
CPR: Classifier-Projection Regularization for Continual Learning

Sungmin Cha¹ Hsiang Hsu² Flavio P. Calmon² Taesup Moon¹

Abstract

We propose a general, yet simple patch that can be applied to existing regularization-based continual learning methods called classifier-projection regularization (CPR). Inspired by both recent results on neural networks with wide local minima and information theory, CPR adds an additional regularization term that maximizes the entropy of a classifier’s output probability. We demonstrate that this additional term can be interpreted as a projection of the conditional probability given by a classifier’s output to the uniform distribution. By applying the Pythagorean theorem for KL divergence, we then prove that this projection may (in theory) improve the performance of continual learning methods. In our extensive experimental results, we apply CPR to several state-of-the-art regularization-based continual learning methods and benchmark performance on popular image recognition datasets. Our results demonstrate that CPR indeed promotes a wide local minima and significantly improves both accuracy and plasticity while simultaneously mitigating the catastrophic forgetting of baseline continual learning methods.

1. Introduction

Catastrophic forgetting (McCloskey & Cohen, 1989) is a central challenge in continual learning (CL): when training a model on a new task, there may be a loss of performance (e.g., decrease in accuracy) when applying the updated model to previous tasks. At the heart of catastrophic forgetting is the *stability-plasticity dilemma* (Carpenter & Grossberg, 1987; Mermillod et al., 2013), where a model exhibits high stability on previously trained tasks, but suffers

¹Department of Electrical and Computer Engineering, Sungkyunkwan University, Suwon, South Korea ²John A. Paulson School of Engineering and Applied Sciences, Harvard University, MA, USA. Correspondence to: Taesup Moon <tsmoon@skku.edu>.

from low plasticity for the integration of new knowledge (and vice-versa). Attempts to overcome this challenge in neural network-based CL can be grouped into three main strategies: regularization methods (Li & Hoiem, 2017; Kirkpatrick et al., 2017; Zenke et al., 2017; Nguyen et al., 2018; Ahn et al., 2019; Aljundi et al., 2018b), memory replay (Lopez-Paz & Ranzato, 2017; Shin et al., 2017; Rebuffi et al., 2017; Kemker & Kanan, 2017), and dynamic network architecture (Rusu et al., 2016; Yoon et al., 2018; Golkar et al., 2019). In particular, regularization methods that control model weights bear the longest history due to its simplicity and efficiency to control the trade-off for a fixed model capacity.

In parallel, several recent methods seek to improve the generalization of neural network models trained on a single task by promoting *wide local minima* (Keskar et al., 2016; Chaudhari et al., 2019; Pereyra et al., 2017; Zhang et al., 2018). Broadly speaking, these efforts have experimentally shown that models trained with wide local minima-promoting regularizers achieve better generalization and higher accuracy (Keskar et al., 2016; Pereyra et al., 2017; Chaudhari et al., 2019; Zhang et al., 2018), are better calibrated (Pereyra et al., 2017), and can be more robust to weight perturbations (Zhang et al., 2018) when compared to usual training methods. Despite the promising results, methods that promote wide local minima have yet to be applied to CL.

In this paper, we make a novel connection between wide local minima in neural networks and regularization-based CL methods. The typical regularization-based CL aims to preserve *important* weight parameters used in past tasks by penalizing large deviations when learning new tasks. As shown in the top of Fig. 1, a popular geometric intuition (as first given in EWC (Kirkpatrick et al., 2017)) for such CL methods is to consider the (uncertainty) ellipsoid of parameters around the local minima. When learning new tasks, parameter updates are selected in order to not significantly hinder model performance on past tasks. Our intuition is that promoting a wide local minima—which conceptually stands for local minima having a *flat*, rounded uncertainty ellipsoid—can be particularly beneficial for regularization-based CL methods by facilitating diverse update directions for the new tasks (*i.e.*, improves plasticity) while not hurting the past tasks (*i.e.*, retains stability). As shown in the bottom of Fig. 1, when the ellipsoid containing the parameters with

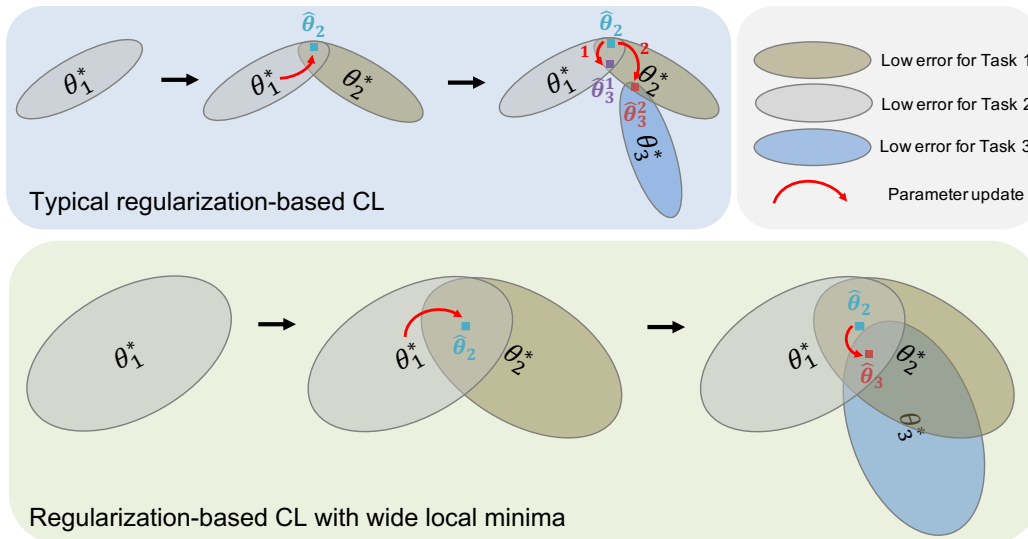


Figure 1: In typical regularization-based CL (top), when the low-error ellipsoid around local minima is sharp and narrow, the space for candidate model parameters that perform well on all tasks (*i.e.*, the intersection of the ellipsoid for each task) quickly becomes very small as learning continues, thus, an inevitable trade-off between stability and plasticity occurs. In contrast, when the *wide local minima* exists for each task (bottom), it is more likely the ellipsoids will significantly overlap even when the learning continues, hence, finding a well performing model for all tasks becomes more feasible.

low-error is wider, *i.e.*, when the wide local minima exists, there is more flexibility in finding a parameter that performs well for all tasks after learning a sequence of new tasks. We provide further details in Section 2.1.

Based on the above intuition, we propose a general, yet simple patch that can be applied to existing regularization-based CL methods dubbed as *Classifier-Projection Regularization* (CPR). Our method implements an additional regularization term that promotes wide local minima by maximizing the entropy of the classifier’s output distribution. Furthermore, from a theory standpoint, we make an observation that our CPR term can be further interpreted in terms of information projection (I-projection) formulations (Cover & Thomas, 2012; Murphy, 2012; Csiszár & Matus, 2003; Walsh & Regalia, 2010; Amari et al., 2001; Csiszár & Matus, 2003; Csiszár & Shields, 2004) found in information theory. Namely, we argue that applying CPR corresponds to projecting a classifier’s output onto a Kullback-Leibler (KL) divergence ball of finite radius centered around the uniform distribution. By applying the Pythagorean theorem for KL divergence, we then prove that this projection may (in theory) improve the performance of continual learning methods.

Through extensive experiments on several benchmark datasets, we demonstrate that applying CPR can significantly improve the performance of the state-of-the-art regularization-based CL: using our simple patch improves *both* the stability and plasticity and, hence, achieves better average accuracy almost uniformly across the tested algorithms and datasets—confirming our intuition of wide local

minima in Fig. 1. Furthermore, we use a feature map visualization that compares methods trained with and without CPR to further corroborate the effectiveness of our method.

Related work Several methods have been recently proposed to reduce catastrophic forgetting (see (Parisi et al., 2018) for a survey). In this paper, we mainly focus on regularization-based CL methods (Li & Hoiem, 2017; Kirkpatrick et al., 2017; Aljundi et al., 2018a; Chaudhry et al., 2018; Zenke et al., 2017; Nguyen et al., 2018; Ahn et al., 2019; Aljundi et al., 2018b). Broadly speaking, the motivation behind regularization-based CL is to measure the importance of model parameters in previous tasks. This measure is then used in a regularization term for overcoming catastrophic forgetting when training for new tasks. Consequently, the main research focus of regularization-based CL is creating metrics for quantifying weight importance on previous tasks (e.g., (Kirkpatrick et al., 2017; Aljundi et al., 2018a; Chaudhry et al., 2018; Zenke et al., 2017; Nguyen et al., 2018; Ahn et al., 2019)). In contrast, here we focus on developing a general method for augmenting regularization-based CL instead of proposing (yet another) new metric for weight importance. The method introduced here, namely CPR, can be applied to *any* regularization-based CL method to simultaneously improve both plasticity and stability.

The work closest to ours is (Aljundi et al., 2018b), which encourages sparsity of representations for each task by adding an additional regularizer to regularization-based CL methods. Note that the motivation of (Aljundi et al., 2018b)—imposing sparsity of neuron activations—is considerably different from ours, which is to promote wide local minima.

Moreover, whereas (Aljundi et al., 2018b) focuses on average accuracy, we carefully evaluate in our experiments the advantage of the added CPR regularization in terms of increasing both plasticity and stability of CL *in addition to* accuracy.

Several papers have recently proposed methods that promote wide local minima in neural networks in order to improve single-task generalization, including using small mini-batch size (Keskar et al., 2016), regularizing the output of the softmax layer in neural networks (Szegedy et al., 2016; Pereyra et al., 2017), using a newly proposed optimizer which constructs a local-entropy-based objective function (Pereyra et al., 2017) and distilling knowledge from other models (Zhang et al., 2018). We expand upon this prior work and investigate here the role of wide local minima in CL. Our objective is to train neural networks that converge to wide local minima for each task, and subsequently benchmark the advantage of wide local minima in CL through numerous experiments. To the best of our knowledge, this is the first paper to study the role of wide local minima in CL.

