

Deep Learning 435: Final Report

Tomer Butbul
Rutgers University
May 2025

Abstract—This project evaluates the performance of a modified AlexNet architecture on the CIFAR-10 image classification dataset. Through a targeted subset of 48 experiments, we explored the effects of learning rate, batch size, dropout, batch normalization, and optimizer choice on accuracy and training efficiency. The results offer practical insights into deep learning configuration strategies under computational constraints.

I. INTRODUCTION

The objective of this project was to explore how architectural modifications and hyperparameter tuning affect the performance of convolutional neural networks, using AlexNet as the base model. By leveraging PyTorch, we implemented a series of controlled experiments on the CIFAR-10 dataset. This dataset consists of 60,000 32x32 color images categorized into 10 classes, with 50,000 images used for training and 10,000 for testing. Compared to larger datasets like ImageNet, CIFAR-10 poses a unique challenge due to its low resolution and limited per-class examples, necessitating adjustments to the standard AlexNet architecture.

Our implementation features a modular AlexNet variant that allows toggling batch normalization and dropout. A training pipeline was developed to evaluate different optimizers and hyperparameters, and a script was written to run 48 pre-defined experiments automatically. These experiments incorporate early stopping and log results such as training duration, validation accuracy, and convergence behavior, offering a robust framework for deep learning experimentation.

Before diving into our custom model design, it is important to understand the origin and core contributions of the original AlexNet architecture. Introduced by Krizhevsky et al. in 2012, AlexNet was a breakthrough in deep learning, demonstrating the power of large-scale convolutional neural networks on the ImageNet dataset, which contains over a million high-resolution labeled images across 1,000 categories. The architecture includes five convolutional layers and three fully connected layers, uses ReLU activation for faster convergence, employs dropout to reduce overfitting, and takes advantage of GPU parallelism for efficient training. AlexNet achieved top-1 and top-5 error rates of 37.5% and 17.0%, respectively, significantly outperforming prior approaches. Its innovations in architecture, normalization, and regularization formed the foundation for modern deep learning models [1].

II. MATHEMATICAL FOUNDATIONS

The backbone of convolutional neural networks like AlexNet lies in several fundamental operations. These operations not only enable deep feature extraction but also influence

how efficiently and accurately the model converges during training.

A. Convolutional Layer

The convolution operation extracts spatial features from images by applying a learnable filter over the input. Mathematically, for input I and filter K , the convolution output O at position (i, j) is:

$$O(i, j) = (I * K)(i, j) = \sum_m \sum_n I(i + m, j + n) \cdot K(m, n)$$

This allows the network to detect edges, textures, and shapes. Deeper layers learn increasingly abstract patterns. In our experiments, convolutional depth played a critical role in generalization capacity.

B. ReLU Activation

After convolution, AlexNet uses the Rectified Linear Unit:

$$f(x) = \max(0, x)$$

ReLU introduces non-linearity while avoiding the vanishing gradient problem common in saturating activations like sigmoid or tanh. Its simplicity accelerates convergence and improved training dynamics, especially evident in the performance gains observed in Figure 1 and Figure 2.

C. Dropout Regularization

Dropout randomly deactivates neurons during training:

$$\tilde{h}_i = h_i \cdot z_i, \quad z_i \sim \text{Bernoulli}(p)$$

where z_i is a binary mask and p is the keep probability. This prevents co-adaptation of neurons and acts as an implicit ensemble method. Our results show that dropout can increase overfitting risk when combined improperly (e.g., with SGD), as seen in overfitting clusters in Figure 3.

D. Batch Normalization

Batch normalization stabilizes training by normalizing activations:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}, \quad y^{(k)} = \gamma \hat{x}^{(k)} + \beta$$

where μ_B and σ_B^2 are batch mean and variance, and γ, β are learnable scale and shift parameters. BN reduces internal covariate shift, enabling higher learning rates and faster convergence. This directly explains the more stable and higher test accuracy of BN-based configurations.

E. Loss Function and Backpropagation

The model minimizes cross-entropy loss:

$$\mathcal{L} = - \sum_{i=1}^C y_i \log(\hat{y}_i)$$

where y_i is the true label (one-hot) and \hat{y}_i is the softmax output. Gradients of this loss are propagated backward using chain rule:

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial z} \cdot \frac{\partial z}{\partial w}$$

Efficient backpropagation enables deep networks like AlexNet to learn from large-scale data. We observed faster and smoother convergence with Adam, attributed to its adaptive learning rate computed as:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t, \quad v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\hat{w}_t = \frac{m_t}{\sqrt{v_t} + \epsilon}$$

where g_t is the gradient and m_t, v_t are first and second moment estimates.

III. MODEL DESIGN

The core model architecture includes five convolutional layers followed by three fully connected layers, using ReLU activations and max pooling. The final layer outputs 10 logits, one for each class in CIFAR-10. Optional components such as batch normalization and dropout (at a rate of 0.5) were selectively integrated to study their individual impact. All trainable weights were initialized using Kaiming Normal initialization, which is well-suited for ReLU-based networks.

