

Incremental Few-Shot Learning with Attention Attractor Networks

Seungmin Lee (profile2697@gmail.com; 2013-11420), Dept. of Computer Science and Engineering, Seoul National University

1. Introduction

This paper tried to solve *incremental few-shot learning*. In incremental few-shot learning, we assume there are base classes and novel classes which are disjoint each other. A model is initially trained using examples from base classes with sufficient labels. Then without re-training on the base classes, we train the model on few-shot labeled novel classes. After trained on novel classes, performance of the model is evaluated on both base and novel classes. If we train the model naively, the performance on base classes is largely reduced (*Catastrophic forgetting*). To avoid this problem, the authors proposed an **Attention Attractor Networks Regularizer** that can be interpreted as learned memory for base classes. They also showed the regularizer can be trained using recurrent back-propagation.

2. Methods

In this method, a model trained using examples from base classes in **Pretrain Stage 2.1**. Then the model is trained by repeating **Incremental Few-Shot Episodes Stage 2.2** and **Meta-Learning Stage 2.3**.

2.1. Pretraining Stage

This step aims to get a good feature extractor f and a base classifier parameterized with W_a before training on novel classes. Pretaining stage using typical cross entropy for the base dataset $\{(x_{a,i}, y_{a,i})\}_{i=1}^{N_a} \in \mathcal{D}_a$ where $y_{a,i} \in \{1 \dots K\}$ and $x_{a,i}$ are i -th label and example from the base dataset \mathcal{D}_a , respectively.

2.2. Incremental Few-Shot Episodes

In this step, we will do a few-shot training on the episode ϵ from the few-shot dataset \mathcal{D}_b . At this time, the episode is composed of support set S_b and query set Q_b , which play the same role as training set and validation set in supervised learning. As we progress through each episode, we learn W_b , which is called fast weights, where W_a is fixed, and the loss to W_b consists of:

$$L_s(W_b, \theta_{aE}) = \text{cross_entropy}(W_b, S_b) + R(W_b, \theta_{aE})$$

Where R is the attention attractor networks regularizer and parameterized by θ_{aE} . θ_{aE} is learned in the meta-learning stage, which is fixed at this stage. As explained earlier, the performance of the model should eventually be good for both base classes and novel classes, which means that the performance for the union $Q_a + b$ of the query set Q_a for \mathcal{D}_a and the query set Q_b for \mathcal{D}_b should

be high. R is used to alleviate the problem of poor Q_a performance due to catastrophic forgetting when only cross entropy is used. In the Incremental Few-Shot Episodes phase, the meta parameters θ_{aE} are fixed. This regularizer term and θ_{aE} are described in the Attention Attractor Networks Regularizer.

2.2.1 Attention Attractor Networks Regularizer

This section describes the Regularizer R . $R(W_b, \theta_{aE})$ is defined as

$$R(W_b, \theta_{aE}) = \sum_k' 1, K' \text{ squared_mahalanobis_distance}(W_b, k', u_k)$$

In this case, γ is a learnable parameter and attractor u_k is defined as follows.

$$u_k' = \sum_k a_k', k U_k + U_0$$

The regularizer term induces the encoding of information used to classify base classes due to attractors when W_b encodes information about a novel class. This is because U_k is the memory encoding base class k , and attractor is the weighted sum of U_k . This term prevents catastrophic forgetting and learns about novel classes.

2.3. Meta-Learning Stage

In this step, you learn the meta parameters θ_{aE} for the episode ϵ . This θ_{aE} is trained to minimize the expected cross entropy loss for $Q_a + b$ generated through multiple episodes. At this time, Recurrent Back-propagation can be used to learn θ_{aE} . The reason for using this is that BPTT takes too long to learn, and T-BPTT is prone to bias in performance depending on the T value.