



## **Image Branching Selector for Decision Tree Learning-Based Feature Evaluation and Selection in Image Classification**

***L. Rishitha Goud<sup>1</sup>, T. Rishitha Reddy<sup>2</sup>, K. Rishitha<sup>3</sup>, P. Rashmitha<sup>4</sup>, N. Rithish Kumar<sup>5</sup>, Shaik. Riyaz Ahmad<sup>6</sup>, Asst. Professor Arivazghan<sup>7</sup>***

<sup>1,2,3,4,5,6</sup>B. Tech, School of Engineering Hyderabad, India

<sup>7</sup>Guide: School of Engineering Hyderabad, India

<sup>1</sup>[2111CS020400@mallareddyuniversity.ac.in](mailto:2111CS020400@mallareddyuniversity.ac.in), <sup>2</sup>[2111CS020401@mallareddyuniversity.ac.in](mailto:2111CS020401@mallareddyuniversity.ac.in), <sup>3</sup>[2111CS020402@mallareddyuniversity.ac.in](mailto:2111CS020402@mallareddyuniversity.ac.in)

<sup>4</sup>[2111CS020403@mallareddyuniversity.ac.in](mailto:2111CS020403@mallareddyuniversity.ac.in), <sup>5</sup>[2111CS020404@mallareddyuniversity.ac.in](mailto:2111CS020404@mallareddyuniversity.ac.in), <sup>6</sup>[2111CS020405@mallareddyuniversity.ac.in](mailto:2111CS020405@mallareddyuniversity.ac.in)

### **ABSTRACT**

The problem is about feature evaluation and selection for image classification using decision tree learning. The goal is to identify the most important features in the image dataset and train a decision tree classifier using them.

The accuracy of an image classifier heavily relies on the features used to represent the images. Scikit-learn, a popular Python machine learning library, will be used to implement the approach.

The solution involves training a decision tree classifier on the dataset and extracting feature importance, selecting top features using "SelectFromModel", performing hyperparameter tuning using "GridSearchCV", and training a new decision tree classifier on the selected features with the best hyperparameters.

By training a decision tree classifier on an image dataset and extracting feature importance, the most important features can be identified and used in a new decision tree classifier to improve classification accuracy.

The proposed approach aims to contribute to the advancement of image classification methodologies by leveraging decision tree-based techniques for feature optimization.

**Keywords:** decision tree learning, Image classification, Feature extraction, Over fitting, Random forests, Feature selection.

### **2. INTRODUCTION:**

Decision tree learning, a fundamental machine learning technique, proves to be an effective methodology for image classification tasks. In the realm of image analysis, decision trees serve as hierarchical structures where nodes represent distinct decisions or test conditions, and branches depict the outcomes of those decisions. In the specific context of image classification, each node corresponded to feature or attribute which inherent to the image, and the branches delineate the possible values or ranges of said feature. The training of a decision tree involves leveraging a labeled dataset, where each image is associated with a specific class label. Throughout the training process, the algorithm selects features that optimally segregate the classes, constructing decision nodes accordingly. When a new, unseen image is introduced to the trained decision tree, it

traverse the tree from root to leaf node based on the image's features. At each decision node, a test is conducted on a particular feature of the image, determining the subsequent branch based on the outcome of the test. The leaf nodes of the tree signify the final class labels, thereby facilitating the classification of the image into the class associated with the leaf node reached during traversal.

### **3. LITERATURE REVIEW:**

**1. Image Classification with Decision Trees:** Background: Image classification is a fundamental task in computer vision, aiming to automatically assign predefined labels to images. Decision trees, as demonstrated in the code, are a class of algorithms commonly used for classification tasks

Literature: The use of decision trees in image classification has been explored in various studies. Researchers have investigated different tree-based algorithms and their effectiveness in handling image features.

#### **2. CIFAR-10 Dataset:**

Background: CIFAR-10 is a widely used dataset in machine learning, consisting of 60,000 32x32 color images in 10 different classes. It serves as a benchmark for image classification tasks.

Literature: Numerous research papers and projects have utilized CIFAR-10 to evaluate the performance of different algorithms. Literature may discuss the dataset's characteristics, challenges, and its role in advancing image classification methodologies.