2. CPR: Classifier-Projection Regularization for Wide Local Minimum

In this section, we elaborate in detail the core motivation outlined in Fig. 1, then formalize CPR as the combination of two regularization terms: one stemming from prior regularization-based CL methods, and the other that promotes a wide local minima. Moreover, we provide an information-geometric interpretation (Csiszár, 1984; Cover & Thomas, 2012; Murphy, 2012) for the observed gain in performance when applying CPR to CL.

We consider continual learning of T classification tasks (with known task boundaries), where each task contains N training sample-label pairs $\{(\mathbf{x}_n^t, y_n^t)\}_{n=1}^N$, $t \in [1, \dots, T]$ with $\mathbf{x}_n^t \in \mathbb{R}^d$, and the labels of each task has M_t classes, *i.e.*, $y_n^t \in [1, \dots, M_t]$. We denote $f_\theta : \mathbb{R}^d \rightarrow \Delta_M$ as a neural network-based classification model with softmax output layer parametrized by θ .

2.1. Motivation: Introducing wide local minima in continual learning

Consider the setting of typical regularization-based CL (top of Fig. 1). We denote θ_i^* as a local minima of task i . From the shape of the low-error ellipsoids, after learning task 2, an appropriate regularization strength can make the parameter update from θ_1^* to $\hat{\theta}_2$ since it can achieve low-errors on both tasks 1 and 2. However, while learning task 3, the ellipsoids may not overlap enough, and it becomes infeasible to obtain a parameter that performs well on all three tasks. In this case, regularization-based CL determines the trade-off between stability and plasticity in terms of its regularization strength; namely, the larger strength (direction 1) results in

a parameter with more stability, $\hat{\theta}_3^1$, so that less forgetting on tasks 1 and 2 is achieved, whereas the smaller strength (direction 2) leads to more plasticity so that the updated parameter $\hat{\theta}_3^2$ performs well on more recent tasks (2 and 3) at the cost of compromising the performance for task 1.

In contrast, when the wide local minima exists for each task (bottom of Fig. 1), it is more likely that the ellipsoids will have non-empty intersections. A regularization-based CL may therefore more easily find a parameter, $\hat{\theta}_3$, that is simultaneously close to the the local minimas for each task, *i.e.*, $\{\theta_i^*\}_{i=1}^3$. This intuition suggests that once we promote the wide local minima of neural networks during continual learning, both the stability and plasticity will improve and result in higher accuracy—which is precisely what we verify in our experimental results for CPR (see Sec. 3).

2.2. Classifier projection regularization for continual learning

Regularization-based continual learning Typical regularization-based CL methods attach a regularization term that penalizes the deviation of important parameters learned from past tasks in order to mitigate catastrophic forgetting. The general loss form for these methods when learning task t is

$$L_{\text{CL}}^t(\theta) = L_{\text{CE}}^t(\theta) + \lambda \sum_i \Omega_i^{t-1} (\theta_i - \theta_i^{t-1})^2, \quad (1)$$

where $L_{\text{CE}}^t(\theta)$ is the ordinary cross-entropy loss function for task t , λ is the dimensionless regularization strength, $\Omega^{t-1} = \{\Omega_i^{t-1}\}$ is the set of estimates of the weight importance, and $\{\theta_i^{t-1}\}$ is the parameter learned until task $t-1$. A variety of previous work, *e.g.*, EWC (Kirkpatrick et al., 2017), SI (Zenke et al., 2017), MAS (Aljundi et al., 2018a), and RWalk (Chaudhry et al., 2018), proposed different ways of calculating Ω^{t-1} to measure weight importance.

Single-task wide local minima Several recent schemes have been proposed (Pereyra et al., 2017; Szegedy et al., 2016; Zhang et al., 2018) to promote wide local minima of a neural network for solving a single task. These approaches can be unified by the following common loss form

$$L_{\text{WLM}}(\theta) = L_{\text{CE}}(\theta) + \frac{\beta}{N} \sum_{n=1}^N D_{\text{KL}}(f_\theta(\mathbf{x}_n) \| g), \quad (2)$$

in which g is some probability distribution in Δ_M that regularizes the classifier output f_θ , β is a trade-off parameter, and $D_{\text{KL}}(\cdot \| \cdot)$ is the KL divergence (Cover & Thomas, 2012). Notice, for example, when g is the uniform distribution P_U in Δ_M , the regularization term corresponds to maximizing the entropy as in (Pereyra et al., 2017), and when g is another classifier’s output $f_{\theta'}$, then (2) becomes equivalent to the loss function proposed in (Zhang et al., 2018).

CPR: Achieving wide local minima in continual learning Combining the above two regularization terms, we propose the CPR as the following loss form for learning task t :

$$L_{\text{CPR}}^t(\theta) = L_{\text{CE}}^t(\theta) + \frac{\beta}{N} \sum_{n=1}^N D_{\text{KL}}(f_{\theta}(\mathbf{x}_n^t) \| P_U) + \lambda \sum_i \Omega_i^{t-1} (\theta_i^t - \theta_i^{t-1})^2, \quad (3)$$

where λ and β are the regularization parameters. The first regularization term promotes the wide local minima while learning task t by using P_U as the regularizing distribution g in (2), and the second term is from the typical regularization-based CL. Note that this formulation is oblivious to Ω_{t-1} and, hence, it can be applied to *any* state-of-the-art regularization-based CL methods. In our experiments, we show that the simple addition of the KL-term can significantly boost the performance of several representative state-of-the-art methods, confirming our intuition on wide local minima for CL given in Section 2.1 and Fig 1. Furthermore, we show in the next section that the KL-term can be geometrically interpreted in terms of I-projections (Csiszár, 1984; Cover & Thomas, 2012; Murphy, 2012), providing an additional argument (besides promoting wide local minima) for the benefit of using CPR in continual learning.

2.3. Interpretation by information projection

Given a distribution P and a convex set of distributions \mathcal{Q} in the probability simplex $\Delta_m \triangleq \{\mathbf{p} \in [0, 1]^m \mid \sum_{i=1}^m \mathbf{p}_i = 1\}$, information projection (I-projection) aims to find P^* in \mathcal{Q} such that the KL divergence between P^* and P is minimized, *i.e.*,

$$P^* = \arg \min_{Q \in \mathcal{Q}} D_{\text{KL}}(Q \| P). \quad (4)$$

The above quantity has several operational interpretations in information theory (e.g., in universal source coding (Cover & Thomas, 2012)). The I-projection enables a “geometric” interpretation of KL divergence, where $D_{\text{KL}}(Q \| P)$ behaves as the squared Euclidean distance, (Q, P^*, P) form a “right triangle,” and the following lemma resembles the KL divergence equivalent of the Pythagorean theorem (not satisfied in general by the KL divergence) (Cover & Thomas, 2012).

Lemma 1. *Suppose $\exists P^* \in \mathcal{Q}$ such that $D_{\text{KL}}(P^* \| P) = \min_{Q \in \mathcal{Q}} D_{\text{KL}}(Q \| P)$, then*

$$D_{\text{KL}}(Q \| P) \geq D_{\text{KL}}(Q \| P^*) + D_{\text{KL}}(P^* \| P), \quad \forall Q \in \mathcal{Q}. \quad (5)$$

A natural extension of the I-projection is to seek the conditional distribution $Q_{Y|X}$ in a set \mathcal{C} that is closest (measured by the KL divergence) to a given conditional distribution

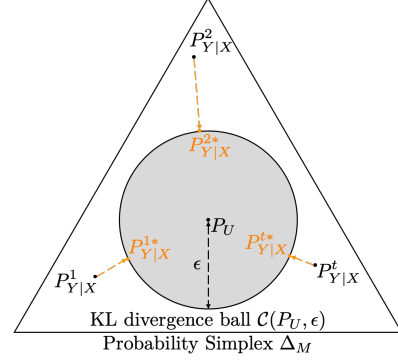


Figure 2: CPR can be understood as a projection onto a finite radius ball around P_U .

$P_{Y|X}$. Viewing a classifier (e.g., a neural network with a softmax output layer) as a conditional probability distribution $P_{Y|X}$, where Y is the class label and X is the input, we call this extension as the *classifier projection*.

Formally, given a convex set \mathcal{C} of conditional distributions, the classifier projection is defined as

$$P_{Y|X}^* = \arg \min_{Q_{Y|X} \in \mathcal{C}} \mathbb{E}_{P_X} [D_{\text{KL}}(Q_{Y|X}(\cdot | X) \| P_{Y|X}(\cdot | X))]. \quad (6)$$

We consider a simple CL setting with single head and fixed number of classes. Then, we pick the set of possible classifiers \mathcal{C} to be a KL divergence ball centered at the uniform distribution P_U , *i.e.*,

$$\mathcal{C}(P_U, \epsilon) \triangleq \{Q_{Y|X} \in \Delta_M \mid \mathbb{E}_X [D_{\text{KL}}(Q_{Y|X} \| P_U)] \leq \epsilon\}.$$

We select P_U since it is the centroid of Δ_M and, hence, the worst-case divergence between any distribution and P_U is at most $\log M$. The following proposition is a direct consequence of Lemma 1.

Proposition 1. *For any classifier $P_{Y|X}^{t-1*} \in \mathcal{C}(P_U, \epsilon)$ for task $t-1$ with data distribution P_X^{t-1} , and let any classifier for task t be $P_{Y|X}^t \notin \mathcal{C}(P_U, \epsilon)$ and $P_{Y|X}^{t*}$ be the projected classifier by (6), then*

$$\mathbb{E}_{P_X^{t-1} P_X^t} \left[-\log P_{Y|X}^t P_X^{t-1} \right] \geq \mathbb{E}_{P_X^{t-1} P_X^t} \left[-\log P_{Y|X}^{t*} P_X^{t-1} \right]. \quad (7)$$

Proposition 1 indicates that when evaluated on the previous task, the classifier of the current task is more similar (in terms of cross-entropy) to each other after projection, thus guaranteeing a smaller change in training loss and accuracy. From the vantage point of classifier projection, the CPR regularization term in (3) can be viewed as the Lagrange dual of the constraint $Q_{Y|X} \in \mathcal{C}(P_U, \epsilon)$ —the term that projects the classifier of individual tasks towards the uniform distribution in order to minimize changes when training sequential tasks (See Fig. 2).

