

Neural Computing Assignment

Food Classification with CNN

Group 41
Enzo Perk (s3894347)
Sem Broere (s3410552)
Bing Steens (s3901149)
Thijmen Bouman (s3543404)

Leiden University
May 16, 2025

1 Introduction

The primary goal of this task was to design, implement, and test a deep model for food image classification into 91 food categories. Food images are difficult to classify due to large intra-class variability and inter-class similarity, meaning that pictures in the same class can have very different appearances, and some different food types look very similar. We approached the task by developing an in-house Convolutional Neural Network (CNN). This report presents our work, from the initial setup and data manipulation to design, training, and evaluation of the model. The most accurate model reached a test accuracy of 62.36% following 50 epochs of training.

2 Implementation Details

2.1 Initial Setup: Imports, GPU Verification, and Reproducibility

The first code cell initializes the project environment. It loads the necessary libraries like `'torch'` for deep learning, `'torchvision'` for data manipulation, `'numpy'` for numerical computation, and `'matplotlib'` for plotting. It also verifies if a CUDA-compatible GPU is present. A key function, `'set_seed'`, is defined and called, with a seed of 42. This function standardizes random number generation across the whole PyTorch, NumPy, and Python `'random'` module, and initializes CUDA for deterministic use with the same seed, which is necessary to ensure reproducible experimental outcomes.

2.2 Data Preprocessing: Augmentation and Loading

This cell is dedicated to data preparation of images. Two chains of transformations are defined by the use of `'transforms.Compose'`. For train data, `'transform_train'` does resizing to 144x144, random crop to 128x128, random horizontal flip, and color jittering (contrast and brightness by 0.1), followed by tensor conversion and normalization to ImageNet standard. These are the augmentations done to increase dataset diversity and model robustness. For test data, `'transform_test'` does the same resizing, center crop to 128x128, tensor conversion, and normalization, but without random augmentations in order to provide consistent evaluation. The data is loaded from `'train'` and `'test'` directories through `'datasets.ImageFolder'`. `'DataLoader'` instances, `'train_loader'` and `'test_loader'`, are established with batch size 32 and 4 worker processes, with shuffling enabled for the training loader and shuffling disabled in the test loader.

2.3 Model Architecture: FoodCNN Definition

A personalized Convolutional Neural Network, `'FoodCNN'`, is defined here, inherited from `'nn.Module'`. Its structure consists of a feature extractor (`'self.features'`) and a classifier (`'self.classifier'`). The feature extractor consists of four sequential blocks. The first block is made up of two convolutional blocks (3 input channels to 64, then 64 channels to 64, kernel size 3, padding 1), both before batch normalization and ReLU activation, and finishing in a max pooling layer. Subsequent blocks progressively deepen channels (64 to 128, 128 to 256, then 256 to 512 channels), each employing a convolutional layer (kernel size 3, padding 1), batch normalization, ReLU, and max pooling (for blocks 2 and 3, kernel

size 2) or adaptive average pooling (for block 4, output size (1,1)). The classifier then flattens these 512 feature channels and inserts a dropout layer with rate 0.4 to regularize, and finally passes the output to a linear layer to produce 91 output logits, one per food class. The 'forward' method handles data flow throughout these components.

2.4 Core Logic: Training and Evaluation Functions Definition

This cell defines the 'train_and_evaluate' and 'evaluate' functions, the core of the learning and evaluation of the model. The 'train_and_evaluate' function organizes the entire training process. It initializes the Adam optimizer with a given learning rate ('lr') and weight decay of 0.0001, and a 'OneCycleLR' learning rate scheduler which adjusts the learning rate dynamically between epochs (with 'max_lr=0.003', 'epochs=num_epochs', 'steps_per_epoch=len(train_loader)', 'pct_start=0.3', 'div_factor=10', 'final_div_factor=100'). The loss function employed is Cross-Entropy Loss with label smoothing factor 0.05. The function iterates over a number of epochs ('num_epochs'). For every epoch, it loops through training data in batches: it performs a forward pass, calculates loss, backpropagates the gradients ('loss.backward()'), and updates model weights ('optimizer.step()') and scheduler ('scheduler.step()'). It calls 'optimizer.zero_grad()' before gradient calculation. At the end of every training epoch, it calls 'evaluate' function to calculate performance on the test set. If test accuracy ('test_acc') is higher than the best so far ('best_acc'), model state dictionary is saved to 'best_model.pth'. The 'evaluate' function is simpler: it sets the model to evaluation mode ('model.eval()'), iterates over the provided data loader without computing gradients (with 'torch.no_grad()'), and calculates the accuracy by comparing predictions ('outputs.max(1)') with ground truths.

