

A CNN Based Approach To Fruit Classification

Jonathan Day, Jimmy Gore

December 18, 2024

1 Introduction

The fruit classification problem focuses on identifying and categorizing fruits based on visual data, such as images. This involves distinguishing various fruit types, such as **apples**, **bananas**, and **mangos**, using their unique features like shape, color, and texture. The challenge lies in handling variations in lighting, orientation, size, and background that can affect image quality and the consistency of fruit appearance.

To address this, *convolutional neural networks* (CNNs) are used due to their ability to automatically extract specific features from images. A CNN processes input images through:

- **Convolutional layers:** Identify patterns like edges and textures.
- **Pooling layers:** Reduce dimensionality and retain important features.
- **Fully connected layers:** Perform the final classification.

By training the CNN using python's torch [1] and torchvision libraries [2], on a labeled dataset of fruit images [3], it learns to recognize distinguishing characteristics for each fruit class. This approach is widely applied in areas such as **food quality assurance**, **automated sorting systems**, and **educational tools**. We will use a CNN-based approach to classify fruits in images, focusing on **accuracy**, **efficiency**, and **scalability**. Our goal is to develop a robust fruit classification system that can handle various fruit types, image conditions, and dataset sizes. By leveraging the power of CNNs, we aim to achieve high classification accuracy and real-time performance, making our system suitable for practical applications in **agriculture**, **retail**, and **research**.

2 Related Work

Classification tasks using convolutional neural networks (CNNs) have been explored quite a lot, providing a large amount of existing work to guide use in developing our own techniques. While large-scale models like those trained on ImageNet [4] may be overly complex for this task, they offer valuable insights that can assist in our approach.

ImageNet

ImageNet is arguably the most comprehensive hand-labeled image dataset and has been a staple in the computer vision community since its introduction in 2009. Its vast scale and diversity have made it a benchmark for image classification tasks.

AlexNet

AlexNet, a groundbreaking CNN architecture, achieved a top-5 error rate of 15.3% on the ImageNet dataset. Its success paved the way for further research into CNN architectures and their applications in image classification.

Cat and Dog Classification

The task of classifying cats and dogs has become iconic in CNN research. This problem has inspired countless tutorials and serves as a foundational example for extending CNN-based approaches to other classification tasks, including fruit classification.

3 Proposed Method

The backbone of our approach is a custom CNN model, which incorporates attention mechanisms to enhance classification performance. The architecture consists of multiple convolutional layers, designed to extract key features such as texture, shape, and color. To further improve the model’s focus on relevant image regions, we introduce two attention mechanisms: Channel Attention (CAM) and Spatial Attention (SAM).

3.0.1 CNN Architecture

The proposed CNN model consists of the following components:

- **Convolutional Layers:** Four convolutional layers are used to extract hierarchical features from the input images, progressively capturing patterns like edges, textures, and object details.
- **Pooling Layer:** A max-pooling layer is applied after each convolutional block to reduce spatial dimensions while retaining essential information, minimizing computational overhead.
- **Channel Attention Mechanism:** Enhances feature extraction by emphasizing the most relevant channels, helping the model focus on critical features specific to each fruit class.
- **Spatial Attention Mechanism:** Captures spatial dependencies within feature maps, improving the model’s ability to localize distinguishing areas in the images.

- **Fully Connected Layers:** Two fully connected layers are employed for final classification, transforming the learned feature representations into fruit class probabilities.

This architecture leverages a combination of convolutional operations, attention mechanisms, and fully connected layers to achieve high accuracy and robust performance for the fruit classification task.

3.0.2 Channel Attention (CAM)

The Channel Attention mechanism computes the importance of individual feature channels. By applying adaptive pooling and fully connected layers, it identifies and emphasizes the most relevant channels, improving the network’s ability to find small variations between fruit categories.

3.0.3 Spatial Attention (SAM)

The Spatial Attention mechanism enhances the model’s focus on specific regions within the input image. It utilizes convolution operations over pooled feature maps to concentrate on important areas, such as fruit boundaries, making sure significant features are prioritized during classification.

3.1 Training Techniques

Our training loop ensures that our model learns effectively from the dataset and is robust and accurate within the variations in the dataset. Key techniques include:

- **Loss Function:** Cross-entropy loss is used for multi-class classification.
- **Optimizer:** Adam optimizer ensures efficient convergence.
- **Data Augmentation:** Techniques such as random horizontal flips, rotations, and normalized pixel values simulate diverse real-world scenarios to ensure robustness.

By integrating attention mechanisms into our CNN architecture, we achieve a robust and scalable fruit classification system, capable of handling the diverse dataset and real-world scenarios with the highest accuracy possible.

4 Experiments

4.1 Initial Experimentation

Our initial experiment focused on training a CNN model on a diverse dataset of fruit images. The dataset included a variety of classes such as apples, bananas, strawberries, mangos, and grapes. Additionally, the dataset contained both real

and cartoon images of fruits, providing a challenging classification task due to variations in appearance, style, and complexity.