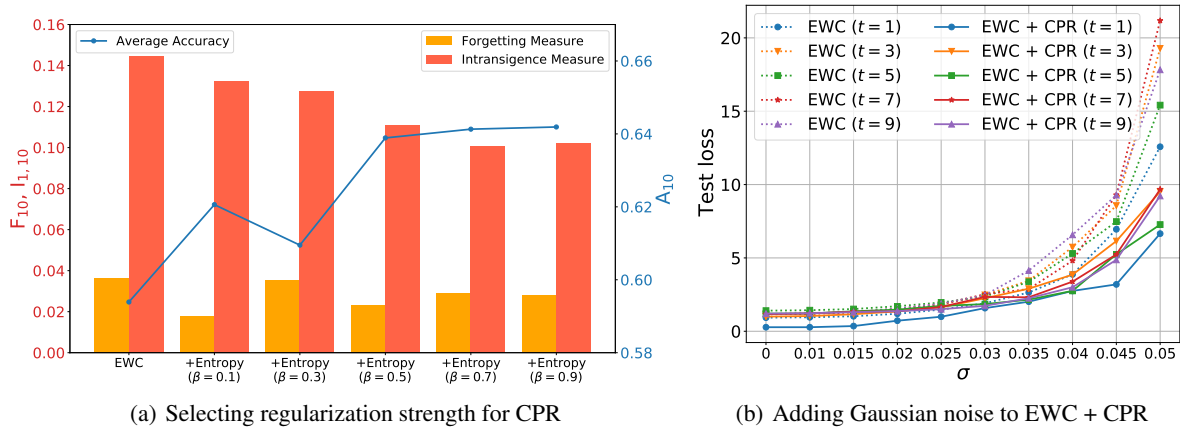


Figure 3: Verifying the regularization for wide local minima

3. Experimental Results

We apply CPR to four regularization-based supervised CL methods: EWC (Kirkpatrick et al., 2017), SI (Zenke et al., 2017), MAS (Aljundi et al., 2018a), and RWalk (Chaudhry et al., 2018), and further analyze CPR via ablation studies and feature map visualizations.

3.1. Data and evaluation metrics

We select CIFAR-100 (Krizhevsky et al., 2009), CIFAR-10/100 (Krizhevsky et al., 2009), Omniglot (Lake et al., 2015), and CUB200 (Welinder et al., 2010) as benchmark datasets. Note that we ignore the permuted-MNIST dataset (LeCun et al., 1998) since most state-of-the-art algorithms can already achieve near perfect accuracy on it. CIFAR-100 is divided into 10 tasks where each task has 10 classes. CIFAR-10/100 additionally uses CIFAR-10 for pre-training before learning tasks from CIFAR-100. Omniglot has 50 tasks, where each task is a binary image classification on a given alphabet. For these datasets, we used a simple feed-forward convolutional neural network (CNN) architecture. For the more challenging CUB200 dataset, which has 10 tasks with 20 classes for each task, we used a pre-trained ResNet-18 (He et al., 2016) as the initial model. Training details, model architectures, hyperparameters tuning, and source codes are available in the Supplementary Material (SM).

For evaluation, we first let $a_{k,j} \in [0, 1]$ be the j -th task accuracy after training the k -th task ($j \leq k$). Then, we used the following three metrics to measure the continual learning performance:

- **Average Accuracy (A)** is the average accuracy A_k on the first k tasks after training the k -th task, i.e., $A_k = \frac{1}{k} \sum_{j=1}^k a_{k,j}$. While being a natural metric, Average Accuracy fails to explicitly measure the plasticity and stability of a CL method.
- **Forgetting Measure (F)** evaluates stability. Namely,

we define the forgetting measure f_k^j of the j -th task after training k -th task as $f_k^j = \max_{l \in \{j, \dots, k-1\}} a_{l,j} - a_{k,j}, \forall j < k$, and the *average forgetting measure* F_k of a CL method as $F_k = \frac{1}{k-1} \sum_{j=1}^{k-1} f_k^j$.

- **Intransigence Measure (I)** measures the plasticity. Let a_j^* be accuracy of a model trained by fine-tuning for the j -th task *without* applying any regularization. The intransigence measure $I_{s,k}$ is then defined as $I_{s,k} = \frac{1}{k-s+1} \sum_{j=s}^k i_j$, where $i_j = a_j^* - a_{j,j}$.

The F and I metrics were originally proposed in (Chaudhry et al., 2018), and we slightly modified their definitions for our usage. Note that a **low** F_k and $I_{1,k}$ implies high stability (low forgetting) and high plasticity (good forward transfer) of a CL method, respectively.

3.2. Quantifying the role of wide local minima regularization

We first demonstrate the effect of applying CPR with varying trade-off parameter β in (3) by taking EWC (Kirkpatrick et al., 2017) trained on CIFAR-100 as a running example. Fig. 3(a) shows how the aforementioned metrics varies as β changes over $[0.1, \dots, 1]$. First, we observe that A_{10} certainly increases as β increases. Moreover, we can break down the gain in terms of $I_{1,10}$ and F_{10} ; we observe $I_{1,10}$ monotonically decreases as β increases, but F_{10} does not show the similar monotonicity although it also certainly decreases with β . This suggests that enlarged wide local minima is indeed helpful for improving both plasticity and stability. In the subsequent experiments, we selected β using validation sets by considering all three metrics; among the β 's that achieve sufficiently high A_{10} , we chose one that can reduce F_{10} more than reducing $I_{1,10}$, since it turns out improving the stability seems more challenging. (In fact, in some experiments, when we simply consider A_{10} , the chosen β will result in the lowest $I_{1,10}$ but with even higher F_{10} than the case without CPR.) For comparison purposes,

Table 1: Experimental results on CL benchmark dataset with and without CPR. Blue color denotes the case which CL method is positively affected by CPR and red color represents a negative impact of CPR.

| Dataset | Method | Average Accuracy (A_{10}) | | | Forgetting Measure (F_{10}) | | | Intransigence Measure ($I_{1,10}$) | | |
|-----------------------------|--------|-------------------------------|-----------------|------------------|---------------------------------|------------------|------------------|--------------------------------------|------------------|------------------|
| | | W/o CPR | W/ CPR | diff (W-W/o) | W/o CPR | W/ CPR | diff (W/-W/o) | W/o CPR | W/ CPR | diff (W-W/o) |
| CIFAR100 ($T = 10$) | EWC | 0.6002 | 0.6328 | +0.0326 (+5.2%) | 0.0312 | 0.0285 | -0.0027 (-8.7%) | 0.1419 | 0.1117 | -0.0302 (-21.3%) |
| | SI | 0.6141 | 0.6476 | +0.0336 (+5.5%) | 0.1106 | 0.0999 | -0.0107 (-9.7%) | 0.0566 | 0.0327 | -0.0239 (-42.2%) |
| | MAS | 0.6172 | 0.6510 | +0.0338 (+5.5%) | 0.0416 | 0.0460 | +0.0044 (+10.6%) | 0.1155 | 0.0778 | -0.0377 (-32.6%) |
| | Rwalk | 0.5784 | 0.6366 | +0.0581 (+10.0%) | 0.0937 | 0.0769 | -0.0169 (-18.0%) | 0.1074 | 0.0644 | -0.0430 (-40.0%) |
| CIFAR10/100 ($T = 11$) | EWC | 0.6950 | 0.7055 | +0.0105 (+1.5%) | 0.0228 | 0.0181 | -0.0048 (-21.1%) | 0.1121 | 0.1058 | -0.0062 (-5.5%) |
| | SI | 0.7127 | 0.7186 | +0.0059 (+0.8%) | 0.0459 | 0.0408 | -0.0051 (-11.1%) | 0.0733 | 0.0721 | -0.0012 (-1.6%) |
| | MAS | 0.7239 | 0.7257 | +0.0017 (+0.2%) | 0.0479 | 0.0476 | -0.0003 (-0.6%) | 0.0603 | 0.0588 | -0.0015 (-2.5%) |
| Omniglot ($T = 50$) | Rwalk | 0.6934 | 0.7046 | +0.0112 (+1.6%) | 0.0738 | 0.0707 | -0.0031 (-4.2%) | 0.0672 | 0.0589 | -0.0084 (-12.5%) |
| | EWC | 0.6632 | 0.8387 | +0.1755 (+26.5%) | 0.2096 | 0.0321 | -0.1776 (-84.7%) | -0.0227 | -0.0239 | -0.0012 (-5.3%) |
| | SI | 0.8478 | 0.8621 | +0.0143 (+1.7%) | 0.0247 | 0.0167 | -0.0079 (-32.0%) | -0.0258 | -0.0282 | -0.0065 (-25.3%) |
| | MAS | 0.8401 | 0.8679 | +0.0278 (+3.3%) | 0.0316 | 0.0101 | -0.0215 (-68.0%) | -0.0247 | -0.0314 | -0.0067 (-27.1%) |
| CUB200 ($T = 10$) | Rwalk | 0.8056 | 0.8497 | +0.0440 (+5.5%) | 0.0644 | 0.0264 | -0.0380 (-59.0%) | -0.0226 | -0.0294 | -0.0068 (-30.1%) |
| | EWC | 0.5363 | 0.5864 | +0.0501 (+9.3%) | 0.0437 | 0.0494 | +0.0058 (+13.3%) | 0.1155 | 0.0580 | -0.0575 (-49.8%) |
| | SI | 0.5457 | 0.5627 | +0.0170 (+3.1%) | 0.0531 | 0.0471 | -0.0060 (-11.3%) | 0.0954 | 0.0838 | -0.0116 (-12.2%) |
| | MAS | 0.5857 | 0.5952 | +0.0096 (+1.6%) | 0.0690 | 0.0626 | -0.0065 (-9.4%) | 0.0411 | 0.0373 | -0.0037 (-9.0%) |
| Rwalk | 0.5261 | 0.5567 | +0.0306 (+5.8%) | 0.0544 | 0.0431 | -0.0113 (-20.8%) | 0.1158 | 0.0934 | -0.0225 (-19.3%) | |

