implement. Its simplicity makes it an excellent choice for introductory learning and quick prototyping.
- It is effective in scenarios where decision boundaries are complex and not easily defined by mathematical equations.
- KNN is a non-parametric model, meaning it does not make assumptions about the underlying data distribution. This makes it suitable for situations where the data distribution is not well-known or may change over time.
- KNN relies on local information, considering only nearby data points for predictions. This makes it robust to outliers and less sensitive to the global structure of the data.
- Unlike many other machine learning models, KNN does not require a lengthy training period. The model "learns" during the prediction phase by comparing new instances to the stored training data.

- KNN is often employed as a base learner in ensemble methods like bagging and boosting, contributing to the diversity of the ensemble.
- The predictions made by KNN are interpretable and can be easily explained. This is beneficial in scenarios where model interpretability is crucial.
- KNN does not make assumptions about the underlying distribution of the data, making it applicable to a wide range of problems without the need for extensive preprocessing.

4.5 Random Forest

Random Forest is a popular machine learning algorithm that belongs to the supervised learning technique.

4.4.1 Introduction to the Model

A random forest algorithm consists of many decision trees. The ‘forest’ generated by the random forest algorithm is trained through bagging or bootstrap aggregating. Bagging is an ensemble meta-algorithm that improves the accuracy of machine learning algorithms.

The (random forest) algorithm establishes the outcome based on the predictions of the decision trees. It predicts by taking the average or mean of the output from various trees. Increasing the number of trees increases the precision of the outcome.

A random forest eradicates the limitations of a decision tree algorithm. It reduces the overfitting of datasets and increases precision. It generates predictions without requiring many configurations in packages (like scikit-learn).

4.4.2 Working of Random Forest

Here's a concise explanation of how the Random Forest algorithm works:

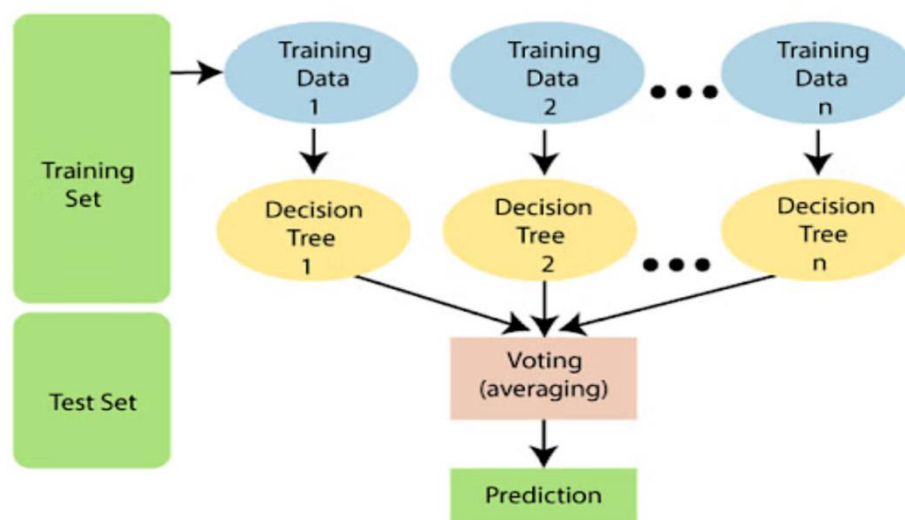


Figure 12: Working of the Random Forest

- **Ensemble of Decision Trees:** Random Forest builds an ensemble of decision trees during the training phase. Each tree is constructed independently and, importantly, based on a random subset of the training data.
- **Random Subset Selection:** For each tree in the ensemble, a random subset of the training data is selected with replacement (bootstrap sampling). This introduces diversity among the trees.
- **Random Feature Selection:** When splitting nodes in each decision tree, a random subset of features is considered. This prevents specific features from dominating the tree-building process and promotes robustness.
- **Tree Construction:** Each decision tree is grown by recursively splitting nodes based on the selected subset of features. The splitting criteria typically involve maximizing information gain (for classification) or minimizing variance (for regression).
- **Voting or Averaging:** For classification tasks, the final prediction is determined by a majority vote among the ensemble of trees. For regression tasks, the predictions are averaged.
- **Reduced Over-fitting:** The randomness introduced during subset selection and feature selection helps to de-correlate the individual trees, reducing over-fitting and improving the generalization performance of the model.
- **Feature Importance:** Random Forest can quantify the importance of each feature in the classification or regression process, providing insights into which features contribute most to predictive accuracy.
- **Parallelization:** The construction of individual trees in the ensemble is often parallelized, making Random Forest a computationally efficient algorithm.

4.4.3 Essential Features of Random Forest

- **Miscellany:** Each tree has a unique attribute, variety and features concerning other trees. Not all trees are the same.
- **Immune to the curse of dimensionality:** Since a tree is a conceptual idea, it requires no features to be considered. Hence, the feature space is reduced.
- **Parallelization:** We can fully use the CPU to build random forests since each tree is created autonomously from different data and features.
- **Train-Test split:** In a Random Forest, we don't have to differentiate the data for train and test because the decision tree never sees 30% of the data.
- **Stability:** The final result is based on Bagging, meaning the result is based on majority voting or average.

4.4.4 Significance of the model

- Random Forest typically delivers high accuracy in classification tasks. By aggregating predictions from multiple decision trees, it mitigates the risk of over-fitting and generalizes well to new, unseen data.
- Random Forest can effectively capture complex relationships in the data, making it suitable for tasks where the decision boundaries are intricate and nonlinear.
- Random Forest provides an internal evaluation mechanism through out-of-bag samples. Since each tree is trained on a different subset of the data, out-of-bag instances serve as a built-in validation set, offering an estimate of the model's performance without the need for a separate validation set.
- The ensemble of decision trees in Random Forest is diverse due to the random subset selection during both instance and feature sampling. This diversity contributes to the model's stability and ability to capture intricate patterns within the data.
- It can handle large datasets efficiently.

CHAPTER 5: METHODOLOGY

This study aims to evaluate various machine learning and deep learning methods for accurately classifying diverse cattle breeds using image data. The challenge lies in differentiating breeds based on their varied appearances and features. By employing algorithms like SVM, KNN, DT, Random Forests, and CNN, this research targets to explore their effectiveness, strengths, and limitations in this classification task. The methodology encompasses several key stages, including data collection and pre-processing, feature extraction, CNN model architecture design, model training and evaluation, hyper-parameter tuning, and considerations regarding limitations and assumptions.

5.1 Dataset

5.1.1 Dataset Description

The dataset employed in this study comprises a diverse collection of images featuring multiple cattle breeds, pivotal for conducting tasks related to cattle breed classification. We have collected 10,000 data. Out of the 10,000 physical data points we initially collected, a portion had to be excluded due to inaccuracies or limitations in specifying a single cattle breed. The curating process has involved carefully filtering out instances with improper information, ensuring the final dataset maintains a high standard of accuracy and relevance for our cattle breed classification model. Essential dataset details are as follows:

Total Files: Within this dataset, there exists a total of 5025 image files, each distinctly portraying specific cattle breeds.

Total Cattle Breed: The total number of cattle breeds mentioned in the provided list is 5: Holstein Friesian, Black Angus, Red Chittagong, Shahiwal, White Native.

5.1.2 Objectives of Collecting the Data

The primary objective of collecting the extensive dataset of cattle is to facilitate the development of a robust and accurate cattle breed classification model. The goal is to encompass a diverse range of cattle images representing various breeds found in different environments, including dairy farms, cow huts, random streets, and fields. By capturing images from these diverse locations, we aim to ensure that our model would be well-equipped to handle the variability in cattle appearances across different settings.

5.1.3 Source of Data

We have collected the dataset from various sources. We have sourced images directly from renowned dairy farms such as Iqbal Agro, Saara Agro, Asian Agro, A.U.R Agro farm, Green Harvest Agro and Bhuiyan Agro. Additionally, we ventured into cow huts, capturing images from locations like Bibirhat, 1 Kilometer, Chowdhury Hut, Shagorika, Ilias Brothers Hut, Qurbani Hut, and Moijjartek Gorur Bazar. In order to account for the diversity of cattle appearances, we have also included images from streets and fields.

5.1.4 Data Collection Method

Over the course of an extensive 20 to 25 days, our data collection efforts have involved physically visiting numerous locations. The data collection process has encompassed a combination of on-site photography and systematic image gathering. Visiting each location, we employed digital cameras and smartphones to capture high-resolution images of cattle. This approach has allowed us to cover a wide range of environmental conditions and ensured that the dataset is reflective of the diverse contexts in which cattle are found. The use of multiple devices and careful documentation of the collection process has helped maintain data consistency and integrity.

5.1.5 Image Collection and Labeling Process

Images collected from the various sources are carefully labeled to provide the necessary ground truth for training the classification model. Each image is assigned a label corresponding to its respective cattle breed. This labeling process has involved a meticulous review and categorization, ensuring the accuracy of the ground truth.

To enhance image clarity, we employed spot fixing and PhotoRoom, an AI photo editing tool, to remove extraneous objects such as additional cow legs or isolated body parts in some images. This meticulous editing process contributes to a cleaner and more focused dataset for accurate cattle breed classification.



Figure 13: Cattle image: Before and After Object Removal

5.1.6 Challenges Faced

While the data collection process has been extensive, it has not been without challenges. Some notable hurdles include the need to navigate varying lighting conditions, the unpredictability of cattle behavior, and the occasional difficulty in obtaining permission to photograph in certain locations. These challenges underscore the importance of adaptability and meticulous planning during the data collection phase. Despite these obstacles, the diversity and scale of our dataset contribute significantly to the robustness of our cattle breed classification model.

5.2 Data Pre-processing

5.2.1 Image Resizing

The batch size parameter, set to 12, defines the number of images processed simultaneously during each training iteration. Meanwhile, the image dimensions, specified as 180 pixels in both height and width (180x180), indicate the standardized size to which the images will be resized or adjusted for uniformity across the dataset. This configuration ensures that the model processes batches of 12 images at a time, with each image resized to a consistent dimension of 180x180 pixels for effective training and analysis.

5.2.2 Dataset Splitting

This division allows the model to grasp patterns from a substantial part of the data while gauging its adaptability to new, unseen images. Typically, about 70% of the dataset is used for training purposes, while the remaining 30% is set aside to test the model's accuracy and effectiveness. This method minimizes the risk of over-fitting, thereby enhancing the model's reliability when handling new, real-world images.

5.2.3 Extraction of Training Data

The process initiates by creating empty containers, `x_train` and `y_train`, designed to store image data and their respective labels. It iterates through the training dataset (`train_ds`), fetching batches of images (referred to as `images`) along with their corresponding labels (denoted as `labels`). These acquired image batches and their associated labels are sequentially appended to their designated containers (`x_train` for images and `y_train` for labels), thereby forming the complete training datasets, `x_train` and `y_train`. This procedure is fundamental for assembling the requisite data used to train machine learning models.

5.2.4 Extraction of Testing Data

Creating a separate set of data for testing a model's performance mirrors the method used for training data. It involves initializing containers (`x_test` and `y_test`) to hold image data and labels. The process loops through the testing dataset (`test_ds`), gathering image batches and their corresponding labels. These batches are added to `x_test` for images and `y_test` for labels, forming the complete testing datasets. This separation allows evaluating the model's performance on new, unseen data.

5.2.5 Optimization

Through optimization, we aim to reduce the errors or losses in the model, enhancing its capability to make precise predictions on fresh, unseen data. Optimizations, such as caching, shuffling, and prefetching, are applied to both training and testing datasets to enhance their processing efficiency and overall performance during model training and evaluation. This ultimately leads to a more effective utilization of computational resources and potentially expedites the model's training and evaluation processes. For this, we used AUTOTUNE. AUTOTUNE is utilized for optimization purposes, dynamically fine-tuning Tensor Flow's operations to maximize performance based on the current computing resources available.

5.2.6 Normalization

Our primary purpose is to rescale numerical values to a standard range, typically $[0, 1]$ or $[-1, 1]$, to ensure consistent and efficient model training. The normalization process involves adjusting these pixel values to fit within a specific range. Often, normalization rescales pixel values to fall between 0 and 1. This is commonly achieved by dividing each pixel value by the maximum value it can have. Normalization standardizes the input data, preventing certain features (like pixel intensity) from dominating the learning process, ensuring that all features contribute more equally. It can accelerate model convergence during training by aiding in the optimization process.

