



Continual Learning with Guarantees via Weight Interval Constraints

***Maciej Wołczyk, Karol Piczak,**
Bartosz Wójcik, Łukasz Pustelnik, Paweł Morawiecki,
Jacek Tabor, Tomasz Trzciński, Przemysław Spurek*

Background and motivation

- **Continual learning** – learning from a stream of data while:
 - reducing forgetting,
 - maximizing knowledge transfer,
 - and satisfying other desiderata on computation and memory.
- The no-forgetting is usually a soft constraint, but sometimes remembering might be critical!
- In this work, we devise a method which can put guarantees on forgetting.
- **Main idea: define a region of viable parameters for the current task, and stay within it for the rest of the training.**

Definition

- In this more restrictive setting, the learning objective for task T_j can be written as [1]:

$$\arg \min_{\theta} \ell(T_j, \theta) \text{ satisfying } \ell(T_m, \theta) \leq \ell(T_m, \theta_m^*)$$

for all $m = 1, \dots, j - 1$

where θ_m^* are the parameters obtained directly after learning task m

and ℓ is the cross-entropy loss over the whole task dataset.

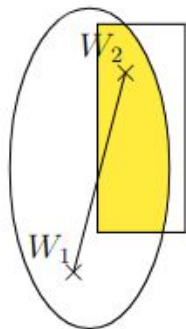
- Alternatively, we look for parameter of the new task θ within the set of parameters Θ_m that satisfies the inequality condition.

Perfect CL is NP-hard

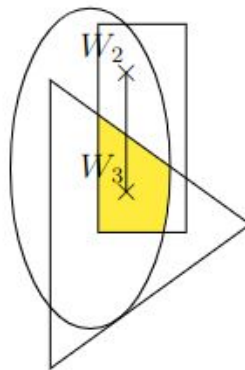
TASK 1



TASK 2

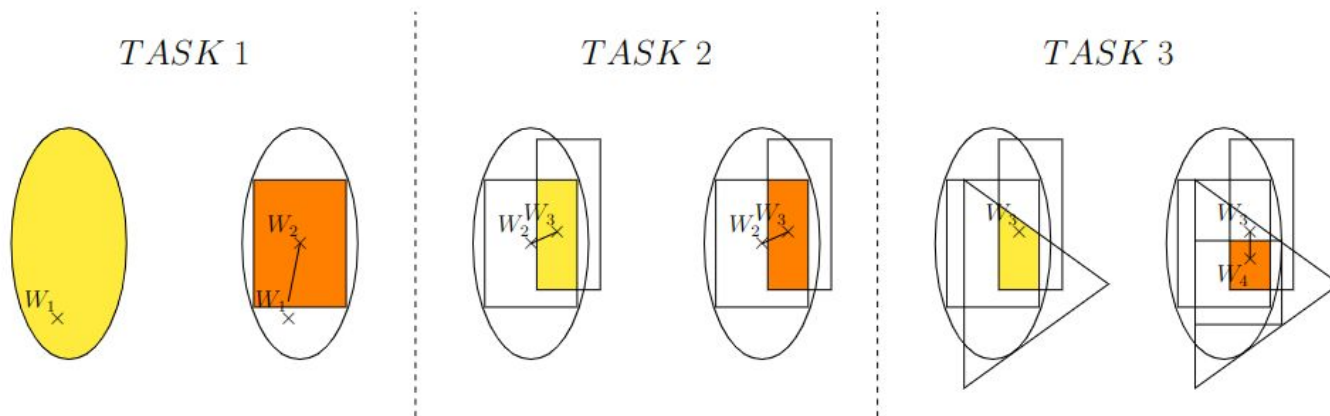


TASK 3



- Problem: the viable parameter region Θ_m is highly irregular
- Previous work [2] showed that even if we assume that the regions are polytopes, it's still NP-hard.
- In our work, we try to solve this by considering a simpler setting with additional assumptions.

Intersection of parameter regions



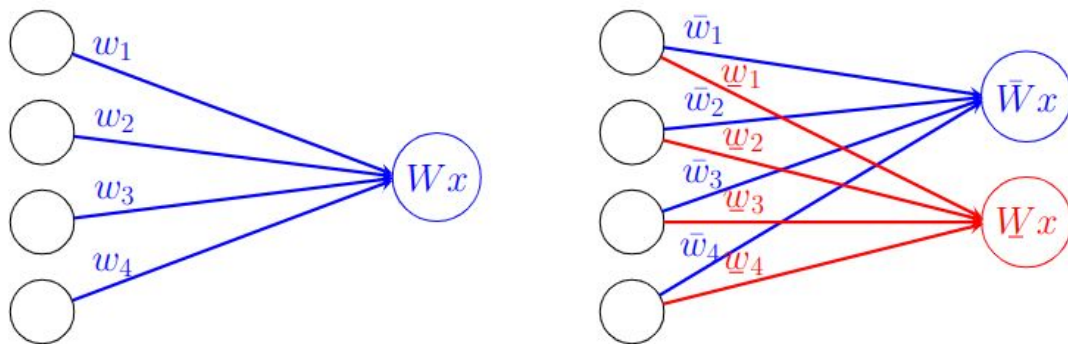
- At each point, we try to find the largest hyperrectangle that fits in the viable region Θ_m .
- Finding intersections between hyperrectangles can be done in polynomial time.
- In our work, we show how to implement this idea efficiently through interval arithmetic.