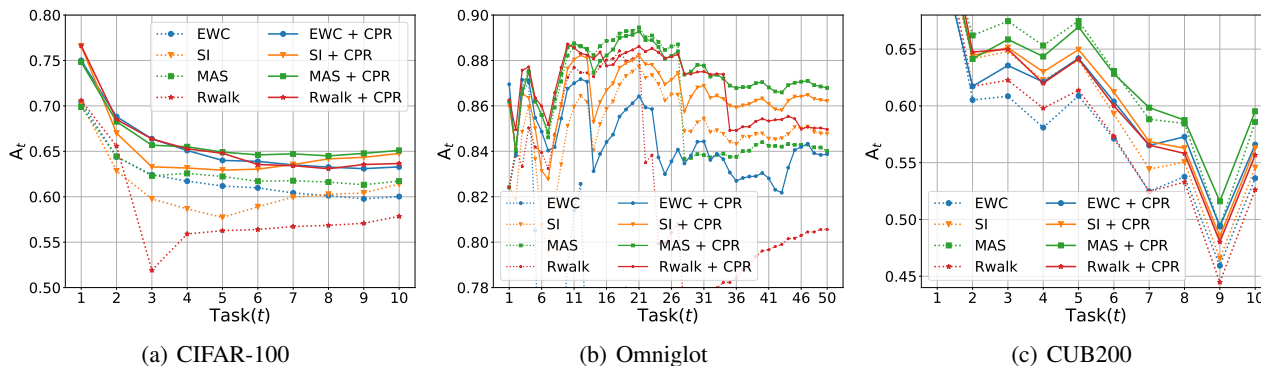


Figure 4: Experimental results on CL benchmark dataset

we also provide experiments using Deep Mutual Learning (Zhang et al., 2018) and Label Smoothing (Szegedy et al., 2016) regularizer for achieving the wide local minima in the SM; their performance was slightly worse than CPR.

With the best β in hand, Fig. 3(b) experimentally verifies whether using CPR indeed makes the local minima wide. Following the methodology in (Zhang et al., 2018), we perturb the network parameters, after learning the final task, of EWC and EWC+CPR by adding Gaussian noise with increasing σ , then measure the increase in *test* loss for each task. From the figure, we clearly observe that EWC+CPR has a smoother increase in test loss compared with EWC (without CPR) in each task. This result empirically confirms that CPR indeed promotes wide local minima for each task in CL settings and validates our initial intuition given in Sec. 2.1. In the SM, we repeat the same experiment with MAS (Aljundi et al., 2018a).

3.3. Comparison with state-of-the-art

Next, we apply CPR to the state-of-the-art regularization-based CL on the benchmark datasets and measure the performance improvement with the three metrics in Section 3.1.

For the regularization strengths, we first select the best λ without CPR, then choose β according to the procedure in Section 3.2. The results in Table 1 are averaged over 10 repeated experiments with different random initialization and task sequence using the chosen (λ, β) . The hyperparameters are reported in the SM.

CIFAR-100 and CIFAR-10/100 In Table 1 and Fig. 4(a), we observe that CPR consistently improves *all* regularization-based methods for *all* tested datasets in terms of increasing A_{10} and decreasing $I_{1,10}$, and also consistently decreases F_{10} except for MAS in CIFAR-100. Additionally, we find that for CIFAR-10/100, the orders of the 10 tasks in CIFAR-100 and CIFAR-10 affect the performance of the CPR; namely, in the SM, we show that when CIFAR-10 tasks are positioned in different positions rather than at the beginning, the gain due to CPR got much bigger.

Omniglot This dataset is well-suited to evaluate CL with long task sequences (50 tasks). In Table 1, it is clear that the CPR considerably increases both plasticity and stability in long task sequences. In particular, CPR significantly decreases F_{10} for EWC and leads to a huge improvement in A_{10} . Interestingly, unlike the previous datasets, $I_{1,10}$ is

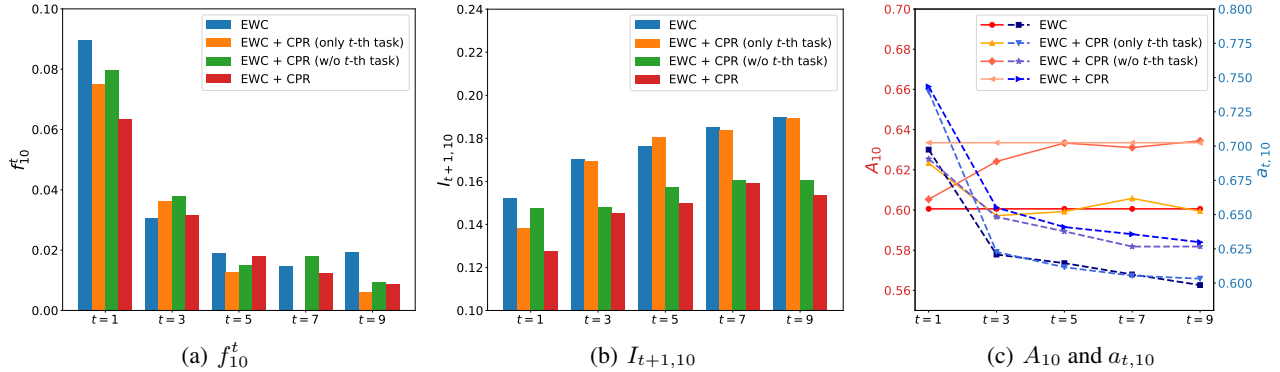


Figure 5: Ablation studies on CL with wide local minima

negative, implying that past tasks help in learning new tasks for the Omniglot dataset; when applying CPR, the gains in $I_{1,10}$ are even better. Furthermore, Fig. 4(b) indicates that applying CPR leads to less variation in A_t curves.

CUB200 The results in Table 1 and Fig. 4(c) show that CPR is also effective when using a pre-trained ResNet model for all methods and metrics, except for EWC. Here, CPR significantly increases A_{10} and reduces $I_{1,10}$ when compared to EWC, whereas F_{10} is slightly increased for EWC + CPR.

3.4. Ablation study

We study the ablation of the CPR on the regularization-based methods using CIFAR-100 with the best (λ, β) found previously, and report the averaged results over 5 random initializations and task sequences in Fig. 5. The ablation is performed in two cases: (i) using CPR only at task t , denoted as EWC + CPR (only t -th task), and (ii) using CPR except task t , denoted as EWC + CPR (w/o t -th task). Fig. 5(a) shows f_{10}^t , the amount of forgetting for task t after learning the task 10, and Fig. 5(b) shows $I_{t+1,10}$, the amount of gap with fine-tuning after task t . In Fig. 5(a), we observe that CPR helps to decrease f_{10}^t for each task whenever it is used (except for task 3), but f_{10}^t of EWC + CPR (w/o t -th task) shows a more random tendency. On average, EWC + CPR does reduce forgetting in all tasks, demonstrating the effectiveness of applying CPR to all tasks. Notably in Fig. 5(b), $I_{t+1,10}$ of EWC + CPR (only t -th task) is lower than that of EWC + CPR (w/o t -th task) only when $t = 1$; this indicates that CPR is most beneficial in terms of plasticity when CPR is applied as early as possible to the learning sequence. EWC + CPR again achieves the lowest (i.e., most favorable) $I_{t+1,10}$. Fig. 5(c), as a further evidence, also suggests that applying CPR for $t = 1$ gives a better accuracy. Moreover, the accuracy of EWC + CPR (w/o t -th task) gets closer to the optimal EWC + CPR, which is consistent with the decreasing difference of $I_{t+1,10}$ between EWC + CPR (w/o t -th task) and EWC + CPR in Fig. 5(b). The EWC + CPR still gives the best A_{10} and individual $a_{t,10}$

accuracy. We emphasize that model converging to a wide local minima from the first task onwards considerably helps the training of future tasks as well, i.e., a significant increase in the plasticity can be achieved. By using this finding, we conducted an experiment on the case where CPR have to learn unscheduled additional tasks and got the impressive experimental result which is reported in SM.

3.5. Feature map visualization using UMAP

We present next two-dimensional UMAP (McInnes et al., 2018) embeddings to visualize the impact of CPR on learnt representations. We compare representations produced by models trained on CIFAR-100 in two cases: (i) an oracle model which learns from the first and the t -th task at training time t , and (ii) sequential CL using EWC and EWC + CPR. We sample 30% of the test data for producing the visualization. Details and parameters for UMAP are provided in the SM.

We first visualize $O_{t,1}$, defined as the output feature map of the first output layer given the first task’s test data after training the t -th task. The first row of Fig. 6 displays the respective embeddings, where c_t corresponds to the center point of the cluster for the t -th task. In the ideal case (in terms of stability), there would be little to no change in $O_{t,1}$ during CL. This is evident in the embeddings for the joint model, which show that each cluster $O_{t,1}$ is almost perfectly centered. In contrast, the resulting embedding from EWC has a slightly scattered c_t when compared to the joint (oracle) model. This indicates that, whenever the model is trained on a new task, feature maps of the output layer may drift despite EWC’s regularization for previous task parameters. EWC + CPR, in turn, display more centered c_t than EWC, indicating that by applying CPR to EWC model parameters become more robust to change after training future tasks.

