

# AUXILIARY MODELS AND DEPENDENT LEARNING

HYEONGJUN JIN

ABSTRACT. We define what is an auxiliary model for a pre-trained model and present a schema for training such a model in an unsupervised, but dependent manner.

## CONTENTS

### 1. AUXILLIARY MODELS

Assume  $\tilde{\theta}$  is a pre-trained parameter for a parametrized model  $\mathcal{M}$  with respect to an admissible loss function  $\lambda$ .

**Definition 1.1.** Consider a pair  $(\mathcal{N}, g)$  where  $\mathcal{N} \in ((Y \models_{\Phi} Z))$  is a model parametrized by  $\Phi$  and  $g : Z \rightarrow \mathbb{R}$  is a function bounded below such that it is smooth on a open dense subset of  $Z$ . It is called an **auxiliary model for**  $(\mathcal{M}, S, \lambda, \tilde{\theta}, V_{\tilde{\theta}})$  if the loss function

$$\begin{aligned} \nu_g : ((X \models Z)) &\rightarrow \mathbb{R} \\ \mathcal{L} &\mapsto \nu_g(\mathcal{L}) = \sum_{(x,y) \in S} g(\mathcal{L}(x)) \end{aligned}$$

is admissible for the composite model  $\mathcal{N}_{(-)} \circ \mathcal{M}_{(-)} : V_{\tilde{\theta}} \times \Phi \rightarrow ((X \models_{V_{\tilde{\theta}} \times \Phi} Z))$ .

Note that the loss  $\nu_g$  makes the training set  $S$  unsupervised.

### 2. DEPENDENT LEARNING

The learning scheme of an auxiliary model is as follows:

- Do gradient descent to find a minimizer  $\tilde{\phi}$  of the function below while keeping the argument  $\tilde{\theta}$  fixed:

$$\begin{aligned} \mu &= \nu_g \circ \mathcal{N}_{(-)} \circ \mathcal{M}_{(-)} : V_{\tilde{\theta}} \times \Phi \rightarrow \mathbb{R} \\ (\theta, \phi) &\mapsto \sum_{(x,y) \in S} g(\mathcal{N}_{\phi}(\mathcal{M}_{\theta}(x))) \\ (\tilde{\theta}, \phi_{old}) &\mapsto (\tilde{\theta}, \phi_{new}). \end{aligned} \tag{1}$$

Here we update the variable  $\phi$  only by considering the projected gradient  $\nabla \mu(-|\tilde{\theta})$ :

$$\begin{aligned} d\mu_{(\theta, \phi)} : T_{\theta} V_{\tilde{\theta}} \times T_{\phi} \Phi &\rightarrow \mathbb{R} \rightsquigarrow \nabla \mu : V_{\tilde{\theta}} \times \Phi \rightarrow TV_{\tilde{\theta}} \times T\Phi \\ \Rightarrow \nabla \mu(-|\tilde{\theta}) &= \text{proj}_{T\Phi} \nabla \mu(\tilde{\theta}, -) : \Phi \xrightarrow{(\tilde{\theta}, -)} V_{\tilde{\theta}} \times \Phi \xrightarrow{\nabla \mu} TV_{\tilde{\theta}} \times T\Phi \rightarrow T\Phi \end{aligned} \tag{2}$$

As is  $\nabla \mu$ , this gradient is carrying information of the direction of steepest change for  $\mu(\tilde{\theta}, \phi)$  for a given  $\phi$  according to a change  $(\tilde{\theta}, \phi) \mapsto (\tilde{\theta} + \Delta\theta, \phi + \Delta\phi)$ .

(Of course we assume  $\Theta$  and  $\Phi$  to have Riemannian metrics so that  $d\mu$  can be transformed into the gradient vector field  $\nabla \mu$  on  $V_{\tilde{\theta}} \times \Phi$ .)

Note the dependence of the objective function  $\mu$  in the direction of  $\theta$ , as the term  $\mathcal{M}_{\tilde{\theta}}$  can vary according to an infinitesimal change on  $\tilde{\theta}$ , making the update of  $\phi$  to depend on the pre-trained model's current state, whence the name **dependent learning**.

- After the dependent learning of  $\mathcal{N}$  with respect to  $\nu_g$ , we make the composite model learn simultaneously with respect to some loss function  $h(\lambda, \nu_g)$ , making the model parameters of  $\mathcal{M}$  and  $\mathcal{N}$  interdependent. We denote simply by  $\mathcal{M} \bowtie \mathcal{N}$  a model incorporated in this way.

(Simplest such  $h$  is  $h(x, y) = x + y$ ,  $h(x, y) = x^y$ , etc.)

GitHub of example models(updating): [Application to Several Models](#)

(In the GitHub, training of the evaluator is solely done with the reproduced pictures of the VAE. It's remarkable to note that the evaluator seems to value the original data as better photo even in this bias.)