Practical implementation – intervals

- In practice, we implement the hyperrectangle region by representing each weight as an interval rather than a point in space:

$$[\underline{W}_k, \bar{W}_k] = [W_k - \varepsilon_k, W_k + \varepsilon_k]$$

- We use the tools from interval arithmetic to implement interval neural networks:



Interval upper bound

- We can find an upper bound $\hat{\ell}(T_j, \Theta)$ on the worst-case loss for task T_j on the parameter region Θ represented by the hyperrectangle:

$$\hat{\ell}(T_j, \Theta) \geq \max_{\theta \in \Theta_m} \ell(T_j, \theta)$$

- As long as we stay within Θ during the training, the loss won't be higher than $\hat{\ell}(T_j, \Theta)$.
- This is easy to compute and differentiable wrt. parameters.
- We can minimize it to find a region of parameters with low loss everywhere

InterContiNet

Algorithm 1 InterContiNet training procedure for a given task

Input: model trained on the previous task (weights W^* , radii ε^*), current task T_j .

Reparameterize W, ε using Eq. (5).

for epoch in $1 \dots \text{center_epochs}$ **do**

 Update μ by minimizing $\ell(T_j, \Theta)$ {Train centers}

end for

Initialize ε as largest possible within the previous interval.

for epoch in $1 \dots \text{radii_epochs}$ **do**

 Update ν by minimizing $\hat{\ell}(T_j, \Theta)$ {Train radii}

if $\hat{\text{acc}} \geq \text{acc} \cdot \text{acc_thresh}$ **then**

 break

end if

end for

return W, ε

Our method works in two stages:

1. **Find the centers W** that minimize cross-entropy loss while staying within the parameter region.
2. **Shrink the radii ε** to minimize the worst-case loss.

Experiments

Table 1. The average accuracy across all five tasks of the split MNIST protocol, evaluated after learning the whole sequence. Each value is the average of five runs (with standard deviations).

Method	Incremental task	Incremental domain	Incremental class
SGD	96.27 \pm 0.38	64.58 \pm 0.26	19.01 \pm 0.04
Adam	95.53 \pm 3.16	59.32 \pm 1.08	19.74 \pm 0.01
L2	96.31 \pm 0.41	72.21 \pm 0.17	18.88 \pm 0.18
EWC	97.01 \pm 0.13	76.90 \pm 0.41	18.90 \pm 0.06
oEWC	97.01 \pm 0.13	77.02 \pm 0.51	18.89 \pm 0.07
SI	96.19 \pm 0.63	80.62 \pm 0.17	17.94 \pm 0.57
MAS	96.52 \pm 0.14	84.41 \pm 0.42	17.38 \pm 4.19
LwF	97.03 \pm 0.05	82.76 \pm 0.17	49.37 \pm 0.68
InterContiNet	98.93 \pm 0.05	77.77 \pm 1.24	40.73 \pm 3.26
Offline	99.74 \pm 0.03	99.03 \pm 0.04	98.49 \pm 0.02

Table 2. The average accuracy across all five tasks of the split FashionMNIST protocol, evaluated after learning the whole sequence. Each value is the average of five runs (with standard deviations).

Method	Incremental task	Incremental domain	Incremental class
SGD	92.22 \pm 3.06	82.77 \pm 0.44	19.91 \pm 0.01
Adam	88.13 \pm 5.64	78.48 \pm 0.47	19.96 \pm 0.01
L2	97.36 \pm 0.17	92.65 \pm 0.09	26.91 \pm 1.23
EWC	97.53 \pm 0.15	92.12 \pm 0.18	19.90 \pm 0.01
oEWC	96.70 \pm 0.57	88.83 \pm 0.30	19.87 \pm 0.02
SI	97.00 \pm 0.25	91.45 \pm 0.07	19.97 \pm 0.34
MAS	97.43 \pm 0.14	91.74 \pm 0.19	10.00 \pm 0.00
LwF	98.10 \pm 0.07	88.63 \pm 0.12	39.51 \pm 1.45
InterContiNet	98.37 \pm 0.06	92.65 \pm 0.40	35.11 \pm 0.02
Offline	97.98 \pm 0.05	96.39 \pm 0.06	82.54 \pm 0.13

Ablation studies

- By tuning the hyperparameters we can balance **plasticity** (accuracy on the current task, top row) and **stability** (accuracy on the first task, bottom row).
- **Check out our paper for more details!**