To accelerate the training process and improve performance, we utilized a GPU-accelerated environment. We began with a simple CNN architecture and progressively introduced attention mechanisms, such as channel and spatial attention, to enhance the model’s ability to distinguish between fruit classes. later, we added additional convolutional layers and increased the number of filters to improve the models feature extraction ability.

4.2 Hyperparameter Tuning and Data Augmentation

To optimize the performance of our model, we conducted experiments with various hyperparameters, including:

- **Learning Rate:** Adjusted to achieve convergence without overshooting.
- **Batch Size:** Smaller batches were used for memory efficiency, while larger batches were tested for stability.
- **Number of Epochs:** Experimented with longer training periods to balance underfitting and overfitting.

These experiments were crucial in identifying the best configuration for our CNN-based fruit classification task.

4.3 Comparative Analysis with Established Architectures

To evaluate the performance of our model against well-established architectures, we cloned and trained models such as AlexNet and ResNet50 using the same dataset and hyperparameters for consistency. The results are summarized below:

- **AlexNet:** Achieved a test accuracy of 74.90% and a validation accuracy of 76.90%.
- **ResNet50:** Achieved a test accuracy of 75.30% and a validation accuracy of 75.60%.
- **Our Model:** Achieved a test accuracy of 71.40% and a validation accuracy of 71.70%.

These results demonstrate that our custom CNN architecture, while simpler, performs competitively against advanced architectures when trained on the same dataset.

4.4 Summary

Our experiments highlight the iterative process of model design and refinement, emphasizing the value of incremental improvements, hyperparameter tuning, and comparative analysis with established models to develop an effective fruit classification system.

5 Results

5.1 Optimal Model Configuration

The optimal configuration for our model was determined through experimentation and automated tuning scripts. The key hyperparameters for the best-performing model are as follows:

- **Number of Epochs:** 30
- **Learning Rate:** 0.0001
- **Batch Size:** 16

5.2 Performance Metrics

The model achieved the following accuracy metrics:

- **Test Accuracy:** 71.40%
- **Validation Accuracy:** 71.70%

5.3 Training Time

The training time for each epoch was recorded under multiple computing configurations. The results are as follows:

- **CPU:** 45–60 seconds per epoch
- **CUDA:** Approximately 30 seconds per epoch

5.4 Conclusion

Our model demonstrated promising results with efficient training times when leveraging GPU acceleration. This configuration serves as a baseline for further improvements and experimentation.

6 Conclusion

In this work, we presented a CNN-based approach to fruit classification, taking advantage of custom architecture designs, attention mechanisms, and data augmentation techniques to improve accuracy and robustness. The combination of channel and spatial attention mechanisms proved to be effective in enhancing the model’s ability to extract and focus on relevant features, leading to comparable results when compared to established architectures like AlexNet and ResNet50.

Through experimentation, we optimized hyperparameters and demonstrated the importance of effective refinement in achieving robust performance. Our

model achieved a test accuracy of 71.40% and a validation accuracy of 71.70% with relatively efficient training times, particularly when utilizing GPU acceleration. These results validate the capability of our CNN architecture to handle the challenges of diverse and complex datasets.

Additionally, our comparative analysis with AlexNet and ResNet50 [5] highlighted the strengths of our lightweight yet effective model. While advanced architectures offer higher accuracy, our approach balances computational efficiency and classification performance, making it a viable solution for real-world applications in agriculture, retail, and more.

Overall, our work underscores the potential of convolutional neural networks in addressing practical image classification challenges, offering a solid foundation for ongoing research and development in this domain.

7 Work Cited

References

- [1] A. Paszke, S. Gross, F. Massa, A. Lerer, M. Beaufait, Z. Lin, Y. Wu, L. Zhang, E. Fey, B. Chai, T. K. L., and K. He, “Pytorch: An imperative style, high-performance deep learning library,” 2019, accessed: 2024-12-18. [Online]. Available: <https://pytorch.org>
- [2] P. Contributors, “Torchvision,” 2024, accessed: 2024-12-18. [Online]. Available: <https://pytorch.org/vision/stable/index.html>
- [3] U. Saxena, “Fruits classification dataset,” 2020, accessed: 2024-12-18. [Online]. Available: <https://www.kaggle.com/datasets/utkarshsaxenadn/fruits-classification>
- [4] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Neural Information Processing Systems (NeurIPS)*, vol. 25, pp. 1097–1105, 2012. [Online]. Available: <https://doi.org/10.1145/3065386>
- [5] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016. [Online]. Available: <https://doi.org/10.1109/CVPR.2016.90>

8 Evaluation Form

Name	Average
Jonathan Day	100%
Jimmy Gore	100%

Table 1: Both participants had equal contributions to the project and worked effectively to achieve the highest overall accuracy out of all teams that took on this project.