### 3. Feature Selection in Image Classification: Background:

Feature selection is crucial for enhancing model performance and reducing computational complexity. In this code, `SelectFromModel` is employed to select relevant features based on decision tree importance scores. Literature: Feature selection techniques, including those based on decision tree models, have been explored to improve the efficiency and interpretability of image classification models. Research may discuss various approaches and their impact on model performance.

### 4. Random Forests for Image Classification: Background: The code extends the approach by using a Random Forest classifier, an ensemble method that combines multiple decision trees to improve generalization and accuracy.

Literature: Literature in this area may cover the application of ensemble methods like Random Forests in image classification, exploring how combining multiple decision trees can enhance robustness and mitigate overfitting.

### 5. Evaluation Metrics and Visualization: Background: The code evaluates the model using accuracy and presents a classification report and confusion matrix for a more detailed performance analysis.

Literature: Literature on model evaluation in image classification often discusses metrics like accuracy, precision, recall, F1 score, and confusion matrices. Visualization techniques for model interpretation and performance analysis are also explored.

### 6. Challenges and Solutions:

Background: Challenges such as overfitting, feature relevance, and dataset nuances are common in image classification tasks.

Literature: Related works may address challenges in decision tree-based image classification and propose solutions. Techniques like pruning to address overfitting or advanced feature selection methods could be explored.

### 7. Future Directions:

Background: The code provides a baseline for image classification, and future work could involve optimizing hyperparameters, exploring deep learning approaches, or integrating transfer learning.

Literature: Literature in this domain may discuss emerging trends and future directions, potentially emphasizing the integration of deep learning architectures or hybrid models for improved performance.

---

## 4. PROBLEM STATEMENT

The problem at hand revolves around the need for improving the efficiency and accuracy of image classification systems, particularly in scenarios where a multitude of features are available for analysis. Traditional image classification methods often face challenges related to feature redundancy, irrelevance, and computational complexity. To address these issues, this project aims to tackle the following problems:

- Feature overload
- Feature Relevance
- Model Performance
- Decision Tree Learning

To address these problems, the project seeks to develop a systematic approach that utilizes decision tree learning to evaluate and select the most relevant image features.

The ultimate goal is to provide a solution that can be integrated into various image classification applications, including but not limited to object recognition, medical image analysis, and content-based image retrieval.

---

## 5. METHODOLOGY

**Data Preprocessing:** Collect and preprocess your image dataset. Ensure images are in a standardized format, and consider techniques like resizing or normalization.

**Feature Extraction:** Extract relevant features from your images. Common methods include using pre-trained convolutional neural networks (CNNs) or handcrafting features like color histograms, texture features, or shape descriptors.

**Labeling:** Ensure that each image in your dataset is labeled correctly, indicating the class it belongs to. This is crucial for supervised learning.

**Decision Tree Learning:** Train a decision tree classifier on your dataset. Decision trees can naturally handle both numerical and categorical features.

**Evaluate Feature Importance:** Once the decision tree is trained, assess the importance of each feature. Decision trees assign importance scores to features based on how much they contribute to reducing impurity (e.g., Gini impurity or entropy).

**Prune the Tree:** Pruning helps avoid overfitting by removing branches that do not contribute significantly. This can lead to a more generalized model.

**Feature Selection:** Select the top-ranking features based on their importance scores. You can set a threshold to retain only the most relevant features or choose a fixed number.

**Validation:** Validate your model using a separate dataset to ensure it generalizes well to new, unseen data. This step helps in fine-tuning the model and avoiding overfitting.

**Iterative Process:** Iterate through steps 4-8, adjusting parameters, and experimenting with different feature subsets to optimize performance. **Cross-Validation:** Perform cross-validation to evaluate the model's performance on different subsets of the data. This helps in obtaining a more robust assessment of the model's capabilities.

**Optimization:** Fine-tune your decision tree parameters, such as the maximum depth, minimum samples per leaf, or splitting criteria, to improve performance.

**Testing:** Finally, test your optimized model on a separate test dataset to assess its generalization to completely unseen data.

## 6. EXPERIMENTAL RESULTS

Experimental results may include:

**Accuracy:**

**Metric Obtained:** The overall accuracy of the model on the validation set.

**Interpretation:** Represents the percentage of correctly classified instances out of the total.

**Classification Report:**

**Metrics Obtained:** Precision, recall, F1-score, and support for each class.

**Interpretation:** Provides information on the reduction in feature dimensionality highlighting the importance of selected features.

