EECS 182	Deep Neural Networks
Spring 2023	Anant Sahai

## 1. Catastrophic Forgetting

The neural networks are vulnerable to the distributional shift. Many questions in AI/ML are related to the distributional shift: out-of-distribution (OoD), domain adaptation/generalization, meta-learning, and so on. In this discussion, we study one of those problems, catastrophic forgetting, alternatively called catastrophic interference. The catastrophic forgetting is the tendency of neural networks to lose information about previously learned tasks when learning the new one. This is also referred to as stability-plasticity dilemma<sup>1</sup>. One potential drawback of a model that is too stable is that it will not be able to consume new information from the future training data. Conversely, a model with sufficient plasticity may suffer from large weight changes and forgetting of previously learned representations.

Neural network models are trained on the assumption of i.i.d, meaning that data points are sampled from a mutually independent and identical distribution. However, this assumption does not apply to real-world applications, such as sequential data stream training settings, which can lead to catastrophic forgetting.

Let's dissect the underlying mechanisms of catastrophic forgetting. Three factors cause catastrophic forgetting mainly: parameter shift, activation shift, and inter-task confusion<sup>2</sup>.

- (a) **Parameter shift** when the gradient descent step updates the neural networks, the parameters drift to the region where the previous tasks' error is high in the parameter space.
- (b) Activation shift Activation shift is the direct ramification of parameter shift. However, focusing on activation can relax catastrophic forgetting as long as activations change minimally when the neural networks are trained with the new task.
- (c) **Inter-task confusion** Since all the tasks are not jointly trained, the outputs and intermediate activations lead to inter-task misclassification.

A cartoon depiction of catastrophic forgetting is in Figure 1. Without any implicit or explicit methods that prevent the model from catastrophic forgetting, the model would forget the useful features to discriminate samples in task 1 when it learns the new task. The right image in Figure 1 is the ideal case in which the model can capture the discriminative features for both previous and new tasks.

(a) What happens in parameter space

Figure 2 is the schematic parameter space of the given model and tasks. The orange oval is the simplified low-loss region of task 0, and the blue oval is that of task 1. The model is trained on the task 0, and  $\theta_0^*$  is the trained parameter of the model, which minimizes the loss for that task. We now train the model with task 1 data.

- i. Mark the  $\theta_1^*$  on the Figure 2, which is the trained parameter for task 1 if the model is too plastic.
- ii. Mark the  $\tilde{\theta_1}$  on the Figure 2, which is the trained parameter for task 1 if the model is too stable.
- iii. What are the problems if the model is too plastic or stable?

<sup>&</sup>lt;sup>1</sup>https://www.sciencedirect.com/science/article/pii/S1364661399012942
<sup>2</sup>https://arxiv.org/pdf/2010.15277.pdf



Figure 1: The binary classification model  $f(x; \theta_0)$  is initially trained with task 1. If the task shifts from 1 to 2, the model may suffer catastrophic forgetting. https://arxiv.org/pdf/1903.06070.pdf



Figure 2: The schematic parameter space.

## (b) Dealing with catastrophic forgetting

We consider three traditional approaches for learning the new tasks: feature extraction, fine-tuning, and joint training.

- i. feature extraction The model trained with the previous tasks is frozen. The output of this model with the new task input is used for to train the new classifier for the new task.
- ii. fine-tuning The model trained with the previous tasks is trained with the new task. To prevent the large shift, the learning rate is typically low.
- iii. joint training The model is trained with both previous task data and new task data

Let's study pros and cons of those methods. Fill in the table below.

Category	Feature Extraction	Fine tuning	Joint training
New task performance			Good
Old task performance	Good		Good
Data storage requirment		Low	
Storing old task data	No	No	

## 2. How to read research papers

One critical skill for improving yourself in deep learning is to read papers. Being able to read papers efficiently and to identify the key contributions regardless of what the authors claimed are some critical skills as you work in a field where some fundamental theory is still an open problem. In this question, you will be reading Learning without Forgetting (Li et al.), a paper that proposes a simple but effective strategy to alleviate the problem of Catastrophic Forgetting.

But first, let's talk about how to read a deep learning paper. Every researcher has a different strategy, but one recurring pattern is to *read with multiple passes*.

- First pass: Read the **title**, **abstract** and **figures**
- Second pass: Read the introduction, conclusion and figures
- Third pass: Read the **method** (skip or skim the math) and skim over the **results** (skip ablation study)
- Fourth pass: Read everything else but skip parts that do not make sense (unless you are doing research in that particular field)

Now spend some time to reading the Learning without Forgetting paper with this multiple-pass strategy, and discuss with your classmates for the following questions.

- (a) What did the authors try to accomplish?
- (b) How is Learning without Forgetting (LwF) different from standard fine-tuning and joint training (multitask learning)?
- (c) Is LwF better than feature extraction in both new and old tasks? How does LwF compare with joint training in terms of efficiency and accuracy?
- (d) Does the new task dataset size affect LwF's effectiveness? What about the old task dataset size?
- (e) What are the key takeaways that you can use yourself? Share a potential project idea where LwF can be helpful.