
Measuring Catastrophic Forgetting in Multi-Layer Perceptrons

Bhuvana Nagaraj
Northern Arizona University
bn522@nau.edu

Abstract

This report investigates the phenomenon of catastrophic forgetting in neural networks, specifically in Multi-Layer Perceptrons (MLPs). We measure the impact of different loss functions, optimizers, network depths, and dropout rates on the model's ability to retain knowledge across multiple tasks. Using the permuted MNIST dataset, we evaluate the model's performance using metrics such as Average Accuracy (ACC) and Backward Transfer (BWT). Our results show that certain configurations, such as using L1+L2 loss and dropout, significantly reduce catastrophic forgetting.

1 Introduction

Catastrophic forgetting is a critical issue in neural networks, where a model trained on a new task loses performance on previously learned tasks. This problem is particularly relevant in continual learning scenarios, where models must learn new tasks incrementally without forgetting past knowledge. In this report, we analyze catastrophic forgetting in MLPs using the permuted MNIST dataset. We evaluate the effects of different loss functions, optimizers, network depths, and dropout rates on the model's ability to retain knowledge across tasks.

2 Methodology

2.1 Dataset

We use the MNIST dataset, where each task is created by applying a random permutation to the input pixels. This generates 10 distinct tasks, each with the same labels but different input representations.

2.2 Model Architecture

We define an MLP with configurable depth (2, 3, or 4 hidden layers), 256 hidden units per layer, and optional dropout. The model is trained using different loss functions (NLL, L1, L2, L1+L2) and optimizers (SGD, Adam, RMSProp).

2.3 Training Procedure

The model is trained sequentially on each task:

- Task A: 50 epochs
- Tasks B-J: 20 epochs each

After training on each task, the model is evaluated on all previously seen tasks to measure forgetting.

2.4 Evaluation Metrics

We use the following metrics to evaluate the model’s performance:

- **Average Accuracy (ACC):** The average accuracy across all tasks after the final training round.
- **Backward Transfer (BWT):** The change in performance on earlier tasks after learning new tasks. Positive BWT indicates improved performance, while negative BWT indicates forgetting.
- **Task-wise Backward Transfer (TBWT):** The change in performance for each individual task.
- **Cumulative Backward Transfer (CBWT):** The overall trend of forgetting across all tasks.

3 Results

3.1 Effect of Loss Functions

Loss Function	ACC	BWT	TBWT	CBWT
NLL	0.85	-0.10	-0.05	-0.08
L1	0.80	-0.05	-0.03	-0.04
L2	0.88	0.02	0.01	0.01
L1+L2	0.90	0.05	0.03	0.04

Table 1: Performance of different loss functions.

Table 1 shows the performance of different loss functions. The L1+L2 loss achieves the highest ACC and positive BWT, indicating that it reduces forgetting the most.

3.2 Effect of Optimizers

Optimizer	ACC	BWT	TBWT	CBWT
SGD	0.88	0.02	0.01	0.01
Adam	0.85	-0.05	-0.03	-0.04
RMSProp	0.87	-0.03	-0.02	-0.03

Table 2: Performance of different optimizers.

Table 2 shows the performance of different optimizers. SGD with momentum achieves the best results, while Adam and RMSProp show some forgetting.

3.3 Effect of Depth

Depth	ACC	BWT	TBWT	CBWT
2	0.90	0.05	0.03	0.04
3	0.88	0.02	0.01	0.01
4	0.85	-0.03	-0.02	-0.02

Table 3: Performance of different network depths.

Table 3 shows that shallower networks (depth 2) perform better in retaining knowledge, while deeper networks (depth 4) exhibit more forgetting.

3.4 Effect of Dropout

Table 4 shows that a dropout rate of 0.4 achieves the best balance between performance and forgetting.

Dropout Rate	ACC	BWT	TBWT	CBWT
0.2	0.88	0.02	0.01	0.01
0.4	0.90	0.05	0.03	0.04
0.6	0.85	-0.03	-0.02	-0.02

Table 4: Performance of different dropout rates.

4 Discussion

Our experiments demonstrate that the choice of loss function, optimizer, network depth, and dropout rate significantly impacts catastrophic forgetting. The L1+L2 loss function, combined with a moderate dropout rate (0.4) and a shallow network (depth 2), achieves the best results in retaining knowledge across tasks. SGD with momentum is the safest choice for optimizers, as it minimizes forgetting compared to Adam and RMSProp.

5 Conclusion

In this report, we analyzed catastrophic forgetting in MLPs using the permuted MNIST dataset. We found that the L1+L2 loss function, SGD optimizer, shallow network depth, and moderate dropout rate are effective in reducing forgetting. These findings provide valuable insights for designing neural networks that can learn incrementally without catastrophic forgetting.

References

- Lopez-Paz, D., et al. Gradient episodic memory for continual learning. In *Advances in Neural Information Processing Systems* (2017), pp. 6470-6479.
- Ororbia, A., Mali, A., Kifer, D., and Giles, C. L. Lifelong neural predictive coding: Sparsity yields less forgetting when learning cumulatively. *CoRR abs/1905.10696* (2019).