

# EWC-detection instantiation

Peng YUN

20200921

We want to find the most probable parameters given the data of old and new tasks  $\mathcal{D} = \mathcal{D}_A \cup \mathcal{D}_B$ , where  $\mathcal{D}_A$  denotes the data of the old task-A, and  $\mathcal{D}_B$  denotes the data of the new task-B. The goal can be achieved by minimizing the negative logarithm of the posterior distribution  $-\log p(\boldsymbol{\theta}|\mathcal{D}_A \cup \mathcal{D}_B)$ . According to the bayesian rule, we can write the posterior term into:

$$-\log p(\boldsymbol{\theta}|\mathcal{D}_A \cup \mathcal{D}_B) = -\log p(\mathcal{D}_B|\boldsymbol{\theta}) - \log p(\boldsymbol{\theta}|\mathcal{D}_A) + \log p(\mathcal{D}_B) \quad (1)$$

To get the optimal parameters  $\boldsymbol{\theta}$ , we solve the following optimization problem:

$$\boldsymbol{\theta}_{AB}^* = \arg \min_{\boldsymbol{\theta}} -\log p(\mathcal{D}_B|\boldsymbol{\theta}) - \log p(\boldsymbol{\theta}|\mathcal{D}_A) \quad (2)$$

It is noted that we get rid of the term  $\log p(\mathcal{D}_B)$  in this optimization problem since the data distribution prior can be considered as a constant.

For the objective function, we can easily compute the likelihood term  $-\log p(\mathcal{D}_B|\boldsymbol{\theta})$  with the task-B dataset at hand, and it is exactly the linear combination of classification and regression loss for object detection, which we denote it as  $\mathcal{L}_{det}(\boldsymbol{\theta}, \mathcal{D}_B)$ .

The term  $-\log p(\boldsymbol{\theta}|\mathcal{D}_A)$  is intractable, since we do not have  $\mathcal{D}_A$  at hand at the time of incrementally training task-B. The mechanism behind the prior-based method EWC is to restore the prior information of old task and then use the priors instead of old-task data to preserve learned knowledge. What we have at hand is  $\boldsymbol{\theta}_A^*$  which minimizes the posterior term  $-\log p(\boldsymbol{\theta}|\mathcal{D}_A)$ , and therefore  $\frac{\partial}{\partial \boldsymbol{\theta}} \log p(\boldsymbol{\theta}|\mathcal{D}_A)|_{\boldsymbol{\theta}_A^*} = 0$ . We unfold the  $\log p(\boldsymbol{\theta}|\mathcal{D}_A)$  at  $\boldsymbol{\theta}_A^*$  with Taylor expansion:

$$\log p(\boldsymbol{\theta}|\mathcal{D}_A) \approx \log p(\boldsymbol{\theta}_A^*|\mathcal{D}_A) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_A^*)^T \frac{\partial^2}{\partial^2 \boldsymbol{\theta}} \log p(\boldsymbol{\theta}|\mathcal{D}_A)|_{\boldsymbol{\theta}_A^*} (\boldsymbol{\theta} - \boldsymbol{\theta}_A^*) \quad (3)$$

We denote the  $\frac{\partial^2}{\partial^2 \boldsymbol{\theta}} \log p(\boldsymbol{\theta}|\mathcal{D}_A)|_{\boldsymbol{\theta}_A^*}$  with  $\mathbf{H}(\mathcal{D}_A, \boldsymbol{\theta}_A^*)$ . It demonstrates the logarithm posterior distribution  $\log p(\boldsymbol{\theta}|\mathcal{D}_A) \sim \mathcal{N}(\log p(\boldsymbol{\theta}_A^*|\mathcal{D}_A), -\mathbf{H}(\mathcal{D}_A, \boldsymbol{\theta}_A^*)^{-1})$  according to Laplacian approximation.