In order to provide further evidence that CPR provides better plasticity on new tasks, we visualized h_t , defined as the embedding for the feature map of the last hidden lay-

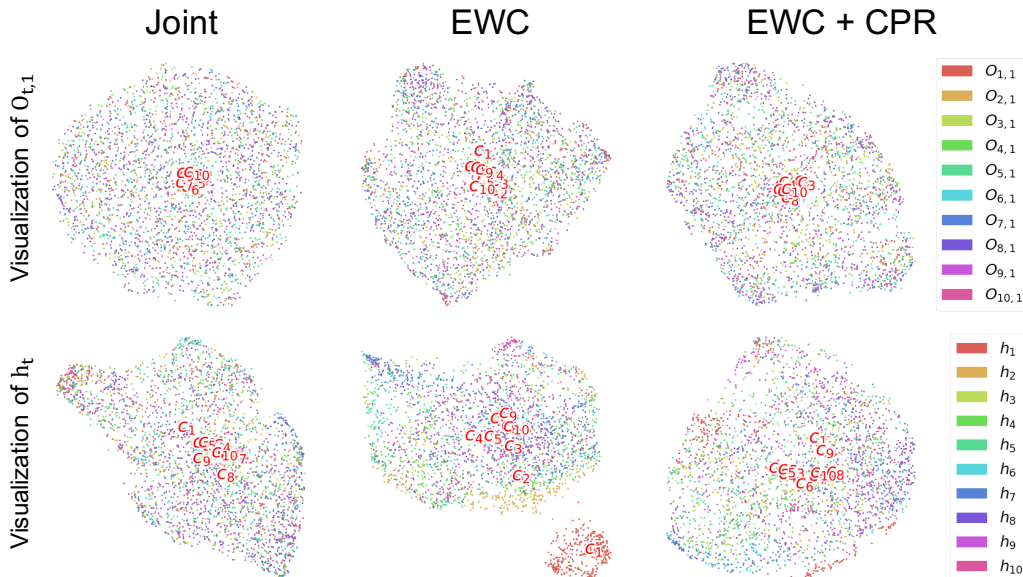


Figure 6: Feature map visualization using UMAP

ers given t -th test data after training the t -th task. In the second row of Fig. 6, Joint and EWC + CPR show closer feature embeddings. EWC, in turn, has a first and second task feature maps divided from other tasks. Strikingly, the feature embeddings for the first task are completely separated. Therefore, we believe that CPR helps the model share feature representations from the start of training, potentially explaining the improvement of the intransigence measure observed in Sec 3.4. We are unaware of prior work that makes use of feature embedding to identify reasons for catastrophic forgetting and limited plasticity of CL methods, and hope that such feature map visualizations become a useful tool for the field. Additional visualizations on different random initializations, different task sequences and MAS (Aljundi et al., 2018a) are reported in the SM.

4. Conclusion

We proposed a simple classifier-projection regularization (CPR) which can be combined with *any* regularization-based continual learning (CL) method. Through extensive experiments, we demonstrated that, by converging to a wide local minima at each task, CPR can significantly increase the plasticity and stability of CL. These encouraging results indicate that wide local minima-promoting regularizers have a critical role in successful CL. Moreover, we observed the impact of CPR through feature map visualizations—a practice that we hope will become more common in future analysis of CL methods. As a theoretical interpretation, we argue that the additional term found in CPR can be understood as a projection of the conditional probability given by a classifier’s output onto a ball centered around the uniform distribution.

References

- Ahn, H., Cha, S., Lee, D., and Moon, T. Uncertainty-based continual learning with adaptive regularization. In *Advances in Neural Information Processing Systems*, pp. 4394–4404, 2019.
- Aljundi, R., Babiloni, F., Elhoseiny, M., Rohrbach, M., and Tuytelaars, T. Memory aware synapses: Learning what (not) to forget. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 139–154, 2018a.
- Aljundi, R., Rohrbach, M., and Tuytelaars, T. Selfless sequential learning. *arXiv preprint arXiv:1806.05421*, 2018b.
- Amari, S.-i., Ikeda, S., and Shimokawa, H. Information geometry of-projection in mean field approximation. *Advanced Mean Field Methods*, pp. 241–258, 2001.
- Carpenter, G. A. and Grossberg, S. Art 2: Self-organization of stable category recognition codes for analog input patterns. *Applied Optics*, 26(23):4919–4930, 1987.
- Chaudhari, P., Choromanska, A., Soatto, S., LeCun, Y., Baldassi, C., Borgs, C., Chayes, J., Sagun, L., and Zecchina, R. Entropy-sgd: Biasing gradient descent into wide valleys. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):124018, 2019.
- Chaudhry, A., Dokania, P. K., Ajanthan, T., and Torr, P. H. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 532–547, 2018.
- Cover, T. M. and Thomas, J. A. *Elements of information theory*. John Wiley & Sons, 2012.
- Csiszár, I. Sanov property, generalized i-projection and a conditional limit theorem. *The Annals of Probability*, pp. 768–793, 1984.
- Csiszár, I. and Matus, F. Information projections revisited. *IEEE Transactions on Information Theory*, 49(6):1474–1490, 2003.
- Csiszár, I. and Shields, P. C. *Information theory and statistics: A tutorial*. Now Publishers Inc, 2004.
- Golkar, S., Kagan, M., and Cho, K. Continual learning via neural pruning. *arXiv preprint arXiv:1903.04476*, 2019.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Kemker, R. and Kanan, C. Fearnnet: Brain-inspired model for incremental learning. *arXiv preprint arXiv:1711.10563*, 2017.
- Keskar, N. S., Mudigere, D., Nocedal, J., Smelyanskiy, M., and Tang, P. T. P. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.
- Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., Hassabis, D., Clopath, C., Kumaran, D., and Hadsell, R. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017. ISSN 0027-8424. doi: 10.1073/pnas.1611835114. URL <https://www.pnas.org/content/114/13/3521>.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Lake, B. M., Salakhutdinov, R., and Tenenbaum, J. B. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Li, Z. and Hoiem, D. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(12):2935–2947, 2017.
- Lopez-Paz, D. and Ranzato, M. A. Gradient episodic memory for continual learning. In *Advances in Neural Information Processing System (NIPS)*, pp. 6467–6476. 2017.
- McCloskey, M. and Cohen, N. J. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pp. 109–165. Elsevier, 1989.
- McInnes, L., Healy, J., and Melville, J. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- Mermillod, M., Bugaiska, A., and Bonin, P. The stability-plasticity dilemma: Investigating the continuum from catastrophic forgetting to age-limited learning effects. *Frontiers in Psychology*, 4:504, 2013.
- Murphy, K. P. *Machine learning: a probabilistic perspective*. MIT press, 2012.
- Nguyen, C. V., Li, Y., Bui, T. D., and Turner, R. E. Variational continual learning. In *International Conference on Learning Representations (ICLR)*, 2018. URL <https://openreview.net/forum?id=BkQqq0gRb>.

- Parisi, G. I., Kemker, R., Part, J. L., Kanan, C., and Wermter, S. Continual lifelong learning with neural networks: A review. *CoRR*, abs/1802.07569, 2018. URL <http://arxiv.org/abs/1802.07569>.
- Pereyra, G., Tucker, G., Chorowski, J., Kaiser, Ł., and Hinton, G. Regularizing neural networks by penalizing confident output distributions. *arXiv preprint arXiv:1701.06548*, 2017.
- Rebuffi, S.-A., Kolesnikov, A., Sperl, G., and Lampert, C. H. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 2001–2010, 2017.
- Rusu, A. A., Rabinowitz, N. C., Desjardins, G., Soyer, H., Kirkpatrick, J., Kavukcuoglu, K., Pascanu, R., and Hadsell, R. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- Shin, H., Lee, J. K., Kim, J., and Kim, J. Continual learning with deep generative replay. In *Advances in Neural Information Processing System (NIPS)*, pp. 2990–2999, 2017.
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2818–2826, 2016.
- Walsh, J. M. and Regalia, P. A. Belief propagation, dykstra’s algorithm, and iterated information projections. *IEEE Transactions on Information Theory*, 56(8):4114–4128, 2010.
- Welinder, P., Branson, S., Mita, T., Wah, C., Schroff, F., Belongie, S., and Perona, P. Caltech-ucsd birds 200. 2010.
- Yoon, J., Yang, E., Lee, J., and Hwang, S. J. Lifelong learning with dynamically expandable networks. In *International Conference on Learning Representations (ICLR)*, 2018. URL <https://openreview.net/forum?id=Sk7KsfW0->.
- Zenke, F., Poole, B., and Ganguli, S. Continual learning through synaptic intelligence. In *International Conference on Machine Learning (ICML)*, pp. 3987–3995, 2017.
- Zhang, Y., Xiang, T., Hospedales, T. M., and Lu, H. Deep mutual learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4320–4328, 2018.

CPR: Classifier-Projection Regularization for Continual Learning

Sungmin Cha¹ Hsiang Hsu² Flavio P. Calmon² Taesup Moon¹

In this supplementary material, we give proofs of the lemma and proposition omitted from Sections 2 , and also provide further details about experiment setups in Section 3.1 , additional experiments on wide local minimum as well as Deep Mutual Learning (Zhang et al., 2018) and MAS (Aljundi et al., 2018) in Section 3.2 . We also report the best regularization strength λ and β in the proposed CPR, and additional experiments to compare with the state of the art on different task arrangements in CL in Section 3.3 . Finally, we provide the hyperparameter settings and additional visualization results for UMAP in Section 3.5 .

1. Mathematical Proofs

1.1. Lemma 1 [Cover & Thomas, 2012, Theorem 11.6.1]

If $D_{KL}(Q\|P)$ is unbounded, then the inequality holds. Assume that $D_{KL}(Q\|P)$ is bounded, then it implies $D_{KL}(Q^*\|P) = \min_{Q \in \mathcal{Q}} D_{KL}(Q\|P)$ is also bounded. Since \mathcal{Q} is a convex set, we consider a convex combination Q^θ of Q^* and Q , i.e., $Q^\theta = (1 - \theta)Q^* + \theta Q \in \mathcal{Q}$, where $\theta \in [0, 1]$. Since Q^* is the minimizer of $D_{KL}(Q\|P)$, we have

$$0 \leq \left. \frac{\partial}{\partial \theta} D_{KL}(Q^\theta\|P) \right|_{\theta=0} \tag{S.1}$$

$$= \left. \frac{\partial}{\partial \theta} D_{KL}((1 - \theta)Q^* + \theta Q\|P) \right|_{\theta=0} \tag{S.2}$$

$$= \left. \frac{\partial}{\partial \theta} \int ((1 - \theta)Q^* + \theta Q) \log \frac{(1 - \theta)Q^* + \theta Q}{P} \right|_{\theta=0} \tag{S.3}$$

$$= \left. \int \frac{\partial}{\partial \theta} \left[((1 - \theta)Q^* + \theta Q) \log \frac{(1 - \theta)Q^* + \theta Q}{P} \right] \right|_{\theta=0} \tag{S.4}$$

$$= \left. \int \left[(-Q^* + Q) \log \frac{(1 - \theta)Q^* + \theta Q}{P} + ((1 - \theta)Q^* + \theta Q) \frac{P}{((1 - \theta)Q^* + \theta Q)} \left(\frac{-Q^* + Q}{P} \right) \right] \right|_{\theta=0} \tag{S.5}$$

$$= \int (-Q^* + Q) \log \frac{Q^*}{P} - Q^* + Q \tag{S.6}$$

$$= \int Q \log \frac{Q^*}{P} - Q^* \log \frac{Q^*}{P} \tag{S.7}$$

$$= \int Q \log \frac{Q}{P} - Q \log \frac{Q^*}{Q} - Q^* \log \frac{Q^*}{P} \tag{S.8}$$

$$= D_{KL}(Q\|P) - D(Q\|Q^*) - D(Q^*\|P), \tag{S.9}$$

where the facts that the exchange of derivatives and integrals is guaranteed by the dominated convergence theorem and that the integrals $\int Q^* = \int Q = 1$. Therefore, we have $D_{KL}(Q\|P) \geq D(Q\|Q^*) + D(Q^*\|P)$, the desired result.