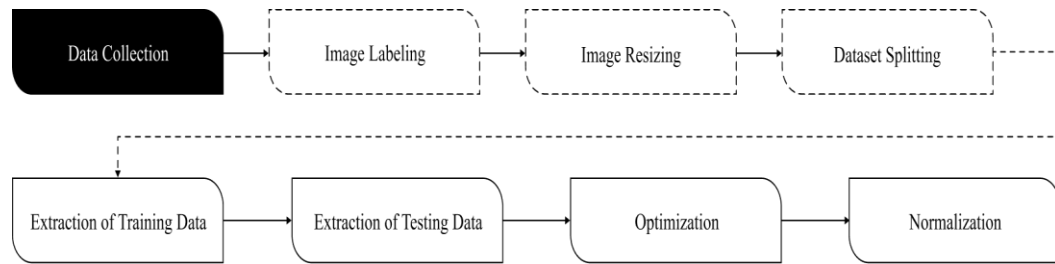


Figure 14: Work Flow of Data Pre-Processing

5.3 Working Process of Convolutional Neural Network (CNN)

We have created a sequential CNN model tailored to handle images of cattle breeds. We trained the initial model using our dataset and have observed its performance across training epochs through visualizations. To enhance diversity and robustness, we have augmented the dataset, incorporating these techniques into the CNN structure. The model has been further trained using data augmentation layers, with a focus on achieving improved performance metrics and potential retraining. The trained model has then been employed to process new cattle images, providing predictions on their breed and displaying both the predicted class and its accuracy. To evaluate the model's effectiveness in predicting various cattle breeds, we have constructed a confusion matrix. The resulting comprehensive report details precision, recall, F1-score, and support metrics for each predicted cattle breed class, allowing for a meaningful comparison against the true labels based on the methodology we have followed in our work.

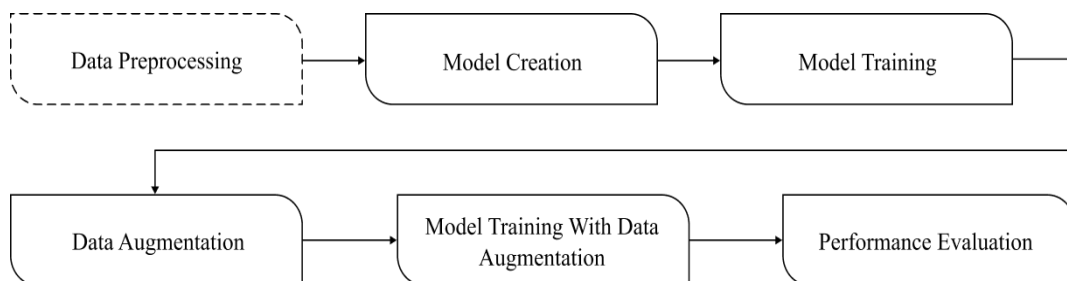


Figure 15: CNN Work Flow

5.3.3 Model Definition

- Convolutional and pooling layers have been included in our sequential model to effectively handle photos of cattle breeds.
- For training, the model has been assembled using the proper loss function and optimizer during our work.

5.3.2 First Model Training

- The training dataset has been utilized to train the first CNN model.
- Understanding model performance has been achieved by visualizing the accuracy/loss trends during training and validation across epochs in our work.

5.3.3 Data Augmentation

- To increase the diversity of the training dataset and the robustness of the model, data augmentation techniques (random flip, rotation, zoom) have been applied.
- This model incorporates the data augmentation pipeline (data_augmentation) as the first layer within the sequential model. After data augmentation, it follows a similar architecture to the previous model: normalization, convolutional layers, max-pooling layers, dropout, and dense layers for classification in our work.

5.3.4 Second Model Training with Data Augmentation

- Data augmentation layers have been added to the CNN architecture from the start.
- The model is compiled with the Adam optimizer, Sparse Categorical Cross entropy loss function, and accuracy as the evaluation metric. Then, it's trained the training dataset and test the testing dataset for 15 epochs.
- For possible gains, the enhanced model has been retrained using the training dataset in our work.

5.3.5 Prediction on New Image

- A fresh breed of cattle image has been preprocessed and loaded.
- To predict the class of the new image, the trained CNN model has been applied. The loaded image has been converted into an array, representing the image in a format suitable for model prediction during our work.
- The pre-trained model has been used to predict the class probabilities for the given image. The predicted class along with the associated accuracy has been displayed.

5.3.6 Confusion Matrix Visualization

- With the trained CNN model, predictions have been created for the test dataset. The trained model has been utilized to make predictions on the test dataset, retrieving the predicted classes by selecting the index with the highest probability during our work.
- The test dataset has been applied to determine the true labels. True labels have been collected from the test dataset to compare them with the predicted classes.
- The confusion matrix has counted the occurrences of each class's true and predicted labels, constructing a matrix that showcases the model's performance in a tabular format during our work.

5.3.7 Classification Report

- Metrics such as precision, recall, F1-score, and support for each class have been calculated and generated in a report.

5.4 Working Process of Support Vector Machine (SVM)

The workflow demonstrates the entire process of training an SVM model for image classification, evaluating its performance on a testing dataset, and generating a

comprehensive report to assess its accuracy and class-wise performance.

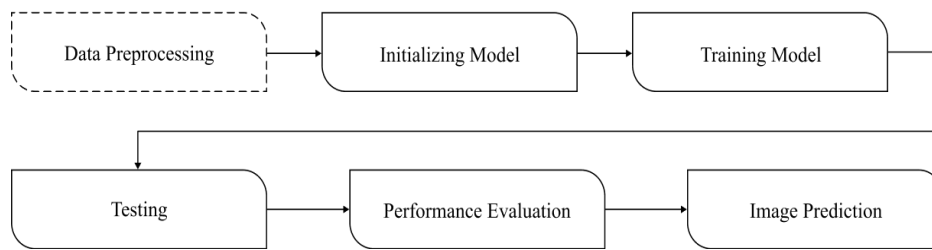


Figure 16: SVM Work Flow

5.4.1 Data Pre-processing for Training & Testing

This processes prepare the image data for optimal utilization within the SVM model, aligning them with the model's specifications and ensuring standardized and optimized training. Training images and labels are retrieved from the dataset. Images are converted into NumPy arrays, conforming to the SVM model's requirements. Reshaping ensures that the arrays match the expected input format of the SVM model. Pixel values in the images are scaled between 0 and 1. This normalization process standardizes pixel values, benefiting the SVM model's training. Similar pre-processing is done for the testing dataset, extracting images and labels and preparing them for prediction.

5.4.2 Initializing and Training the Model

The step initiates an SVM model with a linear kernel to perform classification tasks. It commences by extracting and preparing the training dataset, segregating images and their associated labels. Following the flattening of images and normalization of pixel values within the 0 to 1 range, the code instantiates the SVM model utilizing a linear kernel. Later, the model is trained using the pre-processed training images and their corresponding labels. This training phase aids the SVM algorithm in establishing distinct boundaries between various cattle breed classes based on the provided training data.

5.4.3 Testing and Performance Evaluation

Testing and performance evaluation are integral steps that gauge the trained SVM model's capability to generalize to new data. The testing phase evaluates the trained model's performance on a completely unseen dataset, distinct from both the training and testing sets. These testing and evaluation steps provide a comprehensive assessment of the trained SVM model's performance in classifying different cattle breeds using the provided testing images. Metrics such as accuracy, confusion matrix, and classification report metrics offer insights into the model's accuracy, misclassifications, and strengths for individual cattle breed classes.

5.4.4 Unknown Image Prediction

- **Loading the Unknown Image:** The image is fetched from a specified file path or source using image-handling libraries during our work.
- **Resizing & Pre-processing:** The image is adjusted to a specific size or dimensions required by the model for consistency in input during our work.

- **Data Conversion for Model Input:** The image content is converted into a structured array, typically a NumPy array. The array values have been reshaped and normalized during our work to fit the format expected by the machine learning model.
- **Prediction using the SVM Model:** The trained Support Vector Machine (SVM) model has been employed to predict the image's classification during our work.
- **Mapping Predicted Index to Class Name:** The predicted index has been used during our work to retrieve the corresponding class name or label from a predefined list or dictionary. This mapping connects the numeric prediction to a human-readable class name.
- **Displaying Predicted Class:** The predicted class name associated with the unknown image has been presented or printed during our work to communicate the model's inference.

5.5 Working Process of Decision Tree

The Decision Tree (DT) classifier is a machine learning technique that constructs a tree-shaped structure by repeatedly splitting data based on features to classify information. Trained on labeled images showing diverse cattle breeds, it learns patterns to categorize new images. To assess its performance, the model uses metrics like accuracy, confusion matrix, and classification report on a separate test dataset. By analyzing pixel values in images, the DT model creates a series of rules to classify images into specific cattle breed categories.

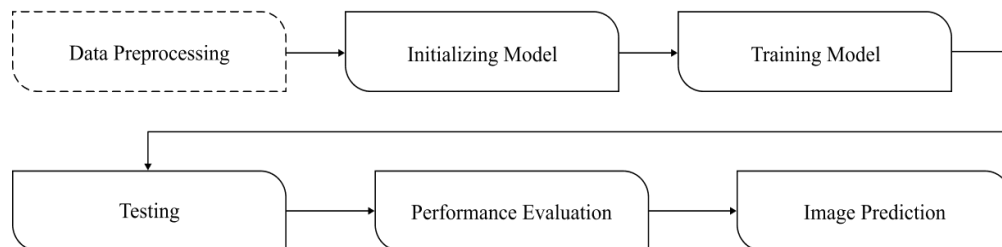


Figure 17: Decision Tree Work Flow

5.5.1 Data Pre-processing for Training

These pre-processing steps are crucial to ensure that the training data is appropriately formatted and scaled before training the Decision Tree classifier, enabling effective learning from the dataset. In this step, collecting images and labels from the training dataset. Converting images to NumPy arrays and normalizing pixel values. Reshaping and formatting the image arrays to match the input format required by the model. Normalizes the pixel values of the images to a range between 0 and 1 by dividing by 255.0. Reshapes the image arrays to a 2D format. This preserves the number of samples while flattening the dimensions into a single vector, maintaining each image's content.

5.5.2 Initializing and Training the Decision Tree Model

Initializing and training the Decision Tree Model involves two main steps:

Model Initialization:

- This step creates a blank canvas for the Decision Tree classifier.
- Initializing the model sets up the structure and parameters necessary for learning from data.

Training Process:

- Leveraging the preprocessed training data comprising images and their respective labels, the model engages in an educative process.
- Throughout this training phase, the model grasps the intricate connections existing between input images and their corresponding labels.
- Through iterative adjustments to its internal parameters, the model aims to minimize inaccuracies, aiming for precise mapping of input images to their relevant labels.

5.5.3 Data Preprocessing for Testing

These steps have collectively ensured that the test/validation data is suitably formatted, pixel values are standardized, and images are shaped appropriately to align with the model's anticipated input structure. This has prepared the test data for the accurate evaluation of the model's performance. Initially, we prepared the test dataset for evaluation by extracting images and labels from the dataset, converting these images into NumPy arrays, and normalizing their pixel values. To ensure uniformity, the pixel values in these images have been normalized to fit within a standardized range (typically 0 to 1) by dividing each value by 255.0. The images have been reshaped to align with the input format used during training.

5.5.4 Evaluation of Decision Tree Model

In the evaluation phase, the Decision Tree model analyzes the test dataset by making predictions for its samples. It measures accuracy by comparing these predictions to the actual true labels. Furthermore, it produces a detailed classification report containing precision, recall, F1-score, and support metrics for each class. This report provides not only overall accuracy but also macro-average and weighted-average metrics. This comprehensive report gives a complete view of the model's performance, highlighting its accuracy, misclassifications, and strengths across different classes in the dataset.

5.5.5 Unknown Image Prediction

The process starts by fetching an image from a specified file path. This image undergoes resizing to specific height and width dimensions, followed by conversion into a NumPy array. To match the model's expected input format, the array's structure is adjusted to a shape defined as (1, -1), ensuring it fits the model's requirements. Additionally, the pixel values in the array are normalized to fall within the range of 0 to 1. Once the image is prepared, the trained Decision Tree model predicts the class label for this single image. The predicted class index is then used to retrieve the corresponding class name. Finally, the predicted class name, indicating the categorization of the unknown image, is displayed

on the console.

5.6 Working Process of K-Nearest Neighbors (KNN)

Using the K-Nearest Neighbors (KNN) algorithm for cattle breed classification with image data appears to be a robust strategy, taking advantage of KNN's simplicity and efficiency in handling complex tasks like breed categorization. The fundamental concept in KNN, employing majority voting among the K most similar neighbors to label an input, is pivotal to its functioning.

Assessing the model's performance through metrics like accuracy, precision, recall, and computational efficiency is essential. Comparing these metrics against benchmarks and alternative classification algorithms is crucial to gauge the effectiveness and efficiency of the KNN model in accurately categorizing various cattle breeds from images.

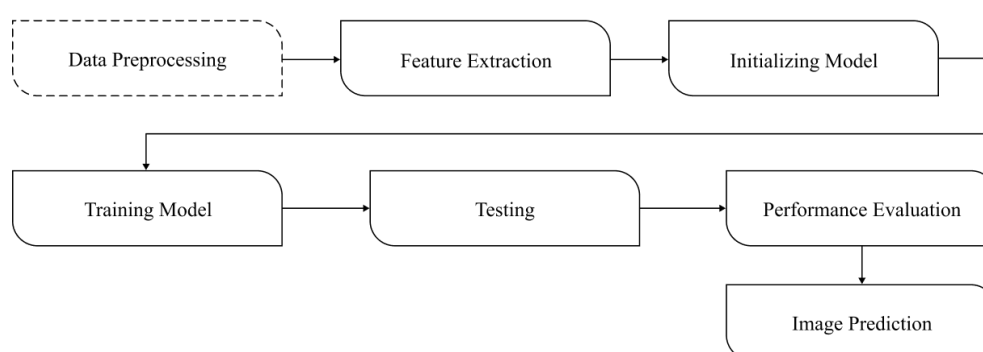


Figure 18: KNN Work Flow

5.6.1 Data Preparation

- Features have been extracted from photos of various cattle breeds using a pre-trained VGG16 Convolutional Neural Network (CNN).
- The dataset has been divided into sets for testing and training during our work.