We can compute  $\mathbf{H}(\mathcal{D}_A, \boldsymbol{\theta}_A^*)$  with empirical fisher information matrix (FIM):

$$\begin{aligned} \mathbf{H}(\mathcal{D}_A, \boldsymbol{\theta}_A^*) &= -\mathbb{F}(\mathcal{D}_A, \boldsymbol{\theta}_A^*) \\ \mathbb{F}(\mathcal{D}_A, \boldsymbol{\theta}_A^*) &= \frac{1}{|S|} \sum_{\tilde{\mathcal{D}}_A \sim \mathcal{D}_A} [(\frac{\partial}{\partial \boldsymbol{\theta}} \log p(\boldsymbol{\theta}|\tilde{\mathcal{D}}_A)|_{\boldsymbol{\theta}_A^*})^T (\frac{\partial}{\partial \boldsymbol{\theta}} \log p(\boldsymbol{\theta}|\tilde{\mathcal{D}}_A)|_{\boldsymbol{\theta}_A^*})] \end{aligned} \quad (4)$$

where  $|S|$  denotes the number of times sampling  $\tilde{\mathcal{D}}_A$  from  $\mathcal{D}_A$ . Consider the task-A is the first task in the task sequence, thus there is no prior information related to  $\theta$ :

$$\begin{aligned}\log p(\theta|\mathcal{D}_A) &= \log p(\mathcal{D}_A|\theta) + \log p(\theta) - \log p(\mathcal{D}_A) \\ &= \log p(\mathcal{D}_A|\theta) + \text{const.}\end{aligned}\quad (5)$$

Therefore, we can compute the empirical FIM right after training task-A with the equation:

$$\mathbb{F}(\mathcal{D}_A, \theta_A^*) = \frac{1}{|S|} \sum_{\tilde{\mathcal{D}}_A \sim \mathcal{D}_A} [(\frac{\partial}{\partial \theta} \mathcal{L}_{det}(\tilde{\mathcal{D}}_A, \theta)|_{\theta_A^*})^T (\frac{\partial}{\partial \theta} \mathcal{L}_{det}(\tilde{\mathcal{D}}_A, \theta)|_{\theta_A^*})] \quad (6)$$

As a result, the objective function of optimizing  $\theta$  for incrementally learning task-B is:

$$\theta_{AB}^* = \arg \min_{\theta} \mathcal{L}_{det}(\theta, \mathcal{D}_B) + \lambda(\theta - \theta_A^*)^T \mathbb{F}(\mathcal{D}_A, \theta_A^*)(\theta - \theta_A^*) \quad (7)$$

where the hyperparameter  $\lambda$  balances the weights between the detection loss and the L2-norm prior constraints.

Then we extend to the third task task-C, and derive the empirical FIM propagation. For three task cases, we want to find the optimal parameter  $\theta$  by minimizing the  $-\log p(\theta|\mathcal{D}_A \cup \mathcal{D}_B \cup \mathcal{D}_C)$ . The optimization problem is seen as:

$$\theta_{ABC}^* = \arg \min_{\theta} -\log p(\mathcal{D}_C|\theta) - \log p(\theta|\mathcal{D}_A \cup \mathcal{D}_B) \quad (8)$$

As before, the likelihood term  $-\log p(\mathcal{D}_C|\theta) = \mathcal{L}_{det}(\theta, \mathcal{D}_C)$  and can be computed easily with the task-C dataset at hand. What we also have at hand is  $\theta_{AB}^*$  which minimizes the posterior term of the two-task case  $-\log p(\theta|\mathcal{D}_A \cup \mathcal{D}_B)$ , thus  $\frac{\partial}{\partial \theta} \log p(\theta|\mathcal{D}_A \cup \mathcal{D}_B) = 0$ . Please recap the equation 1. We can unfold its RHS at  $\theta_{AB}^*$  and get:

