

KU Leuven (Belgium), TU Delft (the Netherlands)

### Abstract

One of the most popular methods for continual learning with deep neural networks is Elastic Weight Consolidation (EWC) [1], which involves computing the Fisher Information. The exact way in which the Fisher Information is computed is however rarely described, and multiple different implementations for it can be found online. Here, I discuss and empirically compare several often-used implementations, which highlights that many currently reported results for EWC could likely be improved by changing the way the Fisher Information is computed.

## **Continual learning**

# Different ways of computing the Fisher



A model  $f_{\theta}$  has been trained on a task by optimizing loss  $\ell_{old}(\theta)$  on training data  $D_{old} \sim \mathcal{D}_{old}$ , resulting in weights  $\hat{\theta}_{old}$ . Goal is to continue training this model on a new task, by optimizing loss  $\ell_{new}(\theta)$  on training data  $D_{new} \sim \mathcal{D}_{new}$ , such that the model maintains its performance on the previous task. Unfortunately, continued training by only optimizing  $\ell_{new}(\theta)$  results in catastrophic forgetting.



**Scope:** continual classification with neural nets (i.e.,  $f_{\theta}$  models the conditional distribution  $p_{\theta}(y|\mathbf{x})$ ).

#### Exact



#### Sampling data points



#### **Sampling labels**

$$F_{\text{old, SAMPLE}}^{i,i} = \frac{1}{|D_{\text{old}}|} \sum_{\boldsymbol{x} \in D_{\text{old}}} \left( \frac{\delta \log p_{\theta}(\boldsymbol{C}_{\boldsymbol{x}} | \boldsymbol{x})}{\delta \theta^{i}} \Big|_{\theta = \hat{\theta}_{\text{old}}} \right)^{2} \text{with } \boldsymbol{C}_{\boldsymbol{x}} \sim p_{\hat{\theta}_{\text{old}}}(.|\boldsymbol{x})$$

#### **Empirical Fisher**

$$F_{\text{old, EMPIRICAL}}^{i,i} = \frac{1}{|D_{\text{old}}|} \sum_{(\boldsymbol{x},\boldsymbol{y})\in D_{\text{old}}} \left( \frac{\delta \log p_{\boldsymbol{\theta}}(\boldsymbol{y}|\boldsymbol{x})}{\delta \theta^{i}} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_{\text{old}}} \right)^{2}$$

#### Batched approximation of Empirical Fisher (e.g., [3, 4])

$$1 \quad \sum \quad \left( \sum \quad \delta \log p_{\theta}(\mathbf{y}|\mathbf{x}) \right)^{2}$$

### Elastic Weight Consolidation (EWC)

When training on a new task, rather than optimizing only  $\ell_{new}(\theta)$ , EWC adds an extra term to the loss that involves the Fisher Information:

$$\ell_{\text{EWC}}(\boldsymbol{\theta}) = \ell_{\text{new}}(\boldsymbol{\theta}) + \frac{\lambda}{2} \sum_{i=1}^{N_{\text{params}}} F_{\text{old}}^{i,i} (\theta^{i} - \hat{\theta}_{\text{old}}^{i})^{2}$$

with  $\lambda$  a hyperparameter and  $F_{old}^{i,i}$  the *i*<sup>th</sup> diagonal element of the model's Fisher on the old data.

### Definition of the Fisher

Following Martens [2], the *i*<sup>th</sup> diagonal element of the model's Fisher on the old data is defined as:

$$F_{\text{old}}^{i,i} := \mathbb{E}_{\boldsymbol{X} \sim \mathcal{D}_{\text{old}}} \left[ \mathbb{E}_{\boldsymbol{y} \sim \boldsymbol{p}_{\hat{\boldsymbol{\theta}}_{\text{old}}}} \left[ \left( \frac{\delta \log \boldsymbol{p}_{\boldsymbol{\theta}} \left( \boldsymbol{y} | \boldsymbol{X} \right)}{\delta \theta^{i}} \right|_{\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}_{\text{old}}} \right)^{2} \right] \right]$$



### **Empirical comparison**



### 

In this definition, there are two expectations:

- 1. an outer expectation over  $\mathcal{D}_{old}$ , which is the (theoretical) input distribution of the old data
- 2. an inner expectation over  $p_{\hat{\theta}_{old}}(y|\mathbf{x})$ , which is the conditional distribution of y given  $\mathbf{x}$  defined by the model after training on the old data

The different ways of computing the Fisher that can be found in the continual learning literature differ in how these two expectations are computed. *Figure 1.* Performance of EWC on Split CIFAR-10 with different ways of computing the Fisher Information for a wide range of hyperparameter values.

#### References & Acknowledgements

- [1] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017.
- [2] James Martens. New insights and perspectives on the natural gradient method. *Journal of Machine Learning Research*, 21(146):1–76, 2020.
- [3] Antonio Carta, Lorenzo Pellegrini, Andrea Cossu, Hamed Hemati, and Vincenzo Lomonaco. Avalanche: A pytorch library for deep continual learning. *Journal of Machine Learning Research*, 24(363):1–6, 2023.
- [4] Da-Wei Zhou, Fu-Yun Wang, Han-Jia Ye, and De-Chuan Zhan. PyCIL: a Python toolbox for classincremental learning. *Science China Information Sciences*, 66:197101, 2023.

Funding: senior postdoctoral fellowship, Resarch Foundation – Flanders (FWO), grant number 1266823N.