¹Department of Electrical and Computer Engineering, Sungkyunkwan University, Suwon, South Korea ²John A. Paulson School of Engineering and Applied Sciences, Harvard University, MA, USA. Correspondence to: Taesup Moon <tsmoon@skku.edu>.

1.2. Proposition 1

Note that $\mathcal{C}(P_U, \epsilon)$ is a convex set by definition since the KL divergence is convex, and hence Lemma 1 applies. By Lemma 1 and the information inequality (i.e., the KL divergence is always non-negative),

$$D_{\text{KL}}(P_{Y|X}^{t-1*} \| P_{Y|X}^t | P_X^{t-1}) \geq D_{\text{KL}}(P_{Y|X}^{t-1*} \| P_{Y|X}^{t*} | P_X^{t-1}), \forall \mathbf{x}_n^1. \quad (\text{S.10})$$

Therefore, we have

$$- \mathbb{E}_{P_{Y|X}^{t-1*} P_X^{t-1}} \left[\log P_{Y|X}^t P_X^{t-1} \right] \quad (\text{S.11})$$

$$= \int P_{Y|X}^{t-1*} P_X^{t-1} \log \frac{1}{\log P_{Y|X}^t P_X^{t-1}} \quad (\text{S.12})$$

$$= \int P_{Y|X}^{t-1*} P_X^{t-1} \log \frac{P_{Y|X}^{t-1*} P_X^{t-1}}{\log P_{Y|X}^t P_X^{t-1}} - P_{Y|X}^{t-1*} P_X^{t-1} \log P_{Y|X}^{t-1*} P_X^{t-1} \quad (\text{S.13})$$

$$= D_{\text{KL}}(P_{Y|X}^{t-1*} \| P_{Y|X}^t | P_X^{t-1}) - P_{Y|X}^{t-1*} P_X^{t-1} \log P_{Y|X}^{t-1*} P_X^{t-1} \quad (\text{S.14})$$

$$\geq D_{\text{KL}}(P_{Y|X}^{t-1*} \| P_{Y|X}^{t*} | P_X^{t-1}) + -P_{Y|X}^{t-1*} P_X^{t-1} \log P_{Y|X}^{t-1*} P_X^{t-1} \quad (\text{S.15})$$

$$= - \mathbb{E}_{P_{Y|X}^{t-1*} P_X^{t-1}} \left[\log P_{Y|X}^{t*} P_X^{t-1} \right], \quad (\text{S.16})$$

where the inequality comes from (S.10).

2. Experimental details of Section 3.1

For training models on CIFAR100, CIFAR10/100 and Omniglot, we used the Adam (Kingma & Ba, 2015) optimizer with initial learning rate 0.001 for 100 epochs. For training CUB200, we set the initial learning rate as 0.0005 and trained the model for 50 epochs. Here we also used the learning rate scheduler which drops the learning rate by half when validation error is not decreased. All experiments was implemented in PyTorch 1.2.0 with CUDA 9.2 on NVIDIA 1080Ti GPU.

Following (Ahn et al., 2019), we use a simple CNN model for training CL benchmark dataset except for CUB200 and details of an architecture is in Table 1 and 2.

Table 1: Network architecture for Split CIFAR-10/100 and Split CIFAR-100

| Layer | Channel | Kernel | Stride | Padding | Dropout |
|--------------------------|---------|--------|--------|---------|---------|
| 32×32 input | 3 | | | | |
| Conv 1 | 32 | 3×3 | 1 | 1 | |
| Conv 2 | 32 | 3×3 | 1 | 1 | |
| MaxPool | | | 2 | 0 | 0.25 |
| Conv 3 | 64 | 3×3 | 1 | 1 | |
| Conv 4 | 64 | 3×3 | 1 | 1 | |
| MaxPool | | | 2 | 0 | 0.25 |
| Conv 5 | 128 | 3×3 | 1 | 1 | |
| Conv 6 | 128 | 3×3 | 1 | 1 | |
| MaxPool | | | 2 | 1 | 0.25 |
| Dense 1 | 256 | | | | |
| Task 1 : Dense 10 | | | | | |
| ... | | | | | |
| Task <i>i</i> : Dense 10 | | | | | |

Table 2: Network architecture for Omniglot

| Layer | Channel | Kernel | Stride | Padding | Dropout |
|------------------------|---------|--------|--------|---------|---------|
| 28×28 input | 1 | | | | |
| Conv 1 | 64 | 3×3 | 1 | 0 | |
| Conv 2 | 64 | 3×3 | 1 | 0 | |
| MaxPool | | | 2 | 0 | 0 |
| Conv 3 | 64 | 3×3 | 1 | 0 | |
| Conv 4 | 64 | 3×3 | 1 | 0 | |
| MaxPool | | | 2 | 0 | 0 |
| Task 1 : Dense C_1 | | | | | |
| ... | | | | | |
| Task i : Dense C_i | | | | | |

3. Additional Experimental Results of Section 3.2

3.1. Experimental Results of Wide Local Minima using Training Data

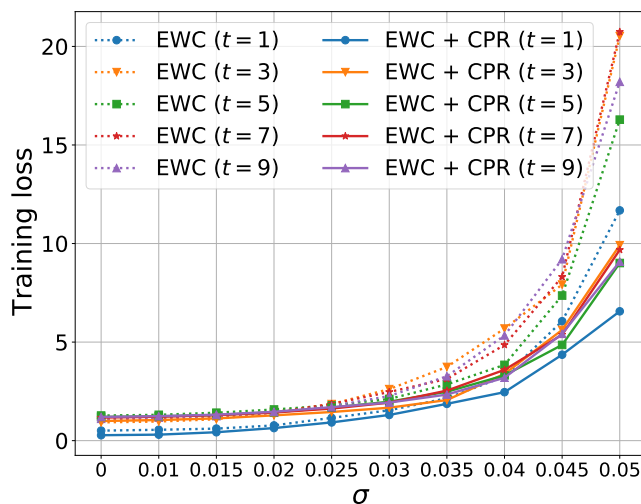


Figure 1: Experimental result of adding Gaussian noise to training data

Figure 1 shows the experimental result of Section 3.2 using training data. We clearly see that training loss of EWC + CPR slowly increases than EWC in all tasks.

3.2. Experimental Results on MAS (Aljundi et al., 2018) and Deep Mutual Learning (Zhang et al., 2018)

We did the same experiments of Section 3.2 using MAS (Aljundi et al., 2018), and Figure 2 shows the results. In Figure 2(a), we observe that MAS shows a clear trade-off between F_{10} and $I_{1,10}$ as β increases, unlike the result of EWC in the manuscript. (We note SI (Zenke et al., 2017) and RWalk (Chaudhry et al., 2018) showed similar trend as EWC (Kirkpatrick et al., 2017) in the manuscript). MAS + CPR achieves the highest accuracy in the range of $0.5 \leq \beta \leq 0.9$ but we can see that $\beta = \{0.7, 0.9\}$ shows a worse F_{10} compared with MAS. Therefore, we can select $\beta = 0.5$ as the best hyperparameter using the criteria for selecting β proposed in Section 3.2 of the manuscript.

We also experimented Deep Mutual Learning (DML) (Zhang et al., 2018) as the regularization for converging wide local minima. We used $\beta = 1$ only because DML reports the best result (with $\beta = 1$) which is converging to a better wide local minima compared to Entropy Maximization (Pereyra et al., 2017). In our experiment, DML shows an increased A_{10} and decreased $F_{10}, I_{1,10}$ but it is not as effective as our CPR. Most decisively, DML requires training at least more than two models so we excluded DML from our consideration.

Figure 2(b) shows the experimental result on adding Gaussian noise to the parameters which is trained on CIFAR-100. We clearly observe that *test* loss of each task more slowly increases by applying CPR to MAS. We believe this is another evidence that CPR can be generally applied to regularization-based CL methods, promoting the wide-local minima.

