

Memory Aware Synapses: Learning what (not) to forget

Rahaf Aljundi¹, Francesca Babiloni¹, Mohamed Elhoseiny²,
Marcus Rohrbach², and Tinne Tuytelaars¹

¹ KU Leuven, ESAT-PSI, IMEC, Belgium

² Facebook AI Research

Abstract. Humans can learn in a continuous manner. Old rarely utilized knowledge can be overwritten by new incoming information while important, frequently used knowledge is prevented from being erased. In artificial learning systems, lifelong learning so far has focused mainly on accumulating knowledge over tasks and overcoming catastrophic forgetting. In this paper, we argue that, given the limited model capacity and the unlimited new information to be learned, knowledge has to be preserved or erased selectively. Inspired by neuroplasticity, we propose a novel approach for lifelong learning, coined Memory Aware Synapses (MAS). It computes the importance of the parameters of a neural network in an unsupervised and online manner. Given a new sample which is fed to the network, MAS accumulates an importance measure for each parameter of the network, based on how sensitive the predicted output function is to a change in this parameter. When learning a new task, changes to important parameters can then be penalized, effectively preventing important knowledge related to previous tasks from being overwritten. Further, we show an interesting connection between a local version of our method and Hebb’s rule, which is a model for the learning process in the brain. We test our method on a sequence of object recognition tasks and on the challenging problem of learning an embedding for predicting <subject, predicate, object> triplets. We show state-of-the-art performance and, for the first time, the ability to adapt the importance of the parameters based on unlabeled data towards what the network needs (not) to forget, which may vary depending on test conditions.

1 Introduction

The real and digital world around us evolves continuously. Each day millions of images with new tags appear on social media. Every minute hundreds of hours of video are uploaded on Youtube. This new content contains new topics and trends that may be very different from what one has seen before - think e.g. of new emerging news topics, fashion trends, social media hypes or technical evolutions. Consequently, to keep up to speed, our learning systems should be able to evolve as well.

Yet the dominating paradigm to date, using supervised learning, ignores this issue. It learns a given task using an existing set of training examples. Once the training is finished, the trained model is frozen and deployed. From then on, new incoming data is processed without any further adaptation or customization of the model. Soon, the model becomes outdated. In that case, the training process has to be repeated, using both the previous and new data, and with an extended set of category labels. In a world