5.6.2 Feature Extraction

- The classification layers of VGG16 have been removed, and it has been used as a feature extractor during our work.
- We have iterated through the training dataset, extracting features using the modified VGG16 model.
- These extracted features, along with their respective labels, have been stored.

5.6.3 KNN Model Training

- A K-Nearest Neighbors (KNN) classifier has been initialized with a specified number of neighbors.
- The KNN model has been fitted using the extracted features and their corresponding labels during our work.

5.6.4 Testing Process

- Features have been extracted from the testing dataset using the same modified VGG16 model during our work.
- These features have been reshaped and concatenated for compatibility with the KNN model.
- The labels for the validation set have been predicted using the trained KNN model.

5.6.5 Evaluation Metrics

- The performance has been visualized using a confusion matrix.
- A classification report containing metrics like precision, recall, and F1-score for each cattle breed has been generated.

5.6.6 Prediction of Unknown Image

- We have loaded new, unseen images of cattle breeds during our work.
- We have utilized the modified VGG16 model to extract features from these images.
- Reshaping the features, we have employed the trained KNN model to predict the class of the unknown images.
- The predicted cattle breeds have been output based on the model's inference.

5.7 Working Process of Random forest (RF)

The dataset has been prepared by utilizing the VGG16 model to extract features from images displaying various cattle breeds during our work. It has been split into sections for training and validation. Features from the training images have been extracted using VGG16, with labels kept alongside these features. Subsequently, a Random Forest classifier has been trained by setting the number of decision trees and using the extracted features and labels. For the testing set, features have been extracted using VGG16, rearranged, and aligned for the Random Forest model. Labels have been determined using the trained classifier. The Random Forest model has been evaluated by comparing predicted labels with the actual ones in the test set. Its performance has been visualized with a confusion matrix, and a detailed classification report, including recall, precision, and F1-score metrics for each cattle breed, has been presented.

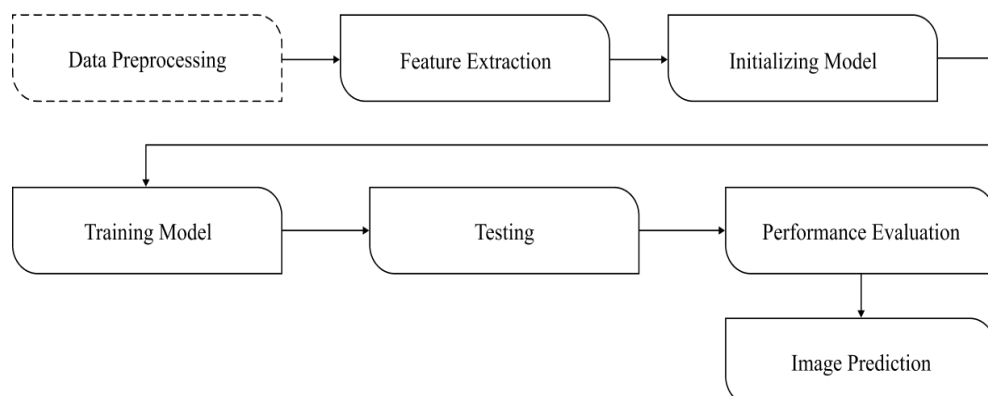


Figure 19: Random Forest Work Flow

5.7.1 Data Preparation

- High-level features have been extracted from pictures of different cattle breeds by using the VGG16 model as a feature extractor during our work.
- The dataset has been divided into sets for validation and training.

5.7.2 Feature Extraction

- Features have been extracted from the images in the training dataset by applying the VGG16 model during our work.
- Labels for these extracted features have been stored with them.

5.7.3 Random Forest Model Training

- The initial number of decision trees has been set in a Random Forest classifier during our work.
- Labels have been applied to the extracted features, and the Random Forest model has been fitted.

5.7.4 Testing Process

- The same VGG16 feature extractor has been used to extract features from the images in the validation dataset during our work.
- To enable the Random Forest model to work with these features, they have been reshaped and concatenated.
- The validation set's labels have been determined by applying the trained Random Forest classifier.

5.7.5 Evaluation Metrics:

- The Random Forest model's accuracy has been determined by contrasting the predicted labels with the actual labels in the testing set during our work.
- A confusion matrix has been used to visualize the model's performance, showcasing the true and predicted labels.
- A classification report for each breed of cattle has been provided, including metrics like recall, precision, and F1-score.

5.7.6 Prediction of Unknown Image:

- A fresh, never-before-seen image of a breed of cattle has been loaded and prepared during our work.
- The VGG16 feature extractor has been utilized to extract features from this picture.
- Using the trained Random Forest model, the features have been reshaped, and the class of the unknown image has been predicted. The anticipated breed of cattle has been provided.

5.8 Performance Evaluation

The performance of our cattle breed classification was assessed using several evaluation metrics, including:

5.8.1 Confusion matrix

Confusion Matrix describes the performance of a model on a set of test data. It gives two types of correct predictions and two types of incorrect predictions for the classifier. A confusion matrix is a table that visualizes the performance of a classification model. Each row of the matrix represents the instances in an actual class, while each column represents the instances in a predicted class.

Table 3 shows the confusion matrix. TP is the predicted output as true positive, TN is the predicted output as true negative, FP is the predicted output as false positive, and FN is the predicted output as a false negative. The accuracy, precision, recall, and f-measure (f-score) are defined in the following:

Table 3: Confusion matrix Predicted

	Predicted Class 0	Predicted Class 1
Actual Class 0	TP	FP
Actual Class 1	FN	TN

5.8.2 Accuracy

Accuracy in the context of a classification system represents the proportion of correct predictions made by the model across all classes. It measures the overall correctness of the model's predictions by considering both true positive and true negative predictions in relation to all predictions made. It shows the performance of the classification system as follows:

$$\text{Accuracy} = \frac{TN + TP}{TN + TP + FN + FP}$$

5.8.3 Precision

Precision is a metric that measures the accuracy of positive predictions made by a model. Specifically, it calculates the ratio of true positive predictions to the total predicted positive instances, indicating how many of the predicted positive cases were actually correct. It is the total number of correctly classified positive divide on the total number of predicted positive examples [4]. The equation of the precision is given as follow:

$$\text{Precision} = \frac{TP}{TP + FP}$$

5.8.4 F-measure

The F-measure, also known as the F1-score, is a combined metric that balances both precision and recall into a single value. It's particularly useful when you want to consider both false positives and false negatives in a classification problem. The equation of the f-measure is given as follows:

$$\mathbf{F\text{-}measure} = \frac{2 * \mathbf{Recall} * \mathbf{Precision}}{\mathbf{Recall} + \mathbf{Precision}}$$

5.8.5 Recall

Recall, also known as sensitivity or true positive rate, measures the proportion of actual positive instances that were correctly predicted by the model. It calculates the ratio of true positives to the sum of true positives and false negatives. The equation of the recall is given as follows:

$$\mathbf{Recall} = \frac{\mathbf{TP}}{\mathbf{TP} + \mathbf{FN}}$$

5.8.6 ROC-AUC Curve

The ROC curve is a graphical representation of a classifier's performance across various thresholds. It illustrates the trade-off between true positive rate (sensitivity) and false positive rate (1 - specificity). AUC quantifies the overall performance of a classifier. It represents the area under the ROC curve, where a higher AUC value (closer to 1) indicates better discriminative ability. AUC serves as a single scalar value summarizing the model's discriminatory power. Higher AUC values signify better classification performance across the entire range of thresholds. It allows for direct comparison between different models. A model with a higher AUC generally has better overall predictive performance.

CHAPTER 6: RESULTS

In this research, the goal was to categorize cattle breeds utilizing a dataset of 5000 images that portrayed five unique breeds: Australian, Black Angus, Red Chittagong, Shahiwal, and White Native. To achieve this, a set of five diverse models including Support Vector Machine (SVM), Decision Tree (DT), K-Nearest Neighbors (KNN), Random Forest, and Convolutional Neural Network (CNN) were employed. The performance of each model underwent evaluation using standard metrics such as accuracy, precision, recall, and F1-score. This comprehensive analysis aimed to understand the capabilities and limitations of these models in accurately classifying the specific cattle breeds showcased in the dataset.

6.1 Individual Classification Model Results

6.1.1 Support Vector Machine (SVM)

Holstein Friesian: The model correctly classified 279 instances of the Holstein Friesian breed, misclassifying 4 as Black Angus and 7 as White Native.

Black Angus: It accurately predicted 306 instances of Black Angus, misclassifying 9 as Holstein Friesian, 1 as Sahiwal, and 3 as White Native.

Red Chittagong: All 286 instances of Red Chittagong were correctly classified.

Shahiwal: 316 instances of Sahiwal were correctly classified, with only 1 misclassified as Holstein Friesian and 1 as Black Angus.

White Native: The model correctly predicted 279 instances of White Native, misclassifying 11 as Holstein Friesian, 1 as Black Angus, 1 as Sahiwal, and 2 as Shahiwal.

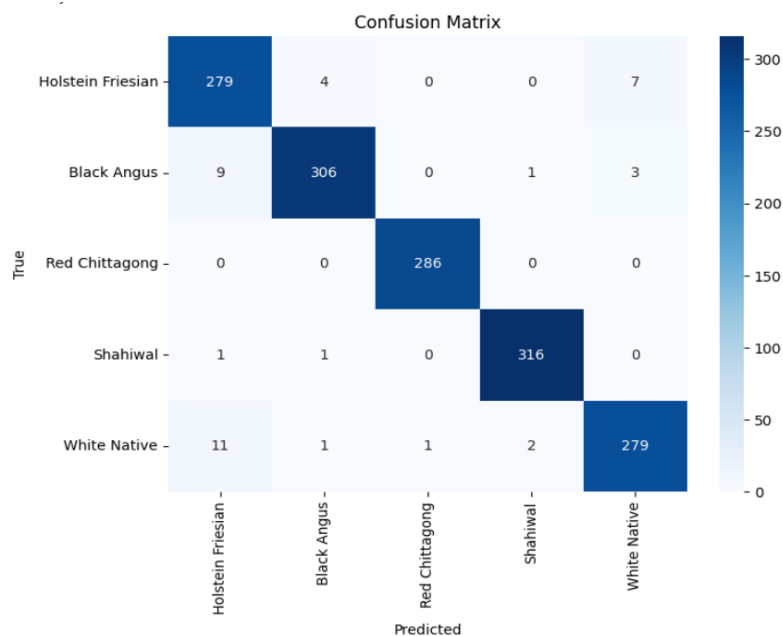


Figure 20 : Confusion Matrix of SVM

These results demonstrate the strengths and weaknesses of the model in distinguishing between specific cattle breeds. While some breeds were classified accurately with minimal errors, others showed higher misclassification rates, indicating potential areas for improvement in the model's predictive capabilities for those breeds.

Classification Report:				
	precision	recall	f1-score	support
Holstein Friesian	0.93	0.96	0.95	290
Black Angus	0.98	0.96	0.97	319
Red Chittagong	1.00	1.00	1.00	286
Shahiwal	0.99	0.99	0.99	318
White Native	0.97	0.95	0.96	294
accuracy			0.97	1507
macro avg	0.97	0.97	0.97	1507
weighted avg	0.97	0.97	0.97	1507

Figure 21: Classification Report of SVM

The SVM model displayed an outstanding overall accuracy of 97% in accurately categorizing cattle breeds. Its precision rates stood out consistently across different breeds, exhibiting high precision for Red Chittagong (100%), White Native (97%), Holstein Friesian (93%), Black Angus (98%), and Sahiwal (99%).

Furthermore, the model showcased impressive recall percentages for Red Chittagong (100%), White Native (95%), Holstein Friesian (96%), Black Angus (96%), and Sahiwal (99%) breeds.

Moreover, the SVM model attained commendable F1 scores for Red Chittagong (100%), White Native (96%), Holstein Friesian (95%), Black Angus (97%), and Sahiwal (99%) breeds.

6.1.2 Decision Tree (DT)

Holstein Friesian: The model correctly classified 271 instances of Holstein Friesian, misclassifying 3 as Black Angus, 1 as Red Chittagong, 4 as Shahiwal, and 11 as White Native.

Black Angus: It accurately predicted 295 instances of Black Angus, but misclassified 14 as Holstein Friesian, 3 as Red Chittagong, 1 as Shahiwal, and 6 as White Native.

Red Chittagong: The model correctly identified 284 instances of 'Red Chittagong', but misclassified 2 as Shahiwal.

Shahiwal: It accurately predicted 317 instances of 'Shahiwal' with 1 misclassification as Holstein Friesian.

White Native: The model accurately predicted 267 instances of 'White Native', but misclassified 4 as Holstein Friesian, 6 as Black Angus, 14 as Red Chittagong, and 3 as Shahiwal.