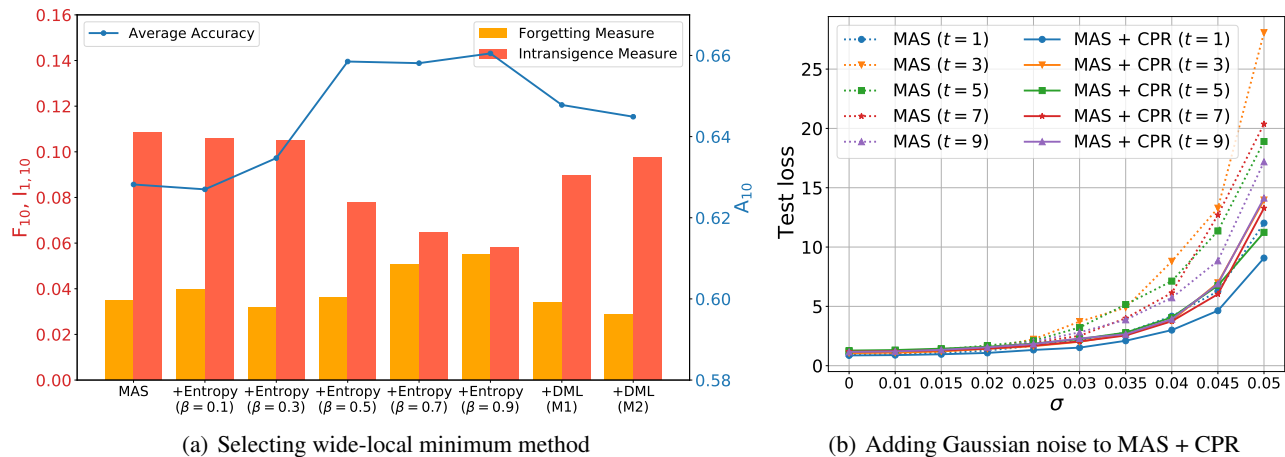


Figure 2: Experiments for selecting the regularization on CIFAR100

4. Selected Best Hyperparameters

Table 3: Best hyperparameters for each regularization-based CL method and CPR

| Best λ / Best β | CIFAR100 | CIFAR10/100 | CIFAR50/10/50 | CIFAR100/10 | Omniglot | CUB200 |
|-------------------------------|--------------|--------------|---------------|--------------|---------------|---------------|
| EWC | 12,000 / 0.5 | 25,000 / 0.4 | 12,000 / 0.8 | 20,000 / 0.6 | 100,000 / 1.0 | 300,000 / 0.4 |
| SI | 1 / 0.8 | 0.9 / 0.2 | 2 / 0.9 | 2 / 0.5 | 8 / 0.7 | 50 / 0.6 |
| MAS | 3 / 0.5 | 1 / 0.2 | 2 / 0.1 | 2 / 0.4 | 10 / 0.6 | 50 / 0.6 |
| RWalk | 8 / 0.9 | 4 / 0.4 | 10 / 0.6 | 10 / 0.8 | 3,000 / 0.6 | 300 / 0.9 |

For each dataset, we firstly searched best λ for each regularization-based CL method and then we selected best β for CPR. All best hyperparameters are proposed in Figure 3.

5. Experimental Results on CIFAR100/10, CIFAR50/10/50

As an additional experiments of Section 3.3 in the manuscript, we experimented on CIFAR100/10 and CIFAR50/10/50, which are the different versions of CIFAR10/100. Namely, we changed the order of the tasks and varied the location for which CIFAR-10 task is inserted. Table 4 and Figure 5 show the results. We can achieve better relative improvements on all metrics compared to CIFAR-10/100.

Table 4: Experimental results on continual learning scenarios with and without CPR. Blue color denotes the case which CL method is positively affected by CPR and red color represents a negative impact of CPR.

| Dataset | Method | Average Accuracy (A_{10}) | | | Forgetting Measure (F_{10}) | | | Intransigence Measure ($I_{1,10}$) | | |
|-------------------------------|--------|-------------------------------|--------|-------------------------|---------------------------------|--------|-------------------------|--------------------------------------|--------|-------------------------|
| | | W/o CPR | W/ CPR | diff (W-W/o) | W/o CPR | W/ CPR | diff (W-W/o) | W/o CPR | W/ CPR | diff (W-W/o) |
| CIFAR50/10/50 ($T = 11$) | EWC | 0.5978 | 0.6346 | +0.0368 (+6.2%) | 0.0288 | 0.0292 | +0.0004 (+1.4%) | 0.1682 | 0.1311 | -0.0371 (-22.1%) |
| | SI | 0.6184 | 0.6468 | +0.0284 (+4.6%) | 0.0598 | 0.0532 | -0.0066 (-11.0%) | 0.1194 | 0.0970 | -0.0224 (-18.8%) |
| | MAS | 0.6172 | 0.6238 | +0.0066 (+1.1%) | 0.0484 | 0.0448 | -0.0036 (-7.4%) | 0.1310 | 0.1277 | -0.0033 (-2.5%) |
| | Rwalk | 0.5697 | 0.6315 | +0.0619 (+10.9%) | 0.0781 | 0.0548 | -0.0233 (-29.8%) | 0.1515 | 0.1109 | -0.0406 (-26.8%) |
| CIFAR100/10 ($T = 11$) | EWC | 0.5808 | 0.6158 | +0.0376 (+6.5%) | 0.0304 | 0.0238 | -0.0066 (-21.7%) | 0.1694 | 0.1378 | -0.0317 (-18.7%) |
| | SI | 0.6116 | 0.6332 | +0.0216 (+3.5%) | 0.0681 | 0.0692 | -0.0011 (-1.6%) | 0.1044 | 0.0832 | -0.0212 (-20.3%) |
| | MAS | 0.6138 | 0.6363 | +0.0214 (+3.5%) | 0.0536 | 0.0532 | -0.0004 (-0.7%) | 0.1153 | 0.0942 | -0.0211 (-18.3%) |
| | Walk | 0.5618 | 0.6113 | +0.0495 (+8.8%) | 0.0924 | 0.0852 | -0.0072 (-7.8%) | 0.1322 | 0.0892 | -0.0430 (-32.5%) |

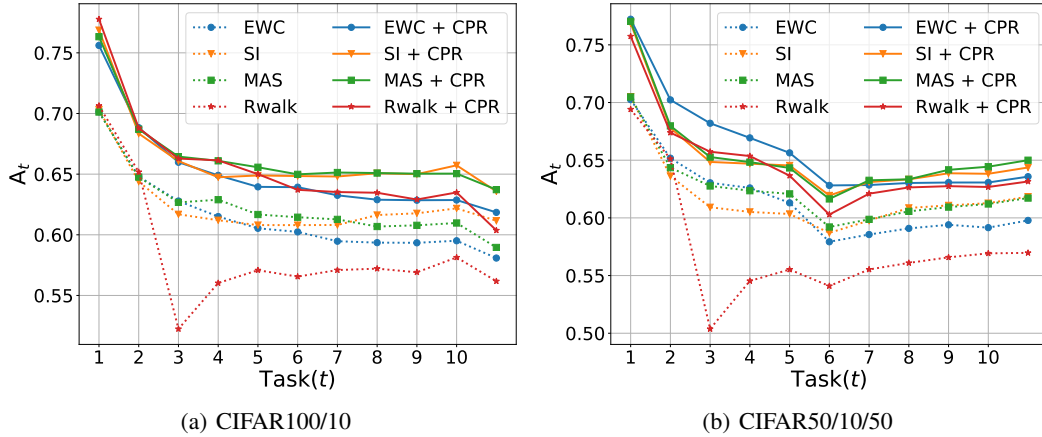


Figure 3: Average accuracy for CIFAR10/100 and CIFAR50/10/50

6. Hyperparameter Settings and Visualization Details of UMAP

From several visualizations, we found out that best hyperparameters for UMAP (McInnes et al., 2018) as $\{n_neighbors = 200, min_dist = 0.1, n_components = 2\}$ and we got all visualization results with these hyperparameters. We used raw features of $O_{t,1}$ as a input of UMAP, however, for visualizing h_t , we reduced the dimension of h_t to 50 by using PCA.

7. Additional Feature Map Visualizations

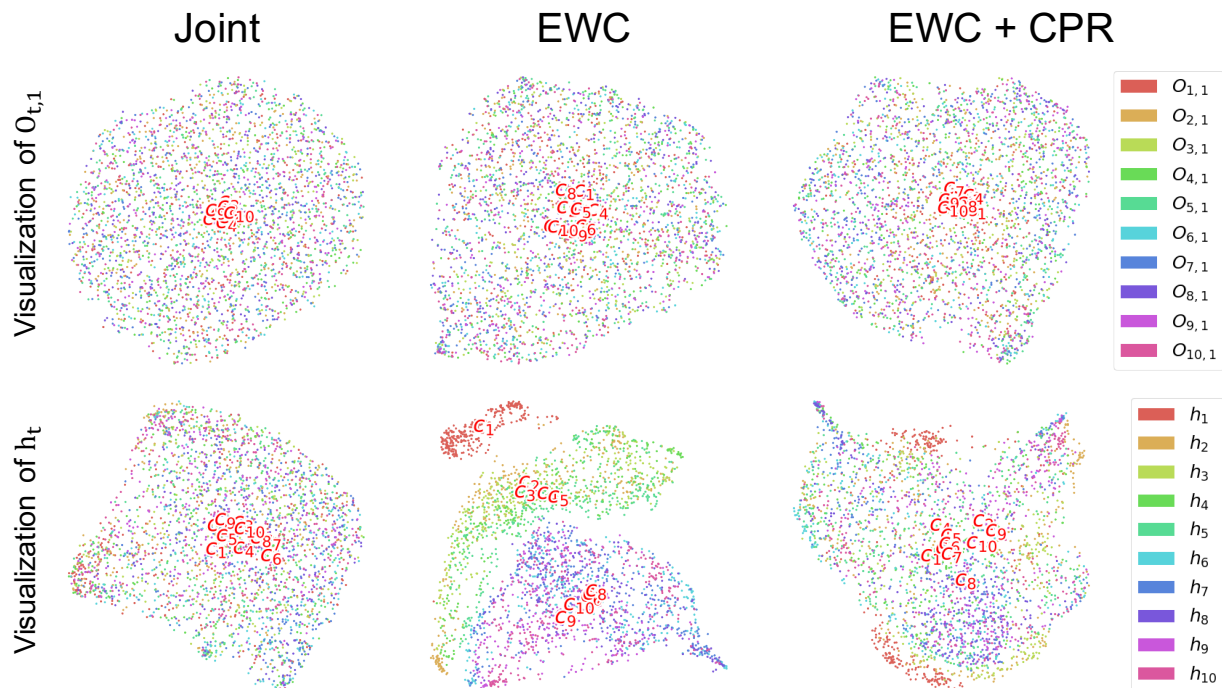
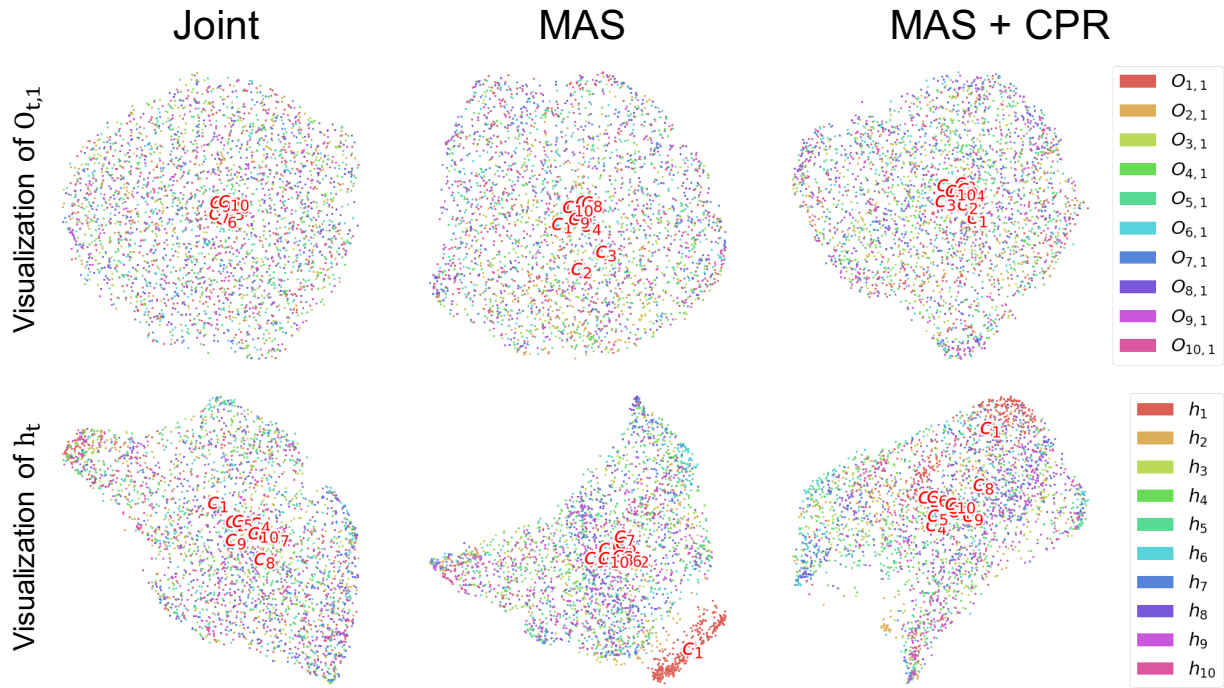