2.5 Model Training Execution

This cell runs the model training. The random seed is fixed to 42 with 'set_seed(42)' for this run. An instance of 'FoodCNN' is created and moved to the GPU with '.to(device)'. The 'train_and_evaluate' function is then called, providing the model and data loaders with 50 epochs of training ('num_epochs = 50'), and the starting learning rate of 0.001 ('learning_rate = 0.001'). The output should report that the model is trained for 50 epochs with an end best accuracy of 62.36%. As we can see when we run the model, the scheduler indeed updates the learning rate in a curve-like fashion.

2.6 Post-Training: Loading Best Model and Final Evaluation

As training is done, this cell loads the best performing model. A new instance of 'FoodCNN' is created, and its weights are loaded from the saved 'best_model.pth' file using 'model.load_state_dict(torch.load("best_model.pth"))'. The 'evaluate' function is invoked on this loaded model and the 'test_loader' to re-ensure its performance.

2.7 Accuracy Validation: Averaged Test Results Calculation

In order to make the stability of the evaluation definite, this cell develops a function 'calculate_test_accuracy'. This function runs the 'evaluate' process five times over the test set and then computes the average and standard deviation of these five accuracy scores using 'np.average' and 'np.std'. As we can see from the output, the model does not change throughout testing, and keeps the same result. This is due to the seed of choosing the dataset. It is possible to change the subset of the dataset, but due to training limitations, we chose not to add any randomness in the loading of datasets into subsets.

2.8 Display Configuration of Hyperparameters Summary

This cell simply prints out a summary of the key hyperparameters in the successful training run. These include things like the "Custom CNN with deeper layers" architecture, 91 classes, learning rates (initial 0.001, max 0.003 for 'OneCycleLR'), Adam optimizer, batch size 32, 50 epochs, 128x128 image input size (after 144x144 cropping), ImageNet normalization, a 0.4 dropout rate, and the OneCycleLR learning rate scheduler.

2.9 Qualitative Analysis: Visualizing Model Predictions

The final implementation block qualitatively checks. The best model is loaded, and 10 random images are selected from the test set. There is an 'imshow' function (which unnormalizes images with 'img = img * 0.45 + 0.22', it is not perfect, as the normalization happens with some degree of randomness, but it is good enough to represent as a proof of concept) to look at these images. For every image, the model predicts its class. The image is then shown with both its label and

the model's predicted label. The sample output we did displayed 7 out of 10 correct predictions on these images selected randomly, displaying both successes and some misclassifications (e.g., "cannoli" being predicted as "foie_gras").

3 Discussion

The developed 'FoodCNN' model achieved a test accuracy of 62.36% on a challenging 91-class food image dataset. This result is way better than random chance (approx. 1.1%) and indicates that the model successfully learned discriminative features from the images.

Several implementation choices contributed to this outcome. The data augmentation strategy (random crops, flips, color jitter) was crucial for exposing the model to varied data and improving generalization. The CNN architecture, featuring multiple convolutional blocks with Batch Normalization and ReLU, followed by an Adaptive Average Pooling layer, allowed for hierarchical feature extraction. Batch Normalization likely aided training stability and speed, while Dropout in the classifier helped mitigate overfitting. The combination of the Adam optimizer and the OneCycleLR scheduler is known for efficient convergence. Label smoothing also acted as a regularizer, potentially preventing overconfident predictions.

Despite the reasonable accuracy, the qualitative evaluation showed some misclassifications. These errors can be attributed to the inherent difficulty of food classification, such as high intra-class variance (e.g., different presentations of the same dish) and inter-class similarity (e.g., visually similar but distinct dishes). The dataset size per class, if imbalanced, could also play a role.

Potential improvements could involve using transfer learning with a pre-trained backbone (e.g., ResNet, EfficientNet), which often provides a stronger feature extraction base. More extensive hyperparameter tuning, perhaps using automated methods, could also yield better results. Exploring more advanced data augmentation techniques or attention mechanisms within the network could further enhance performance. The "Early stopping patience: 5" mentioned in the hyperparameter summary was not explicitly seen in the training loop logic provided (which saved based on best accuracy per epoch); if implemented, it could prevent overfitting by stopping training when validation performance plateaus.

4 Conclusion

This project involved the design, implementation, and training of a custom CNN for food image classification. The model achieved a final test accuracy of 62.36% on a dataset of 91 food categories after 50 epochs of training. Key elements of the approach included robust data augmentation, a multi-layered CNN architecture incorporating Batch Normalization and Dropout, and an effective training regimen using the Adam optimizer and OneCycleLR scheduler. The results demonstrate a good grasp of fundamental deep learning concepts applied to a complex computer vision task. While the current performance is noteworthy, further improvements could be explored through techniques such as transfer learning and more sophisticated regularization methods.