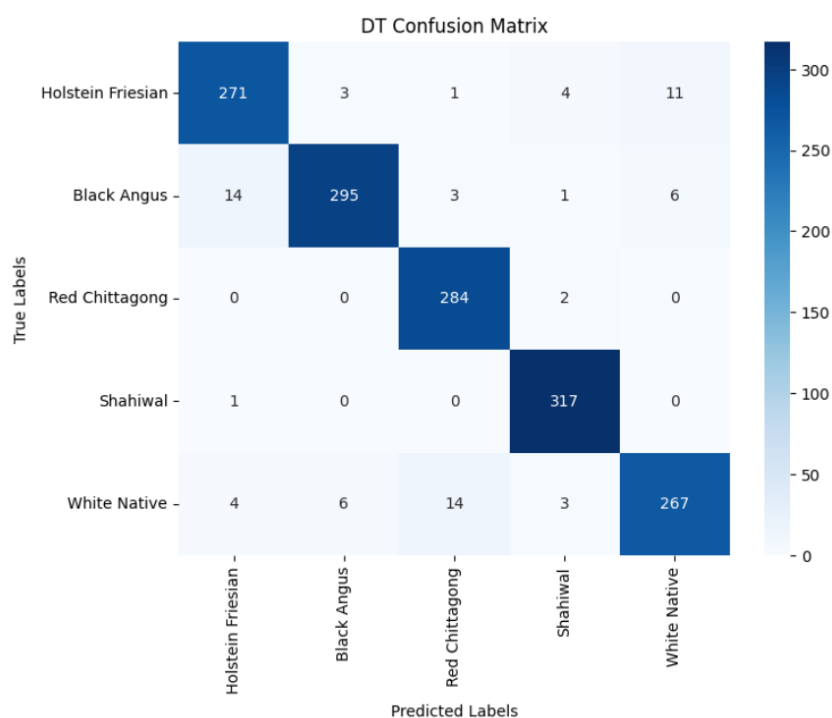


Figure 22: Confusion Matrix of DT

The Decision Tree model demonstrated a strong overall accuracy of 95% in effectively categorizing cattle breeds. Notably, it showcased impressive precision rates across diverse breeds, with high precision for Red Chittagong (94%), White Native (94%), Holstein Friesian (93%), Black Angus (97%), and Sahiwal (97%).

Additionally, the model exhibited notable recall percentages for Red Chittagong (99%), White Native (91%), Holstein Friesian (93%), Black Angus (92%), and Sahiwal (100%) breeds.

Moreover, the Decision Tree model achieved commendable F1 scores for Red Chittagong (97%), White Native (92%), Holstein Friesian (93%), Black Angus (95%), and Sahiwal (98%) breeds.

DT Classification Report:				
	precision	recall	f1-score	support
Holstein Friesian	0.93	0.93	0.93	290
Black Angus	0.97	0.92	0.95	319
Red Chittagong	0.94	0.99	0.97	286
Shahiwal	0.97	1.00	0.98	318
White Native	0.94	0.91	0.92	294
accuracy			0.95	1507
macro avg	0.95	0.95	0.95	1507
weighted avg	0.95	0.95	0.95	1507

Figure 23: Classification Report of DT

6.1.3 K-nearest Neighbors (KNN)

Holstein Friesian: The KNN model correctly classified 257 instances as Holstein Friesian, while misclassifying 15 as 'Black Angus' and 18 as 'White Native'. There were no instances misclassified as 'Red Chittagong' or 'Shahiwal'.

Black Angus: It correctly predicted 296 instances as 'Black Angus', misclassifying 3 as Holstein Friesian, 5 as 'Red Chittagong', 11 as 'Shahiwal', and 4 as 'White Native'.

Red Chittagong: The model accurately classified 254 instances as 'Red Chittagong', misclassifying 2 as Holstein Friesian, 8 as 'Black Angus', and 22 as 'White Native'. No instances were misclassified as 'Shahiwal'.

Shahiwal: It correctly classified all 314 instances as 'Shahiwal', with only 4 instances misclassified as 'Black Angus'.

White Native: The KNN model correctly predicted 248 instances as 'White Native', misclassifying 5 as Holstein Friesian, 18 as 'Black Angus', 17 as 'Red Chittagong', and 6 as 'Shahiwal'.

		KNN Confusion Matrix				
True Labels	Holstein Friesian	257	15	0	0	18
	Black Angus	3	296	5	11	4
	Red Chittagong	2	8	254	0	22
	Shahiwal	0	4	0	314	0
	White Native	5	18	17	6	248
		Holstein Friesian	Black Angus	Red Chittagong	Shahiwal	White Native

Figure 24: Confusion Matrix of KNN

The K-Nearest Neighbors (KNN) model demonstrated a robust overall accuracy of 91% in effectively categorizing cattle breeds. Particularly noteworthy were its precision rates across various breeds, exhibiting high precision for Red Chittagong (92%), White Native (85%), Holstein Friesian (96%), Black Angus (87%), and Sahiwal (95%).

Additionally, the model displayed notable recall percentages for Red Chittagong (89%), White Native (84%), Holstein Friesian (89%), Black Angus (93%), and Sahiwal (99%) breeds.

Furthermore, the KNN model achieved commendable F1 scores for Red Chittagong (90%), White Native (85%), Holstein Friesian (92%), Black Angus (90%), and Sahiwal (97%) breeds.

KNN Classification Report:				
	precision	recall	f1-score	support
Holstein Friesian	0.96	0.89	0.92	290
Black Angus	0.87	0.93	0.90	319
Red Chittagong	0.92	0.89	0.90	286
Shahiwal	0.95	0.99	0.97	318
White Native	0.85	0.84	0.85	294
accuracy			0.91	1507
macro avg	0.91	0.91	0.91	1507
weighted avg	0.91	0.91	0.91	1507

Figure 25: Classification Report of KNN

6.1.4 Random Forest (RF)

Holstein Friesian: The model accurately classified 289 instances of the Holstein Friesian breed. However, it misclassified only 1 instance as Black Angus.

Black Angus: It correctly predicted 305 instances of Black Angus. Yet, there were misclassifications, including 3 instances as Holstein Friesian, 4 as Red Chittagong, 4 as Shahiwal, and 3 as White Native.

Red Chittagong: The model performed well, correctly identifying 284 instances of Red Chittagong. However, it misclassified 2 instances as White Native.

Shahiwal: It accurately predicted 316 instances of Shahiwal with minimal misclassification, only 1 instance each as Black Angus and White Native.

White Native: The model correctly classified 284 instances of White Native. Nevertheless, it misclassified 2 instances as Shahiwal, and 8 instances as Red Chittagong.

Random Forest Confusion Matrix					
True Labels	Holstein Friesian	Black Angus	Red Chittagong	Shahiwal	White Native
	289	1	0	0	0
	3	305	4	4	3
	0	0	284	0	2
	0	1	0	316	1
	0	0	8	2	284
Predicted Labels					

Figure 26: Confusion Matrix of RF

The Random Forest model showcased a robust overall accuracy of 98% in efficiently categorizing cattle breeds. Notably, it demonstrated impressive precision rates across various breeds, with notable precision for Red Chittagong (96%), White Native (98%), Holstein Friesian (99%), Black Angus (99%), and Sahiwal (98%).

The model displayed noteworthy recall percentages for Red Chittagong (99%), White Native (97%), Holstein Friesian (100%), Black Angus (96%), and Sahiwal (99%) breeds.

Moreover, commendable F1 scores were achieved for Red Chittagong (98%), White Native (97%), Holstein Friesian (99%), Black Angus (97%), and Sahiwal (99%) breeds.

Random Forest Classification Report:				
	precision	recall	f1-score	support
Holstein Friesian	0.99	1.00	0.99	290
Black Angus	0.99	0.96	0.97	319
Red Chittagong	0.96	0.99	0.98	286
Shahiwal	0.98	0.99	0.99	318
White Native	0.98	0.97	0.97	294
accuracy			0.98	1507
macro avg	0.98	0.98	0.98	1507
weighted avg	0.98	0.98	0.98	1507

Figure 27: Classification Report of RF

6.1.5 Convolutional Neural Network (CNN)

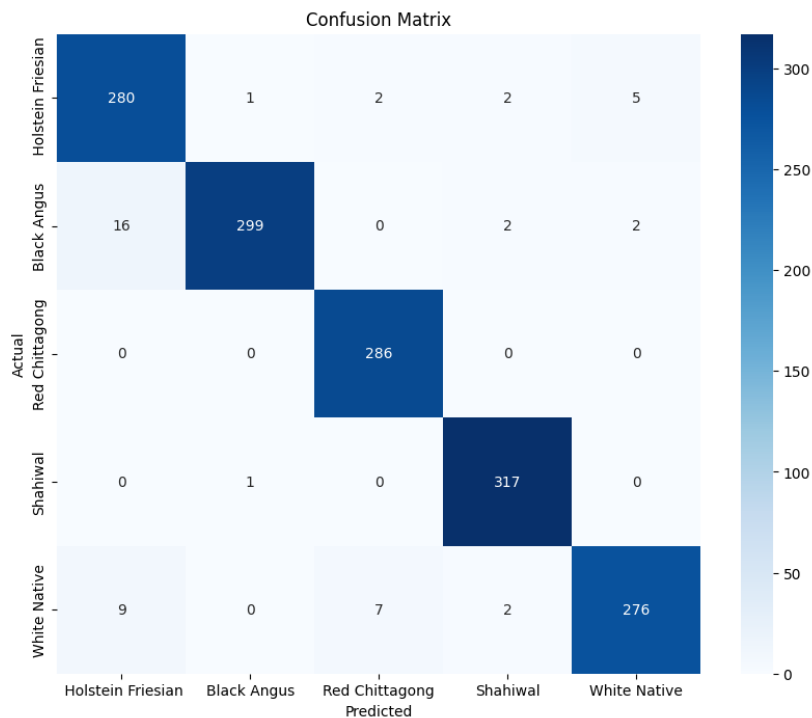


Figure 28: Confusion Matrix of CNN

Holstein Friesian: The CNN correctly predicted 280 instances as Holstein Friesian. However, it misclassified 1 instances as 'Black Angus', 2 instances as 'Red Chittagong', 2 instances as 'Shahiwal' and 5 instances as 'White Native'.

Black Angus: For 'Black Angus', the CNN accurately predicted 299 instances. But, it misclassified 16 instances as Holstein Friesian, 2 instances as 'Shahiwal' and 2 instances as 'White Native'.

Red Chittagong: The CNN performed well, correctly predicting 286 instances of 'Red Chittagong'.

Shahiwal: The model correctly predicted 317 instances as 'Shahiwal'. It misclassified only 1 instances as Black Angus.

White Native: The CNN correctly predicted 276 instances as 'White Native'. Nonetheless, it misclassified 9 instances as Holstein Friesian, 7 instances as 'Red Chittagong' and 2 instances as 'Shahiwal'.

Classification Report:				
	precision	recall	f1-score	support
Holstein Friesian	0.92	0.97	0.94	290
Black Angus	0.99	0.94	0.96	319
Red Chittagong	0.97	1.00	0.98	286
Shahiwal	0.98	1.00	0.99	318
White Native	0.98	0.94	0.96	294
accuracy			0.97	1507
macro avg	0.97	0.97	0.97	1507
weighted avg	0.97	0.97	0.97	1507

Figure 29: Classification Report of CNN

The CNN model demonstrated an exceptional overall accuracy of 97% in efficiently categorizing cattle breeds. Notably, it exhibited remarkable precision rates across various breeds, showcasing high precision for Red Chittagong (97%), White Native (98%), Holstein Friesian (92%), Black Angus (99%), and Sahiwal (98%).

The recall percentages for Red Chittagong (100%), White Native (94%), Holstein Friesian (97%), Black Angus (94%), and Sahiwal (100%) breeds.

F1 scores were achieved for Red Chittagong (98%), White Native (96%), Holstein Friesian (94%), Black Angus (96%), and Sahiwal (99%) breeds.

6.2 Comparative Analysis

- **Random Forest** shows the highest values across all metrics: precision (0.9810), recall (0.9808), F1-score (0.9807), and accuracy (98%). This indicates it performs consistently well across various evaluation aspects.

- **Support Vector Machine (SVM)** and **CNN** both demonstrate strong performances, achieving high accuracy (97%) and competitive precision, recall, and F1-score values.
- **Decision Tree** performs slightly lower than SVM and CNN but still maintains a respectable performance with precision, recall, and F1-score around 95%.
- **K-nearest Neighbor (KNN)** exhibits the lowest performance among the models listed, with precision, recall, and F1-score around 91%. However, it's still reasonably effective, especially depending on the specific context or requirements of the task.

In summary, the **Random Forest** model appears to be the top performer across the board in this evaluation, closely followed by SVM and CNN. The Decision Tree model also shows decent performance, while KNN lags slightly behind in these metrics. The choice of the best model can depend on the specific needs of the problem at hand, considering factors like accuracy, precision, recall, and the trade-offs between these metrics.

Table 4: Comparison of different model

Model	precision	recall	f1-score	Accuracy
Support Vector Machine	0.9731	0.9728	0.9728	97%
Decision Tree	0.9517	0.9516	0.9513	95%
K-nearest Neighbor	0.9095	0.9084	0.9083	91%
Random Forest	0.9810	0.9808	0.9807	98%
CNN	0.9675	0.9677	0.9672	97%

6.2.1 Accuracy Comparison

Random Forest achieves the highest accuracy of 98%. This model demonstrates the highest level of correctness in its predictions compared to the other models evaluated in this comparison. It's particularly effective in classifying instances accurately. Both **Support Vector Machine (SVM)** and **CNN** follow closely with an accuracy of 97%. They perform exceptionally well, almost at par with Random Forest, indicating their robustness in making correct predictions across various instances. **Decision Tree** maintains a slightly lower accuracy of 95%. While still providing commendable results, it falls behind Random Forest, SVM, and CNN in terms of overall correctness in predictions. **K-nearest neighbors (KNN)** shows the lowest accuracy among the models at 91%.