Three optimizers were examined: SGD with momentum, RMSprop, and Adam. To make the experiments both meaningful and computationally feasible, we defined a constrained parameter space. For instance, dropout and batch normalization were not used together; instead, configurations tested one or the other. Adam and RMSprop were evaluated with learning rates of 0.001 and 0.0001, while SGD was tested with higher learning rates (0.01 and 0.001) due to its sensitivity to hyperparameters. The batch sizes were fixed at 32 and 64 to maintain training stability and time efficiency. Our hypothesis was that Adam would provide consistent convergence, batch normalization would improve accuracy, and dropout would serve as a useful regularizer under limited training data.

IV. EXPERIMENT SETUP

The complete hyperparameter space, if fully expanded, would have included 324 combinations. Given the compute time and diminishing returns from certain configurations, we narrowed the scope to 48 well-chosen experiments. The constraints used to trim the space included excluding combinations that applied both dropout and batch normalization simultaneously, limiting learning rate ranges based on optimizer behavior, and retaining only the most informative training/test splits of 70/30 and 90/10.

This streamlined approach maintained the representational diversity necessary to uncover trends while reducing experiment runtime from over 45 hours to a manageable subset. For each configuration, the model was trained from scratch and evaluated independently using consistent scripts and metrics. The collected metrics were stored in a CSV file for post-analysis.

V. RESULTS

Overall, several consistent trends emerged. Batch normalization outperformed dropout across all optimizers in terms of final test accuracy, suggesting it offers better training stability and regularization in the CIFAR-10 context. Adam was the most reliable optimizer, yielding the highest average accuracy and the smoothest convergence curves. In contrast, SGD's performance varied significantly, and one configuration involving SGD with dropout failed catastrophically, resulting in near-random accuracy. Overfitting was evident in multiple dropout experiments, where training accuracy approached 95% but test accuracy stagnated below 70%, indicating poor generalization.

Optimizer sensitivity to learning rate also differed. Adam and RMSprop performed best at a learning rate of 0.001, while SGD required a higher learning rate of 0.01 to reach competitive results. Training time correlated with batch size, with smaller batches (32) consistently taking longer to converge. RMSprop was generally the slowest among the three optimizers. Additionally, early stopping mechanisms triggered later when 90% of the data was used for training, which aligns with expectations that larger training sets require more epochs to overfit or converge.

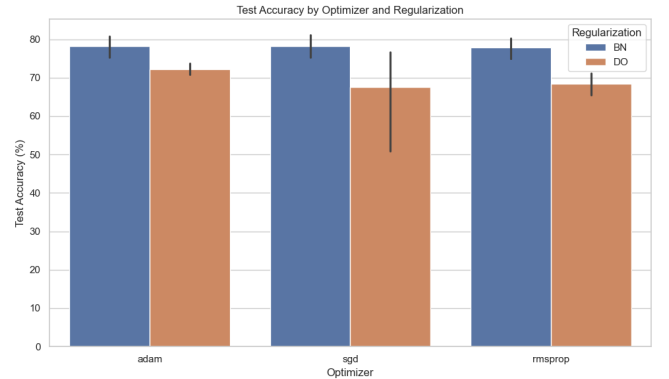


Fig. 1: Test Accuracy by Optimizer and Regularization (BN vs. DO)

Figure 1 illustrates the test accuracy distribution across different optimizers and regularization techniques. Models using batch normalization consistently outperformed those using dropout across all three optimizers. Adam combined with batch normalization achieved the highest accuracy values, reinforcing the stabilizing impact of BN and the adaptability of Adam. Dropout configurations showed higher variance and several low-performing outliers, underscoring its sensitivity and potential for overfitting on smaller datasets.

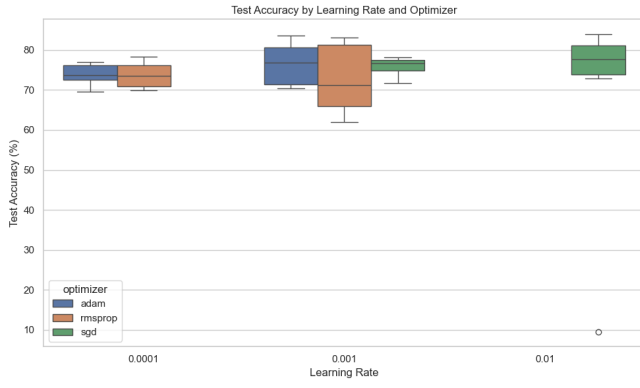


Fig. 2: Test Accuracy Across Learning Rates for Each Optimizer

Figure 2 presents the effect of learning rate on model accuracy, segmented by optimizer. Adam and RMSprop both demonstrated optimal performance at a learning rate of 0.001, while degrading slightly at 0.0001. SGD required a higher learning rate of 0.01 to remain competitive and suffered significant drops at 0.001. This validates the necessity of careful learning rate selection, especially when using more sensitive optimizers like SGD.