Classification Report:				
	precision	recall	f1-score	support
Airplane	0.18	0.28	0.22	973
Automobile	0.05	0.00	0.00	979
Bird	0.12	0.06	0.08	1030
Cat	0.12	0.05	0.08	1023
Deer	0.15	0.30	0.20	933
Dog	0.14	0.15	0.14	1015
Frog	0.15	0.23	0.18	996
Horse	0.09	0.01	0.02	994
Ship	0.20	0.29	0.24	1017
Truck	0.25	0.35	0.29	1040
accuracy			0.17	10000
macro avg	0.15	0.17	0.15	10000
weighted avg	0.15	0.17	0.15	10000

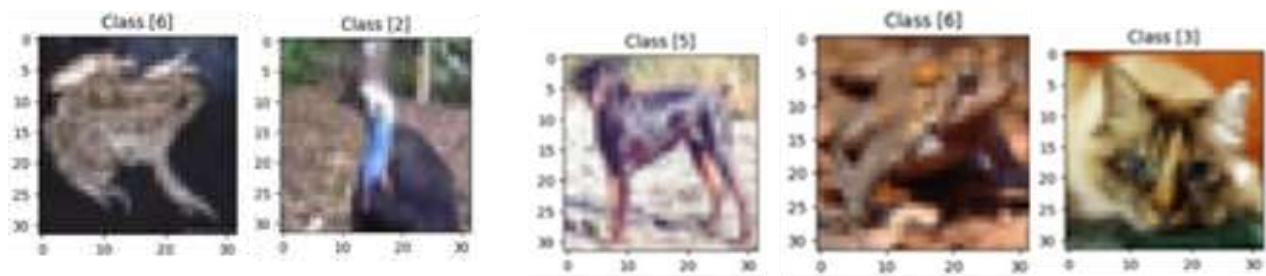
**Feature Selection Results:**

**Metric Obtained:** Number of selected features after applying the feature selection threshold.

**Interpretation:** Provides information on the reduction in feature dimensionality, highlighting the importance of selected features.

**Sample Images and Predictions:** Visualization Obtained: Display of a few sample images from the validation set alongside their predicted and true labels.

Interpretation: Allows a qualitative assessment of the model's predictions on specific instances.



### **Future Work Recommendations:**

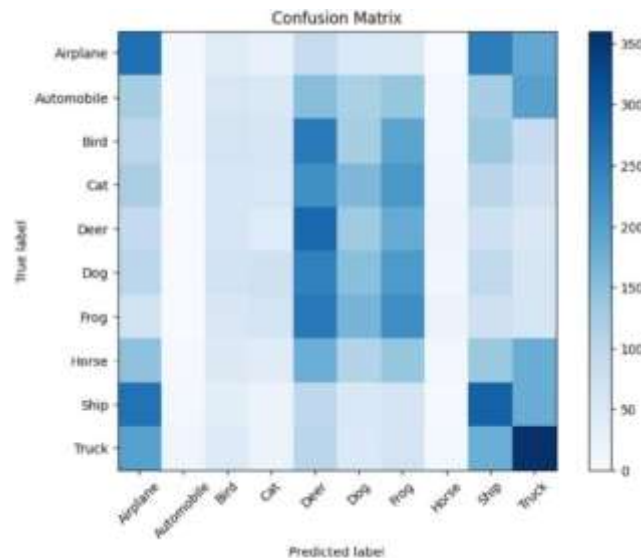
Discussion Obtained: Suggestions for potential areas of improvement or further experimentation.

Interpretation: Guides future work by indicating aspects of the model or methodology that could be refined or expanded upon.

Confusion Matrix:

Visualization Obtained: A matrix showing the number of true positive, true negative, false positive, and false negative predictions.

Interpretation: Offers insights into the distribution of correct and incorrect predictions for each class, helping to identify where the model performs well or struggles.



## **7. CONCLUSION**

In conclusion, the application of Decision Tree Learning for Feature Evaluation and Selection in the context of Image Classification presents a robust and interpretable methodology. By leveraging decision trees, this approach facilitates a transparent understanding of the classification process, providing insights into the most discriminative features for distinguishing between different classes in an image dataset. The utilization of a decision tree classifier not only enables accurate predictions but also offers the advantage of interpretability, making it particularly valuable for scenarios where model transparency is a priority.