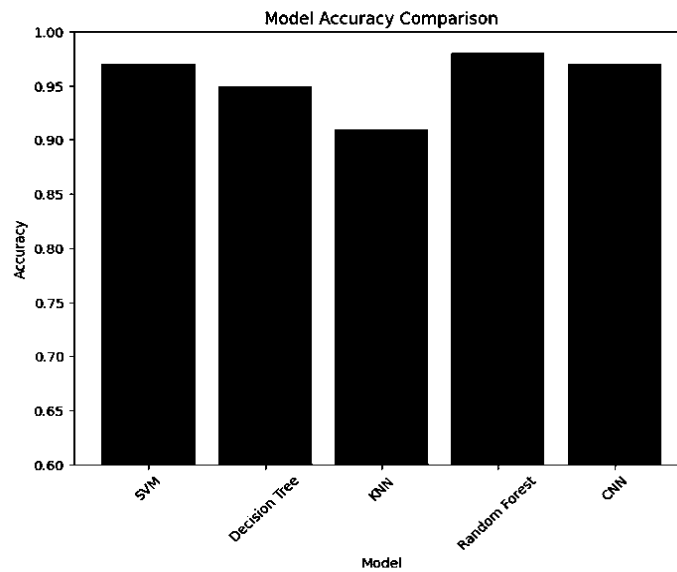


Figure 30: Accuracy Comparison of each model

Random Forest performs best in making correct predictions, closely followed by SVM and CNN. Decision Tree holds a respectable accuracy, though slightly lower, while KNN lags behind in terms of overall correctness in predictions based metric.

6.2.2 Precision Comparison

Random Forest exhibits the highest precision of 0.9810, indicating that when it predicts a positive outcome, it is highly accurate. **Support Vector Machine (SVM)** and **CNN** follow closely with precision values of 0.9731 and 0.9675, respectively. Both models demonstrate high precision in their positive predictions. **Decision Tree** maintains a slightly lower precision of 0.9517, which indicates it's slightly less accurate in correctly predicting positive cases compared to Random Forest, SVM, and CNN. **K-nearest Neighbor (KNN)** shows the lowest precision among these models, with a value of 0.9095.

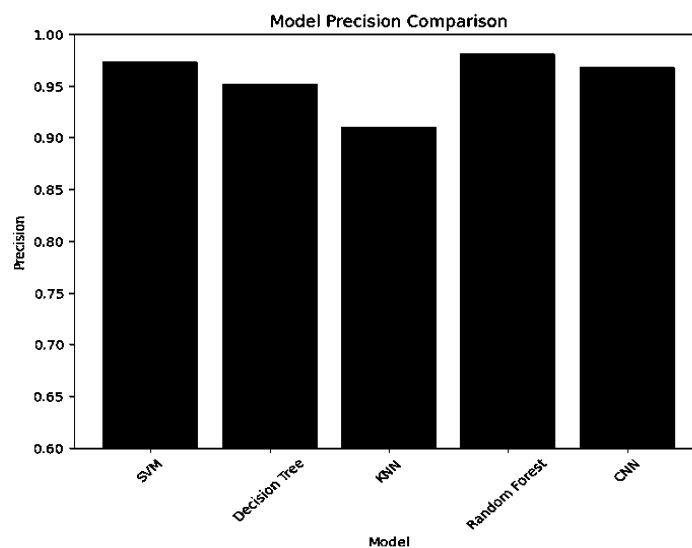


Figure 31: Precision Comparison of each model

Random Forest stands out with the highest precision score, followed by SVM and CNN. Decision Tree follows, while KNN shows the lowest precision among these models.

6.2.3 Recall Comparison

Random Forest demonstrates the highest recall of 0.9808, indicating its ability to effectively capture most of the positive instances within the dataset. **Support Vector Machine (SVM)** closely follows with a recall of 0.9728, showing a high ability to find relevant positive instances, just slightly behind Random Forest. **CNN** exhibits a recall of 0.9677, indicating it performs well in capturing positive instances but slightly lower than SVM and Random Forest. **Decision Tree** maintains a recall of 0.9516, indicating it captures positive instances slightly less effectively compared to the aforementioned models. **K-nearest Neighbor (KNN)** shows the lowest recall among these models with a value of 0.9084.

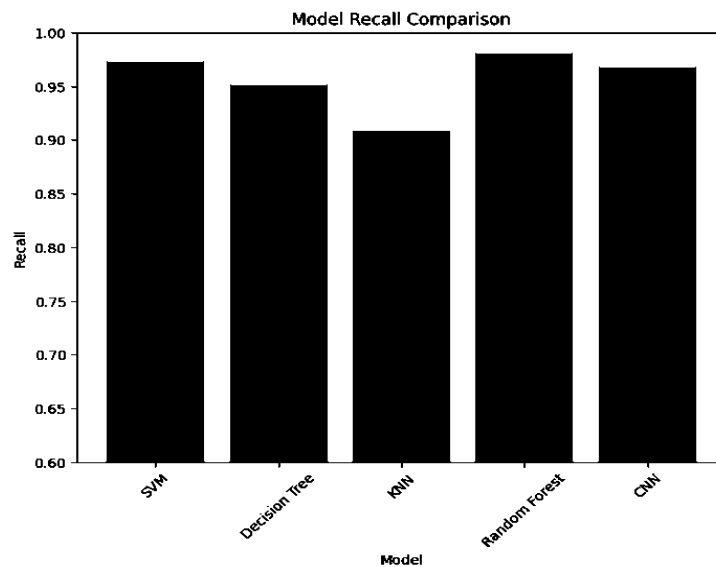


Figure 32: Recall Comparison of each model

In terms of recall, Random Forest performs the best, followed closely by SVM and CNN. Decision Tree follows, while KNN demonstrates the lowest recall among these models.

6.2.4 F1-Score Comparison

Random Forest showcases the highest F1-score of 0.9807, indicating a strong balance between precision and recall, resulting in effective and accurate predictions. **Support Vector Machine (SVM)** follows closely with an F1-score of 0.9728, showing a good balance between precision and recall, although slightly lower than Random Forest. **CNN** demonstrates an F1-score of 0.9672, indicating a strong balance between precision and recall, slightly lower than SVM but still performing well. **Decision Tree** maintains an F1-score of 0.9513, suggesting a reasonable balance between precision and recall but comparatively lower than Random Forest, SVM, and CNN. **K-nearest Neighbor (KNN)**

shows the lowest F1-score among these models with a value of 0.9083, indicating a comparatively lower balance between precision and recall.

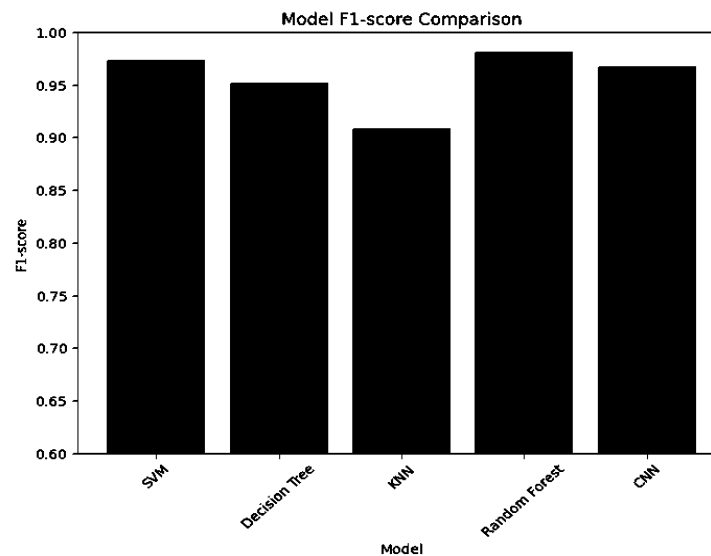


Figure 33: F1-Score Comparison of each model

In terms of F1-score, Random Forest leads the group, closely followed by SVM and CNN. Decision Tree follows, while KNN demonstrates the lowest F1-score among these models, suggesting a less balanced performance between precision and recall.

6.3 ROC-AUC Curve

6.3.1 CNN Curve

Red Chittagong and Sahiwal Breeds (AUC-ROC 1.00): The CNN model achieved a perfect AUC-ROC value of 1.00 for Red Chittagong and Sahiwal breeds. This signifies the model's exceptional ability to accurately differentiate and classify instances of these breeds from others, demonstrating high precision and confidence without any misclassifications.

Black Angus (AUC-ROC 0.98): The CNN model performed well with a high AUC-ROC value of 0.98 for Black Angus. This indicates a strong ability to distinguish instances of this breed with high accuracy and reliability.

Holstein Friesian and White Native Breeds (AUC-ROC around 0.97 and 0.96): These breeds attained AUC-ROC values around 0.97 and 0.96, respectively. While slightly lower than some other breeds, these values still signify good discriminatory capability by the CNN model. It demonstrates the model's ability to effectively differentiate between these breeds with relatively high accuracy.

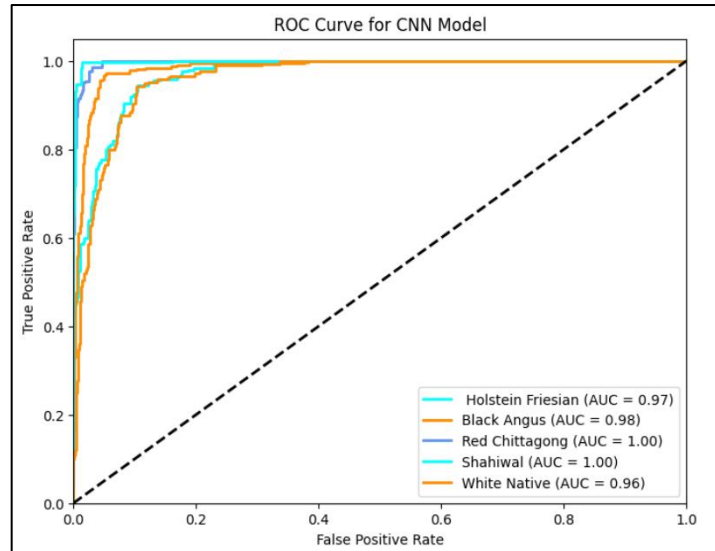


Figure 34: CNN AUC-RUC Curve

6.3.2 SVM Curve

Red Chittagong and Sahiwal Breeds: The SVM model achieved a perfect AUC-ROC value of 1.00 for these breeds. This signifies that the model could effectively classify between instances of these breeds with no misclassifications. It implies an exceptional ability to separate these breeds from others, demonstrating strong predictive power and high confidence in the classification.

Holstein Friesian, Black Angus, and White Native Breeds: These breeds attained a very high AUC-ROC value of 0.99. Though not perfect, this near-perfect score indicates an extremely accurate classification by the SVM model. It showcases the model's excellent discriminatory capability in distinguishing these breeds, with minimal misclassifications or errors.

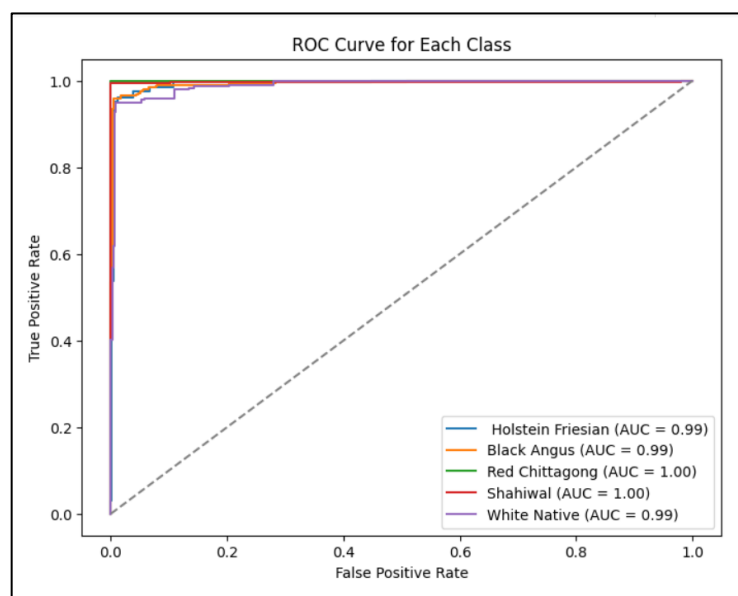


Figure 35: SVM AUC-ROC Curve

6.3.3 Decision Tree Curve

Red Chittagong and Sahiwal Breeds (AUC-ROC 0.99): The Decision Tree model achieved a high AUC-ROC value of 0.99 for these breeds. This indicates the model's excellent ability to distinguish and classify instances of these breeds. It demonstrates strong predictive power and suggests the model's high accuracy in separating these breeds from others.