Figure 4: Visualization Result on EWC (seed = 9)

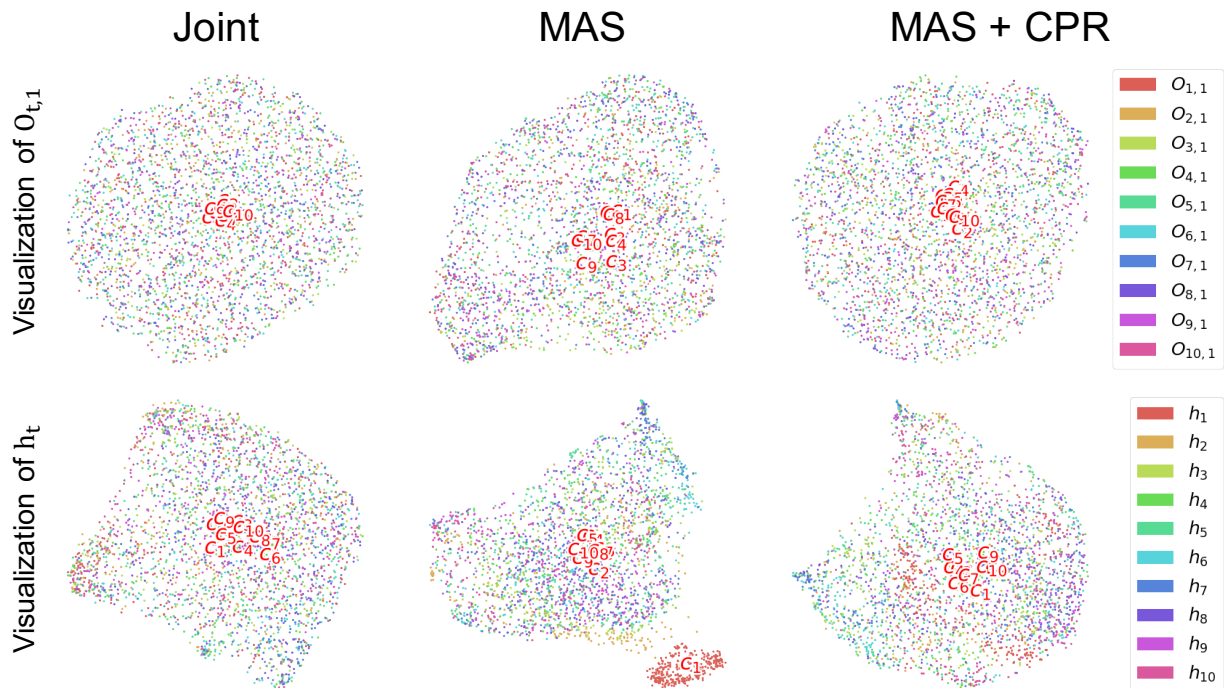
We visualize $O_{t,1}$ and h_t of Joint, EWC (Kirkpatrick et al., 2017), EWC (Kirkpatrick et al., 2017) + CPR with a different seed and visualizations are shown in Figure 7. We hold the experimental settings and we can see the similar pattern of $O_{t,1}$ and h_t , which is already shown in Section 3.5 of the manuscript. Especially, $O_{t,1}$ of EWC showed clearly divided clusters compared with the visualization result in the manuscript, nevertheless, we confirm that the feature maps become to be more

shared and centered by applying CPR to EWC.

We also did same visualization using MAS (Aljundi et al., 2018) and the results are shown in Figure 7. We checked the similar results of $O_{t,1}$ and h_t , and we could see that, by applying CPR to MAS, $O_{t,1}$ and h_t are more centered than before. From these additional visualizations, we want to emphasize that the pattern of $O_{t,1}$ and h_t is a general phenomenon of regularization-based CL methods, and these can show why the typical regularization-based CL methods still suffer from the stability-plasticity dilemma at the feature map level. Also, we could check again that CPR increases the stability and plasticity of the regularizaion-based CL methods by alleviating this phenomenon.



(a) Visualization result on MAS (seed = 0)



(b) Visualization result on MAS (seed = 9)

Figure 5: Feature map visualization of MAS

Table 5: Experimental results on training additional tasks with EWC and EWC + CPR

| | A_{20} | F_{20} | $I_{1,20}$ | $I_{10,20}$ |
|--------------------------|---------------|----------|------------|---------------|
| EWC + CPR (all tasks) | 0.6612 | 0.1229 | 0.1027 | 0.0855 |
| EWC (all tasks) | 0.6195 | 0.1362 | 0.1319 | 0.1156 |
| EWC + CPR (CIFAR-100) | 0.6502 | 0.1486 | 0.0882 | 0.0677 |
| EWC (CIFAR-100) | 0.6143 | 0.1604 | 0.1128 | 0.0870 |

8. Experiments on additional tasks

From Section 3.4 in the manuscript, we demonstrated the critical role of CPR in terms of increasing the plasticity. From this result, we thought that CPR might help to learn additional future tasks well without the hyperparameter search for new whole tasks. To verify our hypothesis, we designed a new task sequence made up of 20 tasks, CIFAR100(10 tasks) + SVHN (Netzer et al., 2011)(5 tasks) + Synthetic MNIST (Roy et al., 2018)(5 tasks) and each task of SVHN and Synthetic MNIST is a binary image classification. Table 5 shows experimental results of EWC (Kirkpatrick et al., 2017) on additional tasks.

We divide the experimental setting as two different cases. The first case is that we newly search the best hyperparameter for all 20 tasks (denoted as all tasks), and in the second case, we just use the best hyperparameter got from CIFAR-100 (denoted as CIFAR-100). EWC + CPR (CIFAR-100) shows a low $I_{10,20}$ compared with EWC (all tasks), as a result, EWC + CPR (CIFAR-100) achieve the higher A_{20} than EWC (all tasks). Also, we observe that, if we find the best hyperparameters (λ, β) for all 20 tasks again, EWC + CPR (all tasks) still achieves the best result in all metrics. In conclusion, we believe that this is a remarkable result, and it shows the effect of wide local minimum in CL continues in additional tasks.

References

- Ahn, H., Cha, S., Lee, D., and Moon, T. Uncertainty-based continual learning with adaptive regularization. In *Advances in Neural Information Processing Systems*, pp. 4394–4404, 2019.
- Aljundi, R., Babiloni, F., Elhoseiny, M., Rohrbach, M., and Tuytelaars, T. Memory aware synapses: Learning what (not) to forget. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 139–154, 2018.
- Chaudhry, A., Dokania, P. K., Ajanthan, T., and Torr, P. H. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 532–547, 2018.
- Cover, T. M. and Thomas, J. A. *Elements of information theory*. John Wiley & Sons, 2012.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.
- Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., Hassabis, D., Clopath, C., Kumaran, D., and Hadsell, R. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017. ISSN 0027-8424. doi: 10.1073/pnas.1611835114. URL <https://www.pnas.org/content/114/13/3521>.
- McInnes, L., Healy, J., and Melville, J. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. Reading digits in natural images with unsupervised feature learning. 2011.
- Pereyra, G., Tucker, G., Chorowski, J., Kaiser, Ł., and Hinton, G. Regularizing neural networks by penalizing confident output distributions. *arXiv preprint arXiv:1701.06548*, 2017.
- Roy, P., Ghosh, S., Bhattacharya, S., and Pal, U. Effects of degradations on deep neural network architectures. *arXiv preprint arXiv:1807.10108*, 2018.
- Zenke, F., Poole, B., and Ganguli, S. Continual learning through synaptic intelligence. In *International Conference on Machine Learning (ICML)*, pp. 3987–3995, 2017.
- Zhang, Y., Xiang, T., Hospedales, T. M., and Lu, H. Deep mutual learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4320–4328, 2018.