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 3. 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 3 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 θ<sub>1</sub><sup>\*</sup> on the Figure 2, which is the trained parameter for task 1 if the model is too plastic.
Solution: If the model is too plastic, the model lose every information that it learned from task 0. Therefore, θ<sub>1</sub><sup>\*</sup> is the middle point of the blue oval.

<sup>&</sup>lt;sup>1</sup>https://www.sciencedirect.com/science/article/pii/S1364661399012942

<sup>&</sup>lt;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

- ii. Mark the  $\tilde{\theta_1}$  on the Figure 2, which is the trained parameter for task 1 if the model is too stable. **Solution:** If the model is too stable, the model cannot learn any new information from the task 1. Therefore,  $\tilde{\theta_1} = \theta_0^*$ . This is called catastrophic remembering. The model suffering from catastrophic remembering cannot transfer any knowledge for learning the new task.
- iii. What are the problems if the model is too plastic or stable?Solution: Plastic: catastrophic forgetting. Stable: catastrophic remembering



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	Solution: Good	Solution: Best in most cases	Good
Old task performance	Good	Solution: Bad	Good
Data storage requirment	Solution: Low	Low	Solution: Large
Storing old task data	No	No	Solution: yes

### Solution:

The solution to the first row "**New task performance**" is more subtle than a universal rule of one being better than the other. Fine-tuning should perform the best in the new task since the model only has to fit the new task data instead of considering *both* the old and new task, which is often a more complex distribution. This narrative is valid when the fine-tuning dataset is large enough (the concept of enough varies among applications). However, when the fine-tuning data is too small, feature extraction can perform better than finetuning if the model capacity is so large that overfitting to the new task is likely. Some examples are when one uses BERT features and CLIP (Contrastive Language Image Pretraining) scores, which are both trained on immense dataset, it is very unlikely to finetune the entire BERT or CLIP model on a reasonable budget.

In conclusion, whether feature extraction or finetuning is better for your application depends on the size of the old and new task dataset.

## 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)

#### **Solution:**

Here are some more concrete guidelines (many of which adapted from Andrew Ng's CS230 lecture at Stanford in 2018):

- Why start with figures? A common pattern in ML/DL papers is that often there is one or two figures that summarizes the entire paper well, either the architecture or method.
- **Skipping** *Related Work*. Some reason to do this are (1) Related Work section is meant to be a place to look for extension of ideas if you're interested in this topic, so chances are unless you are familiar with this topic beforehand, you won't understand the brief description of related papers well (2) due to the nature of peer review, sometimes people cite lots of papers that one believes could potentially be their reviewers, which might not all be helpful material for learning.

- Why skip parts that do not make sense? This might be a seemingly counterintuitive statement. The main reason is that papers are often state-of-the-arts research and sometimes even the authors do not know exactly why a new phenomenon emerges. There are tons of examples of unimportant modules in influential papers: Transducers in LeNet-5 (set many foundations for ConvNets), Memory bank in Contrastive Learning... etc. If one is doing research on the same topic as the paper, then understanding thoroughly can be a good idea. Yet most students will be reading *many more* papers outside of their research / project, in those cases, focusing on the main idea is more effective.
- Why skip the math first? Beginners often get caught up by one or two math equations before understanding the entire picture of the paper. However, whether you understand every single equation does not necessarily indicate your knowledge on the main idea of this paper, and not all papers are well-written. Knowing the high-level goals and methods first can help one understand the math much easier later.
- Why skim over results? Results are important. Yet most of the good papers (either accepted to conferences or being viral on Twitter) are state-of-the-arts in some way already. Focusing on the Method section teaches more for future projects than the results do.

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?

**Solution:** The authors proposed Learning without Forgetting (LwF), which allows a network to learn new tasks without forgetting old tasks. The method leverages knowledge distillation, where the original network is used as a "teacher" to guide the learning of a new network. The LwF approach involves training the network on new tasks while simultaneously trying to maintain its performance on the old tasks by minimizing the differences between the outputs of the original and updated networks.

(b) How is Learning without Forgetting (LwF) different from standard fine-tuning and joint training (multitask learning)?

**Solution:** Standard finetuning is the most commonly used adaptation method, but also the most likely to suffer from catastrophic forgetting among the methods mentioned in Q1. LwF can preserve the model's performance better by using knowledge distillation on the "old task branch". LwF differs from joint training because it allows the model to learn new tasks sequentially, without access to the previous tasks' data.

(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?

**Solution:** LwF performs better in new tasks, but worse in the old tasks. This is reasonable since feature extraction freezes the features which leaves much smaller flexibility for the model to fit the new tasks. On the other hand, LwF slightly outperforms joint training on new tasks, but slightly underperforms on old tasks. In terms of efficiency, LwF will be faster to train since it does not require the old tasks' data.

(d) Does the new task dataset size affect LwF's effectiveness? What about the old task dataset size?Solution: In this figure, we see that all methods suffer from the decrease of new task's data, while



Fig. 5. Influence of subsampling new task training set on compared methods. The x-axis indicates diminishing training set size. Three runs of our experiments with different random  $\theta_n$  initialization and dataset subsampling are shown. Scatter points are jittered

Figure 3: the influence of the dataset size

feature extraction suffers the most. LwF compares favorably with joint training in all sizes of new task's dataset size. As the size of old task's dataset changes, we observe that fine-tuning is affected most drastically, where joint training and feature extraction remain fairly robust. LwF's performance drops as the dataset size decreases but by a much smaller magnitude compare to fine-tuning, and still obtains a comparable performance with joint training when the old task dataset size is cut to 30%.

(e) What are the key takeaways that you can use yourself? Share a potential project idea where LwF can be helpful.

Solution: Any answer where catastrophic forgetting is possible can be valid.

#### **Contributors:**

- Suhong Moon.
- Jerome Quenum.
- Kumar Krishna Agrawal.
- Kevin Li.
- Anant Sahai.