Holstein Friesian, Black Angus, and White Native Breeds (AUC-ROC around 0.96): These breeds attained AUC-ROC values around 0.96. Though slightly lower than Red Chittagong and Sahiwal, these values still signify very good discriminatory capability by the Decision Tree model. It indicates the model's ability to effectively differentiate between these breeds with high accuracy and reliability.

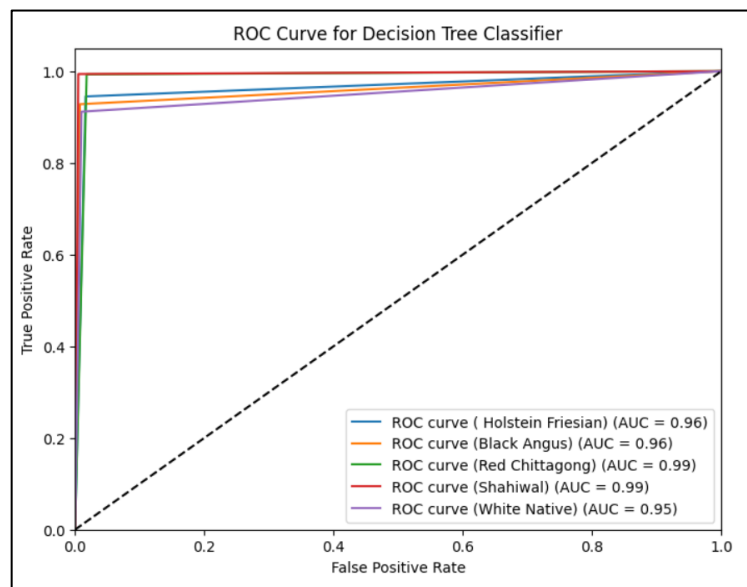


Figure 36: Decision Tree AUC-ROC Curve

6.3.4 KNN Curve

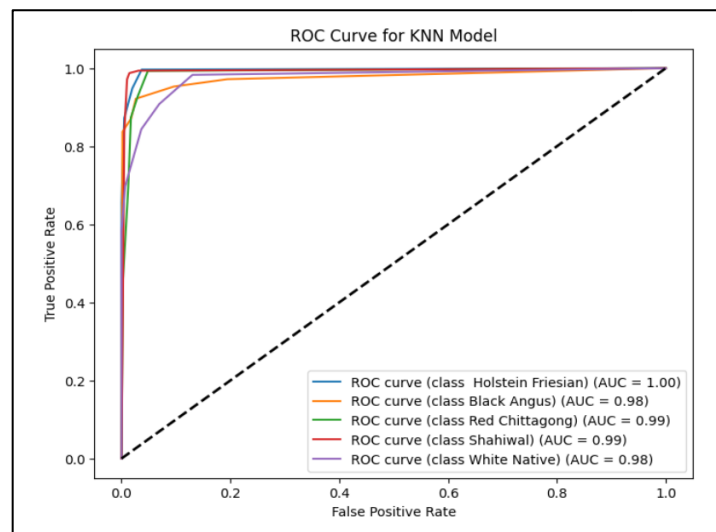


Figure 37: KNN AUC-ROC Curve

Holstein Friesian Breed (AUC-ROC 1.00): The k-NN model achieved a perfect AUC-ROC value of 1.00 for Holstein Friesian. This implies that the model flawlessly separated instances of this breed from others without any misclassifications. It demonstrates an exceptional ability to classify this breed with high precision and confidence.

Red Chittagong and Sahiwal Breeds (AUC-ROC 0.99): Both Red Chittagong and Sahiwal breeds achieved a very high AUC-ROC value of 0.99. This indicates an excellent ability of the k-NN model to distinguish and classify instances of these breeds with high accuracy and reliability.

Black Angus and White Native Breeds (AUC-ROC around 0.98): These breeds obtained AUC-ROC values around 0.98. These values signify a very good discriminatory capability by the k-NN model. It demonstrates the model's strong ability to effectively differentiate between these breeds with high accuracy.

6.3.5 Random Forest Curve

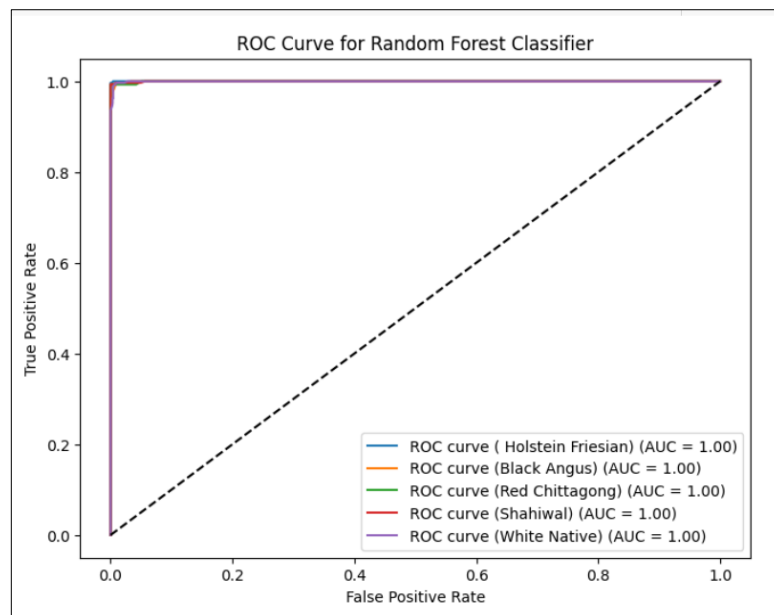


Figure 38: Random Forest AUC-ROC Curve

Each breed (Holstein Friesian, Black Angus, Red Chittagong, Sahiwal, White Native), the Random Forest model achieved a perfect AUC-ROC value of 1.00. This signifies flawless discrimination and classification by the model for all breeds. It demonstrates an exceptional ability to accurately classify each breed with high precision and confidence, without any misclassifications.

6.3.6 Comparative Insights of ROC-AUC Curve

Table 5: Comparative Analysis

Cattle Breed	CNN	SVM	KNN	Decision Tree	Random Forest
Holstein Friesian	0.97	0.99	1.00	0.96	1.00
Black Angus	0.98	0.99	0.98	0.96	1.00
Red Chittagong	1.00	1.00	0.99	0.99	1.00
Sahiwal	1.00	1.00	0.99	0.99	1.00
White Native	0.96	0.99	0.98	0.95	1.00

- **Consistency in High Performance:** Random Forest and some instances of SVM, k-NN, and CNN models consistently achieved perfect or near-perfect scores (1.00 or close to it) across multiple breeds (especially Red Chittagong and Sahiwal), indicating robust performance in breed classification.
- **Variance in Model Performance:** Decision Tree and some instances of SVM and CNN showed slightly lower AUC-ROC values (around 0.95 to 0.99), suggesting relatively strong but not perfect discriminatory power across the breeds.
- **Overall Model Performance:** Random Forest and certain instances of SVM, k-NN, and CNN models showcased exceptional classification abilities across all breeds. However, CNN, while performing well, had slightly lower scores compared to Random Forest, SVM, and k-NN in some cases.
- **Consideration for Specific Breeds:** Different models exhibited varying degrees of performance across breeds. Red Chittagong and Sahiwal breeds were consistently well-classified by most models, while some variability existed for other breeds among the models.

6.4 Observation of the Result

1. **Model Performance Variation:** There's observable variation in the performance metrics among the models. Random Forest stands out with the highest precision, recall, F1-score, and accuracy. SVM follows closely, demonstrating strong precision and overall accuracy. Meanwhile, Decision Tree and KNN exhibit consistent but slightly lower performance across these metrics.
2. **Random Forest Excellence:** The Random Forest model excelled in predicting cattle breeds by combining multiple decision trees to improve accuracy while reducing over-fitting. Its ability to understand intricate data relationships and capture subtle breed patterns, alongside aggregated randomized trees, ensures reliability and precision. By identifying crucial features, it strengthens its exceptional performance in distinguishing between different breeds.

3. **SVM vs. CNN:** SVM, a traditional machine learning method, showed competitive performance similar to the CNN, a deep learning architecture. The expectation might have been that CNN, being a more complex neural network, would significantly outperform SVM in this visual recognition task.
4. **Consistent Performance of Decision Tree and KNN:** Decision Tree and KNN exhibit balanced performance in precision, recall, and accuracy, although they show slightly lower metrics compared to other models. This suggests their stability but might indicate limitations in handling certain breed variations.
5. **High Precision Across Models:** All models demonstrate relatively high precision, indicating their ability to make accurate positive predictions for cattle breeds. This suggests that when these models predict a specific breed, they are often correct.
6. **Consistent Recall Rates:** The recall rates among the models are also generally high and consistent, showing their ability to capture a significant proportion of positive instances for each cattle breed. This indicates that these models effectively identify instances belonging to specific breeds.
7. **Insights for Livestock Management:** Leveraging models like SVM or CNN, which demonstrated competitive performance, could offer valuable insights for automating livestock monitoring systems or breed identification tasks in agricultural settings.
8. **Improved Breed Classification Techniques:** The superior performance of Random Forest in accurately identifying cattle breeds implies its potential as a robust technique for breed classification. Implementing ensemble learning methods like Random Forest could substantially enhance cattle breed classification accuracy, aiding livestock management and breeding programs.
9. **Hybrid Model Approaches:** Investigate the potential of combining the strengths of different models. For instance, blending the interpretability of Decision Trees with the predictive power of neural networks might yield more accurate and interpretable models for cattle breed classification.
10. **Data Augmentation and Pre-processing:** Augment the dataset to include variations in cattle breed images. Robust pre-processing techniques that handle image quality, background variations, or lighting conditions could improve model generalization.
11. **Class Imbalance Handling:** Address potential class imbalances within the dataset. Implement techniques like oversampling minority classes or adjusting class weights to prevent bias in models towards majority classes.

In summary, the results highlight Random Forest as the most effective model, closely followed by SVM and CNN, while Decision Tree and KNN show consistent but slightly lower performance in accurately classifying cattle breeds.

CHAPTER 7: CONCLUSION

In summary, our thesis on "Cattle Breed Classification Using Multiple Approaches" represents a significant stride in the field of livestock classification. With a rich dataset and the utilization of diverse machine learning models, including CNN, SVM, Decision Tree, KNN, and Random Forest, our approach excels in accuracy. Notably, our SVM outperforms comparative studies, and the integration of data augmentation enhances overall model performance. The exceptional accuracy of the Random Forest model further distinguishes our work. This research not only contributes to automated cattle breed identification but also sets a high standard for future endeavors, demonstrating the efficacy of our methodology in advancing precision livestock management.

7.1 Limitations

- The dataset may not fully encompass global cattle breeds, impacting the model's versatility.
- Our models are limited to recognizing only five distinct cattle types.

7.2 Challenges Faced

The pursuit of excellence in cattle breed classification led us through a challenging yet rewarding journey, where we encountered and overcame several hurdles. This section encapsulates the difficulties faced.

- **Data Quality Assurance:** The collection of a substantial dataset, comprising 10,000 cow images, presented a commendable opportunity but also posed challenges in ensuring data quality.
- **Dataset Categories:** Initially categorizing the large dataset(5025) into numerous categories resulted in challenges, particularly with limited training data for each category. To enhance model accuracy, we strategically re-categorized the dataset into five distinct categories, achieving a significant improvement in accuracy.
- **Computational Efficiency:** Execution of our code in the Google Colab environment revealed computational challenges, with code execution requiring an extended period. The prolonged runtime, extending up to 8 hours for such a large dataset, posed concentration challenges.
- **Model Accuracy Enhancement:** Attaining high accuracy in our chosen model proved to be a demanding task, requiring persistent effort and iterative refinement.
- **Exploration of Deep Learning Models:** We were genuinely eager to delve into deep learning models. However, due to our limited knowledge in that domain, we encountered challenges that prevented us from successfully implementing them.

7.3 Future Recommendation

- Continuously update the dataset to encompass a more comprehensive range of cattle variations and environmental conditions.
- Explore the integration of advanced deep learning architectures to enhance classification accuracy.
- Explore the potential for transfer learning to leverage pre-trained models for improved performance.

REFERENCES

1. Briggs, H.M. & D.M. Briggs. Modern Breeds of Livestock. Fourth Edition. Macmillan Publishing Co. 1980 (reprinted with permission from Dr. Briggs).
2. "Breeds of Livestock - Bengali Cattle". Department of Animal and Food Sciences. Oklahoma State University. Retrieved December 22, 2018.
3. Afroz, M. A., Hoque, M. A., and Bhuiyan, A. K. F. H. (2011). Estimation of heritability for growth traits of Red Chittagong cattle in a nucleus herd. *Bang. Veterinar.* 28, 39–46. doi: 10.3329/bvet.v28i1.8812
4. Patel N, Upadhyay S. Study of various decision tree pruning methods with their empirical comparison in WEKA. *Int J Comp Appl.* 60(12):20–25. [Google Scholar]
5. “Buy upgraded Philippine native cattle for sale,” *Alpha Adventure Farms*, 30-Oct-2023. [Online]. Available: <https://www.alphaagventure.com/philippine-native-cattle/>. [Accessed: 19-Dec-2023].
6. *Roysfarm.com*. [Online]. Available: <https://www.roysfarm.com/sahiwal-cattle/>. [Accessed: 19-Dec-2023].
7. G. goelaparna1520 Follow, “Convolutional neural network (CNN) in machine learning,” *GeeksforGeeks*, 25-Dec-2020. [Online]. Available: <https://www.geeksforgeeks.org/convolutional-neural-network-cnn-in-machine-learning/>. [Accessed: 19-Dec-2023].
8. A.Follow, “Support vector machine (SVM) algorithm,” *GeeksforGeeks*, 20-Jan-2021. [Online]. Available: <https://www.geeksforgeeks.org/support-vector-machine-algorithm/>. [Accessed: 19-Dec-2023].
9. *Researchgate.net*. [Online]. Available: https://www.researchgate.net/figure/General-architecture-of-a-support-vector-maching-SVM-model-according-to-55_fig3_348745187. [Accessed: 19-Dec-2023].

