

Energy Minimum Regularization in Continual Learning

Xiaobin Li¹, Lianlei Shan¹, Minglong Li¹, Weiqiang Wang^{1,*}

¹University of Chinese Academy of Sciences
Computer Vision and Multimedia Technology Lab
{lixiaobin161,shanlianlei18, liminglong18, }@mailsucas.ac.cn
wqwang@ucas.ac.cn

December 2, 2020

Background

How to give agents the ability of continuous learning like human and animals is still a challenge. In the regularized continual learning method OWM, the constraint of the model on the energy compression of the learned task is ignored, which results in the poor performance of the method on the dataset with a large number of learning tasks.

Contributions

- We propose an energy minimum regularization(EMR) method to constrain the energy of learned tasks, providing enough learning space for the following tasks that are not learned, and increasing the capacity of the model to the number of learning tasks.
- We propose a new measurement method called AD to measure the anti-degradation degree of model.
- Extensive experiments show the superiority of EMR in learning sequential tasks and EMR can make the model less sensitive to multiple tasks and network size.

Energy Minimum Regularization

We define the energy of a task T_i as

$$\mathbb{E}(T_i) = \sum_{i=1}^{N_l} \sum_{l=1}^L \frac{\|\mathbf{x}_i^l \mathbf{x}_i^{lT}\|_F^2}{\mathbf{x}_i^{lT} \mathbf{x}_i^l}$$

In our EMR model, the loss function is

$$L = L_{cls} + \gamma L_{emr} = L_{cls} + \gamma \mathbb{E}(T) \quad (1)$$

where L_{cls} is classification loss.

Improved Energy Minimum Regularization

In order to reduce the calculation burden, we improve the calculation formula of energy as follow.

$$\mathbb{E}(T_i) = \sum_{i=1}^{N_l} \sum_{l=1}^L \mathbb{E}(\hat{\mathbf{x}}_i^l) \quad (2)$$

For each element $x_{ij}^l \in \mathbf{x}_i^l$,

$$\hat{x}_i^l = \begin{cases} x_i^l & x_i^l < \beta_l \\ 0 & x_i^l \geq \beta_l \end{cases} \quad (3)$$

where β_l is a constant threshold. After modification, we get new feature $\hat{\mathbf{x}}_i^l = [\hat{x}_{i1}^l, \hat{x}_{i2}^l, \dots, \hat{x}_{id}^l]$. The energy of feature $\hat{\mathbf{x}}_i^l$ is

$$\mathbb{E}(\hat{\mathbf{x}}_i^l) = \sum_{j=1}^d |\hat{x}_{ij}^l| \quad (4)$$

Experiments

Experiments and Analysis on MNIST

Methods accuracy comparison on ordered MNIST and 10-shuffled MNIST.
"SFT" denotes sequential fine tuning.

Methods	Ordered MNIST(%)	10-Shuffled MNIST(%)
SFT	10.01	10.03 ± 0.01
EWC	53.6	52.72 ± 1.36
CAB	95.03	94.91 ± 0.30
OWM	96.71	96.30 ± 0.03
EMR	97.89	97.51 ± 0.05

Experiments

Experiments and Analysis on Scene Datasets

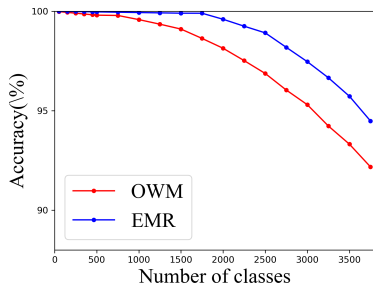
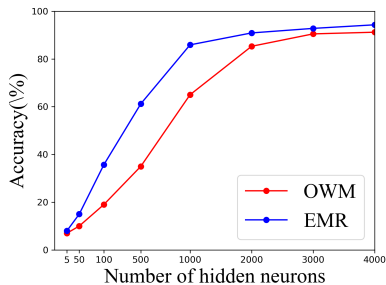
Comparison of performance of different methods in the disjoint CIFAR10 task and ImageNet task. The quantitative metrics is accuracy(%).

Methods	CIFAR10	ImageNet
Pre-train	None	resnet152
SFT	20.14	0.69
EWC	31.09	-
OWM	52.83	73.80
EMR	53.72	76.29

Experiments

Ablation Experiments on CASIA-HWDB

Model Capacity Model capacity experiments on CASIA-HWDB dataset. *Left.* The performance comparison of methods EMR and OWM on tasks with different categories. The hidden neurons is fixed to 4000. *Right.* The performance comparison of methods EMR and OWM on tasks with different neurons. The tasks number is 3755.



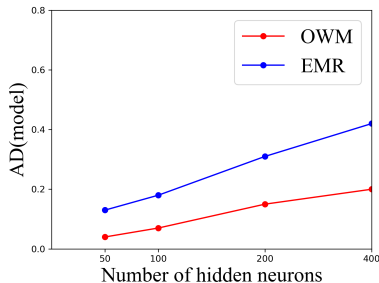
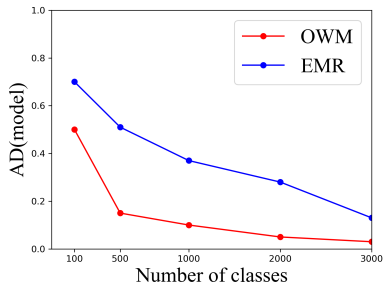
Experiments

Ablation Experiments on CASIA-HWDB

Anti-degradation of Model

We define the anti-degradation degree of the model as

$$AD(model) = \frac{1}{N_c L} \sum_{i=1}^{N_c} \sum_{l=1}^L \frac{\|P^T \mathbf{x}_i^l\|}{\|\mathbf{x}_i^l\|} \quad (5)$$



Conclusion

In this paper, we propose a novel method called energy minimum regularization (EMR) to effectively address the issue of catastrophic forgetting in continual learning and model capacity problem. The proposed EMR has not only a solid theory foundation, but also obtain the support of experimental results. Extensive experiments show the superiority of EMR in learning sequential tasks and EMR can make the model less sensitive to multiple tasks and network size.

Thanks!