$$\begin{aligned}& \log p(\theta|\mathcal{D}_A \cup \mathcal{D}_B) \\ &= \log p(\mathcal{D}_B|\theta_{AB}^*) + \log p(\theta_{AB}^*|\mathcal{D}_A) \\ &= \frac{1}{2}(\theta - \theta_{AB}^*)^T \frac{\partial^2}{\partial^2 \theta} [\log p(\mathcal{D}_B|\theta) + \log p(\theta|\mathcal{D}_A)]_{\theta_{AB}^*} (\theta - \theta_{AB}^*) + \text{const.} \\ &= \frac{1}{2}(\theta - \theta_{AB}^*)^T [\frac{\partial^2}{\partial^2 \theta} \log p(\mathcal{D}_B|\theta)|_{\theta_{AB}^*} + \frac{\partial^2}{\partial^2 \theta} \log p(\theta|\mathcal{D}_A)|_{\theta_{AB}^*}] (\theta - \theta_{AB}^*) + \text{const.}\end{aligned}\quad (9)$$

For  $\frac{\partial^2}{\partial^2 \theta} \log p(\mathcal{D}_B|\theta)|_{\theta_{AB}^*}$ , we can approximate it with the negative empirical FIM as before:

$$\mathbb{F}(\mathcal{D}_B, \theta_{AB}^*) = \frac{1}{|S|} \sum_{\tilde{\mathcal{D}}_B \sim \mathcal{D}_B} [(\frac{\partial}{\partial \theta} \mathcal{L}_{det}(\tilde{\mathcal{D}}_B, \theta)|_{\theta_{AB}^*})^T (\frac{\partial}{\partial \theta} \mathcal{L}_{det}(\tilde{\mathcal{D}}_B, \theta)|_{\theta_{AB}^*})] \quad (10)$$

For  $\frac{\partial^2}{\partial^2 \boldsymbol{\theta}} \log p(\boldsymbol{\theta}|\mathcal{D}_A)|_{\boldsymbol{\theta}_{AB}^*}$ , we have already derived the following equation (equation 3, 6):

$$\log p(\boldsymbol{\theta}|\mathcal{D}_A) \approx \log p(\boldsymbol{\theta}_A^*|\mathcal{D}_A) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_A^*)^T \mathbb{F}(\mathcal{D}_A, \boldsymbol{\theta}_A^*)(\boldsymbol{\theta} - \boldsymbol{\theta}_A^*) \quad (11)$$

Thus, we can get:

$$\frac{\partial^2}{\partial^2 \boldsymbol{\theta}} \log p(\boldsymbol{\theta}|\mathcal{D}_A)|_{\boldsymbol{\theta}_{AB}^*} = \mathbb{F}(\mathcal{D}_A, \boldsymbol{\theta}_A^*) \quad (12)$$

As a result, the objective function for incrementally learning task-C is:

$$\boldsymbol{\theta}_{ABC}^* = \arg \min_{\boldsymbol{\theta}} \mathcal{L}_{det}(\boldsymbol{\theta}, \mathcal{D}_C) + \lambda(\boldsymbol{\theta} - \boldsymbol{\theta}_{AB}^*)^T [\mathbb{F}(\mathcal{D}_A, \boldsymbol{\theta}_A^*) + \mathbb{F}(\mathcal{D}_B, \boldsymbol{\theta}_{AB}^*)](\boldsymbol{\theta} - \boldsymbol{\theta}_{AB}^*) \quad (13)$$

By recursively applying the above equation, the general objective function for incrementally learning task- $\mathcal{T}$  is:

$$\boldsymbol{\theta}_{A..\mathcal{T}}^* = \arg \min_{\boldsymbol{\theta}} \mathcal{L}_{det}(\boldsymbol{\theta}, \mathcal{D}_{\mathcal{T}}) + \lambda(\boldsymbol{\theta} - \boldsymbol{\theta}_{AB}^*)^T \left[ \sum_t \mathbb{F}(\mathcal{D}_t, \boldsymbol{\theta}_{A..t}^*) \right] (\boldsymbol{\theta} - \boldsymbol{\theta}_{AB}^*) \quad (14)$$