10. A.Raj, "An exhaustive guide to Decision Tree classification in python 3.X," *Towards Data Science*, 27-Oct-2021. [Online]. Available: <https://towardsdatascience.com/an-exhaustive-guide-to-classification-using-decision-trees-8d472e77223f>. [Accessed: 19-Dec-2023].
11. "Decision Tree Algorithm in Machine Learning - javatpoint," *www.javatpoint.com*. [Online]. Available: <https://www.javatpoint.com/machine-learning-decision-tree-classification-algorithm>. [Accessed: 19-Dec-2023].
12. O.Mbaabu, "Introduction to random forest in machine learning," *Engineering Education (EngEd) Program / Section*. [Online]. Available: <https://www.section.io/engineering-education/introduction-to-random-forest-in-machine-learning/>. [Accessed: 19-Dec-2023].
13. "K-Nearest Neighbor(KNN) Algorithm for Machine Learning," *www.javatpoint.com*. [Online]. Available: <https://www.javatpoint.com/k-nearest-neighbor-algorithm-for-machine-learning>. [Accessed: 19-Dec-2023].
14. O.Harrison, "Machine learning basics with the K-nearest neighbors algorithm," *Towards Data Science*, 10-Sep-2018. [Online]. Available: <https://towardsdatascience.com/machine-learning-basics-with-the-k-nearest-neighbors-algorithm-6a6e71d01761>. [Accessed: 19-Dec-2023].
15. N.Donges, "Random forest: A complete guide for machine learning," *Built In*, 22-Jul-2021. [Online]. Available: <https://builtin.com/data-science/random-forest-algorithm>. [Accessed: 19-Dec-2023].

APPENDIX

```
from google.colab import drive
drive.mount('/content/drive')
dirname = '/content/drive/My Drive/CattleGroup'

import matplotlib.pyplot as plt
import numpy as np
import os
import PIL
import cv2
import glob
import string
from PIL import Image
from pathlib import Path
import tensorflow as tf
from tensorflow import keras
from tensorflow.keras import layers
from tensorflow.keras.models import Sequential
from tensorflow.keras.preprocessing import image
from mlxtend.plotting import plot_decision_regions
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
import seaborn as sns
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.utils.multiclass import unique_labels
from sklearn import metrics

data_dir = Path(dirname)
dirs = data_dir.glob("*")
labels_dict = { " Holstein Friesian": 0, "Black Angus": 1, "Red Chittagong": 2, "Shahiwal":
3, "White Native": 4 }
labels2breed = {0:" Holstein Friesian", 1:"Black Angus", 2:"Red Chittagong",
3:"Shahiwal", 4:"White Native"}
labels = []
image_data = []
for label in dirs:
    image_label = label.parts[-1]
    cnt = 0
```

```

print(image_label)
for img_path in label.glob("*.jpg"):
    img = image.load_img(img_path, target_size=(100, 100))
    image_array = image.img_to_array(img)
    image_data.append(image_array)
    if image_label in labels_dict:
        labels.append(labels_dict[image_label])
    else:
        print("Label not found:", image_label)
    cnt += 1
print(cnt)

batch_size = 12
img_height = 180
img_width = 180
train_ds = tf.keras.utils.image_dataset_from_directory(
    dirname,
    validation_split=0.3,
    subset="training",
    seed=123,
    image_size=(img_height, img_width),
    batch_size=batch_size)
test_ds = tf.keras.utils.image_dataset_from_directory(
    dirname,
    validation_split=0.3,
    subset="validation",
    seed=123,
    image_size=(img_height, img_width),
    batch_size=batch_size)

class_names = train_ds.class_names
print(class_names)
# Get class names
class_names = train_ds.class_names
# Extract x_train, y_train from train_ds
x_train = []
y_train = []
for images, labels in train_ds:
    x_train.append(images)

```



```

    y_train.append(labels)
x_train = np.concatenate(x_train)
y_train = np.concatenate(y_train)
# Extract x_test, y_test from test_ds
x_test = []
y_test = []
for images, labels in test_ds:
    x_test.append(images)
    y_test.append(labels)
x_test = np.concatenate(x_test)
y_test = np.concatenate(y_test)

plt.figure(figsize=(10, 10))
for images, labels in train_ds.take(1):
    for i in range(9):
        ax = plt.subplot(3, 3, i + 1)
        plt.imshow(images[i].numpy().astype("uint8"))
        plt.title(class_names[labels[i]])
        plt.axis("off")

for image_batch, labels_batch in train_ds:
    print(image_batch.shape)
    print(labels_batch.shape)
    break

AUTOTUNE = tf.data.AUTOTUNE
train_ds = train_ds.cache().shuffle(1000).prefetch(buffer_size=AUTOTUNE)
test_ds = test_ds.cache().prefetch(buffer_size=AUTOTUNE)
normalization_layer = layers.Rescaling(1./255)
normalized_ds = train_ds.map(lambda x, y: (normalization_layer(x), y))
image_batch, labels_batch = next(iter(normalized_ds))
first_image = image_batch[0]
print(np.min(first_image), np.max(first_image))

```

CNN Model

```

num_classes = len(class_names)
model = Sequential([
    layers.Rescaling(1./255, input_shape=(img_height, img_width, 3)),
    layers.Conv2D(16, 3, padding='same', activation='relu'),

```

```

layers.MaxPooling2D(),
layers.Conv2D(32, 3, padding='same', activation='relu'),
layers.MaxPooling2D(),
layers.Conv2D(64, 3, padding='same', activation='relu'),
layers.MaxPooling2D(),
layers.Flatten(),
layers.Dense(128, activation='relu'),
layers.Dense(num_classes)
])

model.compile(optimizer='adam',
              loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True),
              metrics=['accuracy'])
epochs=10
history = model.fit(
    train_ds,
    validation_data=test_ds,
    epochs=epochs
)

acc = history.history['accuracy']
val_acc = history.history['val_accuracy']
loss = history.history['loss']
val_loss = history.history['val_loss']
epochs_range = range(epochs)
plt.figure(figsize=(5, 5))
plt.subplot(1, 2, 1)
plt.plot(epochs_range, acc, label='Training Accuracy')
plt.plot(epochs_range, val_acc, label='Testing Accuracy')
plt.legend(loc='lower right')
plt.title('Training and Testing Accuracy')
plt.subplot(1, 2, 2)
plt.plot(epochs_range, loss, label='Training Loss')
plt.plot(epochs_range, val_loss, label='Testing Loss')
plt.legend(loc='upper right')
plt.title('Training and Testing Loss')
plt.show()

data_augmentation = keras.Sequential(

```

```

[
    layers.RandomFlip("horizontal",
                      input_shape=(img_height,
                                    img_width,
                                    3)),
    layers.RandomRotation(0.1),
    layers.RandomZoom(0.1),
]
plt.figure(figsize=(10, 10))
for images, _ in train_ds.take(1):
    for i in range(9):
        augmented_images = data_augmentation(images)
        ax = plt.subplot(3, 3, i + 1)
        plt.imshow(augmented_images[0].numpy().astype("uint8"))
        plt.axis("off")

model = Sequential([
    data_augmentation,
    layers.Rescaling(1./255),
    layers.Conv2D(16, 3, padding='same', activation='relu'),
    layers.MaxPooling2D(),
    layers.Conv2D(32, 3, padding='same', activation='relu'),
    layers.MaxPooling2D(),
    layers.Conv2D(64, 3, padding='same', activation='relu'),
    layers.MaxPooling2D(),
    layers.Dropout(0.2),
    layers.Flatten(),
    layers.Dense(128, activation='relu'),
    layers.Dense(num_classes)
])

model.compile(optimizer='adam',
              loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True),
              metrics=['accuracy'])
epochs = 15
history = model.fit(
    train_ds,
    validation_data=test_ds,
    epochs=epochs

```

```

)
acc = history.history['accuracy']
val_acc = history.history['val_accuracy']
loss = history.history['loss']
val_loss = history.history['val_loss']
epochs_range = range(epochs)
plt.figure(figsize=(5, 5))
plt.subplot(1, 2, 1)
plt.plot(epochs_range, acc, label='Training Accuracy')
plt.plot(epochs_range, val_acc, label='Testing Accuracy')
plt.legend(loc='lower right')
plt.title('Training and Testing Accuracy')
plt.subplot(1, 2, 2)
plt.plot(epochs_range, loss, label='Training Loss')
plt.plot(epochs_range, val_loss, label='Testing Loss')
plt.legend(loc='upper right')
plt.title('Training and Testing Loss')
plt.subplots_adjust(wspace=0.4)
plt.tight_layout()
plt.show()
plt.show()

img=tf.keras.utils.load_img("/content/asia (477).jpg",target_size=(180,180))
plt.imshow(img)
img_array=tf.keras.utils.img_to_array(img)
img_array=tf.expand_dims(img_array,0)
pred=model.predict(img_array)
score=tf.nn.softmax(pred[0])
print(
    "This image belongs to { } and Accuracy {:.2f} %."
    .format(class_names[np.argmax(score)], 100 * np.max(score))
)

# Get predictions for the test dataset
predictions = model.predict(test_ds)
predicted_classes = np.argmax(predictions, axis=1)
# Get true labels from the test dataset
true_labels = []
for images, labels in test_ds:

```

```

    true_labels.extend(labels.numpy())
# Create confusion matrix
conf_matrix = confusion_matrix(true_labels, predicted_classes)
# Display the confusion matrix using a heatmap
plt.figure(figsize=(10, 8))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues',
            xticklabels=class_names, yticklabels=class_names)
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()
# Generate classification report
report = classification_report(true_labels, predicted_classes, target_names=class_names)
# Print the classification report
print("Classification Report:")
print(report)

from sklearn.metrics import classification_report
# Generate classification report
report = classification_report(true_labels, predicted_classes, target_names=class_names,
                              output_dict=True)
# Calculate overall precision, recall, F1-score
precision = report['macro avg']['precision']
recall = report['macro avg']['recall']
f1_score = report['macro avg']['f1-score']
print(f"Overall Precision: {precision:.4f}")
print(f"Overall Recall: {recall:.4f}")
print(f"Overall F1-Score: {f1_score:.4f}")

from sklearn.metrics import roc_curve, auc
from sklearn.preprocessing import label_binarize
from itertools import cycle
# Binarize the labels
y_test_bin = label_binarize(y_test, classes=np.arange(len(class_names)))
# Compute probabilities of the test set
probabilities = model.predict_proba(x_test)
# Compute ROC curve and ROC area for each class
fpr = dict()
tpr = dict()

```

```

roc_auc = dict()
for i in range(len(class_names)):
    fpr[i], tpr[i], _ = roc_curve(y_test_bin[:, i], probabilities[:, i])
    roc_auc[i] = auc(fpr[i], tpr[i])
# Plotting ROC curves
plt.figure(figsize=(8, 6))
colors = cycle(['aqua', 'darkorange', 'cornflowerblue']) # You may need more colors for
more classes
for i, color in zip(range(len(class_names)), colors):
    plt.plot(fpr[i], tpr[i], color=color, lw=2,
             label='{0} (AUC = {1:0.2f})'
             ".format(class_names[i], roc_auc[i]))
plt.plot([0, 1], [0, 1], 'k--', lw=2)
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve for CNN Model')
plt.legend(loc="lower right")
plt.show()

```

SVM Model

```

from sklearn import svm
# Data Preprocessing for Training
train_images = []
train_labels = []
for images, labels in train_ds:
    train_images.extend(images.numpy())
    train_labels.extend(labels.numpy())
train_images = np.array(train_images)
train_labels = np.array(train_labels)
train_images = train_images.reshape(train_images.shape[0], -1)
train_images = train_images.astype('float32') / 255.0
# Initialize and train the SVM model
svm_model = svm.SVC(kernel='linear')
svm_model.fit(train_images, train_labels)

# Data Preprocessing for testing
val_images = []

```

```

val_labels = []
for images, labels in test_ds:
    val_images.extend(images.numpy())
    val_labels.extend(labels.numpy())
val_images = np.array(val_images)
val_labels = np.array(val_labels)
val_images = val_images.reshape(val_images.shape[0], -1)
val_images = val_images.astype('float32') / 255.0
val_predictions = svm_model.predict(val_images)

# Prediction on a Single Image
image_path = '/content/asia (477).jpg'
img_height = 180
img_width = 180
image = Image.open(image_path)
image = image.resize((img_height, img_width))
image_array = np.array(image)
image_array = image_array.reshape(1, -1).astype('float32') / 255.0
predicted_class = svm_model.predict(image_array)[0]
predicted_class_name = class_names[predicted_class]
print(f"The predicted class for the unknown image is: {predicted_class_name}")

# Calculate accuracy
accuracy = accuracy_score(val_labels, val_predictions)
print(f"Accuracy: {accuracy * 100:.2f}%")
conf_matrix = confusion_matrix(val_labels, val_predictions)
# Plot confusion matrix using seaborn heatmap
plt.figure(figsize=(8, 6))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues',
            xticklabels=class_names, yticklabels=class_names)
plt.xlabel('Predicted')
plt.ylabel('True')
plt.title('Confusion Matrix')
plt.show()
# Generate classification report
class_report = classification_report(val_labels, val_predictions,
target_names=class_names)
# Display classification report with class names
print("Classification Report:")