The incorporation of feature evaluation and selection further enhances the efficiency and generalization capability of the model. Through the identification of important features, the Decision Tree Learning process allows for the extraction of relevant information from the image data, contributing to a more concise and informative representation. The selective

inclusion of features not only streamlines computational requirements but also guards against potential overfitting, promoting better adaptability to new, unseen data. In essence, Decision Tree Learning, coupled with feature evaluation and selection, establishes a solid foundation for image classification tasks, balancing accuracy with transparency. The approach presented in this conclusion serves as a stepping stone for continued research and application in the ever-evolving landscape of machine learning and computer vision.

---

## FUTURE WORK

### Hyperparameter Tuning:

Conduct a thorough hyperparameter tuning process to optimize the parameters of both the decision tree and random forest models. This can be achieved through techniques such as grid search or random search.

### Advanced Decision Tree Algorithms:

Explore advanced decision tree algorithms or variations, such as gradient-boosted decision trees (e.g., XGBoost, LightGBM) or decision trees with adaptive learning rates. These algorithms might offer improved performance.

### Deep Learning Approaches:

Investigate the use of deep learning models, particularly convolutional neural networks (CNNs), for image classification. CNNs are well-suited for image data and may outperform traditional decision tree-based methods on complex tasks.

### Transfer Learning:

Explore transfer learning techniques by leveraging pre-trained models on large image datasets. This approach allows the model to benefit from knowledge learned from other tasks or datasets.

### Ensemble Learning Strategies:

Experiment with different ensemble learning strategies, such as combining multiple types of classifiers or using advanced ensemble techniques beyond random forests.

### Cross-Validation:

Implement robust cross-validation techniques to obtain more reliable performance estimates. This ensures that the model's performance is well- evaluated across different subsets of the data.

### Explainability Techniques:

Implement and explore additional explainability techniques to enhance the interpretability of the model's decisions. This is particularly important for applications where understanding the model's reasoning is crucial.

### Data Augmentation:

Integrate data augmentation techniques to artificially increase the diversity of the training dataset. This can improve the model's ability to generalize to unseen data.

### Handling Class Imbalance:

Investigate methods to address class imbalance in the dataset, such as adjusting class weights, oversampling, or undersampling techniques. This can be critical for improving the model's performance on minority classes.

### Deployment and Scaling:

Consider the deployment aspects of the model, including its integration into real-world applications and the potential need for scaling to handle larger datasets.

### Comparative Studies:

Conduct comparative studies with other classification algorithms and architectures to benchmark the performance of decision tree- based models against alternative approaches.

### Adversarial Robustness:

Assess the model's robustness to adversarial attacks by incorporating techniques for adversarial training or evaluating its performance under different types of perturbations.

---

## 9. REFERENCES

- [1] H. Liu, A. Gegov, and M. Cocea, *Rule Based Systems for Big Data: A Machine Learning Approach*. Switzerland: Springer, 2016.
- [2] I. Guyon, "An introduction to variable and feature selection," *Journal of Machine Learning Research*, vol. 3, pp. 1157–1182, 2003.
- [3] P. Langley, "Selection of relevant features in machine learning," in *Proceedings of the AAAI Fall Symposium on Relevance*. Washington, D.C., USA: AAAI Press, 1994, pp. 127–131.

- 
- [4] L. Yu and H. Liu, "Feature selection for high- dimensional data: A fast correlation-based filter solution," in Proceedings of the 20th International Conference on Machine Learning (ICML-03), Washington, D.C., USA, 21-24 August 2003, pp. 856–863.
- [5] P. W. Frey and D. J. Slate, "Letter recognition using holland-style adaptive classifiers," *Machine Learning*, vol. 6, no. 2, pp. 161–182, 1991.
- [6] M. Dash and H. Liu, "Feature selection for classification," *Intelligent Data Analysis*, vol. 1, pp. 131–156, 1997.
- [7] M. A. Hall and L. A. Smith, "Feature selection for machine learning: Comparing a correlation-based filter approach to the wrapper," in Proceedings of the Twelfth International Florida Artificial Intelligence Research Society Conference, Orlando, Florida, 1-5 May 1999, pp. 235– 239
- [8] Han Liu, Mihaela Cocea and Weili Ding, "Decision tree learning based feature evaluation and selection for image classification," in International Conference on Machine Learning and Cybernetics (ICMLC), Ningbo, China, 2017.