### 3. APPLICATIONS

#### 1. Auto-evaluator for any generative model

1-1. Data selection    1-2. Data augmentation    1-3. Performance improvement

Our auto-evaluator in its simplest form is as follows: We set  $Z = \mathbb{R}^k$  for some  $k$  and use the exponentially decreasing function on positive end as our loss function:

$$g : \mathbb{R}^k \rightarrow \mathbb{R}, \quad g(x_1, \dots, x_k) = \sum_{1 \leq i \leq k} \exp(-ReLU(x_i)).$$

The gradient descent with respect to this function make any parametrized model  $\mathcal{N}_{(-)} \in ((Y \models_{\Phi} Z))$  to diverge at the end when we employ the honest gradient descent as our learning technique. As soon as we employ the dependent learning, however, the parameter  $\phi$  stays in a compact region in  $\Phi$ , meaning its convergence which is the practical part for use.

As our auto-evaluator converges to a trained parameter  $\tilde{\phi}$ , we let  $\mathcal{N}_{\tilde{\phi}}(y) \in Z = \mathbb{R}^k$  the value vector or the values. We can use either the sum of its entries or the sum of the log of its entries as the value of  $y$ .

#### 2. Hyperparameter adjustor

#### 3. Weight scheduler

#### 4. Everything you want

### 4. APPENDIX. MATHEMATICAL DEFINITIONS

We begin by defining a model in mathematical terms.

**Definition 4.1.** *Let  $X$  and  $Y$  be smooth manifolds and  $C(X, Y)$  be the set of all continuous maps from  $X$  to  $Y$ . Then  $C(X, Y)$  is a topological space with the compact-open topology.*

(1) *By a **model  $\mathcal{M}$  for  $Y$  based on  $X$** , we refer an element  $\mathcal{M} \in C(X, Y)$  that is smooth on an open dense subset of  $X$ . Let  $((X \models Y))$  denote the subspace of  $C(X, Y)$  consisting of all models for  $Y$  based on  $X$ .*

(2) *For models  $\mathcal{M} \in ((X \models Y))$  and  $\mathcal{N} \in ((Y \models Z))$  where  $Z$  is another smooth manifold, the composite  $\mathcal{N} \circ \mathcal{M}$  is again a model. We call it the **composite model**.*

(3) A **loss function** on  $((X \models Y))$  is a function  $\lambda : ((X \models Y)) \rightarrow \mathbb{R}$  that is continuous and bounded below.

One typical way to set a loss function is the following.

**Definition 4.2.** (1) A **(supervised) dataset** for  $((X \models Y))$  is a finite subset  $S \subset X \times Y$ .

(2) Given a dataset  $S$  and a function  $f : Y \times Y \rightarrow \mathbb{R}$  which is bounded below and smooth on an open dense subset of  $Y \times Y$ , a **supervised f-loss** (or simply an **f-loss**) with respect to  $S$  is a loss function  $\lambda_f : ((X \models Y)) \rightarrow \mathbb{R}$  defined as follows:

$$\lambda_f(\mathcal{M}) = \sum_{(x,y) \in S} f(\mathcal{M}(x), y).$$

In practice, we set  $\lambda$  as a sum discrepancy measures each representing a training objective. Since the compact-open topology on  $((X \models Y))$  is too coarse to have a differentiable structure that is necessary for training a machine, we shall consider the following concept.

**Definition 4.3.** Let  $X$ ,  $Y$ , and  $\Theta$  be smooth manifolds and  $\lambda : ((X \models Y)) \rightarrow \mathbb{R}$  be a loss function.

(1) A collection of models  $\{\mathcal{M}_\theta\}_{\theta \in \Theta}$  for  $Y$  based on  $X$  is said to be **parametrized by**  $\Theta$  if (a)  $\mathcal{M}_{(-)} : \Theta \rightarrow ((X \models Y))$  is continuous and (b)  $(\theta, x) \mapsto \mathcal{M}_\theta(x) \in Y$  is smooth on an open dense subset of  $\Theta \times X$ . Let  $((X \models_\Theta Y))$  denote the space of models for  $Y$  based on  $X$  parametrized by  $\Theta$ .

(2) We say  $\lambda$  is **admissible** for the parametrized model  $\mathcal{M}_{(-)} \in ((X \models_\Theta Y))$  if the composite continuous function

$$\begin{aligned} \lambda \circ \mathcal{M}_{(-)} : \Theta &\rightarrow \mathbb{R} \\ \theta &\mapsto \lambda(\mathcal{M}_\theta) \end{aligned}$$

is smooth on an open dense subset of  $\Theta$ .

(3) A parametrized model  $\mathcal{M}_{(-)} \in ((X \models_\Theta Y))$  is said to be **pre-trained** with respect to the loss function  $\lambda$  if  $\lambda$  is admissible for  $\mathcal{M}_{(-)}$  and there is an observed  $\tilde{\theta} \in \Theta$  equipped with a neighborhood  $V_{\tilde{\theta}}$  on which  $\lambda \circ \mathcal{M}_{(-)} : V_{\tilde{\theta}} \rightarrow \mathbb{R}$  is smooth and attains a local minimum at  $\tilde{\theta}$ .

Note that, for two parametrized models  $\mathcal{M}_{(-)} \in ((X \models_\Theta Y))$  and  $\mathcal{N}_{(-)} \in ((Y \models_\Phi Z))$ , the composite  $\mathcal{N} \circ \mathcal{M} \in ((X \models_{\Theta \times \Phi} Z))$  is also a parametrized model.

DEPT. OF MATHEMATICS, YONSEI UNIVERSITY  
Email address: jhj941211@gmail.com