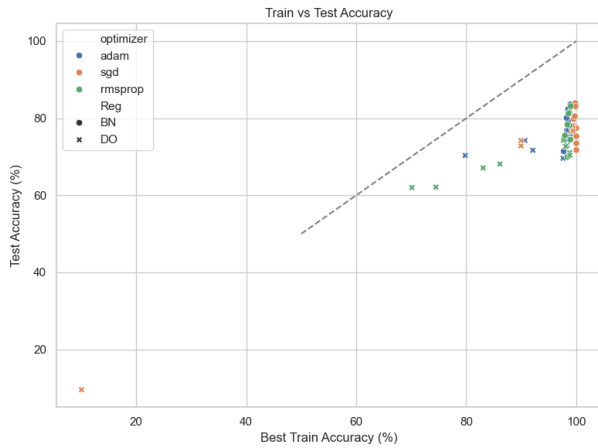


Fig. 3: Scatter Plot of Train vs. Test Accuracy, Showing Overfitting Trends

Figure 3 shows a scatter plot of training versus test accuracy. Points in the upper-left quadrant—high training but low test accuracy—represent overfitted models. This pattern was predominantly observed in dropout-heavy configurations, particularly those using SGD. In contrast, batch normalization configurations were clustered more tightly along the diagonal, indicating better generalization.

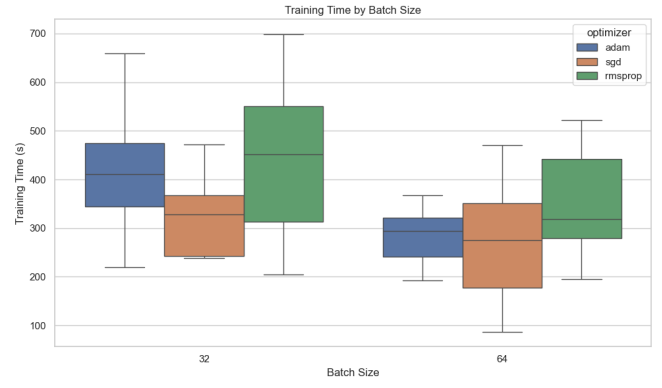


Fig. 4: Training Time Distribution by Batch Size and Optimizer

Figure 4 displays training time as a function of batch size and optimizer. Models trained with batch size 64 were consistently faster, demonstrating the computational benefits of larger batches. Among optimizers, Adam was the fastest on average, followed by SGD and then RMSprop, which took the longest—likely due to its more complex weight updates.

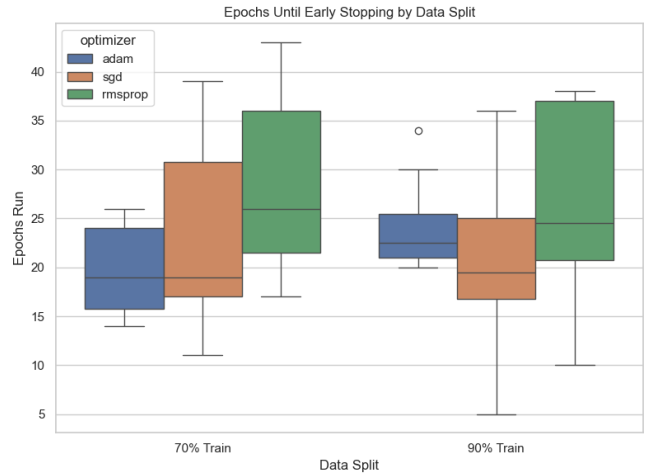


Fig. 5: Number of Epochs Until Early Stopping for Each Train/Test Split

Figure 5 depicts the number of epochs until early stopping for each train/test split. As expected, the 90/10 split required more epochs on average due to the larger training set. However, this also led to marginal improvements in test accuracy for well-regularized models, indicating a trade-off between training time and generalization performance.

VI. CONCLUSION

This project demonstrates how CNN performance on a constrained dataset like CIFAR-10 is highly sensitive to architectural tweaks, regularization choices, and training parameters. The experiments showed that batch normalization is generally more effective than dropout, providing both accuracy gains and training stability. Adam emerged as the most robust optimizer

across all scenarios, and larger batch sizes offered efficiency benefits without compromising model performance.

Dropout should be used cautiously, especially when batch normalization is an available option. Learning rate tuning remains a critical component of model performance, particularly for SGD. Early stopping proved effective for managing overfitting and should be integrated with validation-based metrics for improved reliability.

Future work could extend this analysis by integrating cross-validation and seed averaging to improve statistical significance. Testing modern architectures like ResNet or MobileNet would help compare baseline improvements, and advanced data augmentation techniques could further boost generalization. Overall, this project provided a controlled yet insightful exploration into CNN training dynamics using AlexNet on CIFAR-10.

REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, vol. 25, 2012, pp. 1097–1105.