```

```

print(class_report)

from sklearn.metrics import precision_recall_fscore_support
precision, recall, f1_score, _ = precision_recall_fscore_support(val_labels,
val_predictions, average='weighted')
print(f"Overall Precision: {precision:.4f}")
print(f"Overall Recall: {recall:.4f}")
print(f"Overall F1-score: {f1_score:.4f}")
from sklearn.metrics import roc_auc_score, roc_curve
from sklearn.preprocessing import label_binarize
# Binarize the labels for multi-class AUC-ROC calculation
num_classes = len(class_names)
binarized_val_labels = label_binarize(val_labels, classes=np.arange(num_classes))
val_probabilities = svm_model.decision_function(val_images)
# Calculate AUC-ROC for each class
auc_roc_per_class = []
plt.figure(figsize=(8, 6))
for i in range(num_classes):
    class_auc = roc_auc_score(binarized_val_labels[:, i], val_probabilities[:, i])
    auc_roc_per_class.append(class_auc)
    # Compute ROC curve for each class
    fpr, tpr, _ = roc_curve(binarized_val_labels[:, i], val_probabilities[:, i])
    plt.plot(fpr, tpr, label=f'{class_names[i]} (AUC = {class_auc:.2f})')
# Plot ROC curve for each class
plt.plot([0, 1], [0, 1], color='gray', linestyle='--')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve for Each Class')
plt.legend(loc='lower right')
plt.show()

```

Decision Tree

```

from sklearn.tree import DecisionTreeClassifier
# Data Preprocessing for Training
train_images = []
train_labels = []
for images, labels in train_ds:
    train_images.extend(images.numpy())
    train_labels.extend(labels.numpy())

```



```

train_images = np.array(train_images)
train_labels = np.array(train_labels)
train_images = train_images.reshape(train_images.shape[0], -1)
train_images = train_images.astype('float32') / 255.0
# Initialize and train the Decision Tree model
decision_tree = DecisionTreeClassifier()
decision_tree.fit(train_images, train_labels)

# Data Preprocessing for testing
val_images = []
val_labels = []
for images, labels in test_ds:
    val_images.extend(images.numpy())
    val_labels.extend(labels.numpy())
val_images = np.array(val_images)
val_labels = np.array(val_labels)
val_images = val_images.reshape(val_images.shape[0], -1)
val_images = val_images.astype('float32') / 255.0
# Prediction on a Single Image
image_path = '/content/asia (477).jpg'
img_height = 180
img_width = 180
image = Image.open(image_path)
image = image.resize((img_height, img_width))
image_array = np.array(image)
image_array = image_array.reshape(1, -1).astype('float32') / 255.0
predicted_class = decision_tree.predict(image_array)[0]
predicted_class_name = class_names[predicted_class]
print(f"The predicted class for the unknown image is: {predicted_class_name}")

# Testing Decision Tree model on the test set
val_predictions = decision_tree.predict(val_images)
# Calculate accuracy
accuracy = accuracy_score(val_labels, val_predictions)
print(f"Accuracy: {accuracy * 100:.2f}%")

# Plot confusion matrix
cm_knn = confusion_matrix(val_labels, val_predictions)
plt.figure(figsize=(8, 6))

```

```

sns.heatmap(cm_knn, annot=True, cmap='Blues', fmt='d', xticklabels=class_names,
yticklabels=class_names)
plt.title('DT Confusion Matrix')
plt.xlabel('Predicted Labels')
plt.ylabel('True Labels')
plt.show()
# Plot classification report
report_knn = classification_report(val_labels, val_predictions, target_names=class_names)
print("DT Classification Report:")
print(report_knn)

from sklearn.metrics import precision_recall_fscore_support
precision, recall, f1_score, _ = precision_recall_fscore_support(val_labels,
val_predictions, average='weighted')
print(f"Overall Precision: {precision:.4f}")
print(f"Overall Recall: {recall:.4f}")
print(f"Overall F1-score: {f1_score:.4f}")
from sklearn.metrics import roc_curve, auc
from sklearn.preprocessing import label_binarize

# Convert labels to one-hot encoded format
val_labels_bin = label_binarize(val_labels, classes=np.unique(val_labels))
# Get the decision scores for each class
val_scores = decision_tree.predict_proba(val_images)
# Compute ROC curve and ROC area for each class
n_classes = len(class_names)
fpr = dict()
tpr = dict()
roc_auc = dict()
for i in range(n_classes):
    fpr[i], tpr[i], _ = roc_curve(val_labels_bin[:, i], val_scores[:, i])
    roc_auc[i] = auc(fpr[i], tpr[i])
# Plot ROC curve
plt.figure(figsize=(8, 6))
for i in range(n_classes):
    plt.plot(fpr[i], tpr[i], label=f'ROC curve ({class_names[i]}) (AUC = {roc_auc[i]:.2f})')
plt.plot([0, 1], [0, 1], 'k--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])

```

```

plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve for Decision Tree Classifier')
plt.legend(loc="lower right")
plt.show()

```

KNN Model

```

from sklearn.neighbors import KNeighborsClassifier
# Extract features using a pre-trained CNN (e.g., VGG16)
base_model = tf.keras.applications.VGG16(weights='imagenet', include_top=False,
input_shape=(img_height, img_width, 3))
# Remove the classification layers
feature_extractor = tf.keras.models.Model(inputs=base_model.input,
outputs=base_model.get_layer('block5_pool').output)
# Extract features from the training dataset
train_features = []
train_labels = []
for images, labels in train_ds:
    features = feature_extractor.predict(images)
    train_features.append(features)
    train_labels.append(labels.numpy())
train_features = np.concatenate(train_features)
train_features = train_features.reshape(train_features.shape[0], -1) # Flatten the features
train_labels = np.concatenate(train_labels)
# Train KNN model
knn_model = KNeighborsClassifier(n_neighbors=5)
knn_model.fit(train_features, train_labels)
val_features = []
val_labels = []
for images, labels in test_ds:
    features = feature_extractor.predict(images)
    val_features.append(features)
    val_labels.append(labels.numpy())
val_features = np.concatenate(val_features)
val_features = val_features.reshape(val_features.shape[0], -1) # Flatten the features
val_labels = np.concatenate(val_labels)
val_predictions = knn_model.predict(val_features)
new_img = tf.keras.utils.load_img("/content/asia (477).jpg", target_size=(img_height,
img_width))

```

```

new_img_array = tf.keras.utils.img_to_array(new_img)
new_img_array = np.expand_dims(new_img_array, axis=0)
# Extract features from the new image
new_img_features = feature_extractor.predict(new_img_array)
new_img_features = new_img_features.reshape(1, -1)
# Use the KNN model to predict the class
predicted_class = knn_model.predict(new_img_features)
predicted_class_name = class_names[int(predicted_class)]
print(f"The predicted class for the unknown image is: {predicted_class_name}")

# Calculate accuracy
accuracy = accuracy_score(val_labels, val_predictions)
print(f" Accuracy: {accuracy * 100:.2f}%")
# Plot confusion matrix
cm_knn = confusion_matrix(val_labels, val_predictions)
plt.figure(figsize=(8, 6))
sns.heatmap(cm_knn, annot=True, cmap='Blues', fmt='d', cbar=False,
            xticklabels=class_names, yticklabels=class_names)
plt.title('KNN Confusion Matrix')
plt.xlabel('Predicted Labels')
plt.ylabel('True Labels')
plt.show()
# Plot classification report
report_knn = classification_report(val_labels, val_predictions, target_names=class_names)
print("KNN Classification Report:")
print(report_knn)
from sklearn.metrics import precision_recall_fscore_support
precision, recall, f1_score, _ = precision_recall_fscore_support(val_labels,
val_predictions, average='weighted')
print(f"Overall Precision: {precision:.4f}")
print(f"Overall Recall: {recall:.4f}")
print(f"Overall F1-score: {f1_score:.4f}")
from sklearn.metrics import roc_curve, auc
# Compute probabilities for each class
val_probs = knn_model.predict_proba(val_features)
# Compute ROC curve and AUC for each class
fpr = dict()
tpr = dict()
roc_auc = dict()

```

```

for i in range(len(class_names)):
    fpr[i], tpr[i], _ = roc_curve(val_labels == i, val_probs[:, i])
    roc_auc[i] = auc(fpr[i], tpr[i])
# Plot ROC curve for each class
plt.figure(figsize=(8, 6))
for i in range(len(class_names)):
    plt.plot(fpr[i], tpr[i], label=f'ROC curve (class {class_names[i]}) (AUC = {roc_auc[i]:.2f})')
plt.plot([0, 1], [0, 1], 'k--', lw=2)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve for KNN Model')
plt.legend(loc='lower right')
plt.show()

```

Random Forest

```

from sklearn.ensemble import RandomForestClassifier
# Extract features using a pre-trained CNN (e.g., VGG16)
base_model = tf.keras.applications.VGG16(weights='imagenet', include_top=False,
input_shape=(img_height, img_width, 3))
feature_extractor = tf.keras.models.Model(inputs=base_model.input,
outputs=base_model.get_layer('block5_pool').output)
# Extract features from the training dataset
train_features = []
train_labels = []
for images, labels in train_ds:
    features = feature_extractor.predict(images)
    train_features.append(features)
    train_labels.append(labels.numpy())
train_features = np.concatenate(train_features)
train_features = train_features.reshape(train_features.shape[0], -1) # Flatten the features
train_labels = np.concatenate(train_labels)

# Train Random Forest model
rf_model = RandomForestClassifier(n_estimators=100, random_state=123)
rf_model.fit(train_features, train_labels)
val_features = []
val_labels = []
for images, labels in test_ds:

```

```

features = feature_extractor.predict(images)
val_features.append(features)
val_labels.append(labels.numpy())
val_features = np.concatenate(val_features)
val_features = val_features.reshape(val_features.shape[0], -1) # Flatten the features
val_labels = np.concatenate(val_labels)
val_predictions = rf_model.predict(val_features)
# Load and preprocess the new image
new_img = tf.keras.utils.load_img("/content/asia (477).jpg", target_size=(img_height,
img_width))
new_img_array = tf.keras.utils.img_to_array(new_img)
new_img_array = np.expand_dims(new_img_array, axis=0)

# Extract features from the new image
new_img_features = feature_extractor.predict(new_img_array)
new_img_features = new_img_features.reshape(1, -1)
# Use the Random Forest model to predict the class
predicted_class = rf_model.predict(new_img_features)
predicted_class_name = class_names[int(predicted_class)]
print(f"The predicted class for the unknown image is: {predicted_class_name}")
# Calculate accuracy
accuracy = accuracy_score(val_labels, val_predictions)
print(f"Accuracy: {accuracy * 100:.2f}%")
# Plot confusion matrix
cm_rf = confusion_matrix(val_labels, val_predictions)
plt.figure(figsize=(8, 6))
sns.heatmap(cm_rf, annot=True, cmap='Blues', fmt='d', cbar=False,
xticklabels=class_names, yticklabels=class_names)
plt.title('Random Forest Confusion Matrix')
plt.xlabel('Predicted Labels')
plt.ylabel('True Labels')
plt.show()
# Generate classification report
report_rf = classification_report(val_labels, val_predictions, target_names=class_names)
print("Random Forest Classification Report:")
print(report_rf)

from sklearn.metrics import precision_recall_fscore_support
precision, recall, f1_score, _ = precision_recall_fscore_support(val_labels,

```

```

val_predictions, average='weighted')
print(f"Overall Precision: {precision:.4f}")
print(f"Overall Recall: {recall:.4f}")
print(f"Overall F1-score: {f1_score:.4f}")
from sklearn.metrics import roc_curve, auc
from sklearn.preprocessing import label_binarize
# Binarize the labels for ROC curve calculation
num_classes = len(class_names)
val_labels_bin = label_binarize(val_labels, classes=np.arange(num_classes))
# Calculate probabilities for each class
val_probs = rf_model.predict_proba(val_features)
# Compute ROC curve and ROC area for each class
fpr = dict()
tpr = dict()
roc_auc = dict()
for i in range(num_classes):
    fpr[i], tpr[i], _ = roc_curve(val_labels_bin[:, i], val_probs[:, i])
    roc_auc[i] = auc(fpr[i], tpr[i])
# Plot ROC curve for each class
plt.figure(figsize=(8, 6))
for i in range(num_classes):
    plt.plot(fpr[i], tpr[i], label=f'ROC curve ({class_names[i]}) (AUC = {roc_auc[i]:.2f})')
plt.plot([0, 1], [0, 1], 'k--')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve for Random Forest Classifier')
plt.legend(loc='lower right')
plt.show()

```

